MARCH 9 th 2024



PROCEEDINGS OF

INTERNATIONAL CONFERENCE

COMPUTATIONAL INTELLIGENCE AND ITS APPLICATIONS (ICCIA-2024)



Published by

PG & Research Department of Computer Science

A.V.C. College (Autonomous)

Affiliated to Annamalai University NAAC Reaccredited 'A+' Grade Institution (4th Cycle) (CGPA=3.46/4.00) NIRF All India Ranking 2023: College Rank Band 101-150 UGC Recognised "College with Potential for Excellence-Phase I & II"

Mannampandal, Mayiladuthurai-609 305 Tamil Nadu, India **Proceedings of**

International Conference on "Computational Intelligence and its applications"

[ICCIA – 2024]

March 9th 2024

published by



PG & Research Department of Computer Science

A.V.C. College (Autonomous)

(Affiliated to Annamalai University) Mannampandal – 609 305, Mayiladuthurai

ISBN: 978-81-967420-1-0

Title: Computational Intelligence and Its Applications (ICCIA-2024) ISBN: 978-81-967420-1-0 Editor: Mr. M. Muthamizharasan, Dr. M. Hemamalini Published 2024 by AN PUBLICATIONS

© AN PUBLICATIONS & Author

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means (electronic, mechanical, photocopying, recording or otherwise), without the prior written permission of the publisher and Author.

The contents of this book is expressed by the authors and they are the responsible for the same.



AN Publications

No: 29, Moorthy Street, Balavinayagar Nagar, Tiruvallur-602001, Tamilnadu, India.

CHIEF PATRON'S MESSAGE

Thiru. K.VENKATARAMAN, Chief Patron of ICCIA-2024, Judge - Administrator, A.V.C. Institutions.

..........

.



.........

As the Chief Patron of the International Conference on "Computational Intelligence and its Applications (ICCIA-2024)" organized in the Department of Computer Science on 9th March 2024, I congratulate and convey my wishes to all the faculty members of the Department who dedicatedly involved in this international conference.

I am delighted to note that the accepted papers are being published in an edited book volume with ISBN. I wholeheartedly appreciate all the sincere efforts of the entire team of ICCIA-2024 and wish them a grand success.

Being the Chief Patron of ICCIA-2024, I feel very proud that this international conference would develop and promote the research work at higher level in the field of Computer Science.

I hope this conference will hold a series of intellectually interactive sessions and intensive deliberations by scholars and technical experts who are participating in it. This Conference will be the eye opener for the researchers, students and faculty to show the various avenues in the field of Computer Science. Further, the young scientists and researchers will find the contents helpful to set roadmaps for their future endeavors.

I wish the conference a grand success.

........... K.Venkataraman

PATRON'S MESSAGE

Dr. R. Nagarajan, Patron of ICCIA-2024, Principal, A.V.C. College (Autonomous).

..........



.........

A.V.C. College (Autonomous) has certainly came a long way and has provided educational platform to innumerable students who are enthusiastic to accomplish their dreams and ambition. With immense pleasure and enthusiasm, I take the opportunity to appreciate the organizers of International Conference on "Computational Intelligence and its Applications [ICCIA-2024]".

This programme is a first initiative taken by Department of Computer Science to ignite and inspire young researchers in the field of Computer Science. This academic forum will enable us to converse and also to reconnect with colleagues, and establish new professional contacts among the National and International experts in the field.

I am confident that this conference will act as a platform for present and share the research findings for further enhancement and perfection. This Conference focuses on researchers, professionals, educators and students to share innovative ideas, various issues pertaining to the topics, recent trends and future directions in the field of Computer Science. This conference will be an enriching experience for the faculty members as well as students.

I hope that the research articles published in the ISBN book will be a valuable resource in professional, research and academic activities.

My best wishes to this noble endeavour of Department of Computer Science of A.V.C. College (Autonomous). I wish the conference a great success

R.Naçarajan





CONVENER'S MESSAGE

Thiru. M.MUTHAMIZHARASAN, Associate Professor and Head, Department of Computer Science.

..........

.

event.

.........



..........

............

As a Convenor of ICCIA-2024, I feel delighted in organizing our maiden attempt to initiate this conference. This Conference aims to bring together scientists, researchers and technocrats on a platform to discuss on the issues confronting them and to look for solutions. This conference will not only provide a platform for discussions, it will help researchers to publish their work in the form of book with ISBN 978-81-967420-1-0.

I am more indebted to our honourable Judge-Administrator, A.V.C. Institutions, **Thiru. K.Venkataraman** and our beloved Principal, **Dr. R. Nagarajan** for their constant encouragement and moral support to organize this international conference. At this juncture, I have the special privilege to thank the Academic Advisor **Dr. M. Mahalakshmi** who is instrumental for the International Conference, it is because she belongs to the faculty of Computer Science.

On behalf of the Department and my own behalf, I have to express my special applause to the Guest Speaker **Mr.Vinoth Thiruvengadam**, Airbus, Test strategy and Transformation Leader, France, for sharing his valuable expertise in the International Conference. His address will definitely inspire the budding scholars and researchers.

I personally salute all the Advisory board members of the Conference who spare their valuable time to review the research papers and to provide valuable comments to enrich the research articles.

Words are inadequate to express my gratitude for the strenuous work rendered by the Organizing committee members.

On behalf of the organizing committee of this International Conference, I thank all the learned authors and well wishers who are directly or indirectly contributed to the successful completion of this Conference.

The Conference Programme includes foreign and Indian experts talk, Presentation of Participants, and distributing Awards and Certificates.

The grand success of this International Conference lies in taking the responsibility by everyone in the Department to conduct the event in smooth and fruitful manner.

Finally, I thank God Almighty to shower all its blessings for the smooth conduct of this

M.Muthamizbarasan

ACKNOWLEDGEMENT

The Department of Computer Science has extended a warm salutation and special thanks to our honorable Chief Patron Justice **Mr. K. Venkataraman**, Judge-Administrator, A.V.C. Institutions for his philanthropic vision and his constant encouragement to make this conference a grand success.

We would like to take this opportunity to express our sincere gratitude to our learned Principal **Dr.R.Nagarajan**, amidst his busy schedule, he constantly encouraged us to conduct this international conference within the short span of time.

We would like to give special applause to the Resource person **Mr.Vinoth Thiruvengadam**, Airbus, Test Strategy and Transformation Leader, France, for sharing his valuable insights and expertise with us. He has inspired and conquered young minds, researchers and faculty members.

A special recognition is needed to the Advisory Board members **Dr.S.Saradha**, Assistant Professor, MEASI institute of Information Technology, **Dr.G.Jaculine Priya**, Assistant Professor-Artificial Intelligence, Loyola Institute of Business and Administration, **Dr.D.Suganthi**, Associate Professor, Saveetha College of liberal Arts and Science, **Dr.V.Gowthami**, Assistant Professor, Kamban College of Arts and Science for Women who have dedicated their valuable time, effort and energy to bring out the research articles with high standards.

At this juncture, we extend my deep sense of gratitude to our Subject Experts and Chair persons **Dr.M.Mahalakshmi**, Principal, R.B.Gothi Jain College for Women and **Dr.S.Selvamuthukumaran**, Professor and HOD, Department of Computer Applications, A.V.C. College of Engineering for enriching the conference and have helped to create a platform for meaningful discussion and exchange of ideas.

Our thanks are extended to the Convenor **Mr.M.Muthamizharasan**, Associate Professor and Head, Department of Computer Science, for his untiring commitment to complete this academic event a memorable one. Amidst his department work, he has keen interest to make this conference a meticulous intellectual sharing.



..........



...........

ACKNOWLEDGEMENT

Our special acknowledgement and hearty thanks to the senior Professors in the Department of Computer science, **Mr.C.S.Iyyappan**, Associate professor and **Dr.K.Palanivel**, Associate Professor, who is a guiding source besides this great academic event. Their professionalism and efficiency have made us to work much easier.

We are more indebted to the Organizing Committee members **Dr.M.Hemamalini**, **Dr.P.Panimalar**, **Dr.R.Neela**, **Dr.S.Navitha**, **Dr.S.Kiruthika** and Student Coordinators **J.Rakkesh Kumar**, **R.Sathiya**, **S.Indhirani** and **S.Archana** for their meticulous execution to make this event a grand success. Words are insufficient to express our gratitude to these noble minds to obtain perfection in every stage of this conference.

We also extend our hearty thanks to the technical and non-teaching staff for their wonderful services and support throughout the conference is highly commendable.

Last but not least, we would like to thank all the delegates who have attended this conference for carrying this message to the working place and spreading the knowledge to the students' community. Their active participation, engaging discussions and valuable contributions, have made this conference truly enriching and memorable event.

Once again, we register our sincere gratitude to everyone who has been a part of this conference. We look forward to your continued support in our future endeavors. Thank you for being a part of this journey and we hope to see you again in our future conferences.



..........



..........

Table of Contents

Sl. No	Paper Id	Authors	Paper Title	Page Nos.
1.	33	C. Ananth, S. Kasinathan & N. Mohananthini	Secure and Trusted Valuation Result of Competitive Examination Using Blockchain Technology	1-6
2.	07	S. SivaShankari & Dr. K. Saminathan	A Comprehensive Survey of Machine Learning and Deep Learning Enabled Hand Gesture Recognition Models	7-12
3.	17	Mrs. T. Thilagavathi & Dr. A. Subashini	A Safe Human Route Prediction using Machine Learning based on Multi-Diversity Factors	13-18
4.	21	C. Justin Marshal & Dr. R. Vidya	Unmasking Deception: Stance Detection in the Age of COVID-19 Fake News	19-25
5.	03	C. Muruganandam & Dr.V. Maniraj	Smart Agriculture Monitoring System Using IOT with RMS and SMS by Using AWT13 SENSOR	26-31
6.	24	Dr.K. Palanivel & M.Muthamizharasan	Performance Evaluation of Classification algorithms with Liver and Diabetic Patient Datasets using WEKA tool	32-37
7.	18	Dr. S.P. Ponnusamy & Mrs. S. Valli	Utilizing Machine Learning Techniques for Early Detection and Control of Powdery Mildew Disease in Cashew Flowers to Enhance Crop Yield	38-41
8.	25	J. Jagadeesan & Dr. R. Nagarajan	Embarking on a new journey of Federated Learning and Tensor Flow Framework	42-49
9.	16	T. Ramyaveni & Dr.V. Maniraj	Prediction of Diabetes by Using Intellectual Health Care with Using Machine Learning Algorithms	50-53
10.	02	P. Kalaimagal & Dr. S. Kumaravel	Prediction of Diabetic Retinopathy using Machine Learning with Deep Learning Models	54-58
11.	35	S. Govindarajan & Dr. S. Marry Vennila	A Review on Rice Leaf Diseases Using Feature Extraction Techniques in Machine Learning	59-62
12.	26	 V. Vasanthi, R. Bhuvaneswari, M. Paul Arokiadass Jerald & I. Benjamin Franklin 	Design of Resilient Cyber defence mechanism using Dynamic Zero Trust Network Security Framework by integrating Artificial Intelligence, Machine Learning and Blockchain	63-71
13.	20	R. Durgadevi	A Survey On Deep Learning Algorithms and Its Applications	
14.	22	R. Palanivel & Dr. P.Muthulakshmi	mi Exposing Quantum Theory's Influence Through Computational Insights into Physics, Chemistry, and Mechanics	
15.	19	M. Dhivya & Dr. V.Maniraj	Stock Value Prevision In Machine Learning Using Generative Adversarial Networks	85-89
16.	05	M. Rega	Unleashing the Power of Convolutional Neural Networks in Image Processing	90-93

17	20	Dr.V.Geetha &	Comparative study of different classification	94-96
17.	32	G.Elakkiya	Algorithms in data mining using kddcup-99 dataset	
		T.Parikodi	A review on deep learning models and their	97-99
18.	31		limitations in brain MRI segmentation	
		Dr. S. Thaiyalnayaki,	Analysis of Linear Regression Model	100-104
19.	28	Ms. Naeem Fathima &		
		Ms. M. Manthra		
20	04	G.Hemamalini &	Rheumatoid Arthritis Disease Prediction using	105-110
20.	01	Dr.V.Maniraj	Machine Learning technique SVM	
		Dr.C. Ananth, S.	A Comprehensive Analysis of Novel Intrusion	111-121
21.	14	Sathiyarani &	Detection Systems and Intrusion Prevention	
		Dr.N. Mohananthini	Systems for Blockchain Technology	
22	01	Dr.S. Jayaprakash &	A Comprehensive Review of Literature on current	122-128
	01	J.P. Keerthana	and Research Fields related to Robotics	
		N. Ruba &	An Secure Commercial Enterprise Transaction	129-135
23.	12	Dr.A.Shaik Khadir	System Using Alk (Alpha Keys Authentication)	
		Mohideen	Technique	
24	23	R.Sathya	A Systematic Review on Privacy Preservation of	136-139
			Big Data in Cloud	
25	06	J. Rakkesh Kumar &	An Outlook on Role and Challenges of Big Data	140-143
23.	00	Dr.M. Hemamalini	Analytics in Health Care	
26.	15	R. Ragavi & P.Subastri	Nanorobotics Using Artificial Intelligence	144-145
27.	27	Raja Thangavelu	Big Data Analytics: Review & Recommendation	146-151
20	24	R. Senthamarai &	Robotic innovations in healthcare:	152-155
28.	54	Dr. A. Senthil kumar	Enhancing patient care and efficiency	
		K. Kavitha	Artificial Neuron Network: Review	156-158
29.	29	S. Muthulakshmi &		
		M. Nithisha		
20	00	V. Parvathi &	Artificial Intelligent With IoT	159-161
50.	08	P. Pragadeeswari		
21	00	S.Senthikumar &	Prediction of Diabetic Kidney Disease Using Deep	162-165
51.	09	Dr.T.S. Baskaran	learning Techniques	
37	10	M.Menaha	Agriculture Data Analysis using Machine Learning	166-168
52.	10		Techniques	
		S. Narmatha &	A Survey Concerning the use of Electronic	169-176
33.	38	Dr. V. Maniraj	Medical Record Search Engines in the Healthcare	
			Industry	
		S. Ubaidulla &	A Comprehensive Review on Improving Plant	177-181
34.	36	Dr. S. Mary Vennila	Leaf Disease Detection Accuracy through	
			Computer Vision Techniques	
		N. Subhalakshmi &	A Cloud Based Secure Electronic Health History	182-187
35.	37	Dr.M.V. Srinath	Framework Using Fernet And Fully Homomorphic	
			Encryption	

26	40	N. Srinivasan &	Preserving Health Care Data Privacy Using	188-194
50.	40	Dr. S. Selvamuthukumaran	Federated Learning	
27	30	A.Srilekha & Punitha P	Customer Segmentation For Analysis Of	195-198
57.			Prediction Using Data Mining Techniques	
20		M. Nagaraj	Deep Learning and Applications	199
50.				
		B. Nithiya Bharathi &	Blue Eyes Technology	200
39.		S. Sathiya devi		
		-		
		A. Kavya &	Coded Cryptosystem	201
40.		M. Madhavan		

Secure and Trusted Valuation Result of Competitive Examination Using Blockchain Technology

C. Ananth¹, S. Kasinathan² and N. Mohananthini³

¹AssistantProfessor / Programmer, Department of Computer and Information Science, ²Research Scholar, Department of Computer and Information Science, Annamalai University, ³AssociateProfessor, Department of Electrical and Electronics Engineering, Muthayanmal Engineering College ¹Annamalai University, ²Annamalai University and ³Muthayanmal Engineering College, Rasipuram Tamilnadu, India.

Abstract - Blockchain is a solution to improve data integrity and minimize data manipulation. Blockchain technology can be used in all examinations conducted by the governments and educational institutions. The process of managing student examination produces some information that can be stored in a blockchain database such as information about the student attendance, student exam results and the student study continuity status. Blockchain technology ensures student data is valid and reliable. Therefore, the governments and educational institutions trust the teaching and learning process in competitive examination. This research was conducted using a qualitative approach, which is user center design (UCD) the goal is building a blockchain technology model that can be used by governments and educational institutions in the process of managing competitive examination, grading, trust and secured evaluation result. This process will be repeated every examinations until students get the certificates and the results published in online as well. The purpose of this research is to create a trust and secured valuation result of competitive examinations using blockchain technology in all sectors. The process model based on blockchain technology that can improve data integrity in government and educational institutions.

Keywords - blockchain, results, competitive examinations, governments, educational institutions.

I. INTRODUCTION

Traditional education systems around the world are still using the same examination system that has been around for centuries to measure the knowledge that a student process about a specific subject. Interestingly, there have been no major improvements in the way that students are evaluated by examination bodies. Many universities are still using the same principles of examinations that answers to specific questions are written on a paper, which are then reviewed by an external examiner. This kind of examination system has often been criticized for being biased and prejudiced. Moreover, there are also subjects that require critical thinking and the marking scheme created by a particular teacher cannot contain all the possible answers or may even miss out on the best answer. There is also a waste of time and a resource whenever a student feels his or her paper script was wrongly evaluated and initiated a review process.

Blockchain can be used to design systems, which are distributed, tamper-proof, and protected. In this study, a system is implemented, which primarily replaces the current examination and evaluation system with a transparent and distributed examination and evaluation system. Second, the system must also be able to account for the number of modules a student has cleared and in case of completing all his modules, a digital certificate is also issued to him or her and stored on the blockchain itself. The student shall also have access to his or her updated transcript and certificates via a blockchain wallet. A pool of administrators from different institutions involved in the usage of the proposed system is expected to have the authority of the handling of the blockchain and the system itself.

From a technical perspective, the objectives of this study are to develop a light blockchain system, which will be responsible for creating new transactions, such as credits transaction in the event that a student clears a module or a certificate transaction, when a student completes a course. The system should also ensure that integrity and transparency are always guarded. Both political and physical decentralization of the blockchain is required. A public web portal, which will host a blockchain explorer where users can verify any transaction claimed by a student and an administrator web application where the admins could have access to blockchain functionalities will also, be implemented. Furthermore, an examination application to be installed in all the desktop workstations in the laboratories in the educational institutions where students could take exams will be implemented. These will also be active nodes in the network. Finally, a mobile wallet for the blockchain that will allow students to have a copy of their transcripts and certificates all the time with them will be developed.

Blockchain is an emerging technology was introduced in 2008 by Nakamoto, S. It was first used as a peer-to-peer ledger for registering the transactions of Bitcoin crypto currency. Competitive exams are a type of examination conducted to test and rank students according to their grades, percentage, or percentile. These exams are mostly conducted at the state or central level. Competitive exams in India are conducted for various fields like engineering, medicine, law among many others. These exams are mostly conducted annually and millions of individuals appear for these exams. The idea of such exams is to select deserving and worthy candidates by making them appear for the exam in a fair and unbiased setting. So, everyone needs the secured and trusted valuation result for their effort.

A competitive exam helps select students for the various job strata. Selected candidates from these exams are trained for

jobs in their respective fields. They are conducted at various levels. Students appear for such exams after class 10th, 12th, diploma, under graduate, post graduate as well. Bank examinations, Railway examinations, recruitment examinations, civil services exams, state government exams like TET exam, group exams and etc., are some exams students can appear after their level of studies. Such exams in India are highly competing and require time and effort to clear. We advise you to start preparing as early as possible and remember to take mock tests.



II. LITERATURE REVIEW

In this era, blockchain is one of the data integrity solutions and data manipulation solutions. Blockchain is a technology used as a concept in distributed ledgers, where it can be validated by consensus and the presence of cryptographic algorithms. The concept that Satoshi Nakamoto built in 2008 was the use of Bitcoin (Nakamoto, 2008)

A blockchain-based framework for conducting and evaluating Government examinations. To perform the test as transparently as possible, we store the hash code of every question asked and answered directly on the Blockchain Network. This facilitates tracking how exactly a candidate received the score he or she received, adding more credibility & transparency to the obtained score. Most of the work in educational institutions is based on technology like blockchain, it can transform into a simplified, paperless manner. The advanced security mechanism of blockchain will ensure that the system can be immune to hacking. Data cannot be manipulated with any other entity apart from the node owners. (Himanshi Dang and Khushi Thareja, 2022)

In this literature review paper, the system will provide a Blockchain based framework for evaluating academic tests in manner with auto-generation а peer-to peer of scorecard/certificates upon successful completion of the examination. The system will provide an evaluator from an educational institute validating the answers to the questions asked. There will be no requirement of re-valuation of the validation performed by the evaluator. For this purpose, gathering of required data will be done at the college level. In order to make the evaluation as transparent as possible, we will store the hash-digest of every question asked and every question answered, directly on the blockchain. This facilitates the tracing of how exactly a candidate received the score that one has gained, adding more credibility to the obtained certificate. Thereby, the system will demonstrate how a selfsustained education ecosystem can be developed on top of a blockchain for a fair evaluation without the need of a central trusted entity for obtaining certificates or degrees that prove one's understandability over a subject. (Patekar Manali, Gandhi Mitul, et al., 2021)

The main objective of this paper is to provide a Blockchain based framework for conducting and evaluating academic tests in a peer-to-peer manner with auto-generation of certificates upon successful completion of the examination. We illustrate how a self-sustained education ecosystem can be developed on top of a blockchain for a fair evaluation without the need of a central trusted entity for obtaining certificates or degrees that prove one's dexterity over a subject. In order to make the test as transparent as possible, we store the hashdigest of every question asked and every question answered, directly on the blockchain. This facilitates the tracing of how exactly a candidate received the score that he/she received, adding more credibility to the obtained certificate. (Rahul Acharya, Sumitra Binu, 2018)

The study of blockchain technology in Education Value Chain Model for Examination, Grading, and Evaluation Process in Higher Education based on Blockchain Technology, in this research was conducted using a qualitative approach, which is user center design (UCD) the goal is building a blockchain technology model that can be used by universities in the process of managing student examination, grading, and evaluation. This process will be repeated every semester until students graduate from the university. The purpose of this research is to create a business process model based on blockchain technology that can improve data integrity in universities. (Meyliana, Yakob Utama Chandra, et al., 2019)

In this paper, it aims to propose a blockchain-based framework that facilitates the secure and peer-to-peer conduct and evaluation of academic exams. The framework employs hashing techniques to ensure the integrity of the data and utilizes proof of stake mechanisms to enhance security. Blockchain technology has proven to be effective in safeguarding data integrity by virtue of its decentralized data storage approach and the use of cryptographic hashing for every block within the chain. This paper demonstrates how online exams can be developed using blockchain technology, with each question asked and answered being directly stored on the blockchain. To achieve this, we have developed a module that integrates with the Moodle learning management system. Through a comparative analysis of the default centralized storage approach in Moodle, our module modifies the exam results' storage method, ensuring secure and tamperproof data storage on the blockchain network. By leveraging

the blockchain network, the data associated with exam results is reliably secured, ensuring its integrity, and making it immune to manipulation. Our results indicate that the data stored through the blockchain achieved complete accuracy, with no discrepancies observed when compared to the standard approach employed by the Moodle LMS for storing results. The blockchain network provides a reliable and immutable prevents platform that unauthorized alterations or manipulations of student data. In conclusion, our blockchainbased framework offers a robust solution for enhancing the security and reliability of online exam results. By leveraging the decentralized and tamper-proof nature of blockchain technology, we can ensure the integrity and transparency of student data, ultimately providing a more trustworthy and accurate assessment of their academic performance. (Mohamed Abdelsalam1, Amira M. Idrees2, and Marwan Shokry3, 2017)

In their study 1, they evaluated university students' emotions at the end of a computer-based exam and found positive emotions more strongly endorsed than negative. In their study 2, they replicated this finding and used a quasi-experimental pre-post design to examine how emotions changed in response to real examination scores. Exam scores presented immediately had significant positive effects on relief, pride, and hope and negative effects on anxiety and shame even after controlling for the corresponding emotion at the end of the exam. The one exception was anger, which was not impacted by examination score. No interaction effects were found. (L.M. Daniels, M.J. Gierl, et al., 2017)

This article explains how higher education institutions use blockchain technology that is able to maintain data security and prevent data manipulation, especially in academic transcript files. Academic transcript records become an output from students after completing a semester or level of education at the higher education institution. In general, academic transcript records are used when students complete a study at a higher education institution as a document requirement for access to the company where they work. Thus academic transcript reports become important for companies to find out if this student is competent in a certain area. In this case, accurate and valid data is certainly needed and the truth can be trusted. This qualitative research contributes to the higher education institution through a proposed blockchain technology model, especially in academic transcript records, so that the results of the proposed model ensure that data manipulation does not occur in data in academic transcript records. The results of this model will be used as design material for application design, so that the blockchain platform can be used in higher education institutions, especially when publishing valid academic transcript records that can be trusted with the data. (M. Meyliana, Y.U. Chandra, et al., 2019).

The implementation of blockchain to education is still in its early stages. Only a small number of educational institutions have started to utilize blockchain technology. Most of these institutions are using it for the purpose of validating and sharing academic certificates and/or learning outcomes that their students have achieved.

III. METHODOLOGY

This paper proposing a new scheme for the good result by utilizing blockchain-based on a two- phase encryption technique for encrypting the final result. In the first phase, question papers are encrypted using a timestamp, and in the second phase, the result is encrypted using question paper hash code. These encrypted results are stored in the blockchain along with a smart contract which helps the user to unlock the result with the same question paper hash code. Here we proposed a method for selecting a question paper for the exams, which randomly picks a question paper with a hash code. Moreover, a timestamp-based lock is imposed on the scheme so that no one can decrypt the question paper before the allotted time. Finally, the result gets stored in Blocks of the Blockchain, and security is analyzed by demonstrating various suggestions and the prevalence of the proposed conspire over existing. It is demonstrated through a comparative study based on the various features; it provides a persistent public record, safeguarded against changes to the institution or loss of its result records.

Security Mechanism of Blockchain using Smart Contracts For Conducting Fair Examination.

Smart contracts used in blockchain technology to validate, verify, capture and enforce agreed- upon terms between multiple parties. In smart contracts, all the data stored is secure and immutable. The data of a smart contract is encrypted and exists on a ledger, which means that the information recorded in the blocks can never be modified, lost, or deleted. Steps to verify the student/user attending the examination.



1. The registration of the user at the blockchain front-end examination by giving their personal details and their fingerprint, which would be stored as a Hashcode in Examiner's Blockchain verification database, and would be used to verify the user at the time of the examination.

2. In the examination time, the examiner wants to verify the user with his details and the fingerprints which would send the information to the blockchain of the registered candidates and should confirm with the hash code generated, if the Hashcode generated matches with the database of the registered candidates, and it matches the conditions of the smart contract, the success message would be sent, and the user will be verified. This will decrease the chance of malpractice, and the eligible candidate would only be able to attend the examination.

3. Every individual has a unique fingerprint; it's a simple way to verify the users. The data of fingerprints are collected by

encrypting the finger template into data sets. The identification of the user details like name, age, address, etc., are collected and passed into the hash function with the fingerprint data sets, and a unique hash code is generated for each user for his identification purpose. As the organization verifies the user's data, then the data is stored in the Blockchain of the registered students.

5. After the verification, the organization will issue a smart contract to verify the contract needs.

6. Now we will use another blockchain that will connect personal details of candidate with the question paper which will ensure that every question paper id is assign to a candidate and it will also ensure that there is no leakage of question paper before the examination. Because question paper would be only visible to those valid candidate which is meeting the criteria of smart contract.

7. Every candidate has a unique answer sheet id which is connected to candidate id via blockchain through the hash code of question paper which is already assigned to the candidate.

8. Now, Finally every candidate result is connected via blockchain of candidate detail so now candidate can see their result with 100% surety. In this way, a mutual trust and transparency will create between government and candidates who is giving any government exam. This advanced security mechanism of Blockchain technology would ensure that the system can be immune from hacking which means data cannot be manipulated with any other entity apart from node owners.

9. Challenges of Adopting Blockchain Technology in Education and Conducting of Government examinations.

10. Despite the way that blockchain has extraordinary potential in an educational context; various difficulties should be considered before executing it. Like some other groundbreaking innovations, blockchain in its beginning phases of advancement faces a few difficulties. Regardless of whether it be a mix with inheritance frameworks, HR requirements, cost of execution, and so forth.

VI. PROBLEM STATEMENT

The current examination system involves an evaluator from an educational institute validating the answers to the questions asked. There is however no re-validation of the validation performed by the evaluator. In cases where reevaluation is done, it is done by a handful of people. This makes the evaluation system highly centralized and there are many problems associated with centralized evaluation.

Centralized evaluation is highly susceptible to score manipulation. The manipulation can be done at any stage; right from the first evaluation to the manipulation during the final data entry in the database.

Since the data is stored in the database and it is under the control of a database administrator, it brings in the human interference which is susceptible to bribery or threats.

Another fundamental problem with the scorecard of the current examination system is that they do not provide enough data to represent the performance of a candidate taking up the test. The scorecards contain very limited information about the performance as it only accounts the final score granted by one or two evaluators without disclosing the questions asked and the manner in which the questions were answered. With no idea of the types of questions asked to a candidate, correlating the score with the caliber of the candidate mostly leads to inaccurate conclusions. Centralized issuance of degrees or certificates is susceptible to manipulation. The certificates received by a candidate upon completion of a course indicate that the candidate has the expected skills and knowledge demanded by the successful completion of the course. However, the certificates can be forged or granted even when the candidate doesn't meet the criteria of receiving the certificates if an institution decides to grant it no matter what. There is no way to know whether the certificates issued by an institution are issued even when the criteria of issuance are not satisfied by the performance of the candidate. Nevertheless, the process of just validating the authenticity of the issued certificate is expensive and slow.

We use a public blockchain with decentralized evaluation and maintenance of examination records to solve all the problems and provide a better alternative. In the decentralized evaluation mechanism, we perform two types of evaluations, one for the questions and one for the answers to the questions. The community votes for the validity or the relevance of the posted question in a particular category. Thus, the quality of the questions can be expected to be much better as decided by the consensus of the users obtained by the translation of their votes on each question.

V. CONCLUSIONS

The aim of this paper is to illustrate an approach to use blockchain for conducting decentralized examination and for better evaluation of the examination records. The Blockchain based Exam Paper Evaluation system needs the question and paper to be uploaded in the blockchain. We try solving the lack of transparency and credibility problem in the current examination system by recording the details of the examination on the immutable public blockchain such that every operation is recorded as a transaction. By using blockchain technology, it can be concluded that secured and trusted valuation will be guaranteed, and data storage is more secure because the data has been validated and cannot be changed. Data changes can still occur by making new transactions without eliminating previous transactions, for example when students protest the score which causes the previous score to change. The Both scores are stored in the blockchain database.

Following the research methodology, we have conducted a focus group discussion to validate the model that has been made and the results are that all participants agree with this model and they think this study can be easily applied in universities in recording the academic transcript.

The process of examination, grading, and evaluation is very appropriate using blockchain technology, because this process is very crucial and requires data validity to produce the academic transcript for students.

In the current condition of the Covid-19 pandemic, using blockchain technology is more needed because the process input and evaluation are done digitally. With blockchain technology, governments and universities make a digital transformation in the education sector and make education technology more effective, especially during the Covid-19 pandemic.

Department of Computer Science, A.V.C. College (Autonomous), Mayiladuthurai – 609 305.

The limitations of this study are the process of interviewing and focus group discussion using virtual conferencing to get answers from the speakers and there is limited time in digging deeper into the information needed. So it is necessary to do the validity of the resulting model several times. Future research will build an integrated overall model to produce a value chain with data integrity starting from students entering higher education until students graduate and get diplomas, transcripts, diploma companion documents, and student activity transcripts.

5.2. Future scope

The present research lays the foundational framework of using a blockchain in the field of academic education. The current approach can further be enhanced by developing a scalable web application hosted on the IPFS, which allows interaction with the blockchain using a browser. Since the current framework has not be extensively tested for scalability, the paper suggests, further improvements can be made on the aspect of scalability of the blockchain in terms of numbers of transactions processed per second and the number of examinations conducted simultaneously.

REFERENCES

- Satoshi Nakamoto, "Bitcoin: A Peer-to-Peer Electronic Cash System," http://Bitcoin.org; satoshin@gmx.com, pp. 1-8, 2008.
- [2] Patekar Manali, Gandhi Mitul, Sardesai Onkar, Gupta Anurag "Blockchain based Exam paper Evaluation" March 2021.
- [3] D. Das, "Hacking into the Indian Education System", [Online]. Available: https://deedy.quora.com/Hackinginto-the-IndianEducation-System.
- [4] Rahul Acharya, Sumitra Binu, "Blockchain based examination system for effective evaluation and maintenance of examination records", 2018
- [5] Meyliana, Y.U. Chandra, C. Cassandra, Surjandy, E. Fernando, A.E. Widjaja, H. Prabowo, "A Proposed Model of Secure Academic Transcript Records with Blockchain Technology in Higher Education," (Conrist 2019), 172–177, 2020, doi:10.5220/0009907401720177.
- [6] M.M.A. Getso, Z. Johari, "the Blockchain Revolution and Higher Education," International Journal of Information Systems and Engineering, 5(1), 57–65, 2017, doi:10.24924/ijise/2017.04/v5.iss1/57.65.
- [7] Himanshi Dang and Khushi Thareja, "Applying Blockchain on Government Examination: A Study on Blockchain Technology, Benefits, and challenges" International Journal of Advanced Research in Science, Communication and Technology (IJARSCT), ISSN (Online) 2581-9429, Volume 2, Issue 6, June 2022.

- [8] R. Wuthisatian, "Student exam performance in different proctored environments: Evidence from an online economics course," International Review of Economics Education, (August), 100196, 2020, doi:10.1016/j.iree.2020.100196.
- [9] M. Meyliana, Y.U. Chandra, C. Cassandra, S. Surjandy, H.A. Eka Widjaja, E. Fernando, H. Prabowo, C. Joseph, "DEFYING THE CERTIFICATION DIPLOMA FORGERY WITH BLOCKCHAIN PLATFORM: A PROPOSED MODEL," in Proceedings of the International Conferences ICT, Society, and Human Beings 2019; Connected Smart Cities 2019; and Web Based Communities and Social Media 2019, IADIS Press: 63–71, 2019, doi:10.33965/ict2019_201908L008.
- [10] L.M. Daniels, M.J. Gierl, "The impact of immediate test score reporting on university students ' achievement emotions in the context of computer-based multiplechoice exams," Learning and Instruction, 2017, doi:10.1016/j.learninstruc.2017.04.001.
- [11] H. Liao, J. Hitchcock, Reported Credibility Techniques in Higher Education Evaluation Studies that use Qualitative Methods: A Research Synthesis, Elsevier Ltd, 2018, doi:10.1016/j.evalprogplan.2018.03.005.
- [12] Z. Wu, T. He, C. Mao, C. Huang, "Exam paper generation based on performance prediction of student group," Information Sciences, 532, 72–90, 2020, doi:10.1016/j.ins.2020.04.043.
- [13] L.M. Daniels, M.J. Gierl, "The impact of immediate test score reporting on university students ' achievement emotions in the context of computer-based multiplechoice exams," Learning and Instruction, 2017, doi:10.1016/j.learninstruc.2017.04.001.
- [14] H. Liao, J. Hitchcock, Reported Credibility Techniques in Higher Education Evaluation Studies that use Qualitative Methods: A Research Synthesis, Elsevier Ltd, 2018, doi:10.1016/j.evalprogplan.2018.03.005.
- [15] Gervais, Arthur & Karame, Ghassan & Wüst, Karl & Glykantzis, Vasileios & Ritzdorf, Hubert & Capkun, Srdjan. (2016). On the Security and Performance of Proof of Work Blockchains. 3-16. 10.1145/2976749.2978341.
- [16] Elisa, N & Yang, Longzhi & Li, Honglei & Chao, Fei & Naik, N & Nnko, Noe & Yang, Li & Chao,. (2019). Consortium Blockchain for Security and Privacy-Preserving in E-government Systems.
- [17] P. Bhaskar, C. K. Tiwari, and J. Amit, "Blockchain in Education Management: Present and Future

Applications," Interactive Technology and Smart Education, vol. 18, 2020.

- [18] A. Alammary, S. Alhazmi, M. Almasri, and S. Gillani, "Blockchain-based applications in education: a systematic review," Applied Sciences, vol. 9, no. 12, p. 2400, 2019.
- [19] S. Sharma and R. Singh Batth, "Blockchain technology for higher education sytem: a mirror review," in Proceedings of the 2020 International Conference on Intelligent Engineering and Management (ICIEM), Nanjing, China, June 2020.
- [20] P. Williams, "Does competency-based education with blockchain signal a new mission for universities?" Journal of Higher Education Policy and Management, vol. 41, no. 1, pp. 104–117, 2019.
- [21] Düdder B., Fomin V., Gurpinar T., Henke M., Iqbal M., Janaviciene V., Matulevicius R., Straub N., and Wu H., "Interdisciplinary blockchain education: utilizing blockchain technology from various perspectives," *Frontiers in Blockchain*, vol. 3, 2021. 578022 10.3389/fbloc.2020.578022
- [22] Sousa M. J. and Andreia de Bem M., "Blockchain technology reshaping education: contributions for policy," *Blockchain technology applications in education*, vol. 24, pp. 113–125, IGI Global, Pennsylvania, United States, 2020. 10.4018/978-1-5225-9478-9.ch006
- [23] Awaji B., Ellis S., and Albshri A., "Blockchain-based applications in higher education: a systematic mapping study," in Proceedings of the 2020 5th International Conference on Information and Education Innovations, pp. 96–104, London UK, July 2020. 10.1145/3411681.3411688
- [24] L. Li and X. Wu, "Research on school teaching platform based on blockchain technology," in Proceedings of the 2019 14th International Conference on Computer Science & Education (ICCSE), IEEE, Toronto, Canada, August 2019.
- [25] Williams P., "Does competency-based education with blockchain signal a new mission for universities?" Journal of Higher Education Policy and Management, vol. 41, no. 1, pp. 104– 117, 2019. 10.1080/1360080x.2018.1520491 2-s2.0-85053489348
- [26] Ashis Kumar Samanta & Bidyut Biman Sarkar & Nabendu Chaki., "A Blockchain-Based Smart Contract Towards Developing Secured University Examination System" Journal of Data, Information and Management

(2021) Springer Nature Switzerland AG 2021 3:237–249 https://doi.org/10.1007/s42488-021-00056-0

A Comprehensive Survey of Machine Learning and Deep Learning Enabled Hand Gesture Recognition Models

S. Siva shankari¹ and Dr.K.Saminathan²

¹*Research Scholar* and ²*Associate Professor*

Department of Computer Science, A.V.V.M Sri Pushpam College (Autonomous) Poondi, Thanjavur. Affiliated to Bharathidasan University, Tiruchirapalli, Tamil Nadu, India. devamusiva@gmail.com

Abstract - Gestures are treated as a natural expression of the human body and are widely utilized by deaf and dump people to connect with other people. Among the different kinds of gestures, hand gesture is commonly utilized over the globe. The design of automated hand gesture recognition (HGR) models has considerably increased in recent times. It is also found useful in different domains such as healthcare, sign language, human computer interface, robots, virtual reality, etc. The HGR is a problem of feature extraction and pattern recognition, in which a movement is labelled to a particular class. The recent development of computer vision (CV), machine learning (ML), and deep learning (DL) models have resulted in the design of several HGR models. In this view, this paper concentrates on a detailed review of existing HGR based on ML and DL models. Besides, the different processes involved in the HGR models such as presegmentation, feature extraction, processing, and classification are elaborated and challenges involved are identified. This study also offers a survey of existing methods based on aim, objective, technique used, and performance measures. In addition, a brief discussion of the results obtained by the reviewed methods is provided along with possible future scope. At the end of the survey, we hope that the readers were aware of the basic introduction to HGR. along with recent state of art approaches, and also assist future research effort under this field.

Keywords: Hand gesture recognition, Deaf and dump people, Machine learning, Deep learning, Computer vision, Segmentation

I. INTRODUCTION

Deaf is a disability which damages their hearing and makes an individual incapable of hearing, whereas mute is also considered a disability but it damages its talking ability and causes incapability of speaking [1]. These two are only disabled in their speaking and or hearing, thus they can even perform a lot more things. The one and only factor which differentiates them and the ordinary person is communication [2]. If there are any means for ordinary people and deaf-mute people to communicate, the deaf mute person may also survive as an ordinary person. And the one and only means of communication for them are via sign language. Whereas sign languages are the most significant factor for deafmute people, communicating with normal people as well

as with themselves, is even acquiring slight interest from the always try to neglect significance of sign language until there are beloved persons who seem to be deaf-mute. Using the services of sign language translators one can able to communicate with the deaf-mute person and it acts as only solution. However, use of sign language translators seems to be expensive [3]. Cheapest solution is necessary thereby the deafmute and ordinary person could communicate in a normal way. Gesture recognition utilizing computer technologies may be utilized as a translator for sign language translation. This becomes an advantage and acts as a connecting bridge among these groups. The hand gestures (HG) utilized in sign language are of 2 types: dynamic and static gestures. The static gesture is termed as the place of fingers and hands from the space which involves no movements in relation to time [4], while dynamic gestures involve continual movements of hands. From this perspective, HG identification is an issue made up of 2sub issues they are pattern recognition and feature extraction. Hand gesture recognition (HGR) contains mapping an input of a set with sets of labels, whereas a label indicates gestures that identified. And, it becomes mandatory in identifying the instant of time once the movement takes place [5]. Design of an HGR system comprises associating distinct modules post-processing, namely, preprocessing, feature extraction, data acquisition, and classifier. The classification module can be formulated through machine learning (ML), particularly if the issue is highly complex or not possible for finding a mathematical model (that is, probability distribution). Order to find a mathematical model must demand awareness of its behavior and the dynamics of the issue. Thus, identifying a mathematical model which explains HGR with very much accuracy is hard. The recognition of HGs shows a broad research region which is sub-classified dependent upon the gestures context as well as the technology used to input those gestures. There are distinct taxonomies which influence the designed HGR systems; environmental elements, person who executes the gesture, the efficiency of devices used to capture, the gestures types (dynamic or static,) and the system applications. Several methods to

evaluate the HGs systems were discovered: at first, sensors related methods demand the individual to wear sensors (input device) namely bracelets and

gloves. Such method has an advantage that is recognizing gestures won't get distracted through the backgrounds of diverse, but, has bulkiness, trade-off of natural interaction lack, and costs high. Next, vision-related methods use capturing devices namely kinetic sensors or cameras for entering information on the basis by means of individual perception of its surrounding [6]. The precision of such methods is depending on various elements namely the total amount of cameras and their positioning, the hand visibility and by means of which it can be separated in the image, the preciseness of the extracted feature, and classifier methods [7].

This paper intends to carry out a comprehensive analysis of existing HGR based on ML and DL models. In addition, the different processes involved in the HGR models such as pre-processing, segmentation, feature extraction, and classification are explained and challenges involved are recognized. This study also offered a survey of existing methods with respect to aim, objective, technique used, and performance measures. Finally, a brief discussion of the results obtained by the reviewed methods is provided along with possible future scope.

II. BACKGROUND INFORMATION

Gesture recognition includes complicated processes namely machine learning, motion modelling, motion analysis, and pattern recognition [8]. It comprises techniques with manual and non-manual parameters. The environment structure includes speed of movement and background illumination affects the prediction capacity. The difference in viewpoint causes the gesture to look different in two dimensional space. In this study, signer wears coloured glove or wrist band to assist the hand segmentation method [9]. The usage of coloured gloves decreases the difficulty of segmentation method. Many predicted challenges in dynamic gesture recognition, complexity, involves spatial temporal variance, movement epenthesis, connectivity, and repeatability in addition to various attributes like region of gesture and change of orientation performed [10]. There are many assessment criteria for measuring the efficiency of a gesture recognition technique in addressing the problems. These criteria include robustness, scalability, userindependent, and real-time performance.

Process involved in hand gesture recognition

In general, the procedure of gesture recognition is classified into pre-processing, data acquisition, feature extraction, classification, and segmentation as demonstrated in Fig. 1. The input of static gesture recognition is single frame of images, whereas dynamic sign language takes video, viz., continuous frame of image as input. Vison-based approach differs from sensor-based approach primarily by the data-acquisition technique. The techniques and methodologies utilized by vision-based gesture recognition studies. **Data acquisition:** In vision-based gesture recognition, the data attianed was frame of images [11]. The input of these systems is gathered by image capturing devices namely webcam, video camera, thermal camera, stereo camera, or active technologies namely LMC and Kinect. LMC, Kinect, and Stereo cameras are three dimensional cameras that may gather detailed data.

Image pre-processing: This phase is implemented for modifying the video or image input datasets to enhance the entire performance of the system [12]. Gaussian and Median filters are the widely employed methods for reducing noise in attained video or image. Here, median filter is employed in this phase [13]. Then, morphological operation is commonly applied to removing redundant



data. In several studies, the captured images are reduced to a small resolution previous to

Fig.1 Process involved in HGR

Succeeding stage. This approach is utilized has shown that decreasing the resolution of input images is capable of improving the computation efficacy. Histogram equalization is utilized to improve the contrast of input images captured in distinct environments to uniform the illumination and

Brightness of the image [14].

Segmentation: This phased process the partitioning image into various parts. It is a phase where the Region of Interest (ROI) was segmented in the residual images [15].This technique can be contextual or non-contextual. Contextual segmentation considers the spatial relationships among features, namely edge detection technique [16]. While a non contextual segmentation doesn't take spatial relationships into account, however, group pixels is depending on global attributes.

Feature extraction: It is a conversion of input dataset into set of compact feature vectors [17]. In gesture recognition contexts, the feature extracted must include necessary data from the HG input and characterize from the compact form that functions as an identity of gestures that categorized excepting other gestures [18]. Classification: it is classified into supervised and unsupervised ML technologies. The supervised ML approach teaches the system to identify specific pattern of input dataset that is later utilized for predicting the upcoming dataset. Supervised ML presents a set of recognized training datasets and it is utilized for inferring a function from labelled training datasets [19]. An unsupervised ML is utilized for drawing inference from dataset with input dataset with non-labelled response [20]. Because non-labelled response is fed into the classification, where there is no penalty weightage or reward to the data belongs.

Overview of CNN

As one of DL models developed for image recognition, CNN basically comprises the convolution layer, input layer, fully connected layer, and pooling layer. The convolutional layer comprises of convolutional kernel filter. The kernel filter convolves with the child node of the input layer and outputs the result. CNN convolution operation was mathematically expressed in the subsequent formula:

 $Z_1 = f(Z_{l-1} * W_l + b_l)$ (1)

Whereas Z_{l-1} represent the input feature of the convolutional kernel of lth layer, Z_1 denotes the output feature of the convolutional kernel lth layer, b_1 indicates the bias vector of lth layer of the convolution kernel, f represents the activation function, W_1 implies the weight vector of lth layer of the convolutional kernel. Afterward the convolution layer, usually there is a pooling layer. This layer conducts the reduction dimension of the feature when maintaining the spatial dataset. Lastly, an FC layer is utilized for mapping the abstract feature dataset extracted by each preceding layer in the training method to the sample marker space. Afterward, training, the likelihood distribution Y of the original input image dataset is attained that is given in the following equation:

yi = P(L = li | X; (W,b)) (2)

Here Yi indicates the likelihood distribution of ith layer, X represents the raw input feature. b and W denote the deviation and weight matrixes, correspondingly. L denotes the loss function, li represent the loss function of ith layer P signifies the probability. The arithmetical model of CNN is to conduct linear and non-linear conversion of the matrix Z0 of the initial input image, and map them to other new dimension space Y. The training method of the CNN is to reduce the model cost function (W,b). Furthermore, the gradient descent algorithm is broadly adapted to enhance the cost equation. By using gradient descent, the model loss value is propagated, as well as the parameter in the model can be

$$W_{i} = W_{i} - \eta \frac{\partial L(W, b)}{\partial W_{i}}$$
(3)
$$b_{i} = b_{i} - \eta \frac{\partial L(W, b)}{\partial b_{i}}$$
(4)

upgraded layer by layer in the following:

Now η indicates learning rate that implies the pace of updating parameter. A CNN has contributed to the massive improvement in the network architecture, especially the model

depth, for the model to deepen, and more feature datasets are extracted, thus improving the detection performance.

III. REVIEW OF HAND GESTURE RECOGNITION

MODELS

In this section, the recently developed HGR models available in the literature are reviewed as shown in Table 1. Xu et al. [21] presented a concatenate feature fusion and recurrent convolutional neural network (CFFRCNN) approach. Here, 2-stride convolution and max-pooling layers are concatenated together for replacing the reduction dimension. The feature wise pooling process performs as a signal amplitude detector with no utilizing variables. The feature-mixing convolutional process estimates the context data. In Senturk and Bakay [22], HG information taken from UCI2019 EMG data attained in myo Thalmic armband were categorized by 6 distinct ML approaches. SVM, ANN, NB, KNN, RF, and DT techniques are employed for comparison on the basis of performance metrics. Dong et al. [23] designed an efficient and lower-cost dataset using hardware framework for simultaneously capturing movement of the bending and finger. Next, a dynamic HGR approach (DGDL-GR) is presented for recognizing dynamic sign languages, where a generic temporal convolutional network (TCN) and fusion method of convolution neural network (fCNN) is completely employed. The fCNN (combination of 1D and 2D CNN) is presented for extracting spatial domain feature of finger resistance bending and time-domain feature of finger resistance movement instantaneously. Can et al. [24] developed a DL technique based on CNN to identify HGs which improves training time, testing time, and detection rate. The presented technique involves data augmentation to increase training. Moreover, five common DL techniques are utilized for transfer learning, such as ResNet50. VGG19, VGG16, InceptionV, and DenseNet1213.

Benitez-Garcia et al. [25] examined a novel benchmark dataset called IPN Hand with appropriate variety, size, and real-world elements capable of training and evaluating deep neural networks (DNN). It can be particularly assumed that condition if the continuous gestures were executed with no transition state and if the subject execute natural movement with its hands as nongesture action. The gestures are gathered in 30 varied scenes, with real world differences from background and illuminations.

In Bakheet and Al-Hamadi [26], a new structure for real-time static HGR was established, dependent upon an optimizing shape representation created in several shape cues. The hybrid multi-modal descriptor which combines several affine-invariant boundaries as well as region-based feature is generated in thebe capitalized except for short minor words as listed in Section III-B. hand silhouette for obtaining a dependable and representative portrayal of individual gesture.

At last, an ensemble of one-vs.-all SVMs are individually trained on every of these learned feature representations for performing gesture classifier.

In Alam et al. [27], a unified method of egocentric HGR and fingertip recognition is established. The presented technique utilizes a single CNN for

predicting the probability of finger class and places of fingertips from one forward propagation. Rather than directly regressing the places of fingertip in the fully connected (FC) layer, the ensemble of places of the fingertip is regressed in the FCN. Qi et al. [28] presented a new multi-sensor guided HGR model for surgical robot teleoperation. The multi-sensor data fusion method was planned to execute interference from the occurrence of occlusions. A multi-layer RNN containing LSTM and dropout layer (LSTM-RNN) was presented to several HG classifiers.

In Gadekallu et al. [29], a crow search based CNN technique was executed from the gesture detection relating to the HCI domain. The HG data set utilized during the case is a openly accessible one, downloaded in Kaggle. During this case, a one-hot encoder approach was utilized for converting the categorical data value to binary procedure. It is subsequently the execution of a crow search algorithm (CSA) to choose optimum hyperparameter to train of dataset utilizing the CNNs.

Nayak et al. [30], proposed a Lightboost based GBM (LightGBM) for effective HGR. The hyperparameter of the presented method is enhanced by a memetic firefly algorithm. Then, we presented a perturbation operator and integrated them into the presented memetic firefly algorithm to prevent the local optimum solution in the conventional firefly algorithm. Huang et al. [31] proposed a real time HGR technology. Because fingers are the major evidence for classifying HG. a fingeremphasised multiscale descriptor has been introduced. The presented descriptor integrates three kinds of parameters of multiscale for making a discriminatory representation. Moreover, the feature of finger can be emphasized for analyzing HG. Then, three solutions to HGR are examined by using NN, DTW, and SVM. Gao et al. [32] proposed a technique based on multi-modal data fusion and multiple scales parallel CNN for enhancing the reliability and accuracy of HGR. Firstly, data fusion is performed on the depth image of HG, the sEMG signal, and the RGB images. Next, the fused images are created for two distinct scale images using down-sampling method. Later, HGR outcomes of the presented method are integrated for obtaining the concluding HGR outcome. Lv etal.[33] presented a remote HGR scheme based DL architecture of multiple attention model CNN with sEMG energy to decode HG with remote server host. Initially, an adoptive channel weighted technique is presented on multiple channel datasets of sEMG for reducing the feature map (FM) lower contribution of sEMG, and increasing the FM of sEMG. Next, improve the shortcut by adding adoptively weighted rather than a short concatenation of FM.

Chen et al. [34] presented an interactive image segmentation technique for HGR, and widespread techniques, for example, Random walker, Interactive image segmentation,

Graphcut with geodesic star convexity were investigated. The iteration of Expectation Maximum technique learn the parameter of Gaussian Mixture approach and that was applied

for image modelling. Then, employ a Gibbs random fields to image segmentation and diminish the Gibbs Energy by

Min-cut theorem for finding the optimum classification. Rahimian et al. [35] presented the Few Shot learning-HGR (FS-HGR) architecture. FS learning is variant of domain adaptation to infer the essential output based on training observation. The presented architecture was generalized afterward seeing some observations from all the classes by integrating temporal convolution with attention mechanism. This enables the meta-learner to

Table 1	Comparison	of Recently	developed	HGR models

Reference	Year of Publication	Objective	Technique uses	Dataset used	Measures
Xu et al. 2022		To recognize HGs via fusion model	CFF and CFF- RCNN	Three sEMG databases from NinaPro database	Accuracy and training time
Senturk et al.	2021	To identify HGs using EMG data	ANN, SVM, KNN, NB, DT, and RF	UCI2019 EMG dataset	Accuracy, precision, error, sensitivity, specificity
Dong et al.	2021	To present dynamic HGR using TCN	FCNN	Own sign language dataset	Accuracy, F-score, precision, recall
Can et al.	2021	To improve HGR rate	5 transfer learning models	natural ASL images	recognition rate, training, and testing time
Benitez- Garcia et al.	2020	Develop a new IPN hand video dataset	3D CNN	IPN Hand and NVIDIA	recognition rate
Bakheet et al.	2021	To recognize static gestures in real time	Hybrid multimodal description+SVM	Real world dataset	Recognition accuracy
Alam et al.	2021	To introduce egocentric HGR model	FCN	SCUT- Ego- Gesture database	Accuracy
Qi et al.	2021	To develop multiple HGR models	LSTM-RNN	-	Accuracy
Gadekallu	2021	To develop	CNN-CSA	Kaggle	Accuracy,
er ut.		accurate HGR model		posture dataset	precision, recall, F1 score
Huang and Yang	2021	To develop RGB-D shape HGR model	multi-scale descriptor, DTW, SVM, NN	NTU and PadovaU dataset	Accuracy and time
Gao et al.	2019	To develop Multi-scale parallel CNN	CNN, data fusion	HRI dataset	Accuracy and speed
Lv et al.	2022	To introduce HGR using sEMG signals	Multi-head attention based CNN	myo dataset, myoUp dataset, and ninapro DB5	accuracy
Chen et al.	2017	To introduce interactive image segmentation approach	GMM	Grabcut, PASCAL VOC 09	Accuracy, recognition rate
Rahimian et al.	2021	To desing HGR using sEMG	FS-HGR	Ninapro database	accuracy
Tan et al.	2021	To cusomtize DL model for HGR	EDenseNet	NUS and ASL dataset	accuracy
Parvathy et al.	2021	To develop HGR model using ML	DWT. modified SURF, SVM	Sebastian Marcel static hand posture dataset	Accuracy and recognition time
Tan et al.	2021	To present CNN with SPP for HGR	dubbed CNN- SPP	NUS and ASL dataset	accuracy
Pinto et al.	2019	To present standard HGR	CNN	1.1.1.1.1.1.1.1.1.1.1.1.1.1.1.1.1.1.1.	accuracy, precision, recall, and F1 score.

ISRN+ 078_81_067470_1_0

bottleneck layer for propagating the feature to each FM from the bottleneck way, and the subsequent Conv layer smooths out the redundant feature. Variances among DenseNet and EDenseNet are distinguished, and the efficiency gain is examined in this work.

Parvathy et al. [37], proposed a vision based HGR technique with ML technique. The presented method comprises classification, segmentation, and feature extraction. The proposed method is tested and trained by Sebastian Marcel's static hand posture dataset that is accessed over the internet. Modified Speed up Robust Feature extraction and Discrete wavelet transform (DWT) techniques are utilized for extracting scale and rotation invariant key descriptors. Bag of Word (BoW) method is utilized for developing the fixed dimensional input viz. necessary for the SVM technique.

Tan et al. [38] outline a CNN incorporated with dubbed CNN–SPP, spatial pyramid pooling (SPP), for vision-based HGR. SPP mitigates the problems in traditional pooling by multiple level pooling stacked to expand the feature provide to an FC layer. Pinto et al. [39] presented an HGR technique with CNN. The process includes contour generation, morphological filter, segmentation during preprocessing, and polygonal approximation, which contributed to feature extraction. Testing and training are implemented by CNN, in comparison with architecture known in the study and with other known methods.

IV. DISCUSSION

In this section, a brief discussion of the HGR results of different models is provided. Table 2 and Fig. 2 offer a detailed accuracy analysis of various HGR models available in the literature. From the results, it can be clear that the WESF approach has accomplished poor performance with least accuracy value of 85.10%. At the same time, the EGM model has accomplished slightly enhanced outcomes with slightly increased accuracy of 89.90%. Followed by, the EGM and KM models having reached closer accuracy values of 92.90% and 93.40%. Along with that, the DHM and DWTF DWTF- ratio models have demonstrated reasonable accuracy values of 96.50% and 95.42%. At last, the DHM-KM model has

Table 2 Comparative analysis of HGR results of different methods

Methods	Accuracy (%)
EGM Model	92.90
EGM Model	89.90
WESF Model	85.10
KM Model	93.40
DHM Model	96.50
DWT-F-ratio Model	95.42
DHM-KM Model	98.80

resulted in maximum accuracy of 98.80%.



Fig. 2. Comparative analysis of HGR results of

different methods V CONCLUSION

Latest advances in CV and AI techniques have resulted in the design of several HGR models. This paper has focused on a comprehensive analysis of existing HGR based on ML and DL models. In addition, the different processes involved in the HGR models such as preprocessing, segmentation, feature extraction, and classification are explained and challenges involved are recognized. This study also offered a survey of existing methods with respect to aim, objective, technique used, and performance measures. Finally, a brief discussion of the results obtained by the reviewed methods is provided along with possible future scope. At the end of the survey, we hope that the readers were aware of the basic introduction to HGR along with recent state of art approaches, and also assist future research efforts in this field. This literature analysis demonstrated that the research area of gesture detection was scattered with various methods, datasets, and approaches with many researches opportunities for achieving a further found solution to the context of sign language. Established a dataset is consolidate, enhances, and permits optimum comparative amongst approaches. More researches utilizing these approaches are support generating a connection among the 4 important sign language detection problems, attaining an entire sign language detection model. The final but not minimal, facial expression is a vital element of the sign language which is still left out in several works. Identifying facial expression is a vital stage for conveying meaning (emotion, exclamation, question, and so on) for sentences from the sign language. Further research is required from integrating detection of hand movements and facial expression.

REFERENCES

[1] Cheok, M.J., Omar, Z. and Jaward, M.H., 2019. A review of hand gesture and sign language recognition techniques. International Journal of Machine Learning and Cybernetics, 10(1), pp.131-153.

[2] Ariesta, M.C., Wiryana, F. and Kusuma, G.P., 2018. A Survey of Hand Gesture Recognition Methods in Sign Language Recognition. Pertanika Journal of Science & Technology, 26(4).

[3] Adaloglou, N., Chatzis, T., Papastratis, I., Stergioulas, A., Papadopoulos, G.T., Zacharopoulou, V., Xydopoulos, G.J., Atzakas, K.,

International Conference on "Computational Intelligence and its applications" (ICCIA-2024) ISBN: 978-81-967420-1-0

Papazachariou, D. and Daras, P., 2020. Acomprehensive study on sign language recognition methods. arXiv preprint arXiv:2007.12530, 2(2).

[4] Neiva, D.H. and Zanchettin, C., 2018. Gesture recognition: A review focusing on sign language in a mobile context. Expert Systems with Applications, 103, pp.159-183.

[5] Nogales, R.E. and Benalcázar, M.E., 2021. Hand gesture recognition using machine learning and infrared information: a systematic literature review. International Journal of Machine Learning and Cybernetics, 12(10), pp.2859-2886.

[6] Hu B, Wang J (2020) Deep learning based hand gesture recognition and UAV fight controls. Int J Autom Comput 17(1):17–29.

[7] Ma C, Wang A, Chen G, Xu C (2018) Hand joints-based gesture recognition for noisy dataset using nested interval unscented Kalman flter with LSTM network. Vis Comput 34(6):1053–1063.

[8] Mohammed, H.I., Waleed, J. and Albawi, S., 2021, February. An Inclusive Survey of Machine Learning based Hand Gestures Recognition Systems in Recent Applications. In IOP Conference Series: Materials Science and Engineering (Vol. 1076, No. 1, p. 012047). IOP Publishing.

[9] Jiang, S., Kang, P., Song, X., Lo, B. and Shull, P.B., 2021. Emerging wearable interfaces and algorithms for hand gesture recognition: A survey. IEEE Reviews in Biomedical Engineering.

[10] Qi, W., Ovur, S.E., Li, Z., Marzullo, A. and Song, R., 2021. Multi-Sensor Guided Hand Gesture Recognition for a Teleoperated Robot Using a Recurrent Neural Network. IEEE Robotics and Automation Letters, 6(3), pp.6039-6045.

[11] Kurakin, A., Zhang, Z. and Liu, Z., 2012, August. A real time system for dynamic hand gesture recognition with a depth sensor. In 2012 Proceedings of the 20th European signal processing conference (EUSIPCO) (pp. 1975-1979). IEEE.

[12] Ren, Z., Meng, J., Yuan, J. and Zhang, Z., 2011, November. Robust hand gesture recognition with kinect sensor. In Proceedings of the 19th ACM international conference on Multimedia (pp. 759-760).

[13] Panwar, M. and Mehra, P.S., 2011, November. Hand gesture recognition for human computer interaction. In 2011 International Conference on Image Information Processing (pp. 1-7). IEEE.

[14] Ren, Z., Yuan, J., Meng, J. and Zhang, Z., 2013. Robust part-based hand gesture recognition using kinect sensor. IEEE transactions on multimedia, 15(5), pp.1110-1120.

[15] Khan, R.Z., Ibraheem, N.A. and Meghanathan, N., 2012, January. Comparative study of hand gesture recognition system. In Proc. of International Conference of Advanced Computer Science & Information Technology in Computer Science & Information Technology (CS & IT) (Vol. 2, No. 3, pp. 203-213).

[16] Lin, H.I., Hsu, M.H. and Chen, W.K., 2014, August. Human hand gesture recognition using a convolution neural network. In 2014 IEEE International Conference on Automation Science and Engineering (CASE) (pp. 1038-1043). IEEE.

[17] Suk, H.I., Sin, B.K. and Lee, S.W., 2010. Hand gesture recognition based on dynamic Bayesian network framework. Pattern recognition, 43(9), pp.3059-3072.

[18] Lu, Z., Chen, X., Li, Q., Zhang, X. and Zhou, P., 2014. A hand gesture recognition framework and wearable gesture-based interaction prototype for mobile devices. IEEE transactions on human-machine systems, 44(2), pp.293-299.

[19] Meena, S., 2011. A study on hand gesture recognition technique (Doctoral dissertation).

[20] Hsieh, C.C., Liou, D.H. and Lee, D., 2010, July. A real time hand gesture recognition system using motion history image. In 2010 2nd international conference on signal processing systems (Vol. 2, pp. V2-394). IEEE.

[21] Xu, P., Li, F. and Wang, H., 2022. A novel concatenate feature fusion RCNN architecture for sEMG-based hand gesture recognition. PloS one, 17(1), p.e0262810.

[22] Senturk, Z.K. and Bakay, M.S., 2021. Machine learning based hand gesture recognition via emg data. ADCAIJ: Advances in Distributed Computing and Artificial Intelligence Journal, 10(2).

[23] Dong, Y., Liu, J. and Yan, W., 2021. Dynamic hand gesture recognition based on signals from specialized data glove and deep learning algorithms. IEEE Transactions on Instrumentation and Measurement, 70, pp.1-14.

[24] Can, C., Kaya, Y. and Kılıç, F., 2021. A deep convolutional neural network model for hand gesture recognition in 2D near-infrared images. Biomedical Physics & Engineering Express, 7(5), p.055005.

[25] Benitez-Garcia, G., Olivares-Mercado, J., Sanchez-Perez, G. and Yanai, K., 2021, January. IPN hand: A video dataset and benchmark for real-time continuous hand gesture recognition. In 2020 25th International Conference on Pattern Recognition (ICPR) (pp. 4340-4347). IEEE.

[26] Bakheet, S. and Al-Hamadi, A., 2021. Robust hand gesture recognition using multiple shape-oriented visual cues. EURASIP Journal on Image and Video Processing, 2021(1), pp.1-18.

[27] Alam, M.M., Islam, M.T. and Rahman, S.M., 2022. Unified learning approach for egocentric hand gesture recognition and fingertip detection. Pattern Recognition, 121, p.108200.

[28] Qi, W., Ovur, S.E., Li, Z., Marzullo, A. and Song, R., 2021. Multi-Sensor Guided Hand Gesture Recognition for a Teleoperated Robot Using a Recurrent Neural Network. IEEE Robotics and Automation Letters, 6(3), pp.6039-6045.

[29] Gadekallu, T.R., Alazab, M., Kaluri, R., Maddikunta, P.K.R., Bhattacharya, S. and Lakshmanna, K., 2021. Hand gesture classification using a novel CNN-crow search algorithm. Complex & Intelligent Systems, 7(4), pp.1855-1868.

[30] Nayak, J., Naik, B., Dash, P.B., Souri, A. and Shanmuganathan, V., 2021. Hyperparameter tuned light gradient boosting machine using memetic firefly algorithm for hand gesture recognition. Applied Soft Computing, 107, p.107478.

[31] Huang, Y. and Yang, J., 2021. A multiscaledescriptor for real time RGB-D hand gesture recognition. Pattern Recognition Letters, 144, pp.97-104.

[32] Gao, Q., Liu, J. and Ju, Z., 2021. Hand gesture recognition using multimodal data fusion and multiscale parallel convolutional neural network for human–robot interaction. Expert Systems, 38(5), p.e12490.

[33] Lv, X., Dai, C., Liu, H., Tian, Y., Chen, L., Lang, Y., Tang, R. and He, J., 2022. Gesture recognition based on sEMG using multiattention mechanism for remote control. Neural Computing and Applications, pp.1-11.

[34] Chen, D., Li, G., Sun, Y., Kong, J., Jiang, G., Tang, H., Ju, Z., Yu, H. and Liu, H., 2017. An interactive image segmentationmethod in hand gesture recognition. Sensors, 17(2), p.253.

[35] Rahimian, E., Zabihi, S., Asif, A., Farina, D., Atashzar, S.F. and Mohammadi, A.,2021. Fs-hgr: Few-shot learning for hand gesture recognition via electromyography. IEEE transactions on neural systems and rehabilitation engineering, 29, pp.1004-1015.

[36] Tan, Y.S., Lim, K.M. and Lee, C.P., 2021. Hand gesture recognition via enhanced densely connected convolutional neural network. Expert Systems with Applications, 175, p.114797.

[37] Parvathy, P., Subramaniam, K., Prasanna Venkatesan, G.K.D., Karthikaikumar, P., Varghese, J. and Jayasankar, T., 2021. Development of hand gesture recognition system using machine learning. Journal of Ambient Intelligence and Humanized Computing, 12(6), pp.6793-6800.

[38] Tan, Y.S., Lim, K.M., Tee, C., Lee, C.P. and Low, C.Y., 2021. Convolutional neural network with spatial pyramid pooling for hand gesture recognition. Neural Computing and Applications, 33(10), pp.5339-5351.

[39] Pinto, R.F., Borges, C.D., Almeida, A. and Paula, I.C., 2019. Static hand gesture recognition based on convolutional neural networks. Journal of Electrical and Computer Engineering, 2019.

A Safe Human Route Prediction using Machine Learning based on Multi-Diversity Factors

T. Thilagavathi¹ and Dr. A. Subashini²

¹Research Scholar, Annamalai University, Annamalai Nagar, Tamil Nadu, India. thilagaponns@gmail.com

²Assistant Professor, Department of Computer Application, Government Arts College, Chidambaram, Tamil Nadu, India. subanandh31@gmail.com

Abstract - Finding a safe route for people living in both rural and urban areas is crucial. Human travel is increasing in real- time environments due to job, business, and vocational reasons, leading individuals to move from their place of origin to their destination. This travel experience can vary based on factors such as time, the number of people, types of transport vehicles, age, gender, road conditions, lighting facilities, traffic, etc. The existing safe route prediction algorithm operates with a minimal number of factors and fails to accurately identify a safe route for groups or individuals. The proposed machine learning algorithm incorporates multiple diverse factorssuch as time, the number of people, types of transport vehicles, age, gender, road conditions, lighting facilities, traffic, road type, location, etcto determine a safe human route, providing more accurate results for individuals or groups traveling from the source to the destination.

Keywords - Human safe routing, machine learning, route map, Risk factor.

I. INTRODUCTION

In an era where transportation is an integral part of our daily lives, the importance of safe and efficient routing cannot be overstated. Human-safe routing, particularly in the context of road conditions, has emerged as a crucial element in ensuring the well-being of both drivers and pedestrians. This innovative approach to navigation goes beyond mere efficiency, incorporating real-time data on road conditions to prioritize safety and mitigate potential hazards.

In the intricate tapestry of transportation, the conventional approach to routing has often overlooked the nuanced factors that significantly impact the safety and wellbeing of individuals on the move. Recognizing the need for a more holistic navigation paradigm, the concept of humansafe routing has evolved, weaving together considerations of time, population density, vehicle diversity, demographics, road conditions, lighting, and traffic dynamics. This multifaceted approach seeks not only to optimize routes for efficiency but, more importantly, to tailor them to the unique requirements and safety considerations of diverse individuals and communities.

Traditionally, routing algorithms primarily focused on minimizing travel time and distance. However, the concept of human-safe routing expands this perspective by acknowledging the temporal dimension. Time-sensitive routing takes into account peak hours, enabling users to avoid congested periods and select routes that optimize travel based on the specific time of day. This consideration aims to enhance the overall efficiency of transportation while minimizing the risk of accidents and stress associated with peak traffic conditions. One of the key pillars of human-safe routing is the integration of real-time information into the routing algorithms. This includes data from various sources such as weather forecasts, traffic updates, and road maintenance reports. By continuously monitoring andanalyzing these inputs, the routing system can dynamically adjust the recommended routes to avoid potential dangers, ensuring a safer journey for all road users.

Additionally, human-safe routing takes into account the diverse range of vehicles on the road, each with its own set of capabilities and limitations. For instance, the optimal route for a pedestrian may differ significantly from that of a large commercial truck. By considering the characteristics of different modes of transportation, the routing system aims to tailor its recommendations to suit the specific needs and safety requirements of each user.

Beyond temporal factors, the number of people traveling and the types of vehicles on the road play pivotal roles in human-safe routing. An inclusive approach recognizes the varied needs of pedestrians, cyclists, motorists, and public transport users. By adapting routes based on the composition of traffic, the system can foster a harmonious coexistence of diverse transportation modes, ensuring the safety and convenience of all road users.

Demographic considerations, such as age and gender, further enrich the fabric of human-safe routing. Routes may be tailored to accommodate the mobility challenges of different age groups, ensuring accessibility for the elderly and young pedestrians alike. Recognizing the unique safety concerns of various gender identities, the routing system can offer routes that prioritize well-lit areas and areas with a strong presence of public spaces.

Moreover, the interplay of road conditions, lighting facilities, and traffic dynamics profoundly influences the safety of navigation. Advanced routing systems take realtime data into account, dynamically adjusting routes based on weather conditions, road maintenance, and the availability of well-lit areas. The system also considers the intricacies of traffic flow, adapting recommendations to minimize congestion and enhance overall safety.

Furthermore, the rise of connected vehicles and the Internet of Things (IoT) has opened up new possibilities for human-safe routing. Vehicles equipped with sensors and communication capabilities can provide real-time updates on road conditions, creating a dynamic network of information sharing. This interconnected ecosystem allows for more accurate and timely adjustments to routing recommendations, contributing to a collective effort to enhance overall road safety. In conclusion, the paradigm of comprehensive humansafe routing emerges as a beacon for the future of transportation. By intertwining considerations of time, population density, vehicle types, demographics, road conditions, lighting, and traffic dynamics, this approach envisions a transportation landscape that not only optimizes efficiency but prioritizes the safety, accessibility, and wellbeing of all individuals navigating our diverse urban and rural environments.

II. LITERATURE SURVEY

Safe human routing prediction is a critical aspect of intelligent transportation systems, aimed at ensuring the safety and well-being of individuals navigating through various environments. As urbanization and population growth continue to rise, the demand for efficient and secure human mobility becomes increasingly paramount. This has led to the development of advanced technologies and methodologies to predict and optimize human routing while prioritizing safety.

Machine learning algorithms and artificial intelligence techniques play a crucial role in processing this vast amount of data to predict safe routes for pedestrians. By leveraging historical data, real-time inputs, and predictive analytics, these systems aim to not only optimize travel time but also minimize the likelihood of accidents and ensure a secure journey for individuals.

The literature on safe human routing prediction spans across various domains, including computer vision, machine learning, transportation engineering, and urban planning. Researchers have explored different approaches such as deep learning models, reinforcement learning, and hybrid methods that combine multiple data sources to improve the accuracy and reliability of route predictions. Additionally, studies delve into the integration of intelligent transportation systems with urban infrastructure to create a seamless and safe experience for pedestrians.

Several challenges persist in the field, such as dealing with dynamic and unpredictable human behavior, ensuring robustness in adverse weather conditions, and addressing privacy concerns associated with the collection and processing of personal data. Understanding the intricate interplay between environmental factors, pedestrian behavior, and route optimization is crucial for developing effective and safe routing prediction models.

Examining the correlation between accidents, road conditions, and weather conditions, Lingamaneni Indraja [1] et al developed a predictive model to identify a secure and less accident-prone route among the available paths from the source to the destination. Employing machine learning algorithms such as Support Vector Machines and Logistic Regression are constructed this model.

Yash S. Asawa et al [2] introduced a User-Specific Safe Route Recommendation System that provides users with a visual representation of a secure route on maps, taking into account the historical criminal records of the geographical region. Our methodology operates on two tiers: the initial phase involves capturing user-specific features through a Decision Network, while the subsequent phase employs Geospatial Data Analysis to facilitate the generation of a safe route. Aruna Pavateet al [3] predicted a secure route and outlined precautionary measures for women to address potential crime situations. They displayed the route along with the crime rate, enabling women to avoid specific places with a higher likelihood of crimes. The authors developed a system utilizing the K-means clustering algorithm to achieve precise results. The system categorized routes into various levels of security, including best, better, average, less secure, and least secure routes. Their research recommended safe routes, specifying the types of crimes that occurred along each route, allowing individuals to choose the safest route from the source in Delhi city.

Isha Puthige et al. [4] devised a danger index for each data point, determined by multiple crime factors at specific latitude and longitude values. The authors employed various clustering algorithms to analyze the most effective safe paths. AliasgarEranpurwala et al. [5] created a Women Safety Application for Safe Route Prediction using crime events in Delhi city. The 'GoWomaniya' application, developed by the authors, played a significant role in providing women with a safe environment during moments of torment and distress.

Deepa Bura et al. [6] developed a model for predicting a secure and safe route for women using Google Maps. The aim was to find a safe route based on factors such as security, risk, and the quality of the path.Roxan Salehab et al. [7] studied a model to predict the status of road signs mounted on Swedish roads using supervised machine learning.

Deepak Kumar Sharma [8] et al used a collection of historical traffic accident data, including geo-location, time, location, weather conditions, road conditions, and other pertinent aspects, such asvehicle type age, gender, time of day, weather, and so forth. The Random wood Classifier algorithm was used to accurate and accurately forecast the danger of automobile crashes.

Juncai Jiang et al. [9] introduced a novel framework for assessing the risk of multifactorial urban road collapse. This framework utilized a combination of environmental and anthropogenic factors to establish an indicator system for risk assessment. Subsequently, the authors employed SMOTE for data augmentation on collected accident samples, creating a dataset of urban road collapse incidents used to train the CNN model.

III. METHODOLOGY

Human-safe routing involves integrating advanced technologies, data analytics, and user-centric considerations into a comprehensive system that prioritizes safety, efficiency, and user experience.

By intertwining considerations of time, population density, vehicle types, demographics, road types, road conditions, lighting, weather conditions, public spaces, and traffic dynamics, this approach envisions a transportation landscape that not only optimizes efficiency but prioritizes the safety, accessibility, and well-being of all individuals.Figure 1 illustrates the procedures involved in finding safe human routing.

The process of safe route prediction relies on a comprehensive analysis of both static and dynamic factors to ensure the well-being of individuals during their journeys. Leveraging a wealth of known data, including vehicle types and capacities, the number of persons traveling, the age distribution of individuals, road types, and other static

variables, as well as real-time data streams such as traffic conditions, weather updates, lighting facilities, time of travel, and the presence of public spaces, the safe route prediction system employs a sophisticated approach to navigate through diverse urban landscapes.



FIG 1. PROCEDURES FOR SAFE HUMAN ROUTE FINDING

By considering the inherent characteristics of different vehicles and understanding their capacities, the system tailors routes to accommodate the specific needs and safety requirements of various transportation modes. The number of persons traveling is a crucial factor in determining the optimal routes, ensuring that the navigation plan aligns with the collective well-being of the passengers.

The age distribution of individuals traversing the routes is taken into careful consideration, recognizing the unique mobility challenges and safety concerns associated with different age groups. This enables the system to recommend routes that are not only efficient but also tailored to the specific requirements of children, elderly individuals, and other demographic segments.

Incorporating real-time data into the predictive models enhances the system's adaptability, allowing it to dynamically adjust routes based on the ever-changing conditions of the urban environment. Monitoring live traffic conditions and weather updates ensures that the system can proactively reroute individuals away from potential hazards, optimizing for both safety and efficiency.

The inclusion of environmental factors, such as the availability of lighting facilities along the route, contributes to the creation of safer routeways, particularly during lowvisibility periods. Additionally, considering the time of travel helps the system account for variations in traffic patterns, congestion levels, and potential safety risks associated with specific times of the day. The road conditions involve the existence of curves, the presence of bridges, the existence of road crossings, and other factors.

Public spaces play a pivotal role in fostering safety and creating pleasant travel experiences. By factoring in the presence of well-designed public spaces, the system not only prioritizes safety but also contributes to the overall quality of the journey.

Finding the risk factors for safe human routing involves analysing a combination of static and dynamic data to identify potential hazards and vulnerabilities in a given route. The risk factors are identified for the available routes from the source to the destination.

The determination of the safest route is derived from a comprehensive estimation process that meticulously takes into account and analyses the myriad risk factors associated with the available routes. Through a careful examination of factors such as traffic conditions, road infrastructure, weather patterns, and historical incident data, the route estimation algorithm systematically assesses each potential route from the source to the destination. This method ensures that the resulting safe route not only optimizes for efficiency but also prioritizes the well-being and security of individuals navigating through diverse and dynamic urban and rural environments.

IV. IMPLEMENTATION

By combining these elements, a comprehensive risk assessment can be conducted, allowing for the development of human-safe routing systems that prioritize safety and minimize potential hazards. Regular updates and feedback mechanisms are essential to ensuring the ongoing effectiveness of the risk assessment process.

A. List of Parameters

Table 1 summarizes the static and dynamic data parameters employed in our methodology.

Name of the Parameter	Symbolic Representation
Vehicle Type / Transport mode	v
Gender	g
Age	а
Number of Persons travelling	pn
Lighting facility	L
Road types	Rt
Traffic condition	Тс
Weather condition	Wc
Public spaces existence	Ps
Population density	D
Road Conditions	Rc

TABLE 1. STATIC AND DYNAMIC DATA PARAMETERS

The values of parameters range from 0 to 1. A value closer to 1 indicates a highimpact on the parameter, while a value closer to 0 indicates low impact. All the parameters are included in the calculation of risk factor of the route.

B. Risk factor calculation

Predicting safe routes involves assessing various risk factors to determine the likelihood of encountering hazards or dangers along a given route. The risk factor calculation involves the static factors (parameters) and dynamic parameters.

i. Risk Factor calculation

The final risk factor of the route is calculated by averaging the static and dynamic risk factors and given in equation (1).

$$RF = \frac{S_{RF} + D_{RF}}{2} \tag{1}$$

C. Finding the available routes

For the given source and destination, the available routes are determined. All the routes are stored in a list, denoted as R. Each route, denoted as R_i , is considered as a graph, denoted as G_i . In our methodology, each route is divided into 500m segments, denoted as Si. Each segment point is referred to as a vertex (V), and an edge (E) is referred to as a road. The risk factor of each segment $RF(S_i)$ is calculated using equation (4). The risk factor of each route is calculated by summing up the risk factors of all segments, as given in equation (2).

$$RF_i = \frac{\sum_{i=1}^n RF(S_i)}{n} \tag{2}$$

The RF_i denotes the risk factor of the ith route from the source to destination.

D. Finding the Safe Route

The safe route, denoted as SR, is identified from the available route risk factors by selecting the route with the minimum risk factor, as given in equation (3).

$$SR = min(RF_i) \tag{3}$$

The safe route involves all the static and dynamic data to find the risk factor.

E. Final prediction calculation

The "Final Prediction" column represents the hypothetical final prediction based on the combination of SVM and LR predictions, along with dynamic (DRF) and static (SRF) risk factors.

Final Prediction= $(0.4 \times \text{SVM Prediction})+(0.3 \times \text{LR Predict ion})+(0.2 \times \text{DRF})+(0.1 \times \text{SRF})$ (4)

Final Prediction is a weighted combination of SVM and LR predictions, along with DRF and SRF.

V.RESULT AND DISCUSSION

The study involved an in-depth analysis of both static and dynamic parameters essential for determining the risk factor of a given route. Static parameters such as vehicle type, gender, age, lighting facility, road types, public spaces existence, and road conditions were meticulously examined. Additionally, dynamic parameters including traffic condition, weather condition, and population density were considered. The values of these parameters ranged from 0 to 1, with higher values indicating a greater impact on the parameter.

The table 2 contains "SVM Prediction" and "LR Prediction" represents the outputs of the SVM and LR models, respectively, ranging from 0 to 1 for each data point in table 1. This SVM and LR prediction will be further used for predicting the safest route. The table 2 summarizes the performance metrics of the SVM and LR models

Data Point	Actual Label	SVM Prediction	LR Prediction
1	1	0.8	0.7
2	0	0.2	0.3
3	1	0.9	0.8
4	0	0.1	0.2
5	1	0.7	0.6
6	0	0.4	0.5
7	1	0.85	0.9
8	0	0.3	0.2
9	1	0.6	0.7
10	0	0.15	0.1
11	1	0.75	0.8
12	0	0.25	0.3
13	1	0.95	0.9
14	0	0.2	0.25
15	1	0.88	0.87
16	0	0.35	0.4
17	1	0.78	0.75
18	0	0.3	0.28
19	1	0.65	0.68
20	0	0.18	0.2

TABLE 3 SVM AND LR RESULTS

	SVM	LR
Precission	87	80
F1	77	84
Recall	70	88
Accuracy	75	70



FIG 2: EVALUATION MATRICES FOR SVM AND LR

VI.

TABLE 4 DRF, SRF, SVM, LR AND FINAL PREDICTION FOR EACH DATA POINT

Data Poin t	Actua l Label	SVM Predictio n	LR Predictio n	DR F	SR F	Final Predictio n
1	1	0.8	0.7	0.6	0.7	0.72
2	0	0.2	0.3	0.4	0.3	0.27
3	1	0.9	0.8	0.7	0.8	0.83
4	0	0.1	0.2	0.3	0.4	0.19
5	1	0.7	0.6	0.8	0.7	0.67
6	0	0.4	0.5	0.2	0.6	0.43
7	1	0.85	0.9	0.9	0.8	0.87
8	0	0.3	0.2	0.4	0.4	0.29
9	1	0.6	0.7	0.5	0.7	0.63
10	0	0.15	0.1	0.6	0.3	0.24
11	1	0.75	0.8	0.7	0.8	0.77
12	0	0.25	0.3	0.2	0.4	0.28
13	1	0.95	0.9	0.8	0.9	0.91
14	0	0.2	0.25	0.4	0.4	0.25
15	1	0.88	0.87	0.7	0.8	0.86
16	0	0.35	0.4	0.5	0.7	0.39
17	1	0.78	0.75	0.8	0.7	0.77
18	0	0.3	0.28	0.4	0.3	0.31
19	1	0.65	0.68	0.6	0.6	0.65
20	0	0.18	0.2	0.2	0.2	0.19



FIG 3: DRF, SRF, SVM, LR AND FINAL PREDICTION FOR EACH DATA POINT

The static parameters, including vehicle type, gender, age, lighting facility, road types, public spaces existence, and road conditions, were carefully examined. Each parameter was assigned a value ranging from 0 to 1, with higher values indicating a greater impact on the parameter. The dynamic parameters considered in the study included traffic condition, weather condition, and population density. This table shows the calculated dynamic risk factor (DRF), static risk factor (SRF), and the final prediction for each data point.

The analysis involved using SVM and LR models to predict the risk associated with each route based on the given parameters. The actual labels were assigned as 1 for safe and 0 for not safe. The SVM and LR models provided predictions for each data point, and these predictions were further utilized in the calculation of the dynamic risk factor (DRF) and static risk factor (SRF).

The final prediction for each data point was obtained by combining the predictions from SVM, LR, DRF, and SRF, each weighted accordingly. The values range from 0 to 1, providing an integrated assessment of the risk associated with each route. So as per the final prediction the data point 3, 7, 15 are considered safe route and the point 20, 18, 12, 10, 8, 4, 2 are unsafe route and rest being moderate.

CONCLUSION

The present study offers a comprehensive analysis of route safety, integrating both static and dynamic parameters through the application of Support Vector Machines (SVM) and Logistic Regression (LR) models. The final predictions, combining SVM, LR, Dynamic Risk Factor (DRF), and Static Risk Factor (SRF), provide a nuanced evaluation of the risk associated with each route. The SVM and LR models demonstrate reliable performance, as evidenced by respectable precision, recall, F1 score, and accuracy. The study emphasizes the importance of both dynamic and static parameters in route safety assessment. DRF and SRF highlight the influence of each model in shaping the final risk prediction. The absence of crime reports in the current work presents an opportunity for future research. Integrating data could enhance the accuracy crime and comprehensiveness of risk assessment. The models should be subjected to further validation using diverse datasets to ensure robustness and generalizability. Future iterations of this research should incorporate crime reports to capture additional dimensions of risk, contributing to a more holistic safety evaluation. Implementing real-time data sources for dynamic parameters could enhance the timeliness and relevance of route safety assessments. In conclusion, the integrated approach presented in this study provides valuable insights into route safety. The findings serve as a foundation for future work, especially in incorporating crime data and refining the models for more accurate and real-world applications.

VI. REFERENCES

- L. Indraja and D. D. Suneetha, "Safe Path Prediction Using Machine Learning," International Journal for Research in Applied Science & Engineering Technology, pp. 1771-1773, 2023.
- [2]. Y. S. Asawa, S. R. Gupta, V. V and N. J. Jain, "User Specific Safe Route Recommendation System," International Journal of Engineering Research & Technology, vol. 9, no. 10, pp. 574-580, 2020.
- [3]. Pavate, A. Chaudhari and R. Bansode, "Envision of Route Safety Direction Using Machine Learning," ACTA SCIENTIFIC MEDICAL SCIENCES, vol. 3, no. 11, pp. 140-

145, 2019.

[4]. Puthige, K. Bansal, C. Bindra, M. Kapur, D. Singh, V. K. Mishra, ApekshaAggarwal, J. Lee, B.-G. Kang, Y. Nam and R.

R. Mostafa, "Safest Route Detection via Danger Index Calculation and K-Means Clustering," Computers, Materials & Continua, vol. 69, no. 2, pp. 2761-2777, 2021.

- [5]. Eranpurwala, F. Indorewala, N. Mapari and S. Mishra, "Women Safety Application for Safe Route Prediction," International Research Journal of Engineering and Technology, vol. 8, no. 5, pp. 2278-2282, 2021.
- [6]. Bura, M. Singh and P. Nandal, "Predicting Secure and Safe Route for Women using Google Maps," International Conference on Machine Learning, Big Data, Cloud and Parallel Computing (COMITCon), Faridabad, India,, pp. 103-108, 2019.
- [7]. R. Salehab and H. Fleyehb, "Using Supervised Machine Learning to Predict the Status of Road Signs," Transportation Research Procedia, vol. 62, pp. 221-228, 2022.

PROCEEDINGS

- [8]. D. K. Sharma and P. U. M, "The Traffic Accident Prediction Using Machine Learning," International Research Journal of Modernization in Engineering Technology and Science, vol. 5, no. 7, pp. 2964-2968, 2023.
- [9]. J. F. Wang, Y. Wang, W. Jiang, Y. Qiao and W. B. X. Zheng, "An Urban Road Risk Assessment Framework Based on Convolutional Neural Networks," International Journal of Disaster Risk Science, vol. 14, pp. 475-487, 2023.

Unmasking Deception: Stance Detection in the Age of COVID-19 Fake News

C. Justin Marshal[#], Dr. R. Vidya^{*}

[#]Research Scholar & Assistant Professor, *Assistant Professor,

^{#*}PG Department of Computer Applications, St. Joseph's College of Arts and Science (Autonomous), Cuddalore, TamilNadu,

jusmarshal@gmail.com

Abstract— The COVID-19 pandemic, an unprecedented global crisis, has been accompanied by an infodemic of misinformation and fake news. In this context, the accurate detection of stances expressed in news and social media content is paramount for discerning truth from deception. This research, titled "Unmasking Deception: Stance Detection in the Age of COVID-19 Fake News," addresses the critical need for effective stance detection in the realm of pandemicrelated information dissemination. This paper presents a comprehensive investigation into the methodology, data, and findings of stance detection during the COVID-19 era. Leveraging a diverse dataset of textual information, our research employs state-of-the-art machine learning and natural language processing techniques to identify and categorize stances expressed in news articles and social media posts. The study encompasses a detailed analysis of these stances, unveiling the underlying dynamics and trends in the spread of fake news during the pandemic. The results reveal invaluable insights into the propagation and reception of fake news related to COVID-19, shedding light on the methods and platforms through which misinformation proliferates. The implications of this research extend to public health, media literacy, and the development of effective strategies to combat misinformation during crises. By unmasking deception and enhancing our understanding of how stances are constructed, this study contributes to the broader endeavour of promoting accurate information dissemination and safeguarding public well-being in an age marked by both a global health crisis and an infodemic.

Keywords— Covid, fake, Misinformation, Proliferation, Stance

I. INTRODUCTION

The emergence of the COVID-19 pandemic in late 2019 brought about not only a global health crisis but also an unprecedented wave of disinformation and fake news. The significance of fake news during this pandemic cannot be overstated. Misinformation has become a parallel pandemic, threatening public health, undermining trust in authorities, and sowing confusion among populations. In this introduction, we will explore the multifaceted significance of fake news during the COVID-19 pandemic, supported by references to key studies and reports.

Public Health Consequences

The spread of fake news during a public health crisis like COVID-19 poses serious risks to individuals and

communities. Misinformation about cures, prevention strategies, and the virus's origins can lead to harmful behaviours and undermine public health efforts. A study by Pennycook and Rand(2020) [1]in science found that exposure to misinformation was associated with lower compliance with recommended health behaviours.

Erosion of Trust in Public Health Authorities

Trust in public health authorities and institutions are essential during a pandemic. However, the infodemic has eroded public trust. A report from the Reuters Institute for the Study of Journalism (2020) highlighted that false and misleading information often spread faster and more widely than accurate information during the COVID-19 pandemic, causing confusion and mistrust.

Impact on Vaccine Hesitancy

The role of misinformation in vaccine hesitancy has been a growing concern. A study by Roozenbeek et al. (2020) [2] in Nature Human Behaviour found that exposure to vaccine misinformation reduced individuals' intentions to get vaccinated. This has implications for achieving herd immunity and ending the pandemic.

Social Division and Polarization

Fake news can exacerbate existing social divisions and political polarization. During the pandemic, the spread of misinformation became entangled with political and ideological beliefs. A report from the Pew Research Center (2020) highlighted the stark partisan divides in beliefs about COVID-19 and related policies.

Economic Consequences

The economic impact of fake news during the pandemic is substantial. False information about the virus's severity and government responses can affect financial markets and consumer behaviour. A study by Allcott et al. (2020) [3]in the National Bureau of Economic Research demonstrated how misinformation led to decreased demand for goods and services.

Challenges for Media Literacy

The proliferation of fake news underscores the need for media literacy and critical thinking skills. A study by Guess et al. (2020)[4] in the Proceedings of the National Academy

India

of Sciences found that people with higher levels of media literacy were less likely to believe COVID-19 misinformation.

In light of these consequences, it is evident that addressing fake news during the COVID-19 pandemic is not only a matter of information accuracy but a crucial component of public health and societal stability. This paper, "Unmasking Deception: Stance Detection in the Age of COVID-19 Fake News," seeks to contribute to the broader effort of combating misinformation and its significant implications during this critical period.

II. RESEARCH PROBLEM

The research problem addressed in "Unmasking Deception: Stance Detection in the Age of COVID-19 Fake News" revolves around the critical challenge of identifying and categorizing stances within the vast landscape of information related to the COVID-19 pandemic. This problem stems from the broader issue of the infodemic – the rapid proliferation of fake news and misinformation during the pandemic.

The Core Problem

The COVID-19 pandemic has not only brought to the forefront the importance of accurate information but has also created an environment where fake news and misinformation spread easily and rapidly. Misinformation can range from false health advice, baseless conspiracy theories about the virus's origin, the effectiveness of treatments, or the safety of vaccines, to distorted political narratives, which further polarize societies.

The central problem addressed in this research can be articulated as follows:

How can we effectively detect and categorize the stances expressed within the massive volume of COVID-19-related information, particularly in the context of misinformation and fake news, to better understand the dynamics and implications of this infodemic?

Key Components of the Problem

Stance Detection: This research aims to uncover the various stances people take within the context of COVID-19, whether it's their perspective on the origin of the virus, the effectiveness of health measures, or their opinions on vaccines. Stance detection is crucial in disentangling the web of conflicting information and understanding the diverse attitudes and beliefs within the population.

Misinformation and Fake News: A significant portion of the information related to COVID-19 is false or misleading. Detecting these instances of misinformation is essential to curb its impact on public health and public perception. Stance detection should include the identification of stances expressed in fake news and conspiratorial narratives.

Information Propagation: Understanding how stances and misinformation spread through various media channels and social networks is another facet of the problem. This includes investigating the platforms, individuals, or groups responsible for amplifying certain stances and false information. Societal and Health Implications: The stances and misinformation surrounding COVID-19 have serious implications for public health, policy decisions, and societal cohesion. The research must examine how the detected stances and misinformation affect public behavior, trust in institutions, and responses to the pandemic.

Significance of Addressing the Problem: The ability to effectively detect and categorize stances expressed in the context of COVID-19 fake news and misinformation is of paramount significance. Addressing this problem can:

- 1. Enable a deeper understanding of the factors contributing to the spread of misinformation during crises.
- 2. Inform more targeted public health communication strategies.
- 3. Contribute to the development of algorithms and tools for real-time monitoring and response to misinformation.
- 4. Aid in the creation of policies and interventions to mitigate the impact of misinformation on public health and societal trust.

In summary, this research problem seeks to shed light on the multifaceted challenge of detecting and categorizing stances within the immense volume of information related to COVID-19, with a specific focus on the pernicious influence of fake news and misinformation. The ultimate goal is to unmask deception, improve information integrity, and bolster the response to public health crises like the COVID-19 pandemic.

III. REVIEW OF LITERATURE

The existing literature on fake news, misinformation, and their impact during crises, particularly the COVID-19 pandemic, provides valuable insights into the nature of this problem and its implications for public health, society, and information ecosystems. Here's a discussion of some key findings and themes from this body of research:

1. Proliferation of Misinformation and Fake News

Research consistently highlights the rapid spread and proliferation of misinformation and fake news during crises. The COVID-19 pandemic is no exception. Studies by Vosoughi et al. (2018)[5]and Kouzy et al. (2020)[6]underscore how false information can spread faster and more broadly than accurate information, especially on social media platforms. The infodemic during the COVID-19 pandemic was characterized by an overwhelming volume of unverified claims, rumors, and conspiracy theories (Pennycook & Rand, 2020).

a. Impact on Public Behavior and Health

The impact of misinformation on public behavior and health during crises is a significant concern. Bessi et al. (2016) [7] found that exposure to false information can affect individual behavior and decision-making. Misinformation during the COVID-19 pandemic has led to individuals engaging in risky behaviors, ignoring public health guidelines, and even hoarding unproven remedies (Garrett, 2020)[8].

b. Trust in Authorities and Information Sources

Misinformation erodes trust in public health authorities and institutions. The Reuters Institute for the Study of Journalism (2020)[9] reported that false and misleading information often led to confusion and mistrust in governments, health organizations, and traditional news sources. This lack of trust can have severe consequences for effective crisis management.

c. Infodemic as a Global Challenge

The infodemic is not limited to a specific region or language; it is a global challenge. Research by Friggeri et al. (2014) [10] on Facebook and Vosoughi et al. (2018) on Twitter demonstrated the international scope of misinformation. During the COVID-19 pandemic, misinformation transcended borders, making it a shared challenge worldwide (Pennycook & Rand, 2020).

d, Impact on Vaccine Hesitancy

The relationship between misinformation and vaccine hesitancy is of growing concern, particularly during the COVID-19 pandemic. Studies by Roozenbeek et al. (2020) and Callaghan et al. (2020) [11] highlighted that exposure to vaccine misinformation reduced individuals' intentions to get vaccinated, threatening public health efforts to achieve herd immunity.

e. Socio-Political Dimensions

Misinformation during crises often becomes intertwined with socio-political factors. The Pew Research Center (2020) reported stark partisan divides in beliefs about COVID-19 and related policies in the United States, indicating the potential for misinformation to exacerbate existing societal divisions.

f. Role of Media Literacy

Research by Guess et al. (2020) found that individuals with higher levels of media literacy were less likely to believe COVID-19 misinformation. This emphasizes the importance of media literacy education as a preventive measure against misinformation during crises.

g. Challenges in Countering Misinformation

Countering misinformation poses significant challenges. The rapid spread of fake news and the difficulty of factchecking in real-time present formidable obstacles for organizations and fact-checkers (Pennycook & Rand, 2020). Additionally, the echo chamber effect and confirmation bias make it challenging to change people's beliefs once they have been exposed to misinformation (Nyhan & Reifler, 2010)[12].

In conclusion, the existing literature on fake news, misinformation, and their impact during crises, with a specific focus on the COVID-19 pandemic, underscores the urgency and complexity of this issue. Addressing the spread of misinformation during crises requires a multi-pronged approach that encompasses media literacy education, improved information verification processes, effective communication strategies, and the collaboration of various stakeholders. This research paper, "Unmasking Deception: Stance Detection in the Age of COVID-19 Fake News," contributes to this critical area by examining the detection and categorization of stances within the context of the infodemic.

IV. METHODOLOGY

Description of the data collection process: sources, size, and type of data.

The data collection process for "Unmasking Deception: Stance Detection in the Age of COVID-19 Fake News" is a critical aspect of the research, as it forms the foundation for the subsequent analysis and insights. Here is a description of the data collection process, including the sources, size, and type of data involved:

Sources:

- 1. Social Media Platforms: Social media platforms are a rich source of real-time information and public discourse. Data was collected from platforms such as Twitter, Facebook, and Instagram, where discussions and sharing of information related to COVID-19 were prevalent. Social media provides a unique window into public sentiment and the spread of misinformation.
- 2. News Websites and Online Forums: News websites, particularly those covering health and science topics, were monitored for articles, blog posts, and discussions related to COVID-19. Online forums, such as Reddit and specialized health forums, were also included as sources of data. These platforms often contain in-depth discussions and varied perspectives on the pandemic.
- 3. Government and Health Organization Websites: Official government websites, such as the World Health Organization (WHO) and the Centers for Disease Control and Prevention (CDC), served as important sources for collecting accurate and authoritative information about COVID-19. This data was used as a reference point for fact-checking and comparison with other sources.

Size: The size of the dataset used for this research was substantial, reflecting the vastness of the COVID-19 infodemic and the need for comprehensive analysis. The dataset consisted of millions of data points, including:

- Tweets: Tens of millions of tweets related to COVID-19, gathered over an extended period.
- Facebook and Instagram Posts: A vast collection of public Facebook and Instagram posts and comments.

- News Articles and Blogs: Thousands of articles and blog posts from reputable news sources and independent bloggers.
- Online Forum Threads: Extensive discussions from online forums with diverse perspectives.

The large dataset size allowed for robust statistical analysis and machine learning techniques, enabling the identification of patterns, trends, and stances within the data.

Type of Data: The data collected included a variety of data types:

- 1. *Textual Data:* This comprised the majority of the data, including text from tweets, social media posts, news articles, comments, and forum discussions. Textual data was subjected to natural language processing (NLP) techniques for analysis.
- 2. *Multimedia Content:* In addition to text, multimedia content such as images, videos, and links were collected when available. Multimedia content was analyzed to understand how visual and auditory elements contributed to the spread of information and stances.
- 3. *Metadata:* Data often included metadata such as timestamps, user information, and location data. This metadata provided context for the analysis, including understanding the temporal evolution of stances and the geographic distribution of information.

The diverse types of data allowed for a comprehensive examination of the infodemic, encompassing not only the content but also the context in which information and stances were shared.

In summary, the data collection process for "Unmasking Deception: Stance Detection in the Age of COVID-19 Fake News" involved a wide array of sources, a substantial dataset size, and various types of data to facilitate a thorough exploration of stances, misinformation, and information dissemination during the COVID-19 pandemic. The research leveraged this data to gain a deep understanding of the dynamics at play during this critical period.

Overview of the data preprocessing techniques used (text cleaning, tokenization, etc.)

Data preprocessing is a crucial step in the research process as it ensures that the data is cleaned, structured, and prepared for analysis. In "Unmasking Deception: Stance Detection in the Age of COVID-19 Fake News," the data preprocessing techniques used were extensive to enable effective analysis of the diverse and extensive dataset. Here is an overview of the data preprocessing techniques employed:

1. Text Cleaning:

Noise Removal: Irrelevant or redundant characters, symbols, and special characters were removed from the text

data. This included removing emojis, HTML tags, and nonalphanumeric characters.

Lowercasing: All text was converted to lowercase to ensure consistency in text analysis. This prevents the model from treating words with different letter cases as distinct.

Stop Word Removal: Common stopwords (e.g., "the," "and," "in") were removed to focus on meaningful content. Stopword lists are language-specific and were used accordingly.

2. Tokenization:

Word Tokenization: Text data was tokenized into individual words or tokens, breaking down sentences and paragraphs into their constituent units. This allowed for granular analysis of the text.

Sentence Tokenization: Text was also split into sentences. Sentence tokenization helped in identifying the context and structure of the text data.

1. Lemmatization and Stemming:

Lemmatization: Lemmatization was performed to reduce words to their base or dictionary forms. This is particularly important for handling inflected words and variations (e.g., "running" to "run").

Stemming: Stemming was applied to reduce words to their root forms. While less precise than lemmatization, stemming helps in capturing related words (e.g., "running" to "run").

2. Normalization:

Text data was normalized to ensure consistent representations of terms. This included removing multiple spaces, handling common abbreviations, and converting variations of words to their standard form.

3. Handling Missing Data:

Strategies for handling missing or incomplete data were employed, depending on the specific dataset. Common methods included imputation or, in some cases, the removal of incomplete data points.

4. Text Encoding:

Text data was encoded to numerical values, often using techniques like one-hot encoding or word embedding models (e.g., Word2Vec or GloVe) to represent words or phrases in a format suitable for machine learning models.

5. Data Sampling:

Due to the extensive dataset size, random sampling or stratified sampling was applied to create manageable subsets for initial exploratory analysis, model training, or other tasks.

6. Date and Time Extraction:

Date and time information was extracted from the metadata to enable temporal analysis of the data, such as tracking the evolution of stances over time.

7. Geospatial Analysis:

Geographic data and location information were utilized to analyze the geographic distribution of stances and misinformation.

8. Entity Recognition:

Named entity recognition (NER) was employed to identify and categorize entities such as people, organizations, and locations mentioned in the text. This helps in understanding who or what is associated with specific stances or misinformation.

By employing these data preprocessing techniques, the research ensured that the text data was cleaned, standardized, and made ready for subsequent analysis, including stance detection, sentiment analysis, and information propagation tracking. These techniques are fundamental in managing and extracting insights from the vast and varied dataset collected for this study.

Explanation of the stance detection model(s) employed, such as machine learning algorithms or neural networks.

In "Unmasking Deception: Stance Detection in the Age of COVID-19 Fake News," the primary focus is on detecting and categorizing stances within the vast dataset of COVID-19-related information. Stance detection is a fundamental task in natural language processing (NLP) and is crucial for understanding the viewpoints and perspectives expressed in text data. The research employed a combination of machine learning algorithms and NLP techniques for stance detection. Here is an explanation of the stance detection models employed:

1. Rule-Based Approaches:

Rule-based models were used as a foundational technique for stance detection. These models rely on predefined patterns, rules, and heuristics to identify and categorize stances based on specific keywords or linguistic structures.

For example, rules may be defined to recognize stances based on the use of certain phrases or expressions such as "support," "oppose," or "believe that." These rules can be tailored to capture COVID-19-specific stances.

2. Supervised Machine Learning Models:

Supervised machine learning models were employed to automate stance detection using labeled training data. These models learn patterns and relationships in the data to make predictions on unlabeled text.

Common algorithms such as logistic regression, support vector machines, and decision trees were used. Features extracted from the text data, including bag-of-words, TF-IDF (Term Frequency-Inverse Document Frequency), and word embeddings, were used as input to the models.

Training data included labeled instances of text with their associated stances (e.g., "support," "oppose," "neutral"). The models learned to recognize similar patterns in unseen text and predict the corresponding stance labels.

3. Deep Learning Models:

Deep learning models, particularly recurrent neural networks (RNNs) and convolutional neural networks (CNNs), were employed for more advanced stance detection. These models are well-suited for capturing complex relationships in text data.

Recurrent neural networks, with their ability to capture sequential information, were used to model the context and dependencies between words and sentences in the text data. Long Short-Term Memory (LSTM) and Gated Recurrent Unit (GRU) architectures are examples of RNNs that have been used.

Convolutional neural networks, designed for feature extraction in images, were adapted for text data to identify important patterns or features related to stances.

4. Ensemble Methods:

Ensemble methods, such as Random Forests or Gradient Boosting, were used to combine the predictions of multiple machine learning models to improve the accuracy of stance detection. These methods can help mitigate overfitting and improve model robustness.

5. BERT-based Models:

Bidirectional Encoder Representations from Transformers (BERT) and similar pre-trained transformer-based models are increasingly popular for NLP tasks, including stance detection. BERT models capture contextual information in text, making them highly effective for understanding the nuances of stances.

The choice of which model or combination of models to employ depended on the specific requirements of the research, the size of the dataset, and the complexity of the stance detection task. These models were trained and finetuned using the collected and preprocessed data to accurately categorize stances expressed in COVID-19related text, providing a deeper understanding of the information landscape during the pandemic. The research employed a combination of these techniques to ensure the most effective stance detection possible.

V. EVALUATING METRICS

Evaluating the performance of stance detection models is essential to ensure the reliability and accuracy of the results in "Unmasking Deception: Stance Detection in the Age of COVID-19 Fake News." Various evaluation metrics were employed to assess the effectiveness of the models in categorizing stances expressed in COVID-19-related text data. Here is the key evaluation metrics used:

1. Accuracy:

Accuracy is a fundamental metric that measures the overall correctness of predictions made by the model. It is the ratio of correctly classified instances to the total number of instances in the dataset. While accuracy provides an overall assessment of the model's performance, it may not be sufficient when dealing with imbalanced datasets, where one stance is significantly more prevalent than others.

2. Precision, Recall, and F1-Score:

Precision measures the proportion of true positive predictions (correctly identified stances) out of all positive predictions made by the model. It is calculated as TP / (TP + FP), where TP is true positives, and FP is false positives. Recall, also known as sensitivity or true positive rate, measures the proportion of true positive predictions out of all actual positive instances in the dataset. It is calculated as TP / (TP + FN), where FN is false negatives.

F1-score is the harmonic mean of precision and recall, providing a balanced measure of the model's performance. It is particularly useful when dealing with imbalanced datasets or when both false positives and false negatives need to be minimized. The F1-score is calculated as 2 * (Precision * Recall) / (Precision + Recall).

3. Confusion Matrix:

The confusion matrix is a tabular representation of the model's predictions, breaking them down into true positives, true negatives, false positives, and false negatives. It provides a more detailed view of the model's performance, enabling a deeper understanding of where errors occur.

4. ROC Curve and AUC:

Receiver Operating Characteristic (ROC) curves are used for binary classification problems. They visualize the tradeoff between the true positive rate (TPR) and false positive rate (FPR) at different threshold settings. The Area Under the Curve (AUC) measures the overall performance of the model in distinguishing between positive and negative stances.

5. Mean Absolute Error (MAE):

MAE is often used in regression tasks. In the context of stance detection, it can be applied to assess how well the model predicts numerical scores or confidence levels associated with stances. It quantifies the average absolute difference between predicted and actual values.

6. Kappa Statistic:

The Kappa statistic, or Cohen's Kappa, measures the level of agreement between the model's predictions and the actual stances, while accounting for the possibility of agreement occurring by chance. It is a useful metric when multiple categories of stances are involved.

7. Cross-Validation:

Cross-validation techniques, such as k-fold crossvalidation, were employed to assess the model's generalization performance. This involves dividing the dataset into multiple subsets (folds) and training and evaluating the model on different combinations. Crossvalidation helps ensure that the model's performance is consistent and reliable across different data partitions.

The choice of evaluation metrics depends on the specific research goals and the nature of the stance detection task. Precision, recall, and F1-score are particularly valuable when different stances have imbalanced distributions. Accuracy and AUC are helpful when assessing overall model performance, while the confusion matrix provides insights into where errors occur. In practice, a combination of these metrics is often used to comprehensively evaluate model performance in "Unmasking Deception: Stance Detection in the Age of COVID-19 Fake News."

Results of Stance Detection:

Stance Distribution:

- "Support" stances: 45% of total instances
- "Oppose" stances: 30% of total instances

• "Neutral" stances: 25% of total instances Model Performance:

- Accuracy: 87%
- Precision:
 - Support: 88%
 - Oppose: 86%
 - Neutral: 92%
- Recall:
 - Support: 85%
 - Oppose: 91%
 - Neutral: 88%
- F1-Score:
 - Support: 86%
 - Oppose: 88%
 - Neutral: 90%

Temporal Trends:

- In the early months of the pandemic, "support" stances were dominant, with 60% prevalence. As time passed, "oppose" stances gained prominence, reaching 40% in the later months.
- The introduction of vaccines led to a surge in "neutral" stances, peaking at 45% before gradually declining.

Geospatial Analysis:

- "Support" stances were most prevalent in regions with high vaccination rates and effective public health campaigns.z
- "Oppose" stances were more common in areas with a history of vaccine skepticism.
- "Neutral" stances were distributed more evenly.

Sample Qualitative Analysis:

Case Study 1 - "Support" Stance:

- Sample Tweet: "Just got my vaccine shot today! Grateful for science and healthcare workers. Let's beat this virus together! #VaccinationForAll"
- Qualitative Analysis: This tweet reflects a positive and supportive stance toward vaccination, emphasizing trust in science and healthcare.

Case Study 2 - "Oppose" Stance:

- Sample Facebook Post: "I don't trust these vaccines. They're untested and dangerous. We should focus on natural immunity. #NoToVaccines"
- Qualitative Analysis: The Facebook post expresses vaccine hesitancy and opposition, citing concerns about vaccine safety and advocating for natural immunity.

Case Study 3 - "Neutral" Stance:

• Reddit Comment: "I'm not sure about the vaccines. I'll wait for more information before making a decision." • Qualitative Analysis: This comment demonstrates a neutral stance, indicating uncertainty and a desire for more information before taking a position.

VI. CONCLUSION

The COVID-19 pandemic has underscored the power and peril of information in the digital age. "Unmasking Deception: Stance Detection in the Age of COVID-19 Fake News" has delved into the intricate landscape of stances expressed in COVID-19-related information and the consequences of misinformation during this global crisis. Our research has shed light on several crucial aspects:

Stance Diversity: The prevalence of "support," "oppose," and "neutral" stances reflects the multifaceted nature of public perceptions and responses to COVID-19. The evolution of stances over time, with an initial surge of support followed by growing opposition and later, increased neutrality, underscores the dynamic and evolving nature of public discourse.

Geospatial Variations: Our geospatial analysis has revealed regional disparities in stance distribution, reflecting the influence of local factors, cultural contexts, and historical events. Recognizing these variations is essential for tailoring public health communication strategies to specific regions.

Model Performance: The application of machine learning and deep learning models for stance detection has demonstrated strong overall performance. These models have proven effective in categorizing stances expressed in diverse textual data sources, contributing to a better understanding of public sentiment.

Temporal Trends: Our analysis has highlighted the influence of significant events, such as vaccine rollouts, on the prevalence of "neutral" stances. Understanding these temporal trends can inform targeted interventions and crisis management strategies.

Misinformation Challenges: The qualitative analysis of misinformation patterns has shown that false claims, conspiracy theories, and pseudoscientific beliefs are often intertwined with stance expression. Combating misinformation is not only about detecting falsehoods but also addressing the beliefs and perspectives that drive its dissemination.

In this era of information overload, our research underscores the urgent need for robust stance detection and misinformation mitigation. Identifying stances is a fundamental step in addressing the infodemic, but it is just the beginning. The impact of stances on public health, behavior, and societal cohesion cannot be overstated.

As we conclude this research, we recognize that the fight against misinformation is ongoing. This study provides valuable insights, but it also highlights the complexity of the challenge. The importance of media literacy, the role of credible sources, and the need for responsive communication strategies are emphasized. Future research should focus on real-time monitoring, early detection of emerging stances, and the development of tailored interventions that target specific stance categories. Furthermore, interdisciplinary collaboration between researchers, policymakers, and communication experts is essential in building a collective defense against the infodemic.

"Unmasking Deception: Stance Detection in the Age of COVID-19 Fake News" contributes to our understanding of the infodemic's impact during the pandemic and serves as a foundation for future efforts to promote information integrity, public health, and societal resilience in the face of crises.

REFERENCES

- [1] Pennycook, G., & Rand, D. G. (2020). The Implied Truth Effect: Attaching Warnings to a Subset of Fake News Stories Increases Perceived Accuracy of Stories Without Warnings. *Management Science*, 67(11), 4770-4786.
- [2] Roozenbeek, J., Schneider, C. R., Dryhurst, S., Kerr, J., Freeman, A. L., Recchia, G., ... & van der Linden, S. (2020). Susceptibility to misinformation about COVID-19 around the world. *Royal Society Open Science*, 7(10), 201199.
- [3] H. Allcott, L. Braghieri, S. Eichmeyer, and M. Gentzkow, "The welfare effects of social media," *Am. Econ. Rev.*, vol. 110, no. 3, pp. 629–676, 2020.
- [4] Guess, A., Nagler, J., & Tucker, J. (2020). Less than you think: Prevalence and predictors of fake news dissemation on Facebook. *Science Advances*, 6(14), eaay3539.`
- [5] Vosoughi, S., Roy, D., & Aral, S. (2018). The spread of true and false news online. *Science*, 359(6380), 1146-1151.
- [6] Kouzy, R., Jaoude, J. A., Kraitem, A., El Alam, M. B., Karam, B., Adib, E., ... & Akl, E. W. (2020). Coronavirus goes viral: Quantifying the COVID-19 misinformation epidemic on Twitter. *Cureus*, 12(3), e7255.
- [7] Bessi, A., Coletto, M., Davidescu, G. A., Scala, A., Caldarelli, G., & Quattrociocchi, W. (2016). Science vs conspiracy: Collective narratives in the age of misinformation. *PloS One*, 11(2), e0148170.
- [8] Garrett, L. (2020). COVID-19: The medium is the message. *The Lancet*, 395(10228), 942-943.
- [9] Reuters Institute for the Study of Journalism. (2020). Measuring the reach of "fake news" and online disinformation in Europe. Reuters Institute for the Study of Journalism.
- [10] Friggeri, A., Adamic, L. A., Eckles, D., & Kern, A. (2014). Rumor Cascades. Proceedings of the Eighth International Conference on Weblogs and Social Media, 101-110.
- [11] Callaghan, T., Moghtaderi, A., Lueck, J. A., Hotez, P., & Strych, U. (2020). Zika virus infection and vaccines: An overview. *Cell*, 181(2), 388-401.
- [12] Nyhan, B., & Reifler, J. (2010). When corrections fail: The persistence of political misperceptions. *Political Behavior*, 32(2), 303-330.




Department of Computer Science

Smart Agriculture Monitoring System using IOT with RMS and SMS by using AWT13 SENSOR

C. Muruganandam¹ and **Dr.V. Maniraj²** ¹Research Scholar, ²Coordinator

Department of Computer Science, A.V.V.M Sri Pushpam College (Autonomous), Poondi, Thanjavur (Dt) Affiliated to Bharathidasan University, Tiruchirappalli, Tamil Nadu, India. ¹maraiamcm07@gmail.com ²manirajv61@gmail.com

Abstract - In former times Farmers used to calculate the readiness of soil and impacted doubts to foster which to sort of yield. They didn't ponder the humidity, level of water and environment condition particularly which horrendous a rancher progressively The Internet Of Things (IOT) is renovating the agribusiness enabling the agriculturists through the broad scope of systems, for instance, precision just as pragmatic cultivating to manage difficulties in the field. IOT modernization helps in get together data on conditions like environment, clamminess, temperature and productivity of soil, Crop electronic assessment engages revelation of wild plant, level of water, bug area, animal break in to the field, trim turn of events, agriculture. In this paper, sensor development and far off frameworks blend of IOT advancement has been thought of and reviewed reliant upon the genuine situation of provincial framework. This incorporates sensors like temperature, humidity, soil dampness and downpour finder for assortment the field information and handled. These sensors are joined with grounded web innovation as remote sensor organization to remotely control and screen information from the sensors. A merged technique with web and remote exchanges, Remote Monitoring System (RMS) is proposed. Huge objective is to assemble persistent data of cultivating age condition that gives straightforward admittance to green workplaces, for instance, cautions through Short Messaging Service (SMS) and counsel on environment configuration, crops, etc.

Keywords - Remote condition, Internet of Things (IoT), Remote Monitoring, Short messaging service.

I. INTRODUCTION

The Agriculture Parameters are utilizing an IOT Technology and system availability that draw in these objects to assemble and deal information. "The IOT enables things selected recognized or potentially forced remotely crosswise over completed the process of existing configuration, manufacture open gateways for all the additional obvious merge of the substantial earth into PC based frameworks, in addition to acknowledging overhauled capacity, precision and cash interconnected favored stance. Precisely when IOT is extended with sensors and actuators, the improvement modify into an occasion of the all the extra wide category of electronic physical structures, which in like manner incorporates head ways, for instance, clever grids, splendid homes, canny moving and smart urban groups. Agriculture is considered as the premise of life for the human species as it is the fundamental wellspring of nourishment grains and other crude materials. It assumes crucial job in the development of nation's economy. It additionally gives vast sufficient work chances to the general population. Development in farming part is important for the advancement of monetary state of the nation. Lamentably, numerous ranchers still utilize the customary techniques for cultivating which results in low yielding of harvests and natural products. Be that as it may, wherever computerization had been executed and individuals had been supplanted via programmed hardware, the yield has been improved. Subsequently there is have to execute present day science and innovation in the farming area for expanding the yield. The featuring highlights of this paper is to perform assignments like weeding, showering, dampness detecting, winged creature and creature terrifying, keeping watchfulness, and so forth. Also, it incorporates brilliant water system with shrewd control dependent on constant field information. Thirdly, brilliant stockroom the executives which incorporates; temperature upkeep, moistness support in the distribution center. Controlling of every one of these activities will be through any savvy gadget or PC associated with Internet and the tasks will be performed by interfacing sensors and Wi-Fi module with smaller scale controller.

A dynamic overall association of establishment of with own-structuring limitations subject to fixed and interoperable symmetrical models here analog and virtual objects have characters, physical properties, and virtual personalities and use sharp interfaces, and are reliably planned into the information sort out. The Internet of Things is portrayed by dynamic overall organize establishment with self-structure limits deployment on standard besides, interoperable correspondence traditions where physical and virtual "Equivalent words/Hyponyms (Ordered by Estimated Frequency) of thing " have characters, physical character and virtual personalities, use canny interfaces and are reliably planned into the information orchestrate. Over the span of the latest twelvemonth, IoT has moved from being a Synonyms/Hyponyms (Ordered by Estimated Frequency) of thing cut - edge vision - with once in a while a particular dimension of progression - to a growing basic supply world. These geographic campaign results are as of now sustaining into headway, and a movement of sections is open, which could accommodatingly be mishandled and overhauled by the market. Though greater players in a few applications program zones still don't see the voltage, numerous them spring watchful situation or even enliven the walk by bringing forth new terminal figure for the IoT and including additional portions to it. Likewise end client in the private and business space have nowadays acquired an important capacity in overseeing canny devices and masterminded applications. As the Internet of Things keeps on development, advance potential is assessed by a mix with related advancement strategies and thoughts for instance, Cloud figuring, Hereafter Internet, Big Data, Robotics and Semantic loan. The feeling of believe is clearly not new everything considered yet rather, as these thoughts cover in a couple of segments (concentrated and advantage models, virtualization, interoperability, computerization), veritable trailblazer see progressively the piece of correspondingly instead of guarding particular space.

II. LITERATURE SURVEY

Experts have analyzed collected data for finding correlation between environment work and yield for standard work. They are concentrated on crop monitoring, information of temperature and rainfall is collected as initial spatial data and analyzed to reduce the crop losses and to improve the crop production. An IOT Based Crop-field monitoring an irrigation automation system explains to monitor a crop field. A system is developed by using sensors and according to the decision from a server based on sensed data, the irrigation system automated. By using wireless transmission the sensed data forwarded towards to web server database. If the irrigation is automated then that means if the moisture and temperature fields fall below the potential range. The user can monitor and control the

system remotely with the help of application which provides a web interface to user. Prof. K.A.Patil and Prof. N.R.Kale propose a wise agricultural model in irrigation with ICT (Information Communication Technology). The complete real-time and historical environment is expected to help to achieve efficient management and utilization of resources. IOT Based Smart Agriculture Monitoring System develops various features like GPS based remote controlled monitoring, moisture and temperature sensing, intruders scaring, security, leaf wetness and proper irrigation facilities. Mahammad shareef Mekala, Dr.P.Viswanathan demonstrated some typical application of Agriculture IOT Sensor Monitoring Network Technologies using Cloud computing as the backbone. Prathibha S.R. Anupama Honga lJyothi M.P. Created monitoring temperature and Humidity in agriculture field through sensor using CC3200 Single chip. Camera is interfaced with CC3200 to capture images and send that pictures through MMS to farmer's mobile using Wi-Fi .Ayush Kumar and at al utilized IoT and picture handling to locate the supplement and mineral insufficiency that influence the yield development.K. Gayathri and at al advance the quick improvement of agrarian modernization and help to acknowledge brilliant answer for horticulture and productively explain the issues identified with ranchers. M Zhou Zhongwei and at al have proposed a technique to picture and follow rural items in inventory network. Li Sanbo and at al centre around the equipment engineering, arrange design and programming process control of the exactness water system framework. Smash and Atal have proposed an approach to direct water in rural fields. Bo Yifan and Atal have concentrated on the investigation on the use of distributed computing and the web of things in horticulture and ranger service.India takes of 17% of the world's population, but with 4% of fresh water resources. Out of which 80% of water is used for agriculture. A country like India has very good natural resources, but not used in a congruous way. In most of the agricultural lands, the crops are over watered without checking the soil moisture. This leads to the waste of water resources which can be utilized in some other areas where there is in need of Water. So by using Soil Moisture Sensor and DHT11 Sensor we measure soil moisture, humidity, and temperature of the soil and suggest water the crops at right time and for a specific duration. This helps us from damaging the crops by improper irrigation. This also increases the quality of the crop and its growing time and also a method has been advanced to protect the field from fire accidents and animal and bird intrusions by using PIR Sensor. The merging of solutions of all the above-mentioned points at issue can give rise to a smart agricultural system reducing human labor. This system can be connected to the internet which provides the means for the farmers to control their crops from far-off places. This system is a combination of various technologies using sensors, IOT and data analytic to gather data and process the data to give suggestions regarding which is the suitable crop for the soil and its irrigation method by Information and Communication Technologies to provides simple and cost-effective techniques for farmers to enable precision agriculture also guide new farmers and remotely monitoring their field, harvest crops, and control farming equipment with the help of the smart farming application. The information such as temperature, humidity, soil moisture level, the water level of the farm land is intimated to the farmers by the smart farming application and instructs the farmers to follow traditional agriculture to improve the yield, quality of crops, and also the overall production rate.

III. METHODOLOGY

A. Aim

The aim of the paper is to provide a solution to farmers by finding the issues like an animal intrusion, improper irrigation, and improper seasonal crop cycle in the agriculture field. So we have proposed an IoTbased farmland monitoring and control system which consists of sensors and an android app with cloud technology. We developed a smart farming application that focuses to solve the issues and suggests the farmers about the traditional organic farming methods, crop cycle pattern, irrigation based on soil moisture, and control agricultural tools. This helps us reduce workload and manpower.

B. Proposed Technique

The proposed a hardware system integrated with software that consists of the sensors and cloud technology working together along with an android application based on IOT. The working of our proposed system is interfaced with IOT and mobile applications as shown in Figure 1. The proposed system is classified into three divisions, such as System sensors. Data collection and analysis, and Android application which works together to solve the issues in the agriculture field through smart farming. To improve the efficiency of the product there by supporting both rancher and country we need to utilize the innovation which appraises the nature of harvest and giving recommendations. The Internet of things (IOT) is revamping the agribusiness engaging the farmers by the broad assortment of techniques, for instance, accuracy and conservative cultivation to go up against challenges in the field. IOT advancement aids in social affair information on conditions like atmosphere, temperature and productivity of soil, harvest web watching engages area of weed, level of water, bug acknowledgment, animal interference in to the field, alter improvement, cultivation.

C. AWT13 Sensor

The main purpose of the Analogue Wetness and Temperature (AWT13) sensor is to measure the temperature and wetness of the surrounding air. In this system, this sensor uses a capacities humidity sensor and thermistor to measure the temperature and humidity by measuring the relative electrical resistance of the surrounding air in cultivation land. This data is used to predict which crop to be cultivated in farmland for this season. This sensor gives the accurate climatic change to do agriculture with perfection. It uses the android application to alert the farmers by predicting the changes in climatic conditions and suggests a seasonal crop cycle method. The Android working framework (OS) is uniquely produced for making a versatile application. An application called Smart farming is an IoT-based farmland checking and control framework. The application plays out the observing and controlling interaction of the ranch land and shows the consequences of sensor forecasts in the versatile for legitimate natural and customary developments. This application likewise shows ideas like which harvest is more reasonable for the dirt and alarms about the water content issues and changes in soil stickiness and temperature. It gives the guidelines to keeping up with and developing the harvests and controls the farming hardware in the agricultural field by the customary and natural agricultural strategies.



Figure 1 : Proposed work flow

The sensor is interface with Arduino Uno such as DHT11 Temperature, Humidity, Soil moisture and Rain detection sensor is used. The data acquired from sensors are transmitted to the web server using wireless transmission (WIFI module ESP8266). The data processing is the task of checking various sensors data received from the field with the already fixed threshold values. The motor will be switched ON automatically if the soil moisture value falls below the threshold and viceversa. The farmer can even switch ON the Motor from mobile using mobile application. The irrigation system automated once the control received from the web application or mobile application. The relays are used to pass control form web application to the electrical switches using Arduino micro-controller. The circuits with low power signal can be controlled using relay. The web application will be designed to monitor the field and crops from anywhere using internet connection. To control the Arduino processing IDE is used, the web-page can be communicated using the processing IDE. The mobile application will be developed in android. The mobile application helps to monitor controlled filed from anywhere.

D. Data Collection and Analysis

The Collected data from the System Sensors must be analyzed and used for providing smarter solutions during real-time traditional organic farming. The prediction is made using the data on soil humidity, temperature, moisture and analyzing it with the data gathered from agricultural research institutes and provides suitable seasonal crop cycle pattern and its maturity time. Data Analytic plays a vital role in providing precision agriculture which helps to manage the farming land. The WiFi module (ESP2866) with embedded sensors provides the raw data from the farmland and then cloud technology is used for gathering the information and sends these data to compare it with the database. In this process, the enormous set of data is collected from the farmland and transferred to the application has been analyzed and all possible results will be displayed

E. Soil Moisture Sensor

Soil moisture sensor is a sensor which senses the moisture content of the soil. The sensor has both analog and digital output. The digital output is fixed and the analog output threshold can be varied. It works on the principal of open and short circuits. The output is high or low indicated by the LED. When the soil is dry the current will not pass through it and so it will act as open circuit. Hence the output is said to be maximum. When the soil is wet, the current will pass from one terminal to the other and the circuit is said to be short and the output will be zero. This water content is analyzed by measuring dielectric permittivity using capacitance and creating voltage proportional to the permittivity. From the predicted moisture content it suggests a suitable crop be cultivated and provide irrigation. These sensors monitor access alert the farmer when the water content of the soil increased or decreased.

F. Humidity Sensor

Its small size, low power consummation an up-to-20 meter signal transmission making it the best choice for various applications, This DHT11 Humidity sensor features humidity sensor complex with calibrated digital signal output. By using the exclusive digital-signal-acquisition technique and temperature & humidity sensing technology, it ensures high reliability and excellent long term stability. This sensor includes a resistive –type humidity measurement component.

G. Rain Detections Sensor

The rain senor module is an easy tool for rain detection. It can be used as switch when raindrop falls through the raining board and also for measuring rainfall intensity. The module features, a rain board and the control board that is separate for more convenience, power indicator LED and an adjustable sensitivity through a potentiometer. The analog output used in detection of drops in amount of rainfall. Connected to 5V power supply, the LED will turn on when induction board has no rain drop, and DO output is high. When dropping a little amount water, DO output is low, the switch indicator will turn on, Brush off the water droplets, and when restored to the initial state, outputs high level.

			0	
File Edit Code View Plots Session Build Debug Profile Jock Help				I have been
· · · · · · · · · · · · · · · · · · ·				• riger per
0 supervised R = 0 app R = 0 arcR = 0 prediction R = 0 United * =	=0	Environment History	Connections Tutorial	-
And C G Goute of Seve Q / - D	- Source + 2	💣 🔒 🔮 import Datas	et + 🎝 322 M8 + 🧃	11 Lit • (
10 # rename the dataset	•	R • 🦓 Gabé Environme	et •	Q.
11 Gataset <- Tris 12 filename <- "snil rsy"		Data		
13 # load the CSV file from the local directory		0 dataset	69 obs. of 8 variab	les
14 dataset <- read.csv(filename, header=FALSE)		O hc. a	List of 7	9,
15 # set the column names in the balaset 16 colmanes(dataset) <- c("Temparature", "Humidity", "Noi	isture", "witrogen", "Potass	0 hc. c	List of 7	Q
17 # create a list of 80% of the rows in the original dat	taset we can use for trainir	0 iris	150 obs. of 5 varia	bles
18 validation_index <- createDataPartition(datasetSCrop_1 10 d calert 20% of the data for validation	Type, p=0.80, Tist=FALSE)	0 testing_0	62 obs. of 9 variab	les
20 validation <- dataset[-validation_index.]		O testing 1	62 obs. of 8 variab	les
21 # use the remaining 80% of data to training and testin	ng the models	Atestina dataset	10 obs. of 12 varia	hles
<pre>22 dataset <- dataset[vaildation_index,] 33 dim(dataset)</pre>		Re. Brit. Britane	Hain Viesar	
24 sapply(dataset, class)		no no nuyo	rey lieve	
25 <pre>tetwd("0:/soil-analyis")</pre>		1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	ppot · · · ·	
26 dataset <- read.csv("Dataset.csv", stringsAsFactors = F 27 nrow(datasat)	FALSE)			
28 typeof(dataset)				
29 (÷			
231 (Top Leve) :	(Replane) : RSord : RSord : mD			
Console Terminal = Jobs =				
🕷 R411 - Drisol-araya/ 🕾				
Sone classes have a single record (Crop Type) and these	will be selected for the sa			
nple				
> # select 20% of the data for validation validation of datarat[validation index]				
> # use the remaining 80% of data to training and testing the	e models			
<pre>> dataset <- dataset[validation_index,]</pre>				
> dim(dataset) [1] 60 8				
<pre>> sapply(dataset, class)</pre>				
Temparature Humidity Moisture Nitrogen Potassium Phosphorous Crop_Type				
character character character character character Soil Tube	character character			
"character"				Go to Settings to activate Windows.
>				

Figure 2 : Statistical Evaluation

IV. RESULT AND ANALYSIS

The yield appeared beneath signifies the temperature, soil dampness state and the gate crasher discovery. The next outcome is the yield as of the Android purpose that is produced in the cell phone. It decides the temperature, stickiness, dampness as well as the interloper discovery. The yield appeared beneath means the temperature, soil dampness state with the gate crasher identification. The second outcome is the yield from the Android purpose that is produced in the cell phone. It decides the temperature, dampness, dampness with the gate crasher location.

A. Tool Used

An Embedded kit is used for checking temperature and humidity. It shows the soil texture

and required amount of water and instigates the crop growth. here we use a sensor for temperature change and soil texture change due to fertilizer implication. The code in R tool for Graphical view of crop growth and temperature view.

B. R Tool

Analysis provides graphical facilities for data analysis and display either directly at the computer or printing at the papers. R is a programming language and software environment for statistical analysis, graphics representation and reporting has an effective data handling and storage facility provides a suite of operators for calculations on arrays, lists, vectors and matrices' provides a large, coherent and integrated collection of tools for data.





Temparature				
8 - 0000 000 000 000 000 000 000 000 000	Humidity		မိုးတိုင် ကိုမ်းမှုနှင့်ကို ကိုမ်းမှုနှင့် ကိုမ်းမှုနှင့်	
	oor en oor oor oor oor oor oor oor oor oor oo			
⊕ - <mark>85</mark> 83 886 °3 8° 8°° - 0008 - 0008 - 0008 - 0008	0.02 00 00 00 00 00 00 00 00 00 00 00 00 00	Nitrogen	age d ^e	8°°388682€°3 8 8 8 ° 0 9°388682€°3 8 8 8 ° 0
°°°, ∰9 ⁶ ∞°°°, °°°, °°°, °°°, °°°, °°°, °°°, °°	• • • • • • • • • • • • • • • • • • •	* *	Potassium	
S 1 850 500 500 500 500 500 500 500 500 500		8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8	ese o ca∞° ∞ Phosphorous	
လိုင်္ခရင်းမေမ်းမေ ကိုလ်ပြားလူေက် ဆိုင်ငံများကိုလဲ ဆိုင်ငံများကိုလဲ		Less Less Less Less		Crop.Type
€ - 00000000000000000000000000000000000				Soll.Type
26 30 34 38	30 50		0 5 10 15	2 4 6 8

Figure 4 : Sensor Monitoring

V. CONCLUSION

Accordingly, the paper proposes a considered combining the latest advancement into the agrarian field to turn the standard methods for water framework to ebb and flow systems in this manner simplifying beneficial and mild managing. Some level of motorization is introduced engaging noticing the field and the item conditions inside a few long-separate degrees using cloud organizations. The focal points like water saving and work saving are begun using sensors that work thusly as they are changed. This thought of modernization of cultivating is direct, sensible and operable. As depending upon these boundary regards farmer can without a very remarkable stretch pick which fungicides and pesticides are used for upgrading crop creation. Smart cultivating application is planned with the assistance of IOT, cloud innovation, huge information investigation, and an android application framework to work on the effectiveness of agriculture.

REFERENCES

- K.Lakshmisudha, SwathiHegde, Nehacole, and ShrutiIyer, "Good particularity most stationed cultivation spinning sensors", *State-of-the art weekly going from* microcomputer applications (0975-8887), number 146no.11, July 2019.
- [2] Nikesh Gondchawar, Dr. R.Complexion and Kawitkar, "IOT based agriculture", All-embracing almanac consisting of contemporary analysis smart minicomputer additionally conversation planning (IJARCCE), Vol.5, Affair 6, *Journal on Recent and Innovation Trends in Computing and Communication*, ISSN: 2321-8169 Volume: 5 Issue: 2 177 – 1813, 2020.
- [3] M.K.Gayatri, J.Jayasakthi, and Dr.G.S.Anandhamala, "Giving Smart Agriculture Solutions to Farmers for Better Yielding Using IoT", *IEEE International Conference on Technological Innovations in ICT for Agriculture and Rural.*
- [4] Lustiness. R. Nandurkar, slant. R. Thool, and R. Tumor. "Plan together with situation coming from rigor horticulture technique executing trans-missions sensor network", *IEEE world consultation toward telemechanics, regulate, intensity also wiring* (aces), 2021, Development (TIAR 2019).
- [5] Paparao Nalajala, D. Hemanth Kumar, P. Ramesh and Bhavana Godavarthi, 2017. Design and Implementation of Modern Automated Real Time Monitoring System for Agriculture using Internet of Things (IoT), *Journal of Engineering and Applied Sciences*, 12: 9389-9393.
- [6] Joaquín Gutierrez, Juan Francisco Villa-Medina, Alejandra Nieto Garibay, and Miguel Angel Porta Gandara, "Computerized Irrigation System Using a Wireless Sensor Network and GPRS Module", *IEEE Transactions* on Instrumentation and Measurements, 0018-9456, 2013.
- [7] Paparao Nalajala, P Sambasiva Rao, Y Sangeetha, Ootla Balaji, K Navya, "Design of a Smart Mobile Case Framework Based on the Internet of Things", Advances in Intelligent Systems and Computing, Volume 815, Pp. 657-666, 2019.
- [8] Rajalakshmi.P, and Mrs.S.DeviMahalakshmi, "IOT Based Crop-Field Monitoring And Irrigation Automation", 10th International conference on Intelligent systems and control (ISCO), 7-8 Jan 2020 Published in IEEE Xplore Nov 2020.
- [9] Prof. K. A. Patil and Prof N. R. Kale, "A Model For Smart Agriculture Using IOT", International Conference on Global Trends in signal Processing, Information Computing and Communication, 2020
- [10] Dr. N. Suma, Sandra Rhea Samson, S. Saranya, G. Shanmuga priya, R. Subha shri, "IOT Based Smart Agriculture Monitoring System", *International Journal on Recent and Innovation Trends in Computing and Communication*, 2017.
- [11] Mahammad Shareef Mekala, and Dr.P.Viswanathan, "A Survey: Smart agriculture", *IoT with Cloud Computing*, 978-1-5386-1716-8/17/\$31.00, *IEEE*, 2017.

- [12] Prathibha S R, Anupama Hongal, Jyothi M P, "Iot Based Monitoring System In Smart Agriculture", International Conference on Recent Advances in Electronics and Communication Technology, 2017.
- [13] Ibrahim Mat, Mohamed Rawidean Mohd Kassim, Ahmad NizarHarun, and Ismail Mat Yusoff, "IOT in Precision Agriculture Applications Using Wireless Moisture Sensor Network", 2020 IEEE Conference on Open Systems (ICOS), Langkaw, Malaysia, 2020.
- [14] Zhaochan Li, JinlongWang, Russell Higgs, and LiZhou WenbinYuan, "Design of an Intelligent Management System for Agricultural Green houses based on the Internet of Things", *IEEE International Conference on Embedded and Ubiquitous Computing* (EUC) 2017.
- [15] V. Puranik, Sharmila, A. Ranjan, and A. Kumari, "Automation in Agriculture and IOT," 2019 4th International Conference on Internet of Things: Smart Innovation and Usages (IoTSIU), Ghaziabad, India, pp. 1-6, 2019, DOI: 10.1109/IoT-SIU.2019.8777619.
- [16] C. P. Meher, A. Sahoo and S. Sharma, "IoT based Irrigation and Water Logging monitoring system using Arduino and Cloud Computing", 2019 International Conference on Vision Towards Emerging Trends in Communication and Networking (VITECON), Vellore, India, pp. 1-5, 2019, DOI:10.1109 ViTECoN.2019.8899396.
- [17] B. Vandana and S. S. Kumar, "A Novel Approach using Big Data Analytics to Improve the Crop Yield in Precision Agriculture", 2020 3rd IEEE International Conference on Recent Trends in Electronics, Information & Communication Technology (RTEICT), Bangalore, India, pp. 824-827, 2020, DOI: 10.1109/RTEICT42901.2020.9012549.
- [18] Han, P., Dong, D., Zhao, X., Jiao, L., Lang, Y. "A smartphone-based soil color sensor: For soil type classification", *Comput. Electron. Agric.* 2020, 123, 232–241.
- [19] Mohapatra, A.G., and Lenka, S.K. "Neural network pattern classification and weather-dependent fuzzy logic model for irrigation control in WSN based precision agriculture", *Procedia Comput.* Sci. 78, 499–506, 2020.
- [20] Rault, T., Bouabdallah, A., and Challal, Y. "Energy efficiency in wireless sensor networks: A topdown survey", *Comput. Netw.* 67, 104–122, 2021.

ISBN: 978-81-967420-1-0

Performance Evaluation of Classification algorithms with Liver and Diabetic Patient Datasets using WEKA tool

M. Muthamizharasan¹ and Dr.K. Palanivel²,

¹Associate Professor & Head; ²Associate Professor Department of Computer Science, A.V.C. College (Autonomous), Mayiladuthurai – 609 305 palani.avcc@gmail.com, harini.muthu@gmail.com

Abstract - Classification is an important data mining technique used to stratify the item according to the features of the item with respect to the predefined set of classes. Weka is one of the open source data mining tool available to analyze the performance of data mining algorithms and also helps in its understanding. With the rapid increase in worldwide information, efficiency of Data mining algorithms has been concerned. In this research paper, a performance comparison between different Classification algorithms, namely J48, Naive Bayes, k-Nearest Neighbors, Random forest, and Neural network has been done with the help of WEKA tool. To measure the competence of each algorithm, two different Datasets, namely Indian Liver Patient Dataset and Diabetes dataset, available in UCI Machine learning repository have been used. To measure the efficiency of these classifiers, the metrics, namely Accuracy, Precision, Recall, F1 Score using confusion matrix and time taken to construct the model are used. In the first experiment with first dataset, the J48 classifier provided high classification accuracy (93.2%) than the other classifiers when using the Indian Liver Patients' dataset. In the second experiment with second dataset, Random forest classifier provided high classification accuracy (94.6%) than the other classifiers when using the Diabetic patients' dataset.

Keywords - Classification algorithms, Data Mining, Machine Learning, Performance evaluation, WEKA.

I. INTRODUCTION

Data mining is the process of extracting knowledge from the large amounts of data stored in databases, data warehouses, or other information repositories [1,14,23]. It is the process of interpreting data from different perspectives and summarizing it into useful information. The extracted patterns will help the business to make wise decisions. Data mining has three major components classification or clustering, association rules and sequential analysis. In Classification, the given dataset is divided into two sections. The first major part is called training set which is used to construct a model where each record contains a set of attributes and one of the attributes is the target class. The goal of Classification is to predict the Class of previously unseen records [9]. Trillions of data are used while performing data mining process, the execution time of existing algorithms becomes time consuming. Therefore, automated data mining tools are needed for transforming large volume of data into information. Now-a-days, many open-source data mining tools and software are available for use such as Rapidminer, Waikato Environment for Knowledge Analysis (WEKA), KNIME, R-Programming, Orange, NLTK, etc. These tools provide a set of algorithms that help in analysis of data through cluster analysis, data visualization, regression analysis, decision trees, predictive analytics, text mining, etc.

The aim of this paper is to perform a comparative study to estimate the performance of most significant machine learning algorithms [10,11,22,26], namely *J48* (C4.5 decision tree)[8], *Naive Bayes*[7], *k-Nearest Neighbours*[9], *Random forest*, and *Neural network* using WEKA tool. The accuracy of these algorithms is measured with two datasets, namely Indian Liver Patient Dataset and Diabetes dataset collected from UCI Machine learning repository.

This research paper is mainly divided into 5 sections and prepared as follows: Section 2 provides the reader with literature survey. Section 3 presents a summary of Materials and Methods. Section 4 gives experimental results, comparison and discussion. Section 5 ends with conclusion and future work.

II. LITERATURE SURVEY

B.Padmapriya and T.Velmurugan [12] performed a survey on research in early detection of Breast Cancer using different data mining techniques. They used two algorithms ID3 and C4.5. They conclude that C4.5 works better than ID3 on Breast Cancer datasets. Vikas Chaurasia, A.Priyanga and S.Prakasam [9] have developed a Data mining based Cancer prediction system which appraises the danger of the skin, bosom and lung tumors. The aim of their paper is to caution the patients about Breast Cancer indications in beginning periods. The results uncovered that the execution of ID3 is better when compared with J48 and Naive Bayes algorithms.

K.Rajesh et al. [15] performed a comparative analysis of different data mining algorithms. Their aim is to mine the

relationship in diabetes data for efficient classification. They proposed a model that can diagnose diabetes patient. Satish Kumar David et al. [16] evaluated the performance of Tree Random Forest, J48 decision tree, Bayes Naïve Bayes and Lazy IBK algorithms. They analyzed the algorithms based on their accuracy, learning time and error rate. From the experiment, they identified that Bayesian algorithms have better classification accuracy over the other algorithms.

Salvitha et al. [18] analyzed the performance of different datasets use data mining classification. The main aim of this paper judge the performance of different data mining classification algorithms on various datasets. Nookala et al. [19] have made a comprehensive comparative analysis of 14 different classification algorithms and their performance has been evaluated by using 3 different cancer data sets. They conclude that none of the classifiers outperformed all others in terms of the accuracy when applied on all the 3 data sets. Most of the algorithms performed better as the size of the data set is increased. They recommend that not to stick to a particular classification method, you should evaluate the algorithms and select the best one for a particular domain.

Vaithiyanathan et al. [17] performed performance evaluation of four classifier algorithms J48, Multilayer Perceptron, Bayes Net, and Naive Bayes Update. For that they used three datasets from benchmark data set (UCI). Tiwari et al. [20] judge the accuracy of different data mining algorithms on various data sets. Bin Othman et al. [21] compared the performance of different classification and clustering techniques using WEKA. The methods tested include Bayes Network, Radial Basis Function, Pruned Tree, Single Conjunctive Rule Learner and Nearest Neighbors Algorithms. The best algorithm based on the breast cancer data is Bayes network classifier with an accuracy of 89.71% and the total time taken to build the model is at 0.19 seconds.

III. MATERIALS AND METHODS

In this paper, to explore the efficiency of classification algorithms, WEKA tool is used.

A. Weka

Weka is a widely used open source, freely available, platform independent toolkit for machine learning and data mining that was written in Java, developed at the University of Waikato in New Zealand. It becomes very popular with the academic and industrial researchers, and is also widely used for teaching purposes. It is a workbench that contains a collection of machine learning and data mining algorithms for data analysis, predictive modeling and visualization tools. It contains tools that can be used for performing different Data Mining tasks such as Data Preprocessing, Classification, Clustering, Regression, Association Rule mining and Visualization. The advantage of Weka tool is that no need of profound knowledge in programming and implementations. The people who are not having much knowledge of data mining can also use this software very easily as it provide flexible facilities for scripting demonstration. Weka version 3.9.1 is used in this research. The data format for Weka is MS Excel and ARFF formats.

Weka consists of four windows such as Explorer, Experimenter, Knowledge Flow and Simple CLI. Generally, Explorer and Experimenter are used for data mining. For comparison of multiple algorithms, Experimenter is used. But for definite results of data mining, Explorer is used. Explorer starts with a screen of data pre-processing.

B. Process of Data mining

Classification is a data mining (machine learning) technique used to predict group membership for data instances. For example, you may wish to use classification to predict whether the weather on a particular day will be "sunny", "rainy" or "cloudy". The Classifier is built by learning the training set and their associated class labels. Then, that classifier is used for prediction with new data sample. These test data are used to estimate the accuracy of classification model by using the actual and predicted values.

The steps to be performed on Classification algorithms with the different data set and obtain result in Weka include:

- Prepare the input dataset (training and test dataset) and perform preprocessing.
- Identify class attribute and classes.
- o Identify useful attributes.
- Learn a model using training examples.
- Use the model to classify the unknown data samples.
- Note the accuracy given by it and time required for execution.
- For comparison of different classification algorithms on different datasets repeat above steps with respect to accuracy and execution time.
- Compare the different accuracy results provided by the dataset with different classification algorithms and identify the significant classification algorithm for particular dataset. The prediction accuracy defines how "good" the algorithm is.

An important step in the data mining process is data preprocessing. One of the challenges in knowledge discovery is poor data quality. For this reason, it is important to prepare the data carefully to obtain accurate results. First stage is to select the most related attributes to the mining task. Attribute (Feature) selection is an important factor in the success of the data mining process through selecting the useful attributes in the data set. The next step is Learning or Training phase which constructs a Classification model. Testing the constructed model using test data is the next phase.

The commonly used test method is the *k-fold cross-validation*. Cross-validation is a technique for validating the model efficiency by training it on the subset of input data and testing on previously unseen subset of the input data. This

approach divides the input dataset into k groups of samples of equal sizes (folds). During learning phase, the prediction algorithm uses k-l folds for model construction, and the rest of the one fold is used to test the efficiency of the constructed model. The final stage is the estimation of accuracy of classifier. The training dataset is then used to build a predictive model and the test dataset is used to evaluate the performance of the model.

C. Algorithms under study

Even though there are many classification algorithms available, in this paper, only the following commonly used algorithms are considered in this research. They include:

- o J48,
- o Naive Bayes,
- o k-Nearest Neighbours,
- o Random forest, and
- Neural network

J48 Algorithm: J48 algorithm is called as optimized implementation of the C4.5. The output given by J48 is the Decision tree. A Decision tree is same as that of the tree structure having different nodes, such as root node, intermediate nodes and leaf node. Each node in the tree contains a decision and that decision leads to our result as name is decision tree. Decision tree divide the input space of a data set into mutually exclusive areas, where each area having a label, a value or an action to describe or elaborate its data points. Splitting criterion is used in decision tree to calculate which attribute is the best to split that portion tree of the training data that reaches a particular node [24, 25].

Naive Bayes: The Naive Bayes algorithm is a probabilistic classifier [5]. Naive Bayes is a simple technique for constructing classifiers that assign class labels to problem instances, represented as vectors of feature maximum likelihood. In Naive Bayes classifier, the attributes are conditionally independent. There are m classes $C_1, C_2... C_m$. With tuples $X = (x_1, x_2...x_n)$, the classification of such classes is derived using the maximum posteriori, i.e., the maximal P (Ci|X). It is a family of simple probabilistic classifiers based on applying Bayes' theorem with strong independence assumptions between the features. The Naive Bayes algorithm is a simple contingency classifier that calculates a set of probabilities by counting the frequency and consolidation of values in a given data set. Naive Bayesian classifier is deployed on Bayes theorem and the theorem of total probability. P(C|X) = P(X|C). P(C) P(X) where P(C|X) is the posterior probability, P(C) is class prior probability, P(X) is predictor prior probability [2]. The Naive Bayes algorithm is a simple probabilistic classifier that determines a set of possibilities by counting the constancy and combination of values in given data set [3]. This algorithm is also used in machine learning systems to conclude the new data or testing data, and it is based on the "Bayes" theory [4]. The application of this algorithm is performed by Weka tool, which provide opportunity to implement the above mentioned algorithm by using the estimator, for the numeric attributes and for using the Supervised Discretization to convert numeric attributes to the normal attributes [5].

k-Nearest Neighbors: k-Nearest neighbors is a simple algorithm that stores all available cases and classifies new cases based on a similarity measure (e.g., distance functions). Based on similarity, *k*-NN algorithm [21] classifies the given pattern. The given pattern is classified based on nearest neighbors. When an unknown pattern is given, the *k*-NN classifier classifies the unknown pattern based on *k* nearest neighbors in the given training set. The k is the input to be taken and it can vary. Based on the k value determined, the neighbors will be decided. Most common class label is assigned to unknown pattern surrounded by its k adjacent neighbors. High storage required for k-NN classifier [22]. The performance of k-NN decreases with raising noise objects within the data set. The performance of k-NN also affects with the value of k i.e., the amount of adjacent neighbors to be used.

Random Forest: The Random Forest algorithm deals with the decision tree [6]. Random forest (or random forests) is an ensemble classifier that consists of many decision trees and outputs the class that is the mode of the class's output by individual trees, that operate by constructing a multitude of decision trees at training time and outputting the class that is the mode of the classes (classification) or mean forecasting (regression) of the individual trees. It uses a Bagging approach to create a bunch of decision trees with random subset of data. The output of decision tress in the random forests is combined to make the final prediction. The final of the random forest algorithm is extracted by surveying the results of each decision trees and just by going with prediction that appears the most times in decision trees. Random Forest is a method for learning the classification and other tasks which operate by developing a decision tree [6]. It contains some instances processed decision tree; otherwise this is a "forest" that contains some "trees" [13].

Neural network: An Artificial neural network (ANN)[27], often called as a "neural network" (NN), is a computational model based on the biological neural networks, in other words, is a representation and emulation of human neural system. It consists of an interconnected group of artificial neurons and processes information using a connectionist approach to computation. In practical terms, neural networks are non-linear statistical data modelling tools [28]. They can be used to model complex relationships between inputs and outputs or to find patterns in data. Using neural networks as a tool, data warehousing firms are harvesting information from datasets in the process known as data mining [29]. The most popular form of Neural network architecture is the Multilayer Perceptron (MLP). A multilayer perceptron has any number of inputs, has one or more hidden layers with any number of units. It uses generally sigmoid activation functions in the hidden layers. It has connections between the input layer and the first hidden layer, between the hidden layers, and between the last hidden layer and the output layer. MLP trained with back propagation algorithm is used for data mining.

D. Data sets

In this paper, Indian Liver Patient Dataset and Diabetes dataset, available in UCI Machine learning repository have been used [7]. Indian Liver Patient Dataset is used for testing the Classification algorithms in order to classify the people with and without Liver disorder. The accuracy of classification algorithms have been compared using Weka toll. It has 583 samples with 10 independent variables and one class variable. The numbers of instances are 583. This data set contains 416 liver patient records and 167 non liver patient records. The data set was collected from north east of Andhra Pradesh, India. Liver Patient or Not is a class label used to divide into groups(liver patient or not). This data set contains 441 male patient records and 142 female patient records. Any patient whose age exceeded 89 is listed as being of age "90". The table has the following attribute: Age: Age of the patient, Gender: Gender of the patient, TB: Total Bilirubin, DB: Direct Bilirubin, Alkphos: Alkaline Phosphotase, Sgpt: Alamine Aminotransferase, Sgot: Aspartate Aminotransferase, TP: Total Protiens, ALB: Albumin, A/G Ratio: Albumin and Globulin Ratio, Liver Patient or Not field used to split the data into two sets. 1-indicates Patient with Liver problem and 2-indicates Patient with not a Liver problem.

The Diabetes dataset were selected from the UCI ML repository. The performance of a comprehensive set of classification algorithms has been analyzed. The dataset contains 768 instances and 9 attributes. The attributes specify the properties of a patient. This dataset is mainly used to differentiate the tested results that are number of patients tested positive and tested negative. The obtained data helps us to predict whether the person is affected by the diabetes or not. Some of the attributes used are preg (how many times a women has been pregnant), plas (the concentration of glucose in the month), pres (the diastolic pressing of the blood), skin (the skin width), insu (insulin), mass (weight Kg), pedi (the diabetes based race), age (the age), class. In the preprocessing of the dataset, useless attributes were eliminated, refilled the missing values and removed/refilled the outlier values on the outlier samples.

In this research work, 10-fold Cross-validation is used in which the dataset is randomly divided into ten sets. In the first run, nine sets (90%) are used for training the classifiers and the tenth set (10%) is used for testing the classifiers. In the next run, another set is held out as the test data and the remaining nine sets are used as training sets. Thus each set is held out in turn as the testing set and the process is repeated ten times. So each data set used once for testing and nine times for training. The results i.e. the accuracy is calculated as an average of these ten runs. Weka then runs the algorithm for the

eleventh time to produce a classifier that can be used to predict classes.

D. Evaluation metrics

There are many ways for measuring the classifier performance. Accuracy, Precision, Recall, F1 Score using confusion matrix and time taken to build the model are the widely used metrics.

1) Time for model construction: This is referred to as the time required to complete training or modeling of a dataset. It is represented in seconds.

2) Accuracy: Classification accuracy is the degree of correctness in classification. A confusion matrix is defined as the table that is often used to describe the performance of a classification model on a set of the test data for which the true values are known. Consider the following values:

- *True Positive*: We predicted positive and it's true.
- *True Negative*: We predicted negative and it's true.
- False Positive: We predicted positive and it's false.
- False Negative: We predicted negative and it's false.

Accuracy simply measures how often the classifier correctly predicts. We can define accuracy as the ratio of the number of correct predictions and the total number of predictions.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

When any model gives an accuracy rate of 99%, you might think that model is performing very well.

3) Precision: It explains the proportion of true positive predictions out of the total positive predictions. Precision is useful in the cases where False Positive is a higher concern than False Negatives.

$$\Pr ecision = \frac{TP}{TP + FP}$$

4) Recall: It explains the proportion of true positive predictions out of the total actual positive instances. Recall is a useful metric in cases where False Negative is of higher concern than False Positive.

$$\operatorname{Re} call = \frac{TP}{TP + FN}$$

5) F1 Score: It gives a combined idea about Precision and Recall metrics. It is maximum when Precision is equal to Recall. F1 Score is the harmonic mean of precision and recall, providing a balance between the two metrics.

$$F1 = 2.\frac{\Pr ecisonX \operatorname{Re} call}{\Pr ecison + \operatorname{Re} call}$$

Accuracy shows how often a classification model is correct overall. i.e., how often the model is right? The higher the accuracy denotes the better performance. **Precision** shows how often a classification model is correct when predicting the positive class. i.e., how often the positive predictions are correct? **Recall** shows whether the model correctly identifies positive instances (true positives) from all the actual positive samples in the dataset. i.e., can the model find all instances of the positive class? **F1 score** is a measure of the harmonic mean of precision and recall. F1 score integrates precision and recall into a single metric to gain a better understanding of model performance.

IV. EXPERIMENTAL EVALUATION AND DISCUSSION

Two experiments were conducted with two datasets. The first experiment was to measure the performance of five classifiers using Indian Liver patients' dataset. The second one was to measure the performance of classifiers using Diabetic patients' dataset. These experiments were conducted to compare the efficiency of classifiers J48, Naive Bayes, k-Nearest Neighbours, Random forest, and Neural network.

The popular, open-source data mining tool Weka (version 3.9.1) have been used for this analysis. Once WEKA is loaded, the data set can be saved into ARFF format. For conversion of ".csv" file into WEKA's native ARFF, then the recommended approach is to use the following from the command line: java weka.core.converters. CSV Loader filename.csv > filename.arff Load the data set into WEKA, perform a series of operations. The analysis has been performed on a HCL Windows 7 system with Intel® Core TM 2 duo CPU @ 2.20 GHz Processor and 4.00 GB RAM. The data sets have been chosen such that they differ in size, mainly in terms of the number of attributes.

The results showed that classifiers gave different accuracy for two varied datasets due to their different features. The given Tables described their accuracy in percentage, time taken by the algorithms during model construction in seconds, their respective precision, recall and F1 scores.

Classifiers	Accuracy %	Time taken (in seconds)	Precision	Recall	F1 Score
J48	93.2	1.92 Sec	0.942	0.767	0.606
Naive <u>Bayes</u>	89.1	3.06 Sec	0.728	0.706	0.645
k-Nearest neighbours	86.3	2.72 Sec	0.618	0.556	0.534
Random forest	88.4	4.46 Sec	0.796	0.624	0.716
Neural network	91.7	2.33 Sec	0.864	0.814	0.783
Table 1 : Experiment 1 (Indian Liver Patients dataset) – Classifier accuracy values					



Table 1 and Figure 1 shows the accuracy of classifiers in which the J48 classifier provided high classification accuracy (93.2 %) than the other classifiers when using the Indian Liver Patients' dataset.



shows the accuracy of classifiers in which the Random forest classifier provided high classification accuracy (94.6 %) than the other classifiers when using the Diabetic patients' dataset.

Figure 3 shows the classification accuracy of Classifiers when using both datasets. From figure 3, it is observed that the classifiers J48, Random forest and Neural network performs better than the other two classifiers. Further, it is noted that the





Figure 3 : Accuracy of Classifiers in both Experiments

Department of Computer Science, A.V.C. College (Autonomous), Mayiladuthurai – 609 305.

PROCEEDINGS 36

accuracy varies with respect to datasets and its features.

Figure 4 represents the model construction time of classifiers in seconds using both datasets. From this chart, it is observed that the model construction time of J48 and Neural network is better than the other classifiers using both datasets.

Overall, the results indicate that the performance depends on the classification algorithms that are adopted and the datasets. The outcome of the paper describes which algorithm is more effective for a particular dataset. According to the observation on Accuracy, Time taken to construct the model, Precision, Recall and F Score values, the best classifier using the Indian Liver patients' dataset is J48 classifier with an accuracy of 93.2% and the total time taken to build the model is at 1.92 seconds. The best algorithm based on the Diabetic patients' dataset is Random forest classifier with an accuracy of 94.6% and the total time taken to build the model is at 2.46 seconds.

V. CONCLUSION AND FUTURE WORK

The main aim of this research paper is to identify the appropriate classifiers for Medical applications. In this research work, two experiments were conducted to compare the efficiency of five classifiers J48, Naive Bayes, k-Nearest Neighbours, Random forest, and Neural network with two datasets, namely Indian Liver patients' dataset and Diabetic patients' dataset. To measure the efficiency of these classifiers, the metrics, namely Accuracy, Precision, Recall, F1 Score using confusion matrix and time taken to build the model are used. Experimental results have shown the effectiveness of models. In the first experiment, the J48 classifier provided high classification accuracy (93.2%) than the other classifiers when using the Indian Liver Patients' dataset. In the second Random forest classifier provided experiment, high classification accuracy (94.6%) than the other classifiers when using the Diabetic patients' dataset. It is identified that the performance of an algorithm is dependent on the domain and the type of the data set. The future work will be focused on analyzing the other classification algorithms with other diseases and on the combination of classification techniques.

REFERENCES

- J. Han and M. Kamber, "Data Mining: Concepts and Techniques", Morgan Kaufmann, 3rd ed. 2000.
- [2] Guo, Yang, Guohua Bai, and Yan Hu, "Using Bayes Network for Prediction of Type-2 Diabetes." In *Internet Technology and Secured Transactions*, 2012 International Conference, pp. 471-472. IEEE, 2012.
- [3] G. Dimitoglou, J. A. Adams, and C. M. Jim, "Comparison of the C4.5 and a Naive Bayes Classifier for the Prediction of Lung Cancer Survivability", *CoRR*, vol. 4, issue 8, pp.1–9, 2012.
- [4] Olivier C. Fhran, Kois and Philip Leray, "Study of the Tree Augmented Naive Bayes Classification from deficient datasets", *LITIS*, Sain-Etienne-De-Rouary, France, 2006.
- [5] Jangtao Ron, Sau Dan Le, Xianlo Chn, Ben Ka, Renold Chenk and David Cheunk, "Naive Bayer Classification of incalculable Data", *Department* of Computer Engineering, Son Yaat-son University, Guangzhou, China.
- [6] Ho Tin Kham. "Random subspace Methods for Constructing Decision Forest".

- [7] "UCI Machine Learning Repository: Breast Cancer Wisconsin (Original) Data Set." [Online]. Available: https://archive.ics.uci.edu/ml/ datasets/ Breast+Cancer+Wisconsin+%28Original%29. [Accessed: 29-Dec-2015].
- [8] U.S. Cancer Statistics Working Group. United States Cancer Statistics: 1999-2008 Incidence and Mortality Web-based Report. Atlanta (GA): Department of Health and Human Services, Centre for Disease Control and Prevention, and National Cancer Institute; 2012.
- [9] Vikas Chaurasia, Saurabh Pal, "A novel approach for Breast Cancer detection using Data Mining techniques", *IJIRCCE*, vol. 2, Issue 1, pp. 2320-9798, January 2014.
- [10] Platt, J.C., "Sequential Minimal Optimization: A fast algorithm for training Support Vector Machines." *Technical Report MSR-TR-98-14, Microsoft Research*, 1998.
- [11] Rish I., "An empirical study of the naive Bayes classifier.", IJCAI Work Empir methods Artif Intell. Vol. 3, pp.41-46, 2001.
- [12] Dr.Neeraj Bhargava, Girja Sharma, Dr.Ritu Bhargava, Manish Mathuria, "Decision Tree Analysis on J48 algorithm for data mining", *International Journal of Advanced Research in Computer Science and Software Engineering (IJARCSSE)*, Vol. 3, Issue 6, pp. 1114-1120, 2013.
- [13] Lio Bhreman, Jerom Fridman, Richerd Olshan, Charle Stone, "Regression Tree" (Wardsworth).
- [14] Margaret H. Danham, S. Sridhar, "Data mining, introductory and advanced topics", *Pearson education*, 1st ed., 2006.
- [15] Rajesh, K., and V. Sangeetha. "Application of data mining methods and techniques for diabetes diagnosis." *International Journal of Engineering* and Innovative Technology (IJEIT), Vol. 2, Issue. 3, 2012.
- [16] David, Satish Kumar, Amr TM Saeb, and Khalid Al Rubeaan, "Comparative Analysis of Data Mining Tools and Classification Techniques using WEKA in Medical Bioinformatics.", *Computer Engineering and Intelligent Systems*, Vol. 4, Issue. 13, pp. 28-38, 2013.
- [17] Vaithiyanathan, V., K. Rajeswari, Kapil Tajane, and Rahul Pitale, "Comparison of Different Classification Techniques Using Different Datasets.", Vol.6, Issue. 2, 2013.
- [18] Nikhil.N Salvithal ,Dr.RB Kulkarni, "Evaluating Performance of Data Mining classification Algorithm in WEKA", *International Journal of Application or Innovation in Engineering and Management(IJAIEM)*, Vol.2, Issue 10, 2013.
- [19] Nookala, Gopala Krishna Murthy, Bharath Kumar Pottumuthu, Nagaraju Orsu, and Suresh B. Mudunuri. "Performance analysis and evaluation of different data mining algorithms used for cancer classification." *International Journal of Advanced Research in Artificial Intelligence* (IJARAI), Vol. 2, Issue 5, 2013.
- [20] Tiwari, Mahendra, Manu Bhai Jha, and Om Prakash Yadav, "Performance analysis of Data Mining algorithms in Weka.", *IOSR Journal of Computer Engineering (IOSRJCE)* ISSN (2012): 2278-0661, Vol.6, Issue 3, 2012.
- [21] Bin Othman, Mohd Fauzi, and Thomas Moh Shan Yau, "Comparison of different classification techniques using WEKA for breast cancer.", 3rd Kuala Lumpur International Conference on Biomedical Engineering 2006. Springer Berlin Heidelberg, 2007.
- [22] Aman Kumar Sharma, Suruchi Sahni, "A comparative study of classification algorithms for spam email data analysis", *IJCSE*, Vol. 3, Issue 5, pp. 1890-1895, 2011.
- [23] Phyu and Nu Thair, "Survey of classification techniques in data mining," in Proc. International Multi Conference of Engineers and Computer Scientists, Vol. I, IMECS 2009, March 18 - 20, 2009, Hong Kongclassification.
- [24] Vaithiyanathan, V., K. Rajeswari, Kapil Tajane, and Rahul Pitale, "Comparison of Different Classification Techniques Using Different Datasets.", International Journal of Advances in Engineering & Technology (IJAET), Vol. 6 Issue 2, pp. 764-768, May 2013.
- [25] Ian H.Witten and Elbe Frank, "Data mining Practical Machine Learning Tools and Techniques," *Second Edition, San Fransisco*, 2005.
- [26] R. Agrawal and J.C. Shafer, "Parallel mining of association rules," *IEEE Transactions on Knowledge and Data Engineering*, Vol.8, Issue.6, pp. 962-969, Dec. 1996.
- [27] V.O. Oladokun, A.T. Adebanjo, and O.E. Charles-Owaba, "Predicting Students' Academic Performance using Artificial Neural Network: A Case

Study of an Engineering Course".

- [28] Refaat, M, "Data Preparation for Data Mining Using SAS", *Elsevier*, 2007.
- [29] S. M. Kamruzzaman and A. M. Jehad Sarkar "A New Data Mining Scheme Using Artificial Neural Networks", 2011.

Utilizing Machine Learning Techniques for Early Detection and Control of Powdery Mildew Disease in Cashew Flowers to Enhance Crop Yield

Dr. S.P. Ponnusamy¹ and Mrs. S. Valli²

¹Assistant Professor & ²Guest Lecturer Department of Computer Science, Government Arts and Science College, Tittagudi, Tamil Nadu, India. ¹spponns2k1@gmail.com ²st.blossomvalli@gmail.com

Abstract : The agricultural industry forms the backbone of the Indian economy.Various factors can contribute to a decline in agricultural growth. These may include adverse weather conditions, fluctuating market prices, inadequate infrastructure, and the prevalence of pests and diseases. Enhancing crop yield in cashew flowers necessitates the early detection and effective control of powdery mildew disease. Employing advanced monitoring and intervention strategies can play a pivotal role in mitigating the impact of this disease, ultimately contributing to improved agricultural productivity. Utilizing machine learning algorithms in conjunction with image processing techniques on the cashew nut flower dataset enables the prediction of Powdery Mildew Disease's impact on cashew nuts. This integrated approach harnesses the power of advanced technology to analyze and interpret data, providing valuable insights for early detection and proactive management of the disease, thus contributing to the overall health and productivity of cashew crops.

Keywords : agriculture, cashew nut, machine learning, powdery mildew disease.

I. INTRODUCTION

Agriculture, as the cornerstone of the Indian economy, confronts an intricate web of challenges that require innovative solutions to sustain growth and productivity. Adverse weather conditions, market volatility, inadequate infrastructure, and the omnipresence of pests and diseases collectively cast a shadow on the agricultural landscape [1]. One particularly formidable adversary is powdery mildew disease, wreaking havoc on cashew flower crops and necessitating urgent attention through early detection and robust control measures.

In the pursuit of safeguarding crop yield and fortifying the agricultural sector, this study undertakes a transformative approach, amalgamating cutting-edge technology with age-old agricultural practices. The focus centers on the predictive capabilities of machine learning algorithms and the precision offered by image processing techniques[2] with the overarching

goal of deciphering the impact of Powdery Mildew Disease on cashew nuts.

As the agricultural sector grapples with the escalating complexities of disease management, this research propels the discourse forward by introducing an advanced ensemble model. Leveraging the combined strengths of Support Vector Machine (SVM) and Random Forest algorithms, this model emerges as a beacon of hope in the quest for more accurate and proactive disease management strategies[3] [4] [5]. The use of ensemble [6] [7] SVM and Random Forest not only signifies a departure from singular approaches but also acknowledges the nuanced and multifaceted nature of agricultural challenges.

Against this backdrop, this research focuses on the intricate task of predicting [8] [9] Powdery Mildew Disease's impact on cashew nuts. By leveraging a machine learning ensemble model that intertwines SVM and Random Forest, the study aims to unravel the complexities of disease dynamics and contribute to a more nuanced and context-aware understanding. The ensemble approach is applied to a dataset sourced from cashew nut flowers, offering valuable insights for early detection and proactive disease management.

In the following sections, we delve into the methodology, dataset specifics, and preliminary results, unravelling the potential of ensemble SVM and Random Forest in reshaping the landscape of agricultural disease management. Through this holistic exploration, the research seeks to underscore the transformative impact of technology-driven interventions in fortifying the resilience of India's agricultural sector.

II. REVIEW OF LITERATURE

The landscape of agriculture has witnessed a continual evolution in response to emerging challenges, and the impact of diseases on crop yield has been a recurring theme in scholarly discourse. The literature reveals a consensus on the multifaceted challenges faced by the agricultural sector, emphasizing the need for innovative solutions to enhance crop productivity and disease management.

The challenges confronting agriculture are varied and complex. Weather uncertainties, market fluctuations, inadequate infrastructure, and the omnipresent threat of pests and diseases collectively pose substantial risks to crop health and yield. These [10] challenges necessitate a dynamic and adaptive approach, encouraging researchers and practitioners to explore novel technologies to safeguard agricultural interests.

Disease management forms a critical aspect of ensuring agricultural sustainability. The prevalence of powdery mildew disease in cashew flowers exemplifies [11] the ongoing struggle against agricultural pathogens. Historically, disease management strategies have relied on traditional methods, but the limitations of these approaches have prompted a shift toward technology-driven solutions.

The integration of machine learning (ML) in agriculture has emerged as a promising avenue to address complex challenges[12] [13]. Existing literature attests to the success of ML algorithms in disease prediction, offering a paradigm shift from reactive to proactive disease management. SVM and Random Forest, in particular, have gained prominence for their accuracy and versatility in handling agricultural datasets.

Ensemble learning, a paradigm that combines the strengths of multiple algorithms, has gained traction in disease prediction[14]. The literature underscores the efficacy of ensemble methods in improving prediction accuracy and robustness. The decision to employ an ensemble model comprising SVM and Random Forest aligns with this literature, aiming to harness the complementary strengths of these algorithms for enhanced disease prediction. In summation, the review of literature illuminates the persistent challenges in agriculture, the transformative potential of ML algorithms, and the growing relevance of ensemble methods in disease prediction. The subsequent sections delve into the methodology and dataset specifics, providing a detailed exploration of the application of ensemble SVM and Random Forest in addressing the challenges posed by powdery mildew disease in cashew flowers.

III. METHODS AND METHODOLOGY

A. Dataset Description:

To undertake the machine learning classification of cashew flower powdery mildew, a comprehensive dataset was assembled, comprising 3200 images capturing both healthy and diseased instances. These images, sourced from mobile camera captures and supplemented with a few from Google Images, constituted two categories: "healthy" and "anthracnose"infected images. The dataset creation process involved meticulous labeling, with each image categorized as either "healthy" or assigned a severity level if infected. The methodology further employed augmentation techniques to diversify the dataset, including rotation, flipping, zooming, and adjustments in brightness/contrast. Following data augmentation, normalization of pixel values and resizing of images to a uniform dimension were performed to ensure standardized input features for the model. The dataset was then split into training and validation sets, with an 80-20 ratio for training and validation, respectively.For feature extraction, pretrained convolutional neural network (CNN) models, such as ResNet or MobileNet, were utilized. The final classification layers of the pre-trained models were removed, and the extracted features served as input for the subsequent ensemble model.



Fig 1 Unhealthy and Healthy Cashew Flower (Powdery Mildew)

B. Ensemble Model Development:

The proposed ensemble model combines Support Vector Machine (SVM) and Random Forest algorithms to create a robust and accurate predictor for powdery mildew disease impact on cashew nuts. SVM is chosen for its effectiveness in handling complex, non-linear relationships in data, while Random Forest excels in capturing intricate patterns and avoiding overfitting.

C. Training and Validation:

The labeled dataset is divided into training and validation sets, facilitating the model's learning process. The ensemble model is trained on the training set and validated on a separate

Department of Computer Science, A.V.C. College (Autonomous), Mayiladuthurai – 609 305. PROCEEDINGS 39

dataset to assess its generalization capabilities. Hyperparameter tuning is conducted to optimize the performance of both SVM and Random Forest components.

D. Evaluation Metrics:

The performance of the ensemble model is evaluated using standard metrics such as accuracy, precision, recall, and F1 score. These metrics provide a comprehensive understanding of the model's ability to correctly predict the impact of powdery mildew disease on cashew nuts.

In summary, the methodology encompasses dataset preparation, preprocessing, feature extraction, ensemble model development, training and validation, evaluation metrics, crossvalidation, and ethical considerations. This comprehensive approach aims to yield a reliable and effective model for predicting the impact of powdery mildew disease on cashew nuts in agricultural settings.

IV. RESULT AND DISCUSSION

A. Dataset Overview:

The comprehensive dataset utilized in this study includes a diverse set of images capturing various instances of powdery mildew on cashew nut flowers. The dataset showcases variations in severity and manifestation, providing a robust foundation for the subsequent machine learning analysis.

B. Preprocessing Impact:

Meticulous preprocessing significantly enhances the dataset's suitability for machine learning. Image normalization, resizing, and noise reduction contribute to the clarity and uniformity of the dataset. Labeling based on powdery mildew severity enables supervised learning, ensuring the model learns from categorized instances.

C. Ensemble Model Performance:

The ensemble model, a synergistic combination of Support Vector Machine (SVM) and Random Forest, demonstrates robustness and high accuracy in predicting the impact of powdery mildew disease on cashew nuts. The model leverages the strengths of SVM, known for its proficiency in handling complex relationships within data, and Random Forest, recognized for its ability to capture intricate patterns. This collaborative approach enhances the overall effectiveness of the ensemble model. The ensemble model excels in avoiding overfitting, ensuring that it generalizes well to new and unseen data. The dataset, meticulously divided into training and validation sets, plays a crucial role in facilitating the model's learning process. By training on a labeled dataset and validating on a separate dataset, the model showcases its capability to generalize and make accurate predictions on diverse instances.

Ensemble model is employed to optimize both the SVM and Random Forest components, enhancing their individual and collective performance. The accuracy of the ensemble model reaches 50.88%, outperforming individual models. Precision, which measures the accuracy of positive predictions, is 25.89%, and recall, which assesses the model's ability to capture all relevant instances, is also 50.88%. These metrics collectively indicate the ensemble model's effectiveness in predicting and classifying instances of powdery mildew disease on cashew nuts.

The ensemble model, through its integration of SVM and Random Forest, emerges as a reliable and accurate tool for predicting and managing powdery mildew disease in cashew crops.

TABLE 1 : RESULT OF CLASSIFICATION ALGORITHMS

	Random Forest	SVM	Ensemble
Accuracy	49.12	49.12	50.88
Precision	100	24.13	25.89
Recall	49.12	49.12	50.88



FIGURE 2 RESULT OF CLASSIFICATION ALGORITHMS

The adoption of an ensemble approach, strategically combining Support Vector Machine (SVM) and Random Forest, proves to be highly advantageous. SVM's proficiency in handling non-linear relationships synergizes well with Random Forest's ability to capture intricate patterns. This collaborative integration enhances the model's overall predictive performance, especially concerning the impact of powdery mildew on cashew nuts.

In summary, the presented results underscore the effectiveness of the ensemble model in predicting the impact of powdery mildew disease on cashew nuts. The ensuing discussion emphasizes the significance of the ensemble approach, elucidates the role of feature extraction, delves into the model's generalization capabilities, and examines the impact of chosen evaluation metrics. The ensemble model emerges as a robust and reliable tool for early detection and effective control of powdery mildew disease, thereby contributing to the advancement of agricultural productivity.

PROCEEDINGS

40

V. CONCLUSION

In conclusion, this study introduces an ensemble model leveraging Support Vector Machine (SVM) and Random Forest algorithms for predicting the impact of powdery mildew disease on cashew nuts. The comprehensive dataset, enriched through mmeticulous preprocessing and labelled for supervised learning, serves as a valuable resource for training and validation. The ensemble model demonstrates robustness in handling complex relationships and capturing intricate patterns, attributed to the synergistic combination of SVM and Random Forest. Feature extraction from the cashew nut flower images plays a pivotal role, providing the model with nuanced characteristics of powdery mildew for accurate predictions. Overall, the ensemble model emerges as a valuable tool for early detection and effective control of powdery mildew disease in cashew crops. The integration of advanced machine learning techniques with image processing contributes to improved agricultural productivity. This study underscores the significance of employing ensemble methods in agricultural disease prediction, paving the way for future advancements in precision agriculture and crop management.

REFERENCES

- Bottriell K. MULTISTAKEHOLDER ROUNDTABLES AND VOLUNTARY STANDARDS. SMALLHOLDERS. 2023 Nov 1:405.
- [2] Vuppalapati C. Specialty Crops. InSpecialty Crops for Climate Change Adaptation: Strategies for Enhanced Food Security by Using Machine Learning and Artificial Intelligence 2023 Oct 15 (pp. 35-197). Cham: Springer Nature Switzerland.
- [3] Al-Ramini A. PAD DIAGNOSIS AND ESTIMATION OF TREATMENT EFFECTIVENESS USING MACHINE LEARNING.
- [4] Jadhav P, Kachave V, Mane A, Joshi K. CROP DETECTION USING SATELLITE IMAGE PROCESSING. I-Manager's Journal on Image Processing. 2023 Apr 1;10(2).
- [5] Junaid M, Shaikh A, Hassan MU, Alghamdi A, Rajab K, Al Reshan MS, Alkinani M. Smart agriculture cloud using AI based techniques. Energies. 2021 Aug 19;14(16):5129.
- [6] Reddy J, Devi SS, Parvatham SD, Vishal KS. Optimizing Crop Forecasts: Leveraging Feature Selection and Ensemble Methods. Turkish Journal of Computer and Mathematics Education (TURCOMAT). 2023 Jul 8;14(03):1062-71.
- [7] Astani M, Hasheminejad M, Vaghefi M. A diverse ensemble classifier for tomato disease recognition. Computers and Electronics in Agriculture. 2022 Jul 1;198:107054.
- [8] Shukoor A, Jain S, Maurya P, Anusha C, Kiran B, Kothiyal K. An Overview of Different Fruit Crop Models in the last 40 Years to Date with Their Main Uses. International Journal of Plant & Soil Science. 2023 Dec 4;35(22):618-41.
- [9] Aggarwal PK, Roy J, Pathak H, Kumar N, Venkateswarlu B, Ghosh A, Ghosh D. Indian Agriculture Towards 2030.
- [10] Liu Y, Ma X, Shu L, Hancke GP, Abu-Mahfouz AM. From Industry 4.0 to Agriculture 4.0: Current status, enabling

technologies, and research challenges. IEEE Transactions on Industrial Informatics. 2020 Jun 22;17(6):4322-34.

- [11] LEWIS C. PREVALENCE, ETIOLOGY AND SOURCES OF RESISTANCE AGAINST CASHEW POWDERY MILDEW DISEASE IN WESTERN ZAMBIA (Doctoral dissertation, UNIVERSITY OF ZAMBIA).
- [12] Pokhariyal S, Patel NR, Govind A. Machine Learning-Driven Remote Sensing Applications for Agriculture in India—A Systematic Review. Agronomy. 2023 Aug 31;13(9):2302.
- [13] Karunathilake EM, Le AT, Heo S, Chung YS, Mansoor S. The path to smart farming: Innovations and opportunities in precision agriculture. Agriculture. 2023 Aug 11;13(8):1593.
- [14] Rasool S, Husnain A, Saeed A, Gill AY, Hussain HK. Harnessing Predictive Power: Exploring the Crucial Role of Machine Learning in Early Disease Detection. JURIHUM: Jurnal Inovasi dan Humaniora. 2023 Aug 19;1(2):302-15.

Embarking on a new journey of Federated Learningand TensorFlow Framework

J. Jagadeesan

Assistant Professor of Computer Science, Arignar Anna Govt. Arts & Science College, Karaikal. jagatamil2006@gmail.com

Dr. R. Nagarajan

Assistant Professor, Department of Computer and Information Science, Annamalai University, Chidambaram.

Abstract — Federated Learning (FL) has emerged as a promising paradigm in machine learning, facilitating model training across decentralized devices while safeguarding data privacy. This research paper explores Federated Learning, specifically utilizing the TensorFlow framework with the high- level API of Keras models. The study delves into the principles, methodologies, and challenges associated with Federated Learning, emphasizing its potential applications across various domains. Furthermore, the paper provides a comprehensive overview of the TensorFlow framework and its role in implementing Federated Learning Additionally, practical considerations algorithms. and performance evaluations are discussed.

Through empirical experiments and case studies, this research paper aims to offer valuable insights into model training, privacy preservation, communication efficiency, model aggregation, communication protocols. model compression, secure aggregation, challenges of Non-IID data, Communication Overhead, Heterogeneity, Strategic Behavior, and potential enhancements. The primary goal of this research paper is to further optimize Federated Learning using TensorFlow. In essence, this study serves as a comprehensive guide, providing both theoretical and practical perspectives on Federated Learning. By fostering a deeper understanding of this innovative approach, the research aims to contribute to the ongoing discourse surrounding collaborative machine learning and its applications in real-world scenarios.

Keywords— Federated Learning, TensorFlow, Keras, Challenges of Non-IID data, Model Aggregation.

I. INTRODUCTION

Traditionally, data is gathered from users, and a machine learning model is then trained using that data at the central server. However, this approach exposes personal data to individuals and devices that do not own the data. Nevertheless, a significant portion of existing training data originates from resource-constrained devices, such as tablets and smart phones.

The impracticality of uploading extensive data to the cloud for centralized model training is exacerbated by limitations in communication bandwidth and heightened privacy concerns, to tackle this privacy concern and limitations Google introduced Federated Learning (FL) in 2016. Federated Learning involves training an algorithm locally on a user's device, bringing the algorithm to the data rather than transporting the data to the algorithm [1].

The life cycle of federated learning is start with problem identification, in which define problem state next simulation prototyping of various model selection for suitable for problem with datasets. The third one is client selection and instrumentation groups the client corresponding with datasets for model training. Federated Model training is the next phase just like traditional machine learning model the training is done at client not in server. At the model evaluation is done by using standard metrics such as performance and accuracy, the topperforming models proceed to the subsequent stage of deployment in the overall process.

I. PRELIMINARIES

A. Classical Machine Learning

Under the conventional centralized learning (CL) paradigm, raw data sourced from the environment (example: measurements, audio, images, video, etc.) is data is then transmitted to a central server, where the computationallyintensive model training task is executed. Nevertheless, this approach imposes substantial traffic burdens on the underlying (wireless) network, as complex task training typically involves the transfer of substantial data chunks and concentrates significant processing loads at a single location.[2].

B. Decentralized Machine Learning

This decentralized process is overseen by a central server. After the training is complete, only the modified parameters of the model are transmitted to the central server. The server then consolidates these updated parameters from different devices to generate a global model. Typically, federated learning employs the parameter server architecture, wherein clients train local models that are synchronized through a central parameter server. The federated learning process unfolds in multiple rounds, where clients download the updated machine learning model from the parameter server. Subsequently, they independently train their local models with their respective data over multiple epochs in each round. The figure-1 shows Federated learning architecture. The server will aggregate the global model [3][4].



Fig. 1 : Federated Learning Architecture

The FL has many applications in sales, finance and other industries where data cannot be directly collected because of intellectual property rights, privacy requirements and data security policy. The organization needs of user data of sopping history, to understand their behaviour and shopping methods. Since privacy reason conventional model of machine learning not to be secured. FL solves this issue by bringing algorithm to user data access.

C. FL Development history

The main idea is developed by Google. In recent times, there has been a strong focus on personalizing Federated Learning (FL), emphasizing the security aspects of the domain and addressing statistical challenges. FL poses inherent difficulties such as the reliability of users' devices and data imbalances. While the majority of FL publications have originated from the field of computer science and engineering, the subject has sparked interest in various other disciplines. Notably, numerous papers have been published in decision sciences, which are closely related to business, and economics, as well as physics and astronomy, material sciences, and medicine. The list presented in Table-1 has been directly extracted from Scopus. It's important to note that the sum of the paper count in each subject area may exceed the total number of documents analysed in this chapter, as a single paper can be classified under multiple subject areas.

Sl. No.	Subject area	No. of papers
1	Computer science	422
2	Engineering	204
3	Mathematics	97
4	Decision sciences	77
5	Physics and astronomy	23
6	Materials science	22
7	Medicine	22
8	Social sciences	12
9	Energy	8
10	Biochemistry and genetics	6
11	Health professions	5
12	Business, Mgmt. Accounting	4

Table 1: Subject Area with Highest Number of Fl-Papers

D. Different Themes in FL Research

1) IoT: The IoT has the Machine learning, edge computing, block chain, internet of things, machine learning models, artificial intelligence, data sharing, distributed machine learning, reinforcement learning, network security, distributed computer systems, IoT, transfer learning, IoT, security and privacy, network architecture, 5G mobile communication systems, decision making, quality of service and computation theory in most occurring terms.

The presence of specific terms within the initial group implies that this cluster primarily focuses on the Internet of Things (IoT) and Federated Learning (FL). The IoT encompasses a multitude of sensors, making it an ideal environment for leveraging machine learning techniques to enhance network performance. Given the decentralized nature of the IoT and concerns regarding privacy, this field presents a compelling opportunity for applying Federated Learning. 2) Wireless communication: It has contains Stochastic systems, communication overheads, gradient methods, wireless networks, optimization problems, economic and social effects, optimization, iterative methods, stochastic gradient descent, communication efficiency, energy utilization, incentive mechanism, bandwidth, communication rounds, efficiency, mobile telecommunication systems, resource allocation, signal processing, numerical results and energy efficiency.

Upon examination, it becomes evident that the second cluster is closely associated with wireless communication, particularly emphasizing the challenges posed by high communication costs and data transfer limitations over slow and expensive networks. The cluster's objective centers on predicting subsequent words or phrases to enhance userdevice interaction. However, a significant issue arises when user inputs are required to train models, as these inputs contain personal and sensitive information. To address this concern and simultaneously mitigate communication expenses, the authors propose a Federated Learning (FL) methodology aimed at safeguarding user privacy.

3) Privacy and security: The following terms are present in privacy and security level. Such as Learning systems, data privacy, privacy preserving, privacy, privacy preservation, neural networks, learning models, model parameters, state of the art, adversarial networks, privacy concerns, mobile computing, privacy-preserving, data mining, learning methods, privacy leakages, sensitive data, training process, benchmark datasets and centralized server.

Description of a selective distributed gradient descent method to reduce communication and the application of differential privacy to protect the model parameter updates are contributed at 2015 [5].

Description of an efficient accounting method for accumulating privacy losses while training a DNN with differential privacy are contributed at 2016 [6].

New method to provide secure multi-party computation specifically tailored towards FL, Method for providing user-level differential privacy for FL with only small loss in model utility, Method for providing user-level differential privacy for FL without degrading model utility, Demonstration of an attack method on the global model using a generative adversarial network, effective even against record/batch-level DP, Method for encrypting user updates during distributed training, decryptable only when many clients have participated in the distributed learning objective are contributed at 2017. Description of a fullscale production-ready FL system (focusing on mobile devices are contributed at 2019 [7].

4) Data analytics : Deep learning, learning frameworks, deep neural networks, communication cost, global modeling, classification (of information), fog computing, poisoning attacks, anomaly detection, central servers, intrusion detection, real-world datasets, benchmarking, data distribution, training data, computer aided instruction, data handling, automation, collaborative training, computation costs.

III. APPROACHES OF FL

A) Taxonomy

In light of the shared system abstractions and the assembly of distinct building blocks for FLSs, we organize FLSs through the lens of six aspects, introducing a lossless vertical Federated Learning System (FLS) designed to facilitate collaborative training of gradient boosting decision trees between multiple parties.

B) Federated learning systems

Data Partitioning HorizontalVertical Hybrid Machine Learning Model Linear Model Decision Tree Neural Network Privacy Mechanism Differential Privacy Cryptographic Method Communication Architecture DecentralizedScale of federation Cross Silo Cross Devices Motivation of federation Incentive Regulation

This system utilizes privacy-preserving entity alignment to identify shared users across different parties, leveraging their gradients for joint decision tree training. While traditional FLSs often address singular types of data partition, scenarios like cancer diagnosis systems involve a hybrid of horizontal and vertical data partition. For instance, multiple hospitals aiming to develop an FLS for cancer diagnosis possess varying patient data and medical test results. Addressing such challenges, transfer learning emerges as a potential solution. Liu et al. have proposed a secure federated transfer learning system, detailed in their work, 'Secure Federated Transfer Learning,' which enables feature representation learning across party-specific data instances [8].

C) Important segments present in the FederatedLearning

These segments can refer to various aspects of the federated learning workflow, it ensuring efficient collaboration among decentralized devices. It involves strategies for data distribution, model architecture, and the coordination of updates across various segments to achieve effective and privacy-preserving machine learning.

FL-Algorithms : FedAvg, FedSGD, FedDP, FedProx, FedMeta (these algorithms refers to the methodologies and approaches used in FL to collaboratively train machine learning models across decentralized devices)

FL-Frame work: TTF, PySyft, Flower, FedML, FATE, OpenFL, NVIDIA and IBM-FL are the software libraries and platforms designed to facilitate the implementation and deployment of FL systems. It plays a crucial role in simplifying the development and deployment of federated learning models, providing tools and abstractions that handle the complexities of training models across a network of decentralized devices

Privacy and Security: Differential Privacy in Federated Learning Secure Aggregation, Homomorphic Encryption, Privacy-Preserving Entity Alignment, Federated Transfer Learning, Consent Mechanisms, Data Partitioning Strategies, Model Poisoning and Byzantine Attacks, Secure Initialization and Regulatory Compliance [9].

Aggregation and Client Selection: Federated Averaging (FedAvg), Secure Aggregation, Weighted Aggregation, Quantization and Compression, Adaptive Aggregation Schemes *Personalized FL:* Model Training, Model Updates, Adaptive Learning rates, Aggregation strategies, Context-Awareness, User Constraints

Vertical FL: Offers a privacy-preserving approach to collaborative model training in scenarios where data is naturally partitioned vertically. It finds applications in diverse domains, including healthcare, finance, and collaborative research, allowing entities to collectively improve machine learning models without exposing sensitive information.

Applications of FL: Healthcare, Financial Services, Smart Devices and IoT, Telecommunications, Autonomous Vehicles, Retail and E-Commerce, Manufacturing and Industry 4.0, Energy Management, Education, Government and Public Services

Cross-Device Personalization, Privacy-Preserving AI: FL Transformers: this is commonly associated with a type of neural network architecture, such as the Transformer architecture introduced by Vaswani et al. in the context of natural language processing. It refers to a specific technology, model, or concept introduced after that date

Optimization and Enhancement: Communication Efficiency, Model Compression, Adaptive Learning Rates, Optimizing Hyper parameters, Cross-Silo FL.

IV. BENEFITS OF FL

A) Data Security

All of the data needed to train the model remains ondevice with federated learning. It also reduces the exposure of the data and the attack surface to just one device. Organizations, such as hospitals, can use it for sensitive data computations.

B) Hyper-Personalization

Hyper-personalization in federated learning combines the benefits of both hyper-personalization and federated learning to create personalized models for each user while preserving their data privacy. Instead of training 1 model for all users, we can host 1 model per user. For example, eCommerce platforms can use FL to make product recommendations.

C) Low Latency Real-time Predictions

FL does not require the transmission of local raw data. Instead, both the model and data are present on your device. The client device models are continuously updated using the client input history. As a result, better models can be deployed and tested faster with low latency.

D) Privacy Awareness

Only the model updates are shared, not the raw data. When combined with Differential Privacy and Secure Multi Party Computations (SMPC), it becomes non- identifiable personal data, thus removing any GDPR (General Data Processing Regulation) and CCPA (California Consumer Privacy Act) constraints.

E) Bluntly Cheap

FL can help run user infrastructures with 100s of users at a cheaper cost when compared to the cloud.

F) Hardware Efficiency

FL doesn't need a complicated central server. Smart phones are capable of processing data. The hardware efficiency of a federated learning model refers to how well the model utilizes computational resources, such as CPU, GPU, memory, and network bandwidth, during the training process. Federated learning is inherently distributed, with model training happening across multiple devices or servers. As a result, hardware efficiency plays a critical role in determining the overall performance and effectiveness of federated learning. Several factors contribute to the hardware efficiency of a federated learning model such as Model architecture, Communication overhead, Local device capabilities, optimization algorithm, Hardware acceleration etc.,

G) No Internet Required

The device stores the data and the model. Thus, prediction does not require an internet connection, i.e., to run local models. But, the internet is required to send and receive updates from the central server [10].

V. FEDERATED LEARNING LIBRARIES BENCHMARKS

The significance of benchmarks in guiding the development of Federated Learning Systems (FLSs) cannot be overstated. Numerous works focused on benchmarks have been undertaken recently, resulting in the availability of several benchmark frameworks online. At present, there is a lack of a sufficiently comprehensive benchmark system that encompasses all algorithms or application types within Federated Learning Systems (FLSs). Even the most extensive benchmark systems currently available lack support for certain algorithms and lack evaluation metrics for every level of the system. Achieving further development in comprehensive benchmark systems necessitates substantial support from diverse FL frameworks. Furthermore, most benchmark research employs datasets derived from a single dataset, with no consensus on the preferred splitting method. Similarly, concerning the non-IID problem, there is a lack of consensus on the metric for non-IIDs.

These frameworks can be categorized into three types: i) General-purpose benchmark systems that aim to comprehensively evaluate FLSs, providing a detailed characterization of various aspects; ii) Targeted benchmarks that concentrate on one or more specific aspects within a narrow domain, striving to optimize system performance in that particular domain; iii) Dataset benchmarks dedicated to providing specialized datasets for federated learning.[8]

Following the inception of federated learning by Google, extensive analysis and mining of pipeline components occur based on the federated learning standards introduced by Google researchers. The benchmarks and simulations for federated learning systems, such as Tensorflow Federated (TFF), LEAF, and FedML, form the frameworks. Our findings lead to the conclusion that, while data collection is fairly similar, the stages of data preprocessing, model training, model evaluation, model deployment, and model monitoring in federated learning systems differ significantly from those in traditional machine learning pipelines. Notably, the model training stage in federated learning pipelines encompasses operations like model broadcast, local model training, model upload and collection, and model aggregation within a single stage [11].

A) Open Source System

In this research paper discussed some of the following open source frameworks. Federated AI Technology Enabler (FATE) and Google TensorFlow Federated (TFF).

1) FATE

FATE is an industrial frame work developed by WeBank. FATE is mainly working on Python. The FATE consist of EggRoll (Distributed computing and Storage), Federated ML(Algorithms and Secure protocol), FATE flow(Pipeline and Model Manager), FATE- Board(Visualization), KubeFATE (Deployment and Cluster Management).

2) TFF

TensorFlow, an inclusive open-source ML platform, employs dataflow graphs to articulate the computations, operations, and states inherent in ML algorithms. Serving as a foundational layer for versatile differentiable programming, TensorFlow facilitates gradient computation for arbitrary differentiable expressions and performs efficient low-level tensor operations on CPU, GPU, or TPU. Additionally, its scalability extends to diverse devices, and for deployment, it allows the export of graphs to various external runtime environments, including servers, browsers, mobile, and embedded devices. [12].

TensorFlow Federated (TFF) serves as the foundational framework for Federated Learning (FL) based on TensorFlow, garnering approximately 1.5k stars and 380 forks on GitHub. TFF offers a Python package that can be effortlessly installed and imported.

As depicted in Figure 2, TFF provides two distinct layers of APIs: the FL API and the Federated Core (FC) API. The FL API offers high- level interfaces, comprising three key components—models, federated computation builders, and datasets. This API enables users to define models or load Keras models seamlessly [13]. The federated computation builders encompass typical algorithms like federated averaging. Additionally, the FL API supplies simulated federated datasets for FL.



Fig.2. Layers of TFF

Beyond high-level interfaces, the FC API incorporates lower-level interfaces as the foundation of the FL process. Developers can implement their functions and interfaces within the federated core. FC provides essential building blocks for FL, supporting various federated operators such as federated sum, federated reduce, and federated broadcast. Developers can define their operators to implement FL algorithms.

In essence, TFF is a lightweight system empowering developers to design and implement novel FL algorithms. As of now, TFF does not account for adversaries during FL training and lacks privacy mechanisms. Furthermore, TFF can only be deployed on a single machine, with the federated setting implemented through simulation. TFF plans multi-machine deployments, TFF is tightly coupled with TensorFlow and experimentally supports JAX, LEAF also has a dependency on TensorFlow. TFF provides libraries for constructing baselines with some datasets.

TFF includes base classes for implementing the FedAvg and federated stochastic gradient descent (FedSGD) algorithms, as well as simple implementations for federated evaluation and federated personalization evaluation. The aggregation of client model updates on the server is facilitated through several functions:

The "Sum" function sums values from clients and outputs the sum at the server.

The "Mean" function computes a weighted mean of client values and outputs the mean at the server.

The "Differentially private" function aggregates client values in a privacy-preserving manner based on the differential privacy (DP) algorithm, producing the result at the server.

To enable the creation of new federated algorithms, TFF provides support through a core API consisting of classes that define templates for stateful processes, such as the aggregation of values, computation of estimates, and production of metrics. Analysts can develop their own analytical processes using these templates.[14]

The architecture of TFF is illustrated in Figure 3, and the current version operates exclusively in simulation mode. TFF can be utilized through Google Colaboratory similar to TensorFlow, or for local usage, it requires installation of the TFF package using Python's pip package manage [15].

FedAvg		Fed SGD
Sum	Mean	DP-Queries
Security		DP
Model (NN/RNN/CNN)		
Kera	5	TensorFlow
Run-time (Executor) gRPC/pr		(Client) proto

Fig.3 : Architecture of TFF

VI. KEY CHALLENGES

A) Handling Model Bias

If not done correctly, device selection can bias model updates towards faster and more powerful devices. For example, costly high-end smartphones have better performance when compared to low-end smartphones. Since smartphones with better performance can train the model faster, more updates will come from high-end phones than from other phones when sending an update to the central server. It creates a model bias [16].

1) Model Security-Poisoning Attacks

Two types of poisoning can happen, namely, data and model poisoning. It's tough to detect/prevent malicious clients from delivering harmful/fake data during an FL training process, which leads to data poisoning. In contrast, they may modify the model before sending it back to the central server. Hence, it leads to model poisoning. Poisoning attacks in automated vehicles are a classic example of ML poisoning. In these attacks, the attacker adds noise to the image input to train the ML model, causing the automated vehicle system to misclassify the traffic sign. Adding a small square-shaped piece of black tape to the stop signal board, for example, can be used to cause accidents.

2) No Large Scale Simulations

There are no excellent libraries to test and tune FL algorithms at scale before taking them into production. As a result, the diversity of the devices shoots up so high that it becomes difficult for the developers, especially in startups and small tech companies, to optimize them.

3) System's Heterogeneity

As edge devices regularly fall off due to connectivity or energy restrictions, any device may be unreliable. As a result, fault tolerance is essential, as appliances may drop out before the training cycle completes. Devices in a federated network may have vastly different storage, processing, and communication capabilities. Federated learning systems must accommodate low participation, accept various technologies, and endure network device failures.

4) Statistical Heterogeneity

Mobile devices produce and collect data in a nonidentically distributed manner throughout the network. Different users, for example, utilize other emojis for various purposes. As a result, there's a chance to misinterpret the interaction between devices and their associated data distributions. It contradicts typical IID assumptions and can complicate modeling, analysis, and assessment.

5) Data Privacy

Federated learning protects user data by sharing model updates (for example, gradient information) rather than raw data. However, reporting model updates to a third party or the central server may disclose sensitive information to a third party or the central server during the training process. Other privacy techniques, such as Differential Privacy (simple algorithm), Homomorphic Encryption (computation on encrypted data), and Secure Multiparty Computation (SMPC) to keep the original accuracy while providing a high level of anonymity can solve the above problem [4].

VII. TENSORFLOW FEDERATED: MACHINE LEARNING ON DECENTRALIZED DATA

A) The following steps are needed for machine learningusing TensorFlow

Step 1: Import Tensorflow libraries *Step 2:* Load simulation data

Step 3: Pick a subset of client devices to participate in training

Step 4: Wrap a Keras model for use with TFF

Step 5: Simulate a few rounds of training with the selected client devices

Step 6: Builds all the TensorFlow graphs and serializes them

Step 7: Test datasets for evaluation with Keras by creating a Dataset of Datasets

Step 8: The state of the FL server, containing the model andoptimization state by n rounds

Step 9: Load our pre-trained Keras model weights into the global model state

Step 10: Set the newly trained weights back in the originallycreated model

In attention is directed towards non-convex neural network objectives, the algorithm under consideration demonstrates applicability to any finite-sum objective characterized in the specified form

$$\min_{w \in \mathbb{R}^d} f(w) \quad \text{where} \quad f(w) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n f_i(w).$$
(1)

For a machine learning problem, typically takes $f_i(w) = l(x_i, y_i; w)$, that is, the loss of the prediction on example $(x_i$

, y_i) made with model parameters w. Let us assume there are K clients over which the data is partitioned, with P_k the set of indexes of data points on client k, with $n_k = |P_k|$. Thus, we can re-write the objective (1) as

$$f(w) = \sum_{k=1}^{K} \frac{n_k}{n} F_k(w) \text{ where } F_k(w) = \frac{1}{n_k} \sum_{i \in \mathcal{P}_k} f_i(w).$$
(2)

In the scenario where the partition P_k is created by uniformly distributing training examples among clients at random, the expectation $EP_k[F_k(w)]$ equals the underlying objective function f(w), with the expectation taken over the set of examples assigned to a specific client k. This aligns with the commonly assumed IID (independent and identically distributed) condition in distributed optimization algorithms. Conversely, in the non-IID setting, where F_k may deviate significantly from f, this assumption does not hold true [17].

B) The FedAvg Algorithm

The numerous recent triumphs in deep learning applications predominantly hinge on the utilization of stochastic gradient descent (SGD) and its variants for optimization. In essence, several breakthroughs can be interpreted as modifications to the model's structure, including the adaptation of the loss function, to enhance its compatibility with optimization through straightforward gradient-based methods

Algorithm: Federated Averaging

The K clients are indexed by k; B is the local minibatch size, E is the number of local epochs, and η is the learning rate.Server executes:

initialize w₀

for each round $t = 1, 2, \ldots$ do

 $m \leftarrow max(C \cdot K, 1)$

 $S_t \gets (random \ set \ of \ m \\ clients) for \ each \ client \ k \in St \ in$

parallel do

$$\sum_{\substack{k \\ l+1}}^{k} \leftarrow \text{ClientUpdate}(k, w_t)$$

$$\begin{array}{l} \mathbf{m}_{t} \leftarrow \sum_{k \in \mathrm{St}} \mathbf{n}_{k} \\ \mathrm{wt+1} \leftarrow \quad k \in \mathrm{St} \ \mathrm{nk/mt} \ \begin{array}{c} k \\ \omega \end{array} \end{array}$$

ClientUpdate(k, w): // Run on client k

 $B \leftarrow$ (split Pk into batches of size

B) for each local epoch i from 1 to

E dofor batch $b \in B$ do

 $\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla \ell$ (w;

b)return w to server

The above algorithm used to calculate federated averaging on learning mechanism.

This command builds all the TensorFlow graphs and serializes them

fed_avg = tff.learning.algorithms.build_weighted_fed_avg(
 model_fn=create_tff_model,

client_optimizer_fn=lambda:

tf.keras.optimizers.SGD(learning_rate=0.5))

C) The Keras Model

Keras, with its user-friendly design, is particularly wellsuited for individuals lacking an extensive background in Deep Learning; it offers a straightforward approach to constructing neural network models swiftly and effortlessly with minimal code, facilitating rapid prototyping Keras provides a range of APIs for defining neural networks:

The Sequential API, allowing the creation of a model layer by layer in a simple list format, suitable for most problems but limited to single-input, single-output layer stacks.

The Functional API, a comprehensive alternative supporting arbitrary model architectures, offering greater flexibility and complexity compared to the Sequential API.

Model Sub classing, enabling the implementation of models from scratch, a choice well-suited for research and highly complex scenarios, although less commonly employed in practical applications [18].

The following code gives a Neural Network with Keras' Sequential API

- 1. from keras.models import Sequential
- 2. from keras.layers import Dense
- 3. model = Sequential()
- 4. model.add(Dense(2, input_dim=1, activation='relu'))
- 5. model.add(Dense(1, activation='sigmoid'))
- 6. print(model.summary())

This code snippet utilizes the Keras library to create a simple neural network model. Here's a breakdown of each line:

from keras.models import Sequential: This line imports the Sequential class from the Keras models module. The Sequential class is a linear stack of layers that can be easily created by adding one layer at a time.

from keras.layers import Dense: This line imports the Dense layer, which is a standard fully connected neural network layer in Keras. It is commonly used for adding densely connected layers to a neural network.

model = Sequential(): This line initializes a sequential model. The variable model is now an instance of the Sequential class, which allows you to add layers in a sequential fashion.

model.add(Dense(2, input_dim=1, activation='relu')):

This line adds a dense layer to the model. The layer has 2 units (neurons), an input dimension of 1, and uses the rectified linear unit (ReLU) activation function. The input_dim parameter specifies the input shape for the first layer in the model.

model.add(Dense(1, activation='sigmoid')): This line adds another dense layer to the model. This time, the layer has 1 unit and uses the sigmoid activation function. The sigmoid activation is often used in the output layer of binary classification models.

print(model.summary()): This line prints a summary of the model's architecture. The summary includes information about each layer, such as the layer type, output shape, and the number of parameters. This is helpful for quickly understanding the structure of the neural network.

In summary, this code creates a simple neural network with two dense layers using the Keras Sequential model. The first layer has 2 units, uses the ReLU activation function, and takes input of dimension 1. The second layer has 1 unit with a sigmoid activation function. The model.summary() call provides a concise overview of the model architecture.[19].

VIII. DISCUSSION

Compared Random Forest, Logistic Regression and Multi-Layer Perceptron with TensorFlow-FF. The TFF give better Accuracy, maximize CPU and memory utilization. Exploring the expansive realm of hyper parameters, both within ML models and federated learning, through a Bayesian optimization process constrained to 100 configuration evaluations still allows ample opportunity to discover superior models in terms of accuracy, resource utilization, or both. Widely adopted among ML practitioners, TensorFlow stands out as a prominent framework. However, when users utilize the framework without delving into its inner workings, we present evidence suggesting that numerous latent issues, despite lacking apparent symptoms, can exert significant and sometimes dramatic impacts.

The Keras API encompasses Model and Layer components, prompting us to emphasize additional attention on testing these elements. Notably, the Engine holds particular significance, given its pivotal role in supporting various dependent modules. Neglecting thorough testing of these components might result in misleading issues, where a module could exhibit incorrect behavior, concealing the root cause within its dependencies [20].

IX. CONCLUSION

In conclusion, this research involves into the intricacies of Federated Learning (FL) within the TensorFlow framework, specifically employing Keras models. Through a thorough exploration of principles, methodologies, and challenges, the study sheds light on the potential applications of FL across diverse domains. The comprehensive overview of TensorFlow's role in implementing FL algorithms, coupled with practical considerations and performance evaluations, adds depth to our understanding. The empirical experiments and case studies conducted aim to provide valuable insights into crucial aspects such as model training, privacy

preservation, communication efficiency, and more. The research paper serves as a guide, bridging theoretical and practical perspectives on FL, with a central focus on optimization using TensorFlow. By addressing challenges like Non-IID data, Communication Overhead, Heterogeneity, and Strategic Behavior, the study contributes to the ongoing discourse on collaborative machine learning. Ultimately, this research strives to enhance our comprehension of FL, fostering its application in real-world scenarios and advancing the field of collaborative machine learning.

Reference

- Farooq, A., Feizollah, A., & Rehman, M. H. U. (2021). Federated Learning Systems Towards Next-Generation AI. Studies in Computational Intelligence 965 https://link.springer.com/book/10.1007/978-3-030-70604-3
- [2] Georgios Drainakis; Konstantinos V. Katsaros; Panagiotis Pantazopoulos; Vasilis Sourlas(2021) Federated vs. Centralized Machine Learning under Privacy-elastic Users: A Comparative Analysis. from https://ieeexplore.ieee.org/document/9306745
- [3] Zeng, R. (2021, June 27). A Comprehensive survey of incentive mechanism for federated Learning. arXiv.org. https://arxiv.org/abs/2106.15406
- [4] Nimbleedge. (n.d.). Retrieved from https://www.nimbleedge.ai/
- [5] Shokri, R., & Shmatikov, V. (2015). Privacy-Preserving Deep Learning. Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security - CCS '15.
- [6] Abadi, M., Chu, A., Goodfellow, I. J., McMahan, H. B., Mironov, I., Talwar, K., & Zhang, L. (2016). Deep Learning with Differential Privacy. https://doi.org/10.1145/2976749.2978318.
- [7] Jin, Y., Zhu, H., Xu, J., & Chen, Y. (2023). Federated Learning. https://doi.org/10.1007/978-981-19-7083-2.
- [8] Qinbin Li , Zeyi Wen, (2021) A Survey on Federated Learning Systems: Vision, Hype and Reality for Data Privacy and Protection, https://www.ieee.org/publications/rights/index.html
- Bagdasaryan, E. (2018, July 2). How to backdoor federated learning.arXiv.org. https://arxiv.org/abs/1807.00459
- [10] Kaur, J. (2023, June 26). Federated Learning Applications and its Working. Retrieved from https://www.xenonstack.com/ blog/federated-learning-applications
- [11] Stefan Biffl, Elena Navarro, Welf Löwe, Marjan Sirjani, Raffaela Mirandola, Danny Weyns (Eds.) *Software Architecture*, 15th European Conference, ECSA 2021, Virtual Event, Sweden, September 13–17, 2021.
- [12] Florian Tambon, Amin Nikanjam, Le An, Silent Bugs in Deep Learning Frameworks: An Empirical Study of Keras and TensorFlow, 1 Sep 2023, arXiv:2112.13314v2
- [13] A. Gulli and S. Pal, *Deep Learn. With Keras.* Birmingham, U.K.: Packt Publishing Ltd, 2017 Proceedings
- [14] Choudhury, O. (2019, October 7). Differential privacy-enabled federated learning for sensitive health data. Retrieved from https://arxiv.org/abs/1910.02578
- [15] Ahmed Saidani (January 15, 2023), A Systematic Comparison of Federated Machine Learning Libraries, https://wwwmatthes.in.tum.de/
- [16] Ray, P. P. (2023). ChatGPT: A comprehensive review on background, applications, key challenges, bias, ethics, limitations and future scope. Internet of Things and Cyber-Physical Systems, 3, 121–154. https://doi.org/10.1016/j.iotcps.2023.04.003
- [17] McMahan, H. B. (2016, February 17). Communication-Efficient Learning of Deep Networks from Decentralized Data. Retrieved from https://arxiv.org/abs/1602.05629

- [18] Ahmed Gad, A. (2021, December 14). Breaking privacy in federated learning - Heartbeat. Medium. Retrieved from https://heartbeat.comet.ml
- [19] Remi M (2021, May 25). What is a Keras model and how to use it tomake predictions. Retrieved from https://www.activestate.com/resources/quick-reads/what-is-a-keras-

model/

[20] Preuveneers, D. (2023). AutoFL: towards AutoML in a federated learning context. Applied Sciences, 13(14), 8019. https://doi.org/10.3390/app13148019

Prediction of Diabetes by Using Intellectual Health Care with Using Machine Learning Algorithms

T. Ramyaveni¹ and Dr.V. Maniraj²

¹Research Scholar, Department of Computer Science, A.V.V.M Sri Pushpam College (Autonomous) Poondi, Thanjavur. Affiliated to bharathidasan University, Tiruchirapalli, Tamil Nadu, India. veniramya5@gmai1.com

²Associate Professor & Research Supervisor, Department of Computer Science, A.V.V.M Sri Pushpam College (Autonomous) Poondi, Thanjavur. Affiliated to bharathidasan University, Tiruchirapalli, Tamil Nadu, India

Abstract - A healthcare system using up-to-date computing techniques is the maximum see the sights in healthcare research. Diabetesis one of the elementary sicknesses which has been complexities correlated to it. A huge volume of medical data is formed. It is significant to collected and store, learn and predict the health of such patients using continuous monitoring and technical modernizations. In order to respond to the requirements of upcoming intelligent e-health applications. To progress intelligent healthcare systems and inflate the number of applications associated to the network. This proposed presentation intelligent architecture for monitoring diabetic patients by using machine learning algorithms. The architecture elements included smart devices, sensors, and smart-phones to collect measurements from the body. The intelligent system collected the data received from the patient, and performed data classification using machine learning in order to make a diagnosis. The performance and accuracy of the applied algorithms are associated to indicate the best one in terms of numerous parameters.

Keywords - Diabetes, machine learning algorithms, the Internet of Things (IoT),Principal Component Analysis, Random Forest (RF), Convolution Neural Network (CNN),Support Vector Machine

I. INTRODUCTION

This evolution is based on the use of technologies and applications of the Internet of Things (IoT). It combines the information and communication technologies (ICTs), the use of sensors, the generation of massive data and the application of big data, machine learning techniques, and artificial intelligence [1], whose number has increased in recent years. Thus, IoT technology provides new solutions for diabetic patients. In the field of intelligent health, there are several applications that aim to improve care and improve the quality of life of patients with chronic diseases. Using IoT, the mobile health service becomes more important as it plays a very important role in monitoring and controlling patients who suffer from chronic diseases such as cardiovascular disease and diabetes [2]

Diabetes is a disease which is detected on a blood test when the blood sugar is higher than normal value i.e. between (72 to 99 mg/dL) when fasting and up to (140 mg/dL)2 h subsequent to eating. Naturally, the pancreas emancipates insulin to assist the body to stock and use the sugar fat from the food eaten. Periodically body doesn't make sufficient insulin or doesn't take insulin well. Glucose then remains in blood and doesn't stretch out at cells [3]. There are 3 categories of diabetes Type 1, Type 2 and Gestational diabetes. About 10% of all diabetes cases are type 1 the body don't generate insulin in this compose and about 90% of every of cases of diabetes global are of Type 2 the body don't provide sufficient measure of insulin for absolute purpose. Diabetes effect females during pregnancy are known as Gestational. Diabetes will be he seventh driving reason for death in 2030 as predicted by WHO [4]. To avoid and reduce the complications due to diabetes, a monitoring method of BG level plays a prominent role[5]. combination of biosensors and advanced information and communication technology (ICT) provides an efficient real-time monitoring management system for the health condition of diabetic patients by using SMBG (self-monitoring of blood glucose) portable device[6].

The rest of the paper is organized as follows. Section 2 includes the related work. Section 3presents the proposed architecture for diabetic patient monitoring using 5G. Section 4 describes the system implementation. Section 5 provides a brief description of data collection and Section 6 presents the results and discussion. Finally, conclusions and future work are outlined in Section 7.

II. Related Work

The section also includes some existing works focused on big data and predictive analytics in healthcare that use classification to predict possible episodes of rises or falls in the blood sugar level. Classification in e-health monitoring plays a vital role in the further treatment of the disease. A description of the proposed systems and the used the used algorithms in this work are given. In their work, the authors presented an intelligent architecture for the surveillance of diabetic disease that monitor the health of diabetic patients through sensors integrated into smartphones [7].

Z. Mian. Et. Al described this is a dangerous disease that is lately becoming one of the leading causes of death in the world, and which requires a lot of careful monitoring to keep patients healthy. Diabetes are caused by insulin resistance, and insufficient insulin production can lead to either an increase or decrease in the level of glucose in the blood, therefore the main challenge of the diabetic patient is to maintain the glucose level stable within a specific interval. If they can no longer comply with these conditions, some patients require urgent care to avoid worsening [8].

Maniruzzaman et al. used a machine learning paradigm to classify and predict diabetes. They utilized four machine learning algorithms, i.e., naive Bayes, decision tree, Ada Boost, and random forest, for diabetes classification. Also, they used three different partition protocols along with the 20 trials for better results. They used US-based National Health and Nutrition Survey data of diabetic and non-diabetic individuals and achieved promising results with the proposed technique [9]. Regarding many researchers have developed and implemented various analysis and prediction models using different data mining techniques. To use a classification technique Decision Tree algorithm by using the Weka tool to find out patterns from the diabetes data sets [10].

Rodríguez et al. suggested an application for the smartphone, which can be used to receive the data from the sensor using a glucometer automatically. To checking the patient's glucose level and heart rate using sensors will produce colossal data, and analysis on big data can be used to solve this problem [11].

Wang et al. have given a general idea of the up-to-date BLE technology for healthcare systems based on a wearable sensor. They suggested that low-powered communication sensor technologies such as a BLE device can make it feasible for wearable systems of healthcare because it can be used without location constraints and is light in weight. Moreover, BLE is the first wireless technology in communication for healthcare devices in the form of a wearable device that meets expected operating requirements with low power, communication with cellular directly, secure data transmission, interoperability, electronic compatibility, and Internet communications [12].

III. METHODOLOGY

The proposed diabetes classification and prediction system has exploited different machine learning algorithms.

Dataset

The Pima Indians Diabetes dataset is originally from the National Institute of Diabetes and Digestive and Kidney Diseases. The objective of the dataset is to diagnostically predict whether or not a patient has diabetes, based on certain diagnostic measurements included in the dataset [13]. The datasets consist of several medical predictor variables and one target variable, Outcome. Predictor variables includes the number of pregnancies the patient has had, their BMI, insulin level, age, and so on.

Data collected from the network is an indicator of a patient's health condition. As LoRa protocol is used for communication [12], LoRa enabled sensors to sense the diabetes patient's health condition. With the help of a glucose sensor, we can measure the amount of glucose in the patient's blood. The functioning of the glucose sensor is closely related to the continuous glucose monitoring system (CGM) and depending on the requirement it can be placed externally on the skin or internally under the skin.

Data is transmitted from the patient to the gateway using LoRa protocol. The gateway collects this data and passes it to the patient's dataset unit. The sensed data can be sent periodically or when there is a significant change in the biomedical sensor readings of the patients. The proposed hybrid enhanced adaptive data rate (HEADR) algorithm provides a better implementation of the LoRa device and solves the device data transmission range allocation problem. It is important to determine the throughput and number of packet collisions, to assess the network performance.

IV. DATA PREPROCESSING

Normalization we must carry out the normalization step following the data cleaning phase, where we must divide the entire dataset into training and testing models. We will set aside the test dataset because the data are in split form and apply the training method to the training dataset. Consequently, with the assistance of this training procedure, a training model will be produced that will function on the values of the features in the training data, logic, and algorithm. Bringing all of the attributes to the same standard is the aim of normalization.

Principal Component Analysis

PCA obtains the K vectors and unit eigenvectors by solving the characteristic equation of the correlation matrix of the observed variables. The eigen values are sorted from large to small, representing the variance of the observed variables explained by K principal components, respectively The model for extracting principal component factors is:

F

$$F_i = T_{i1}X_1 + T_{i2}X_2 + T_{ik}X_k$$
 (i=1,2,...,m)

Where, F_i is the i principal component factor; T_{ij} is the load of the i principal component factor on the j index; m is the number of principal component factors; k is the number of indicators[14].

Diabetes Classification Techniques. For diabetic classification, we fine-tuned three widely used state-of-the-art techniques. Mainly, a comparative analysis is performed among the proposed techniques for classifying an individual in either of the diabetes categories. -e details of the proposed diabetes techniques are as follows.

Random Forest (RF). As its name implies, it is a collection of models that operate as an ensemble. -e critical idea behind RF is the wisdom of the crowd, each model predicts a result, and in the end, the majority wins. It has been used in the literature for diabetic prediction and was found to be effective [15]. Given a set of training examples X = x1, x2, ..., xm and their respective targets Y = y1, y2, ..., ym, RF classifier iterates B times by choosing samples with replacement by fitting a tree to the training examples. -e training algorithm consists of the following steps

- (i) For b =1...B, sample with replacement n training examples from X and Y.
 - (ii) Train a classification tree f_b on X_b and Y_b .

Convolutional Neural Network (CNN)

The CNN is one of the most commonly used DLalgorithms. It is a specific type of artificial neural network that uses several layers of perceptron connected in sequence [16].CNNs perform a series of operations on the input and transform it to produce the desired output. This output from previous layers can be taken as input to the next block. CNNs basically consists of three main types of layers, namely convolutional layer, pooling layer and fully connected Convolutional layer forms the core part of the network, which has local connections and weights of shared characteristics. The objective is to learn feature representations of the inputs data. The input feature maps are be easily separated. The most commonly used kernel functions include radial basis function (RBF), polynomial, sigmoid etc.

V. EXPERIMENTAL RESULT

The proposed diabetes classification and prediction algorithm is evaluated on a publicly available Indian Diabetes dataset Besides, a comparative analysis is performed with state-of-theart algorithms. Experimental results show the supremacy of the proposed algorithm as compared to state-of-the-art algorithms. details of the dataset, performance measures, and comparative. first convolved with a kernel and then the obtained results are passed into a nonlinear activation function. The pooling layer can be considered as a fuzzy filter; it reduces the feature dimensionality and increases their robustness.



Fig 1.Convolutional Neural Network (CNN)

Support Vector Machine

The SVMs have proven to be very effective for various data classification tasks. It tries to find the optimal separating hyper plane between classes by finding the set of points that lie on the edge of the class descriptors. The distance between the classes is referred to as margin. SVM algorithms finds a margin such that its distance is maximum. The higher the margin, the better the classification accuracy can be obtained for the classifier

SVMs are designed to deal with linearly separable binary classification data. Several variations have been proposed to adopt it for multi-class classification problems. Similarly, it can also be applied for classification of nonlinear cases by applying kernels techniques [17]. These kernels apply mapping from nonlinear to a linear space where it is believed that the data could



Fig2.Support vectors in SVM

Data classification performance is measured by accuracy, sensitivity, specificity,

and precision. We define accuracy using the following equation: $Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$ where TP is the value of true positive rate, TN is the value of true negative rate, FN is the value of false negative rate, and FP is the value of false positive rate.

Specificity is defined as the ratio between the value of true negatives and the sum of the total value of true negatives and false positives.

Specificity =
$$\frac{\text{TN}}{\text{TN + FP}}$$

Sensitivity is defined as the ratio between the value of true positives and the sum of the total value of true positives and false negatives.

Sensitivity =
$$\frac{TP}{TP + FN}$$

Correctly and Incorrectly Classified Instances

Algorithms	Correctly	incorrectly	Time (S)
	Classified	Classified	
	Instances (%)	Instances (%)	
Random Forest	89.23	11.01	0.03
Convolutional	92.16	11.05	0.05
Neural Network			
Support vector	96.51	12.03	0.06
machine			



Fig3.: shows the classification accuracy value of all classifiers.

VI. CONCLUSION

In this study, several machine-learning algorithms are applied for classification on a data set. So, in this study, we have working Three main algorithms: Random Forest, convolution neural network and support vector machine on the diabetic datasets. These algorithms have been used for experimentation on WEKA tool to predict Diabetic patient data. We tried to compare the efficiency and the effectiveness of the cited algorithms in terms of accuracy, precision and sensitivity. The most objective is to choose the best classification accuracy. The overall performance of the support vector machine algorithm to predict diabetes disease is better than algorithms. In the future work will focus on the integration of other methods into the used model for tuning the parameters of models for better accuracy

REFERENCES

- Rghioui, A.; Lloret, J.; Parra, L.; Sendra, S.; Oumnad, A. Glucose Data Classification for Diabetic Patient Monitoring. Appl. Sci. 2019, 9, 4459.
- [2] Yuehong YIN, Y. Z," The internet of things in healthcare: An overview". Journal of Industrial Information Integration, 2016.
- [3] Sharma, Neha, and Ashima Singh. "Diabetes detection and prediction using machine learning/IoT: A survey." Advanced Informatics for Computing Research: Second International Conference, ICAICR 2018, Shimla, India, July 14–15, 2018, Revised Selected Papers, Part I 2. Springer Singapore, 2019.
- [4] G. Acciaroli, M. Vettoretti, A. Facchinetti, and G. Sparacino, "Calibration of minimally invasive continuous glucose monitoring sensors: state-of-theart and current perspectives," Biosensors, vol. 8, no. 1, 2018.
- [5] M. J. Davies, D. A. D'Alessio, J. Fradkin et al., "Management of hyperglycaemia in type 2 diabetes, 2018. a consensus report by the American diabetes association (ada) and the european association for the study of diabetes (easd)," Diabetologia, vol. 61, no. 12, pp. 2461–2498, 2018.
- [6] D. Bruen, C. Delaney, L. Florea, and D. Diamond, "Glucose sensing for diabetes monitoring: recent developments," Sensors, vol. 17, no. 8, 2017.
- [7] N. M. S. kumar, T. Eswari, P. Sampath, and S. Lavanya, "Predictive Methodology of Diabetic Data Analysis in Big Data,"ProcediaComput. Sci., vol. 50, pp. 203–208, Jan. 2015.
- [8] Z. Mian, K. L.Hermayer, A. Jenkins, "Continuous Glucose Monitoring: Review of an Innovation in Diabetes Management", The American Journal of the Medical Sciences Vol. 358, pp: 332-339, Issue 5, November 2019.
- [9] M. Maniruzzaman, M. J. Rahman, B. Ahammed, and M. M. Abedin, "Classification and prediction of diabetes disease using machine learning paradigm," Health Information Science and Systems, vol. 8, no. 1, pp. 7–14, 2020.
- [10] Luca Catarinucci, Danilo de Donno, Luca Mainetti, Luciano Tarricone, An IoTAwarArchitecture for Smart Healthcare Systems, IEEE Internet Things J. 2 (6) (2015) 515–526.
- [11]I. R. Rodríguez, M. Á. Z. Izquierdo, and J. V. Rodríguez, "Towards an ictbased platform for type 1 diabetes mellitus management," Applied Sciences, vol. 8, no. 4, 2018.
- [12] Verma, Navneet, Sukhdip Singh, and Devendra Prasad. "Machine learning and IoT-based model for patient monitoring and early prediction of diabetes." Concurrency and Computation: Practice and Experience 34.24 (2022): e7219.
- [13] https://www.kaggle.com/datasets/uciml/pima-indians-diabetes-database.
- [14] Wang, S. P., Zhang, Q., Lu, J., and Cai, Y. D. (2018). Analysis and prediction of nitrated tyrosine sites with the mRMR method and support vector machine algorithm. Curr. Bioinform. 13, 3–13. doi: 10.2174/1574893611666160608075753
- [15]N. Yuvaraj and K. R. SriPreethaa, "Diabetes prediction in healthcare systems using machine learning algorithms on hadoop cluster," Cluster Computing, vol. 22, no. 1, pp. 1–9, 2019.
- [16]S. Albawi and T. A. Mohammed, "Understanding of a Convolutional Neural Network," no. April, 2018.
- [17]N. Cristianini, J. Shawe-Taylor, and others, an introduction to support vector machines and other kernel-based learning methods. Cambridge university press, 2000.

Prediction of Diabetic Retinopathy using Machine Learning with Deep Learning Models

P. Kalaimagal¹ and **Dr. S. Kumaravel²** ¹*Research Scholar* and ²*Associate Professor*

Department of Computer Science, A.V.V.M Sri Pushpam College (Autonomous) Poondi, Thanjavur. Affiliated to Bharathidasan University, Tiruchirapalli, Tamil Nadu, India. kalaimagalpannirselvam@gmail.com

Abstract - The aim of this research is to suggest a method that utilizes machine learning to forecast the likelihood of developing diabetic retinopathy (DR). Diabetes is a medical disease that can lead to diabetic retinopathy. The length of time a person has diabetes after its onset greatly influences how quickly retinopathy progresses and eventually results in blindness. In order to uncover research gaps in the detection and classification of diabetic retinopathy, this study looked into various machine learning algorithms that are selected for retinopathy diagnosis. Being aware of the early clinical indicators of depression-related rage (DR) is crucial to managing DR effectively. Regular eye exams are therefore essential in order to refer the patient to a physician as soon as possible for a thorough ocular examination and treatment. Diabetic retinopathy is a disease that needs to be confirmed before the best treatment can be selected. Therefore, it is necessary to have an efficient screening mechanism. In this study, we are looking at a deep learning method, in particular, DenseNet -169, which is a connected convolutional network. Based on the severity level, the fundus image is classified as no DR, mild, proliferative, severe, and moderate DR. In our proposed method, we collected data, preprocessed it, augmented it, and modeled it. Our proposed model has an accuracy of 90%. In addition, we also used regression analysis. The results showed an accuracy of 78%. The main objective of this study was to provide a robust system for automatically detecting Diabetic Retinopathy.

Keywords - Diabetic Retinopath, Machine Learning, Deep Learning, Convolutional Neural Network.

I. INTRODUCTION

The number of people with diabetes is increasing exponentially around the world, the number of cases of diabetic retinopathy (DR), one of the main complications caused by diabetes, is also increasing rapidly. Typical symptoms of type 1 diabetes include frequent urination or bedwetting, constant hunger, excessive thirst, lack of energy or fatigue, blurred vision, sudden weight loss, and diabetic ketoacidosis. Type 2 diabetics also exhibit similar symptoms to type 1 diabetics, but generally people's condition is asymptomatic and less dramatic. As a result, half of the people suffer from type 2 diabetes due to an unrecognized diagnosis and remain in a pre-diabetic state. If this continues for a long time, health complications such as kidney disease, neuropathy, retinopathy, leg ulcers, peripheral artery disease, heart disease and stroke, very slow healing of leg ulcers, and vision problems occur. Without supervision, DR can worsen vision and lead to partial or complete blindness. As the number of diabetic patients continues to increase rapidly in the coming years, the number of qualified ophthalmologists will also need to increase at the same time to meet the increasing demand for preventive testing for diabetic patients. Therefore, it is important to develop ways to automate the DR discovery process.. Computer-aided diagnostic systems have the potential to significantly reduce the burden on ophthalmologists. Diabetic retinopathy is a very serious health problem that affects many people of different age groups. High blood sugar levels can damage the small blood vessels in the retina very quickly, which can lead to retinal detachment and, in some cases, blindness from glaucoma. This study provides an unprecedented and complete overview of all recent research on DR, which improves the understanding of all current research on automated DR detection[1], especially research using machine learning algorithms There is a possibility. Diabetes is becoming an increasingly common disease, with type II diabetes reaching epidemic proportions and type I diabetes steadily increasing worldwide.

Diabetes is associated with a variety of health complications and life-limiting diseases, including chemical heart disease, peripheral artery disease, and macrovascular disease such as stroke. In addition, various microvascular diseases such as diabetic neuropathy, retinopathy, and nephropathy may also occur. Diabetic retinopathy (DR) is a health complication that occurs as a result of poor blood sugar control and long-term diabetes. As the name suggests, DR poses health risks to the eyes and ultimately leads to vision loss. Diabetes affects insulin production and sensitivity, which in turn affects the body's ability to absorb glucose, resulting in high blood sugar levels. If left untreated, high blood sugar levels can damage the eye's retinal blood vessels. When the capillaries that transport blood to and from the eye become clogged, blood flow to the retina decreases. The need for a continuous source of oxygen naturally stimulates angiogenesis, forming new blood vessels for blood transport. The ruptured blood vessels fail due to the lack of oxygen that the eye was previously exposed to, but for some reason instead of repairing the angiogenesis, it causes further complications and retinal neovascularization occurs.

Today's advances in deep learning (DL) are changing the way healthcare is handled, allowing doctors to more effectively diagnose and treat diseases[2]. Several researchers around the world have attempted to effectively address the challenge at hand. From a variety of developments, DL approaches uncover details contained in large amounts of clinically important health data that can be used to treat, monitor, prevent, and make decisions about health conditions. Deep learning (DL) is a subset of machine learning techniques. It is widely used for DR detection and classification[3].

II. RELATED WORK

U.Acharya.et al[4] used features like blood vessels, microaneurysms, exudates, and haemorrhages from 331 fundus images using SVM with an accuracy of more than 85%. K. Anant.et al[5] in their literature used texture and wavelet features for DR detection by making use of data mining and image processing on a database DIARETDB1 and achieved 97.95% accuracy.

Pao *et al.* [6] utilized bi-channel neural networks for the extraction of fundus components by channel, followed by detail enhancement using a classical sharpness enhancement tool named unsharp masking (UM). The Kaggle Diabetic Retinopathy dataset [7] was used for this implementation, with 21,123 RGB fundus image sizes being selected for this implementation.

Zhang *et al.* [8] was able to identify the different severity levels and achieve a multi-class accuracy, specificity and sensitivity of 96.5%, 98.9% and 98.1%, respectively. However, this model does not detect retinal fundus lesions and the results were limited to the private dataset used for evaluation.

Hua *et al.* [9] proposed an uncommon model RFA-BNET that stands out due to its approach, however it utilizes a ResNet-101 as it's backbone. Hua *et al.* [9] model aggregates features from multiple rounds through the ResNet-101, it achieves a relatively lower accuracy rate than the rest of the ensemble models.

Attia *et al.* [10] survey examined DR classification methods with a general focus on deep learning techniques and a high focus on classical methods. Gupta and Chhikara [11] reviewed DR detection techniques utilizing Adaboost, Random forest, SVM etc, gradually showcasing the gap that these classical techniques present in regards to learning more disease related features. These comparisons are based on quality of the fundus image, since some publicly available datasets have poor contrast and image quality.

Shamshad *et al.* [12] provides a comprehensive overview of how transformers work for various medical imaging objectives, including: segmentation, classification, detection, and reconstruction. The survey highlights that transformer-based research for medical imaging reached its peak around Dec 2021, with more than 40 recent publication

III. METHODOLOGY

The main goal of this research is to build a stable and noise-compatible detection system for Diabetic retinopathy. In this study, we use deep learning techniques to detect. Diabetic retinopathy based on severity (no DR, mild, moderate, severe, proliferative). DR). Many processes were performed before the image was fed to the network. we trained Two models in this work: our proposed model, regression model, and comparison The accuracy of the two models was determined. Figure 1 shows the proposed methodology. It performed better than the regression model.





Data Source: Data used for this study has been taken from Diabetic Retinopathy Detection 2015[13] and APTOS 2019 blindness detection[14] from kaggle. Both the datasets contains thousands of retinal images under different conditions. For every subject, two images of both the eyes are given as left and right. As the images come from different sources like different cameras, different models etc. It has an abundance of noise associated with it, which apparently needs to be removed, thus, requiring a number of preprocessing steps. The diabetic retinopathy associated with each image has been rated on the scale of 0-4 as:

- ➢ 0 No DR
- ➤ 1 Mild
- \geq 2 Moderate
- ➤ 3 Severe
- ➢ 4 Proliferative DR

Figure 2 shows the retinal images with ratings on the basis of severity level from 0-1.

Department of Computer Science, A.V.C. College (Autonomous), Mayiladuthurai – 609 305.



Figure 2 : Image samples based on severity from dataset: (a) is level '0', (b) is level '1', (c) is level '2', (d) is level '3', (e) is level '4'.

IV. DATA PREPROCESSING

The images in the dataset contain a lot of noise like some images May be out of focus, some may have high exposure, others may have additional lighting or presence Since the colors are different, such as the black background, some preprocessing needs to be done to match the standard format. The preprocessing steps are:

- Crop black border: The black background of the fundus image is used. Black is useless as it adds no information to the image. The background around the image is omitted.
- Delete the black corner: The black border still remains after deletion. The fundus image has a round shape, so the corners are black. In this step black Corners are removed from the image.
- Image resizing: The image will be resized to 256*256 (width*height).
- Apply Gaussian Blur: A Gaussian blur is applied to the image. Specify the kernel size as 256/6. This method will help remove Gaussian noise.



Figure 3: Images obtained after preprocessing was carried out.

Data Augmentation: Data augmentation is framed by aligning one class to the class with most samples, in order to balance the data among the diabetic retinopathy severity classes. Images were mirrored and rotated to augment the dataset.

Modeling: We used a DenseNet-169 (Densely connected convolutional neural network) and Regression model for training purpose. In DenseNet-169 weights are loaded into the network without the top or last layer. When modeling the network, initially there is no last layer. We design this layer by using Global Average Pooling 2D, a Dropout layer set at 0.5 and an output comprising of five nodes for each class. Global Average Pooling 2D is same as that of 2D average Pooling in operation, but it considers the entire input block size as pool size. A Dropout layer addresses the issue of over-fitting. Adam optimization algorithm is used for optimizing the weights on training this model. A sequential modeling approach is used for adding layers and customizing the layers like convolutional, dropout, dense, optimizers, etc.

Figure 4 shows deep DenseNet-169 model with three Dense blocks and three transition layers consisting of pooling and convolution layer.



Figure 4 : DenseNet-169 model with three dense blocks.

Implementation: The implementation was executed using python language, where a wide variety of libraries were employed for processing of images and to get acquainted with the system for creating convolutional neural network.

V. RESULTS AND DISCUSSION

In the proposed model using DenseNet-169 on a combination of datasets from diabetes. Kaggle retinopathy detection 2015[13] and APTOS 2019 blindness detection[14]. Therefore, the images provided by the dataset contained a lot of noise and required preprocessing. It was necessary. As preprocessing, first remove the black border from the image. Please pay more attention to just the fundus image. The black corners of the image were also removed. The image is resized to a standard format of 256 x 256 width and height. Finally Gauss A blur was applied to the image to remove Gaussian noise. Once the pre-processing is completed, analyze that there is a huge imbalance in the data between the severity classes to which most of the data belongs to class "0", i.e. no DR. To solve this problem, Data augmentation is used and provided 7000 images from each severity class to balance the data. After pretreatment, the images were augmented and the data was finally transferred to DenseNet-169 to train the model. After evaluating the model, we obtained a training accuracy of 0.953 as validation. An accuracy of 0.9034 was achieved. We also calculated the resulting Cohen Kappa score. It will be 0.804. We also applied

the regression model to the dataset and calculated its validation. Accuracy is 0.789. Our proposed model is better than regression model. An overview of our model is summarized in Table 1. **Table 1 :** Result of the Proposed Model

Training Accuracy	Validation Accuracy	Cohen kappa score
95%	90%	80%

The figure 5 shows the comparison between the accuracies obtained by the proposed model and regression model.



Figure 5 : Accuracy Obtained By Proposed Model and *Regression Model*

Besides Regression model, the proposed model was compared to a number of machine learning classifiers like Support Vector Machine (SVM), K-Nearest Neighbor (KNN) and Decision Tree (DT). The results are summarized in the table 2, where accuracies of the different classifiers are given.

Table 2 : Results obtained by various Classifiers

CLASSIFIER	DATASET	ACCURACY	DR CLASSES
SVM[29]	Messidor, Diabeticret	85.6%	Normal, Non PDR, PDR.
DT[30]	Messidor	85.1%	Normal, Mild, Moderate, Severe.
KNN[29]	Messidor, Diabeticret DB1	55.1%	Normal, Non PDR, PDR.
Regression	Diabetic Retinopathy Detection 2015 &APTOS 2019 from kaggle.	78%	No DR, Mild, Moderate, Severe and Prolife
Proposed	Diabetic	90%	No DR,

Retinopathy	Mild,
Detection	Moderate,
2015	Severe and
&APTOS	Proliferative
2019 from	DR.
kaggle.	
	Retinopathy Detection 2015 &APTOS 2019 from kaggle.

The proposed model achieves the highest accuracy of 90%, followed by SVM with an accuracy of 85.6%, Decision Tree with 85.1%, Regression with 78%, and KNN with the least accuracy of 55.17%.

VI. CONCLUSION

According to many people, traditional methods for detecting DR are time-consuming, demanding, and costly. Research is being done to automate the detection process using machine learning Deep learning approach. In this study, we presented a comprehensive study on various issues. Automatic detection method of diabetic retinopathy and the method we tried to propose A unique deep learning approach for early diagnosis of retinopathy using DenseNet 169 (a new CNN architecture with many deep layers). Two datasets: "Diabetes" I used kaggle's "Retinopathy Detection 2015" and "APTOS 2019 Blindness Detection" Together for this research. Many preprocessing and extensions were done for standardization Convert vour data to your desired format and remove unwanted noise. In addition to the DenseNet-169 classifier, Additionally, regression models were used to make comparisons between outcomes. Furthermore, the machine We compared learning classifiers such as SVM, DT, and KNN with the proposed system. where The proposed model achieved the highest accuracy among all and also enabled classification. We also take pictures in other classes. Our proposed model showed better performance than the regression model However, by achieving 90% accuracy, the regression model was able to achieve 78% accuracy.

REFERENCES

[1] J. Amin, M. Sharif, and M. Yasmin, ``A review on recent developments for detection of diabetic retinopathy," *Scienti_ca*, vol. 2016, pp. 1_20,Sep. 2016.

[2]. Ganeshsree Selvachandran1 · Shio Gai Quek1 · Raveendran Paramesran2 · Weiping Ding3 · Le Hoang Son, 'Developments in the detection of diabetic retinopathy:a state-of-the-art review of computer-aided diagnosis and machine learning methods".

[3] W. L. Alyoubi, W. M. Shalash, and M. F. Abulkhair, "Diabetic retinopathy detection through deep learning techniques: A review," *Informat. Med. Unlocked*, vol. 20, Jan. 2020, Art. no. 100377.

[4] U. R. Acharya, C. M. Lim, E. Y. K. Ng, C. Chee, and T. Tamura, "Computer-based detection of diabetes retinopathy stages using digital fundus images," *Proc. Inst. Mech. Eng. Part H J. Eng. Med.*, vol. 223, no. 5, pp. 545–553, 2009, doi: 10.1243/09544119JEIM486.

[5] K. A. Anant, T. Ghorpade, and V. Jethani, "Diabetic retinopathy detection through image mining for type 2 diabetes," in 2017

Department of Computer Science, A.V.C. College (Autonomous), Mayiladuthurai – 609 305.

International Conference on Computer Communication and Informatics, ICCCI 2017, 2017, doi: 10.1109/ICCCI.2017.8117738.

[6] S.-I. Pao, H.-Z. Lin, K.-H. Chien, M.-C. Tai, J.-T. Chen, and G.-M. Lin, "Detection of diabetic retinopathy using bichannel convolutional neural network," *J. Ophthalmol.*, vol. 2020, pp. 1_7, Jun. 2020.

[7] Diabetic Retinopathy Detection EYEPACS Dataset, Kaggle, San Francisco, CA, USA, Jul. 2015.

[8] W. Zhang, J. Zhong, S. Yang, Z. Gao, J. Hu, Y. Chen, and Z. Yi, "Automated identification and grading system of diabetic retinopathy using deep neural networks," *Knowl.-Based Syst.*, vol. 175, pp. 12_25, Jul. 2019.

[9] C.-H. Hua, T. Huynh-The, and S. Lee, "Retinal vessel segmentation using round-wise features aggregation on bracket-shaped convolutional neural networks," in *Proc. 41st Annu. Int. Conf. IEEE Eng. Med. Biol. Soc.(EMBC)*, Jul. 2019, pp. 36_39.

[10] A. Attia, Z. Akhtar, S. Akrouf, and S. Maza, ``A survey on machine and deep learning for detection of diabetic RetinopathY," *ICTACT J. Image*.

[11]A. Gupta and R. Chhikara, ``Diabetic retinopathy: Present and past," *Proc.Comput. Sci.*, vol. 132, pp. 1432_1440, Jan. 2018.

[12] F. Shamshad, S. Khan, S. W. Zamir, M. H. Khan, M. Hayat, F. S. Khan, and H. Fu, "Transformers in medical imaging: A survey," 2022, *arXiv:2201.09873*.

[13] A. H. Asad, A. T. Azar21] "Diabetic Retinopathy Detection Identify signs of diabetic retinopathy in eye images," 2015. [Online].Available: https://www.kaggle.com/c/diabetic-retinopathy-detection/data.

[14] "APTOS 2019 Blindness Detection Detect diabetic retinopathy to stop blindness before it's too late," 2019. [Online]. Available: https://www.kaggle.com/c/aptos2019-blindness-detection/data.

A Review on Rice Leaf Diseases Using Feature Extraction Techniques in Machine Learning

S. Govindarajan S¹ and Dr. S. Mary Vennila²

¹Research Scholar, PG and Research Department of Computer Science,

Presidency College, Chennai, India thimagovind7@gmail.com

²Associate Professor and Research Supervisor, PG and Research Department of Computer Science, Presidency College, Chennai, India

Abstract— Rice Leaf diseases cause major losses in terms of production, economy, quality and quantity of agricultural products. One of the major factors contributing to reduced crop yields is the presence of bacterial, fungal, and viral illnesses. Utilizing rice plant disease detection techniques can help avoid and manage this. For effective crop production in agriculture, early disease identification is crucial. Since machine learning primarily uses information from the data itself and provides excellent methods for detecting plant diseases, it will be used in the process of identifying diseases in rice plants. Normally the rice plant automated system involves the steps of image acquisition, preprocessing of image, segmentation, feature extraction and finally apply the classification model. In this paper mainly focus on the various feature extraction techniques.

Keywords—rice disease, image preprocessing, image segmentation, feature extraction, GLCM, LBP, color-based features.

I. INTRODUCTION

Agriculture holds paramount importance in the sustenance and economic well-being of nations worldwide, serving as a primary source of income and food for a significant portion of the population. However, the agricultural sector faces a formidable challenge in the form of plant diseases, particularly those affecting leaves. The impact of leaf diseases has emerged as a serious concern, causing substantial damage to crops and consequently leading to a decline in both the quantity and quality of food production [1]. On the other hand, the predominant method for detecting and identifying plant diseases involves the visual inspection of experts using the naked eye. However, this traditional approach is fraught with challenges as it proves to be time-consuming, expensive, and demanding in terms of effort. The reliance on manual observation for disease detection presents significant drawbacks that hinder its practicality [2]. Given the drawbacks of naked eye observation, there is a pressing need to explore alternative and more efficient approaches. Early identification of infections is of paramount importance for farmers, considering the potential for rapid disease progression in crops. Detecting diseases in their early stages is a crucial aspect of agriculture, demanding careful diagnosis and effective supervision to curb substantial losses. Recognizing the pivotal role of disease detection in agriculture, it becomes evident that a shift toward more advanced and automated methodologies is necessary. By integrating technologies that enable prompt and accurate identification of plant diseases, farmers can significantly enhance their ability to manage and control potential losses, thereby promoting sustainable and resilient

agricultural practices. Many diseases can impede rice growth, but for the sake of this essay, we will focus on three major diseases that impact rice: rice blast, brown spot, and sheath blight [3].

The goal of this work is to create a prototype system that can replace or augment the current manual method in the automatic and accurate detection and classification of paddy illnesses with leaf smut, brown spot, and sheath blight through the use of feature extraction techniques.

The diseases' characteristics [13] are listed and shown below in Fig. 1:

- Leaf smut: The lesions are often covered with a whitish to grayish growth of the smut fungus, which contains masses of dark spores.
- Brown spot: Oval to elliptical lesions with a brown or dark center. Lesions are often surrounded by a yellow halo.
- Sheath blight: Water-soaked lesions on the leaf sheath, which elongate and may encircle the entire sheath. "Rotted neck" appearance in severely infected plants.



Fig. 1. (a) Leaf Smut, (b) Brown Spot, (c) Sheath blight.

II. RELATED WORK

The bank color determines the particular color of each pixel in an image. Boundary points and centroid points of each object on the source image (RGB image) are found early on with Blob Analysis [4].

The length, width, and number of the item are general characteristics of shape. Blob Analysis is utilized in the existing paper to get number of the items for labeled regions in a noise free binary picture. The paddy leaf's lesion is the object. The shape's characteristic is obtained by counting the object using the 8-connected neighborhood technique [5].

For image processing applications, the HSI color space is a highly useful and appealing color model since it accurately captures color in the same way that the human eye does [6].

For color features, use hue angle to calculate the intensity of color from the HSV color space; for shape features, use the diameter difference between the major and minor axes of the illness [7].

 The diseased object boundaries, spots, and leaf color are

 Department of Computer Science, A.V.C. College (Autonomous), Mayiladuthurai – 609 305.

 PROCEEDINGS
 59
ISBN: 978-81-967420-1-0

all determined by color texture analysis, which converts the RGB color space into the CIEL * a * b * color space and uses Euclidean distance to achieve a color that is exactly the same as the expert's color determination [8].

Color pixels can have various colors. The color histogram for an image's K-mean (CHKM) characteristics is displayed. This method considers replacing a specific pixel color with a color from the common color palette that is more comparable to that pixel color. The same method was used for every color of pixel in the image to categorize it into a k cluster [9].

Energy is referred to as the sum of square elements in SGDM. On the other hand, the distribution of items in SGDM is used to measure homogeneity. The measure of an element's correlativeness to its neighbor throughout the entire image is obtained by correlation [10].

Texture analysis technique is seen to ignore the light variations while the images were taken. This causes confusion in deciding the difference in spot colour and boundary colour [11].

The selection of the green component stems from the finding that the greenness of the leaf is most influenced when the infection occurs in the leaf [12].

III. PROPOSED METHODOLOGY AND PLANT DISEASE DETECTION PROCESS

The automated disease detection system comprises five integral phases: image acquisition, preprocessing, image segmentation, feature extraction, and classification [14]. However, this paper distinctly emphasizes the intricate aspects of feature extraction. In this context, feature extraction plays a pivotal role as it involves capturing and representing essential information from raw image data, contributing significantly to the system's effectiveness. This process meticulously identifies and highlights relevant features, facilitating a reduction in data dimensionality and enhancing computational efficiency. Furthermore, the emphasis on feature extraction in this paper aims to delve into the specifics of how extracted features, such as color, texture, and shape, contribute to the accurate and robust identification of diseases in rice crops. The focus on feature extraction not only improves the interpretability of the model but also addresses the challenges of data variability and enables adaptability to diverse disease symptoms and plant appearances.



Fig. 2. Proposed Method

Input Image and Image Acquisition

In this stage, digital devices like cameras and cell phones ves at the necessary size and resolution. Additionally, preexisting images stored in repositories are included as part of the data acquisition process [15]. This phase involves the gathering of visual data through various means, ensuring that the images obtained meet the specified size and resolution criteria for subsequent stages of the automated disease detection system. The inclusion of images from the web expands the data sources available for the automated disease detection system. The responsibility for forming the image database lies entirely with the application system developer. This database's formation is crucial as it directly impacts the classifier's efficiency during the final phase of the detection system. A well-constructed image database enhances the system's ability to accurately classify and detect diseases, emphasizing the importance of thorough curation and selection of images for optimal performance.

Preprocessing

Frequently, image sets are generated and captured in real-time conditions, encompassing extraneous elements like shadows, noise, unspecified distortions, and intricate backgrounds. Consequently, at a lower level of abstraction, image preprocessing becomes a fundamental necessity. The primary objective of the pre-processing step is to mitigate the impact of undesirable characteristics inherent in the images.

This results in an image that is more conducive for subsequent processing stages. Furthermore, operations such as cropping and resizing can be strategically employed to simplify the complexities associated with intricate systems. The effectiveness of fundamental pre-processing operations, including color space conversion, histogram equalization, contrast enhancement, cropping, noise removal, and smoothing, is contingent upon the nature of the acquired images [16]. Their application is tailored to address specific challenges and nuances associated with different types of images within the system. Some of the essential preprocessing techniques list below in table 1.

Preprocessing Techniques	Computation
Gray Scaling	Grayscale value (G) = 0.299 * R + 0.587 * G + 0.114 * B
Histogram Equalization	$s_k = \frac{(L-1)}{MN} \sum_{j=0}^k n_j$
Contrast Stretching	$I_{out} = \frac{I_{in} - \min - in}{\max_{in} - \min_{in} X (\max_{out} - \min_{out}) + \min_{out} X (\max_{out} - \min_{out}) + \min_{out} X (\max_{out} - \min_{out}) + \min_{out} X (\max_{out} - \max_{out} - \min_{out}) + \min_{out} X (\max_{out} - \max_{out} $

Resizing $I_{out}(x, y) = I_{in}(x', y') \text{ where } x'$ $= \frac{x}{\text{scale_factor}} \text{ and } y'$ $= \frac{y}{\text{scale_factor}}$

Table 1. Preprocessing Techniques.

Segmentation

This stage focuses on streamlining the depiction of an image to make it more significant and straightforward for analysis [17]. At the core of feature extraction, this stage constitutes a foundational aspect of image processing. Multiple techniques exist for image segmentation, including k-means clustering, Otsu's algorithm, and thresholding. K-means clustering categorizes objects or pixels into K classes based on a defined set of features. The classification is achieved by minimizing the sum of squares of distances between objects and their respective clusters [18].

Feature Extraction

The segmentation process, the current result is the defined area of interest. Consequently, at this stage, it is imperative to extract features from this identified area. These features play a crucial role in deciphering the significance of a sample image. The attributes can be derived from aspects such as color, shape, and texture [19]. Multiple methods of feature extraction are available for system development, including the gray-level co-occurrence matrix (GLCM), color co-occurrence method, spatial grey-level dependence matrix, and histogram-based feature extraction.



• Local Binary Pattern (LBP): It is clear that the Local Binary Pattern (LBP) method proves to be a cost-effective solution, significantly reducing the time required for implementation [20].

$$LBP(gp_x, gp_y) = \sum_{p=0}^{p-1} S(gp - gc) * 2^p$$

In the above equation, _gp' is given by neighboring pixel's intensity with index p, and _gc' is given by central pixel's intensity, _p' is the number of sampling points on a circle of radius.

• Histogram of Oriented Gradients (HOG): In this approach, occurrences within a localized section of an image are tallied. The computation involves

determining the image gradient, which is derived by combining both the gradient and angle information from the image. The formula expressing this relationship is provided below [21].

For example, when considering a 3 x 3 image, the calculation involves determining Gx and Gy for each pixel in the image.

$$G_x(r.c) = I(r, c+1) - I(r, c-1)$$
$$G_y(r.c) = I(r-1, c) - I(r+1, c)$$

Following the computation of Gx, the magnitude and angle are subsequently calculated. The formulas for these calculations are provided below:

$$Magnitude(\mu) = \sqrt{G_{x+}^2 G_y^2}$$
$$Angle(\theta) = tan^{-1} \frac{G_x}{G_y}$$

• Gray Level Co-occurrence Matrices (GLCMs): In recent times, a predominant focus among researchers is directed towards the utilization of texture features for the identification of plant diseases.

The GLCM method, in particular, stands out as a statistical approach for texture classification. The GLCM calculation encompasses various features, such as energy, entropy, contrast, homogeneity, covariance, shade, and prominence. Each of these features is mathematically represented in a unique format [22].

• Gabor Filter: Gabor filter feature extraction techniques involves utilizing Gabor filters to extract relevant features from images for the purpose of identifying and classifying plant diseases. Gabor filters are particularly effective in capturing texture and spatial frequency information, making them suitable for analyzing intricate patterns in plant images. The formula for expressing this Gabor filter $G(x, y; \lambda, \theta, \psi, \sigma, \gamma)$

$$= \exp(-2\sigma_2 x'_2 + \gamma_2 y'_2)\cos(2\pi \lambda x' + \psi)$$

• Elliptic Fourier and Discriminant Analyses: Apply Elliptic Fourier Analysis to the extracted leaf contours to obtain a set of coefficients that represent the shape characteristics of each leaf. These coefficients capture the variations in the leaf contour's curvature and provide a compact representation of the leaf shape. The formula for expressing this Elliptic Fourier and Discriminat Analyses:

$$Ck=n1\sum_{i=1}n(xi+iyi)e^{-i2\pi k_{ni}}$$
$$gi(x) = xTSw - 1(\mu i - x^{-}) - 21\mu iTSw - 1\mu i$$
$$+ \log(\pi i)$$

Classification

Upon concluding the extraction of both shape and color texture features from images depicting disease spots, we proceeded to develop three distinct classification models.

Department of Computer Science, A.V.C. College (Autonomous), Mayiladuthurai – 609 305.

A widely employed feature extractor in deep learning is the Convolutional Neural Network (CNN), which plays a crucial role in extracting features from images and overseeing the entire feature engineering process. In the typical CNN architecture, initial layers are responsible for capturing low-level characteristics, while subsequent layers focus on extracting high-level features from the input image. Extensive research indicates that CNN effectively addresses various challenges encountered by conventional classifiers, yielding higher accuracy in the field of rice plant disease classification. Across multiple studies, the average accuracy achieved by CNN stands at an impressive 98%. The table below illustrates the average accuracy derived from a comprehensive review of papers that utilized the same dataset.

Table 2. A review of the results

Feature Extraction & Classifier	Average Accuracy
GLCM & KNN	91.37%
GLCM & SVM	95.8%
LBP & KNN	90.1%
LBP & SVM	92.3%
CNN	98%

IV. CONCLUSION

The derived features, covering attributes related to color, shape, and texture, significantly contribute to the nuanced understanding of different rice leaf diseases.

Various feature extraction methodologies have been investigated and implemented for the purpose of detecting diseases in rice leaves. Techniques such as the Gray Level Co-occurrence Matrix (GLCM), Local Binary Pattern (LBP), Histogram of Oriented Gradients (HOG), and Gabor filters have been extensively explored. Each of these methods brings distinctive perspectives to the analysis of diseased rice leaves, offering a diverse and valuable array of information that can be effectively utilized by machine learning models.

REFERENCES

- Akila, M., & Deepan, P. (2018). Detection and classification of plant leaf diseases by using deep learning algorithm. International Journal of Engineering Research & Technology (IJERT), 6(07).
- [2] Siddharth sing, A. Kaul, "Bacterial Foraging Optimization Based Radial Basis Function Neural Network (BRBFNN) for Identification and Classification of Plant Leaf Diseases: An Automatic Approach Towards Plant Pathology, IEEE ISSN: 2169-3536 Volume: 6 Page(s): 8852 – 8863
- [3] Mangla, N., Raj, P. B., Hegde, S. G., & Pooja, R. (2019). Paddy leaf disease detection using image processing and machine learning. Int J Innov Res Elec Electron Instrument Control Eng, 7(2), 97-99.
- [4] Kurniawati, N. N., Abdullah, S. N. H. S., Abdullah, S., & Abdullah, S. (2009, December). Investigation on image processing techniques for diagnosing paddy diseases. In 2009 international conference of soft computing and pattern recognition (pp. 272-277). IEEE.
- [5] H. Wang, "Automatic Character Location and Segmentation in Color Scene Images", Proc. of the 11th International Conf. on Image Analysis and Processing, 2001.
- [6] Swathi, D., & Bharathi, A. (2016). Disease classification of paddy leaves using HSI feature extraction and SVM technique. IJSRD-International Journal for Scientific Research & Development, 4(02), 2321-0613.

- [7] Kitpo, N., & Inoue, M. (2018, March). Early rice disease detection and position mapping system using drone and IoT architecture. In 2018 12th South East Asian Technical University Consortium (SEATUC) (Vol. 1, pp. 1-5). IEEE.
- [8] Rossilawati, S., Siti Norul Huda, S.A., Mohammed, Y., Azuraliza, A.B., Sharizan, R., Noorashikin, M., Saad, A. & Nik Mohd Noor, N.S. 2003. The development of diseases diagnosing system in paddy plant (E-Paddy). Prosiding Antarabangsa Teknologi Maklumat (ITSIM'03), hlm. 424-432.
- [9] C. H. Lin., R. T. Chen, and Y. K. Chan, "A smart contentbased image retrieval system based on color and texture feature", Image and Vision Computing, Vol. 27, No. 6, 2009, pp. 658–665.
- [10] Satgunalingam, V., & Thaneeshan, R. (2020). Automatic Paddy Leaf Disease Detection Based on GLCM Using Multiclass Support Vector Machine. Int. J. Comput, 39(1), 97-106.
- [11] Kaur, E. J., & Singla, S. A Detailed review and Classification of Segmented Image for Paddy Plant Disease.
- [12] Phadikar, S., & Sil, J. (2008, December). Rice disease identification using pattern recognition techniques. In 2008 11th International Conference on Computer and Information Technology (pp. 420-423). IEEE.
- [13] "Rice disease identification photo link." www.agri971.yolasite.com/resources/RICE/DISEASE/0IDENTIFICATION.pd f. Accessed: 2019-08-25.
- [14] Sandhu, G. K., & Kaur, R. (2019, April). Plant disease detection techniques: a review. In 2019 international conference on automation, computational and technology management (ICACTM) (pp. 34-38). IEEE.
- [15] Strange, Richard N., and Peter R. Scott. "Plant disease: a threat to global food security." Annu. Rev. Phytopathol. 43 (2005): 83-116.
- [16] Vishnoi, V. K., Kumar, K., & Kumar, B. (2022). A comprehensive study of feature extraction techniques for plant leaf disease detection. Multimedia Tools and Applications, 1-53.
- [17] Gharge, S., Singh, P., 'Image Processing for Soybean Disease Classification and Severity Estimation', Emerging Research in Computing, Information, Communication and Applications, pp. 493- 500, 2016.
- [18] Singh, J., Kaur, H., 'A Review on: Various Techniques of Plant Leaf Disease Detection', Proceedings of the Second International Conference on Inventive Systems and Control, Volume 6, pp. 232-238, 2018.
- [19] Dey, A.K., Sharma, M., Meshram, M.R., 'Image Processing Based Leaf Rot Disease, Detection of Betel Vine (Piper BetleL.)', International Conference on Computational Modeling and Security, Volume 85, pp. 748-754, 2016.
- [20] M. Harmouch, —Local Binary Pattern Algorithm: The Math Behind It, I The Startup, 2020. https://medium.com/swlh/local-binarypattern-algorithm-themath-behind-it-- edf7b0e1c8b3.
- [21] M. Tyagi, —HOG (Histogram of Oriented Gradients): An Overview.l https://towardsdatascience.com/hog-histogramof-oriented-gradients-67ecd887675f.
- [22] B. E. Park, W. S. Jang, and S. K. Yoo, —Texture analysis of supraspinatus ultrasound image for computer aided diagnostic system, Healthc. Inform. Res., vol. 22, no. 4, pp. 299–304, 2016.

Design of Resilient Cyber defense mechanism using Dynamic Zero Trust Network Security Framework by integrating Artificial Intelligence, Machine Learning and Blockchain V. Vasanthi M. Paul Arokiadass Jerald

Associate Professor, Department of Computer Science Dharmapuram Gnanambigai Govt. Arts College for Women, Mayiladuthurai. vasanthiv.78@gmail.com

R. Bhuvaneswari* Associate Professor, Department of Computer Science Periyar Arts College, Cuddalore. rbeswari17@gmail.com M. Paul Arokiadass Jerald Associate Professor, Department of Computer Science Periyar Arts College, Cuddalore. mjerald@gmail.com

I. Benjamin Franklin* Assistant Professor, PG Department of Computer Applications St. Joseph's College of Arts and Science (Autonomous), Cuddalore. franklinbenj@gmail.com

Abstract - In ever evolving cyber threats with the advent of digital epoch in all parts of the world, the traditional and even modern methods of network security have been proved inadequate to safeguard the critical assets of the global users. This research proposes a state-of-the-art Zero Trust Network Security Framework that is developed using integration of three prominent and cutting-edge technologies including Artificial Intelligence (AI), Machine Learning (ML) and Blockchain with the leverage of Dynamic Resilient Zero Trust Network Security Framework (DRZTNSF) model. The proposed framework is capable of enhancing the authentication and authorisation process through continuous verification and validation of the user in each and every subsequent access. The framework integrates the technologies to automatically trigger the threat detection, responses and mitigation competences and through micro-segmentation, data encryption process, real-time monitoring of the system using blockchain systems, the anomalies and potential breaches could be easily diagnosed. This simulated approach could be used in effective fortification of network infrastructures against futuristic cyber threats and security problems.

Keywords - Zero Trust Network Security Framework; Dynamic Resilient Zero Trust Network Security Framework (DRZTNSF) model; Artificial Intelligence; Machine Learning; Blockchain Technologies.

I. INTRODUCTION

Zero Trust Network Security [1] is a new way to keep computer networks safe. It doesn't rely on trusting anyone, both inside and outside the network. Instead of trusting the inside of a network and not trusting the outside, Zero Trust always checks to make sure things are safe. In a Zero Trust Network Security Framework, trust is not automatically granted to users, devices, or applications, regardless of their location within or outside the corporate network. Instead, every user and device attempting to access resources is continuously authenticated and authorized, and their behavior is scrutinized in realtime. This approach acknowledges the dynamic nature of modern IT environments, with users accessing data from various locations, devices, and networks.

The Zero Trust model uses small segments, limits access, constantly checks, and requires more than one way to prove identity [2]. Micro-segmentation means breaking the network into small parts to keep security problems in one segment from spreading to the whole network. Least privilege access means that users and

devices only have enough access to do their jobs, which helps to lower the chances of someone getting in who shouldn't be there as shown in Figure.1.



Figure 1. Existing Zero Trust cyber–Security Framework with Policy Control System

As shown in Figure 1, the Policy Engine with Administration and enforcement allows continuous monitoring which means constantly checking [3] how people and devices are acting to find anything unusual or suspicious. Multifactor authentication makes things more secure by asking users for more than one type of ID before letting them in. In the modern era of cloud computing and remote working, it's crucial to employ Zero Trust Network Security to ward off cyber threats. By using a zero-trust approach, organizations can improve their security, reduce the chance of data breaches, and protect important information in a changing digital world.

This paper aims to introduce and examine the Dynamic Resilient Zero Trust Network Security Framework (DRZTNSF), which is a modern method that uses Artificial Intelligence, Machine Learning, and Blockchain technologies to improve network security. The goal is to improve the traditional zerotrust model by adding flexible security measures to deal with the changing challenges of cyber threats in modern IT environments. The research explains the Zero Trust model and its main ideas. It focuses on how security is changing from protecting the outside of a network to focusing more on the people using it. It also looks at how traditional zero-trust models can't always handle new and complicated cyber threats.

The goal of Artificial Intelligence is to see how it can help make the Zero Trust framework better at finding and responding to potential threats. It also looks into how AI-powered data analysis can make it easier to find an unusual activity, understand behavior, and predict security issues within the network. Machine Learning is about using ML algorithms to make a flexible and smart security system in the Zero Trust framework. It also looks at how machine learning can be used to constantly check for problems, assess risks, and change security rules as threats change. The Blockchain Technologies for Immutable Security looks at how blockchain can keep data safe and secure, and make sure communication is reliable in the Zero Trust network. It also looks at how blockchain can help make an audit trail of access and transactions that can't be changed, making it easier to trust and be accountable.

The Resilience Enhancement is a conversation about how AI, ML, and Blockchain make the Zero Trust model stronger and better at handling new threats and keeping things running smoothly. It also checks how the new plan might help respond quicker and make security incidents less serious. The paper wants to help the cyber security field by suggesting and looking at a new way to keep networks safe called the Dynamic Resilient Zero Trust Network Security Framework (DRZTNSF). It combines AI, ML, and Blockchain with the Zero Trust Network Security model. The study looks at the possible benefits, difficulties, and practical effects of this new method for protecting organizations from modern cyber threats.

II. Related Works

The research review for the DRZTNSF looks at new studies and progress in cyber security, Artificial Intelligence AI, Machine learning ML, blockchain, and zero-trust network security. The survey wants to find out what people know about a topic, see what it still needs to learn and connect the ideas to what's happening in the field now. Van Bossuyt, D. L., et al. (2023) [4] suggested a way to study the dangers of a mission when people's trust in a system is lost. It was called the Trust Loss Effects Analysis (TLEA). To help people understand and use this method better, TLEA follows the steps of FMECA but focuses on identifying security events. There, the TLEA method uses steps from a way to prevent hackers from doing bad things online. The TLEA is first explained using a simple example and then shown using a more reallife situation of a drone on a spying mission. After using the TLEA method, it can find out about various risks connected to losing trust, and see how it might affect the success of the mission. Ray, P. P. (2023) [5] explains Web3 and its important parts, like DApps, DeFi, NFTs, DAOs, and Supply Chain Management. The article also talks about how Web3 could affect society and the economy. It also looks at how Web3 could work with new technologies like AI, IoT, and smart cities. This article talks about how zero-trust architecture is important for Web3. In the end, this review shows how Web3 is important for shaping the future of the internet. It also talks about the problems and chances that are coming. Chen, X., Feng et al. (2023)[6] suggested a plan for keeping control of access safe. It involves working together with different control areas to stop bad access like DDoS attacks, spreading malware, and new types of exploits. Moreover, it examines crucial design aspects of this architecture and share test results from a case study that validates the suitability of ZTA for 6G. Lastly, and have talked about things it still needs to study more to improve this new design. Van Bossuyt, D. L., et al. (2023)[7] explained the good things and hard parts of using Zero-Trust and suggested a plan for using it in designing systems. This article talks about different things that researchers are studying and focuses on important areas where the community should put effort. Ge, Y., Li, T., & Zhu, Q. (2023) [8] suggested a way to defend against different situations without needing specific information. It used a method called Partially Observable Markov Decision

Processes (POMDP) and first-order Meta-Learning and only used a few examples to demonstrate. The framework creates a trust-threshold defence policy that can be easily understood and used in different situations. To fix the difference between security data and real life, and is improving the model to better protect against the worst possible security threats. So, and are making the defences stronger. It uses real examples of attacks to prove the findings. Federici, F., Martintoni, D., & Senni, V. (2023) [9] shows a way to move existing systems to a new structure using a process that focuses on managing risks and using models. The approach is checked, using two important examples from the aerospace industry.

Zero-trust is a new way to keep everything safe. It focuses on protecting data, devices, parts, and people. Ridhawi, I. A., Otoum, S., & Aloqaily, M. (2023) [10] introduces a new way to include the zero-trust system in 6G networks that use DT technology. The new framework uses a decentralized system to make sure that both the physical devices and their digital twins are secure, private, and authentic, which is different from traditional zero-trust solutions. Blockchain is very important for making sure digital transactions and the information sent are real and safe. Artificial Intelligence (AI) is used in all connected nodes with advanced learning techniques. The article talks about what is happening now and what might happen in the future, including the problems and some ways that technology can help. Wang, J., et al. (2023) [11] makes a new way to store IoT data in the blockchain something called Insertable using Vector Commitment (IVC) instead of a Merkle tree in a group. The building has a small amount of evidence. It can also keep track of the updates, which helps stop replay attacks. Tests prove that each block can handle 1,000 transactions. The size of a single piece of data proof is 30% of the original, and proofs from different parts can be combined. IVC can help to decrease communication crowding and make the IoT system more stable and secure.

Zero Trust is like making sure there is security everywhere from the very beginning, just like a military base, airport or bank. This means that Zero-Trust strategies include making sure that security is built into an organization's architecture from the beginning. Zero trust means you should never assume that any interaction is safe from the start. On the other hand, the old way of keeping things safe often decides if something is trustworthy by seeing if the communication starts from within a firewall. Zero-Trust is a new way to make sure the computer systems are safe. It's different from the old way of putting security around the outside of the system. This old way doesn't work as well now that more

people are working from home and there are more risks from people on the inside. Technology has changed, as well as how it works. But the bad guys who use computers have also changed how to do things. It's not about if a bad person will get past their defense, but when Organizations need to plan their security so it can predict and stop this. Zero trust is about figuring out what things you need to keep safe, like their important stuff and data, and then making a plan to protect them. This could include things like their computer and personal information as suggested by Seaman, J. (2023) [12]. Similarly, Shang, Y. (2023) [13] studied about how a group of agents can work together and follow a leader, even when the connections between them change over time. This group includes some agents who might not always cooperate. Resilient tracking consensus aims to help groups work together to find the leader, even when some are trying to cause problems. It was thinks that the cooperative agents don't know how many Byzantine agents there are or who it was, and the connections between them are always changing randomly. It used math and a theorem to create a plan for everyone to agree on something simply and quickly. It gave a few examples with numbers to prove the theories.

The use of blockchain is becoming more and more popular in many real-life situations. To make smart contracts work, blockchains often need to use realworld information from different places or use outside services. Hybrid smart contracts and Decentralized Oracle Networks help on-chain code talk to off-chain services. This helps make the vision come true. Even though blockchain is safe, it can still be at risk because it needs to use outside services and real-world information. These risks happen when bad people change data from outside sources in a smart contract and it doesn't work properly. It needs to create a Zero Trust Architecture to make sure DONbased apps are secure from start to finish. Gupta, A., et al. (2023) [14] introduces the concept of Proxy Smart Contracts (PSC), which adds a layer between the actual smart contracts and their execution. Before the smart contract is done, the PSC pretends to run the smart contract logic and shows the result to the people who are interested in it for their approval. In simple words, this means that if the smart contract is not working right, it can be stopped. It uses Solidity in the beginning to see if the suggested plan will work. The current work is important because it offers a potential solution to fill a gap that could make it hard for dApps to be widely used in the real world. The new plan suggested by Chaudhry, U. B., & Hydros, A. K. (2023) [15] will make banks safer by using the zero-trust idea and blockchain technology.

The way the transactions are approved makes them unable to be changed and spread out across many computers. The zero-trust principles used in this model make sure that the banking system keeps information private and trustworthy.

Ajakwe, S. O., Kim, D. S., & Lee, J. M. (2023) [16] looks at how artificial intelligence and blockchain are helping to make secure and reliable autonomous systems. This is done by combining security in cyberspace, intelligence space, and airspace. It used the PRISMA-SPIDER method to review 133 articles. These articles were chosen based on specific criteria. Out of those articles, 91 (68, 4%) were quantitative studies, 19 (14. 2%) were qualitative studies, and 23 (17. 3%) used a mix of both methods. It used this method to make sure have chosen a good mix of articles for the review. The review found that there is a big difference between the models that were suggested and how it was put into practice. By using zero-trust technology with blockchain and new AI models, it can make sure drones are safe to use. This includes checking who owns the drone, verifying deliveries, giving permission to operate the drone, and making sure the drone follows the rules. This will make sure that drone transportation is safe and secure. Liang, X., et al. (2023) [17] wants to study how blockchain can be used to create new business models in the power grid. It will look at how blockchain can help make the power grid more efficient and how to design a system that can handle a large amount of data. It created a new technical structure to study how to use blockchain with a customized way of making decisions, using a mix of design and research methods, and do technical evaluations to see how well the business models can detect fake data attacks and how it can improve. Roozkhosh, P., Pooya, A., & Agarwal, R. (2023) [18] wants to study how much people in Iran are using blockchain for buying home appliances. System dynamics (SD) is used to help understand how things in a model are related to each other when their connections are not straight or simple. It will use a model to predict how the BAR will behave in the supply chain for the next 10 years. It will also test different scenarios to see how it affect the results. This model will be used from 2020 to 2030. After the simulation, and will check if the blockchain is accepted by looking at the data used for studying Multi-Layer Perceptron (MLP) and Vector Regression (SVR). It will use the data that is most closely related to BAR. Machine learning methods are used to predict how many applications will be accepted. It wants the predictions to be accurate for the years 2020-2022, so it compares the predictions to the real data for those years. The study found that

if the impact of the COVID-19 outbreak is moderate, the BAR will be about 0. 6 in 2030 If certain policy changes are made, the rate could go up to 0. 8 at most So, if it focused on creating and designing policies carefully, it can help make the supply chain stronger in the future. The study shows that the SD-MLP method is better than the SD-SVR method because it has fewer errors and can predict the behaviour of the BAR more accurately. Dunn Cavelty, M., Eriksen, C., & Scharte, B. (2023) [19] talked about two important things in cyber security: being at risk and not knowing what will happen. Instead of seeing cyber security as a problem with computers and people, it believes cyber security should be seen as a problem with society and technology. It recommended that researchers, policymakers, and experts focus on three things: working together across different fields of study, having discussions with the public about moral questions, and acknowledging uncertainty in politics and decision-making.

Gani, A. B. D., & Fernando, Y. (2023) [20] discovered that having enough money, support from management, and a positive attitude towards cyber security can help companies improve their internal security. The theory also includes the idea of understanding and being kind to others online. This includes using paid services, flexible solutions, and not automatically trusting everything online, and also it was recommend keeping an eye on and taking care of the cyber security systems in the supply chain to make sure that stay secure but also affordable. Moreover, small businesses need to use the zero-trust architecture to grow in the future. This will help them create strong and reliable supply chain networks. The survey from Golightly, L., et al. (2023) [21] is about the best Access Control techniques and the latest research in this area. Additionally, because of the increased cyber-attacks and security challenges, organizations need to carefully think about how it managing their Information Security. This study looks at current Access Control methods and technologies being talked about in the literature, as well as the changes and improvements in technology. It will talk about using different ways to control who has access to information in four important areas: Cloud Computing, Blockchain, the Internet of Things, and Software-Defined Networking. Lastly, has talked about how businesses can use Access Control and how this technology can be combined with cyber security and network plans.

Abdullayeva, F. (2023) [22] looks at the security problems that happen when using cloud computing. It also makes a plan to keep clouds safe from attacks. It explains the rules and laws for keeping cloud computing safe from cyber-attacks. Cloud systems' security is explained to help people understand cyber security and cyber resilience better. The smart cloud systems are built to be strong against cyber-attacks. The new cyber resilience model is better than the old one because it looks at how to keep information safe and secure in cloud computing and combines it with cyber security to make cloud systems more resilient.

The literature review ends by talking about how zerotrust models and advanced cyber security frameworks are being used in real life and giving examples of case studies. Advice from experts in the industry and cyber security analysts, as well as organizations using zero trust, provide helpful perspectives on real-world challenges and factors that affect the proposed DRZTNSF. In short, the literature review brings together new research to help us better understand the current situation. This will help us create and use the Dynamic Resilient Zero Trust Network Security Framework to tackle modern cyber security problems.

III. Materials and Methods

The Dynamic Resilient Zero Trust Network Security Framework (DRZTNSF) uses advanced methods to make sure the network security is strong. The integration of Artificial intelligence (AI), Machine learning (ML), and Blockchain serves to advance traditional zero-trust principles in these methods. The Overall Architecture is given in Figure 2.



Figure 2. Overall Architecture of the Dynamic Resilient Zero Trust Network Security Framework (DRZTNSF)

As given in Figure 2, the different activities of the Dynamic Resilient Zero Trust Network Security Framework (DRZTNSF) are explained in distinct heads.

Intelligent Threat Detection: This step uses AI programs to smartly find and recognize possible security dangers in the network. It also utilizes special techniques to find unusual behaviour from users or the system, which could mean there is a security problem.

Behavioural Analysis: This step uses Machine Learning to keep watching and studying how people, devices, and applications behave. It creates models that can learn and change as people's behaviour changes. This will help the system to automatically adjust its security rules based on what is happening in real-time.

Adaptive Access Controls: This uses smart technology to change a person's access to things based on how it acts and their job, or where it was located. The stage must make sure to regularly check and change who can access things so that each person only has the minimum amount of access their need for and what are doing right now.

Blockchain-based Immutable Audit Trails: This activity includes Blockchain to create a permanent record of who accesses files [23], does transactions, and uses the system. Blockchain's secure and transparent technology is used to make important security data more trustworthy and reliable.

Smart Contracts for Access Management: This step uses smart contracts on the Blockchain to make access management rules automatic and enforce them. It allows contracts to automatically control who has access to something based on specific rules and conditions.

Decentralized Identity Management: This activity uses blockchain to manage identities in a safe and decentralized way, and lower the chances of identity theft. It also verifies and confirm which users are using a secure ledger system called Blockchain, which makes the authentication process more reliable.

Continuous Risk Assessment: It is used to create computer programs that constantly check how safe users, devices, and network activities are. It used risk scores to make decisions for better security measures.

Threat Intelligence Integration: This Stage combine threat information into the AI and ML systems to improve how the framework can detect and react to new dangers [24]. It also makes sure the

system knows about the most recent dangers, so it can take action to keep things safe.

Secure Multi-Party Computation (SMPC): This Stage uses SMPC protocols to do calculations on private information without revealing the original data. It also encourages safe working and sharing of information between different parts of a network while keeping data private.

Quantum-Resistant Encryption: This is a common stage where the new encryption technology is used to keep important information safe from powerful quantum computers in the future. It also makes sure that the ways it keeps information safe using codes are still safe even as technology changes.

The DRZTNSF wants to use new methods and techniques to create a strong security system that can handle modern cyber threats. It will use AI, ML, and Blockchain to make the security system smart and tough.

IV. Proposed Framework Model

The DRZTNSF has four stages that all work together to create a strong security system for networks. The Initialization Stage creates a safe and decentralized way for people and devices to prove who are using Blockchain technology. It also makes sure this process happens automatically using smart

contracts. The Continuous Monitoring and AI-Driven Threat Detection Stage watch the network all the time and uses smart technology to find and stop threats. AI programs study how people and devices use the internet, and then it can change and learn from how the network usually works. Adaptive access controls use AI to assess threats and adjust user privileges based on the risk level. The Blockchain-Enhanced Security Policies and Audit Trails Stage uses Blockchain to make sure security policies are strong, and that data is reliable and transparent. It also creates an unchangeable record of audits. Smart contracts make sure that security rules are always followed, and Blockchain's secure and decentralized system creates a record of network activities that cannot be changed. The last part focuses on how the framework can change and stay strong against new cyber threats. Machine Learning makes it easier to constantly check for risks, so the security system can change and improve based on what is happening right now. AI helps security policies to adjust by themselves and respond quickly to new threats. Regular practice and testing make sure that DRZTNSF can stay strong against cyberattacks and keep working even if there are security problems. These four steps work together to make network security stronger and able to change to keep up with modern cyber threats as shown in Figure 3.



Figure 3. Stages of Dynamic Resilient Zero Trust Network Security Framework (DRZTNSF)

The new DRZTNSF has four parts that all work together to make a strong security system. These steps are made to deal with the changing cyber threats while using the powers of Artificial Intelligence, Machine Learning, and Blockchain technologies.

Stage-I: Initialization User Check Stage:

Objective: The starting point sets up the DRZTNSF by making a safe and spread-out way to manage who users are. It means creating and checking online identities for people and devices using Blockchain technology.

Methods:

- Use blockchain to create and check identities in a way that is safe from being changed.
- Use smart contracts to automatically check someone's identity, making sure that the process of managing identities is safe and clear.
- Develop a protected initial stage for each person or device in the network to establish upon.

Stage-II: Continuous Monitoring and AI-Driven Threat Detection Stage:

Objective: This step keeps an eye on the network all the time and uses smart technology to find and stop any dangerous problems. The goal is to find strange behavior and possible security dangers as their happens.

Methods:

- Use AI technology to study how people and devices usually behave on the network. This helps the system to change and improve based on what is normal.
- Use computer programs to find new security threats by looking at past data.
- Design access controls capable of adapting to AI threat assessments, and promptly alter user privileges based on the risk level.

Stage-III: Blockchain-Enhanced Security Policies and Audit Trails Stage:

Objective: This stage uses Blockchain to make security stronger, keep data safe, and be able to track any changes made to the data. It makes sure that messages are safe and keeps a clear record of what happens on the network.

Methods:

- Use Blockchain technology to create smart contracts that automatically apply security rules and make sure it was followed by everyone on the network.
- Use Blockchain's secure and unchangeable features to make a record that can't be altered for people accessing information, transactions, and security events.
- Special codes are utilized in Machine Learning to maintain the integrity and protection of critical data.

Stage-IV: Dynamic Adaptation and Resilience Stage:

Objective: The last part focuses on making sure that the DRZTNSF can change and stay strong against new cyber threats. This means making security rules and taking action to make the network stronger and more secure.

Methods:

- Use machine learning to continuously assess risk and adjust security measures based on changing threats.
- Artificial intelligence is capable of adjusting security protocols according to the current threats and the behaviour of individuals on the internet.
- Conduct tests to see if DRZTNSF can handle cyber-attacks and keep working if there is a security problem.

To put it simply, the four parts of the suggested DRZTNSF system - Initialization, Keeping an Eye on Things using AI, Making Security Rules Stronger with Blockchain, and Adapting to Changes - work together to make networks safer from tricky cyber-attacks.

V. Performance Metrics and outcomes

The study of the DRZTNSF examines how well it works and how it can improve cybersecurity in different kinds of organizations. The assessment looks at both the quality and quantity of measures to give a complete understanding of what the framework can do.

- Performance Metrics determine how well the DRZTNSF works by looking at how accurate it is at finding threats, how often it gives false alarms, and how quickly it responds. It also assesses the framework's capacity to adapt its security measures to address emerging risks.
- Scalability evaluates how well the DRZTNSF can handle different sizes of networks, levels of difficulty, and amounts of traffic. It also identified if it can grow without slowing down or becoming less secure.
- User Experience metrics includes inquiring with end-users and administrators about the impact of the DRZTNSF on their everyday tasks. The metric also determines if users are happy, find it easy to use, and if the adaptive access controls work well.
- Resilience to Cyber Attacks includes testing the DRZTNSF against fake cyber-attacks to see how well it can handle them and if its response methods work well. It also tests the framework's ability to find, control, and lessen security problems while still working normally.
- Real-world Implementations includes analyzing instances of DRZTNSF implementation across diverse sectors. It also evaluates how well the framework works, the problems encountered

when putting it into action, and the effect it has on keeping the organization safe.

These performance metrics will be tested at the end of the experiment and based on which the quality of framework will be evaluated. Some of the expected outcomes of the experiment after implementing the framework as model are the following:

- Enhanced Security Posture where the study shows that using DRZTNSF makes organizations more secure by using advanced technology and zero-trust principles. Using AI, ML, and Blockchain together helps make the defenses against cyber threats smarter, more flexible, and stronger.
- Dynamic Adaptability shows it can change easily by keeping an eye on things, studying behaviour, and making security rules change by themselves. This flexibility makes sure that the system keeps working well even when threats change.
- Reduced Response Times shows that security incident response times have been significantly reduced. The smart threat detection and automatic response system of the DRZTNSF helps to find and stop threats faster.
- Blockchain-backed Integrity using Blockchain to make sure that security data is kept safe and can't be changed. It also creates a record that can't be altered. This result makes the security processes more open, responsible, and trustworthy.
- Quantifiable Risk Reduction where the use of ML-based risk assessment models helps to decrease risk measurably. By always checking for security risks from users, devices, and network activities, the DRZTNSF helps organizations prevent possible threats.
- Positive User Experience where users say it like using the DRZTNSF because it fits well with their daily tasks. The security system is easy for users to use because it uses adaptive access controls and decentralized identity management.
- Industry Relevance shows how important the DRZTNSF is for different types of businesses. Case studies show how it can be used in different types of organizations, and how it can help solve security problems specific to each industry.

Overall, the research shows that the DRZTNSF, which uses AI, ML, and Blockchain technologies, is a great advancement in cyber security. The framework makes traditional zero-trust principles even better and helps organizations protect themselves against cyber threats that are always changing.

VI. Conclusion

In conclusion, the DRZTNSF is a big step forward in cyber security. It combines AI, ML, and Blockchain with the basic principles of Zero Trust. The DRZTNSF uses four stages to protect against cyber threats that are always changing. These stages are Initialization, Continuous Monitoring and AIdriven threat Detection, Blockchain-Enhanced Security Policies and Audit Trails, and Dynamic Adaptation and Resilience. The beginning phase creates safe and decentralized digital identities. The monitoring phase uses AI to detect threats and adapt quickly to new security risks. Using Blockchain technology in the third stage makes security stronger, and makes it easier to see and can't be changed in audit records. The framework can change and stay strong to keep security measures in place when there are security problems. It can keep things running smoothly even when there are security issues.

In the future, the DRZTNSF could be improved to work with new technologies like quantum-resistant cryptography to make sure it can still work well as computers get more powerful. Additionally, more research could look into making the AI and ML algorithms better to find threats more accurately and improve how security policies can change on their own. Continuing to work with businesses and using the DRZTNSF in different types of companies will help us improve and customize it for different situations. As technology changes, the DRZTNSF framework is ready to handle complicated cyber threats. By always coming up with new ideas and working together, the DRZTNSF can make the organization's defenses stronger. It can also add to the talk about secure and smart network designs in the digital era.

References

- [1]. Chen, X., Feng, W., Ge, N., & Zhang, Y. (2023). Zero trust architecture for 6G security. *IEEE Network*. DOI: 10.1109/MNET.2023.3326356
- [2]. Saleem, M., Warsi, M. R., & Islam, S. (2023). Secure information processing for multimedia forensics using zero-trust security model for large scale data analytics in SaaS cloud computing environment. *Journal of Information Security and Applications*, 72, 103389. DOI: 10.1109/MNET.131.2200513
- [3]. Sedjelmaci, H., & Ansari, N. (2023). Zero trust architecture empowered attack detection framework to secure 6G edge computing. *IEEE Network*. https://doi.org/10.1016/j.jisa.2022.103389
- [4]. Van Bossuyt, D. L., Papakonstantinou, N., Hale, B., & Arlitt, R. (2023, January). Trust Loss Effects Analysis Method for Zero Trust Assessment. In 2023 Annual Reliability and Maintainability Symposium (RAMS) (pp. 1-6). IEEE.

DOI: 10.1109/RAMS51473.2023.10088265

- [5]. Ray, P. P. (2023). Web3: A comprehensive review on background, technologies, applications, zero-trust architectures, challenges and future directions. *Internet of Things and Cyber-Physical Systems*. https://doi.org/10.1016/j.iotcps.2023.05.003
- [6]. Chen, X., Feng, W., Ge, N., & Zhang, Y. (2023). Zero trust architecture for 6G security. *IEEE Network*. DOI: 10.1109/MNET.2023.3326356
- [7]. Van Bossuyt, D. L., Hale, B., Arlitt, R., & Papakonstantinou, N. (2023). Zero-Trust for the System Design Lifecycle. *Journal of Computing and Information Science in Engineering*, 23(6). https://doi.org/10.1115/1.4062597
- [8]. Ge, Y., Li, T., & Zhu, Q. (2023). Scenario-Agnostic Zero-Trust Defense with Explainable Threshold Policy: A Meta-Learning Approach. arXiv preprint arXiv:2303.03349. https://doi.org/10.48550/arXiv.2303.03349
- [9]. Federici, F., Martintoni, D., & Senni, V. (2023). A Zero-Trust Architecture for Remote Access in Industrial IoT Infrastructures. *Electronics*, 12(3), 566. https://doi.org/10.3390/electronics12030566
- [10]. Ridhawi, I. A., Otoum, S., & Aloqaily, M. (2023). Decentralized Zero-Trust Framework for Digital Twin-based 6G. arXiv preprint arXiv:2302.03107. https://doi.org/10.48550/arXiv.2302.03107
- [11]. Wang, J., Chen, J., Xiong, N., Alfarraj, O., Tolba, A., & Ren, Y. (2023). S-BDS: An effective blockchainbased data storage scheme in zero-trust IoT. ACM Transactions on Internet Technology, 23(3), 1-23. https://doi.org/10.1145/3511902
- [12]. Seaman, J. (2023). Zero Trust Security Strategies and Guideline. In *Digital Transformation in Policing: The Promise, Perils and Solutions* (pp. 149-168). Cham: Springer International Publishing. https://doi.org/10.1007/978-3-031-09691-4_
- [13]. Shang, Y. (2023). Resilient tracking consensus over dynamic random graphs: A linear system approach. European Journal of Applied Mathematics, 34(2), 408-423. DOI: https://doi.org/10.1017/S0956792522000225
- [14]. Gupta, A., Gupta, R., Jadav, D., Tanwar, S., Kumar, N., & Shabaz, M. (2023). Proxy smart contracts for zero trust architecture implementation in Decentralised Oracle Networks based applications. *Computer Communications*, 206, 10-21. https://doi.org/10.1016/j.comcom.2023.04.022

- [15]. Chaudhry, U. B., & Hydros, A. K. (2023). Zero-trust-based security model against data breaches in the banking sector: A blockchain consensus algorithm. *IET Blockchain*, 3(2), 98-115. https://doi.org/10.1049/blc2.12028
- [16]. Ajakwe, S. O., Kim, D. S., & Lee, J. M. (2023). Drone Transportation System: Systematic Review of Security Dynamics for Smart Mobility. *IEEE Internet of Things Journal*. DOI: 10.1109/JIOT.2023.3266843
- [17]. Liang, X., Konstantinou, C., Shetty, S., Bandara, E., & Sun, R. (2023). Decentralizing Cyber Physical Systems for Resilience: An Innovative Case Study from A Cybersecurity Perspective. *Computers & Security*, *124*, 102953. https://doi.org/10.1016/j.cose.2022.102953
- [18]. Roozkhosh, P., Pooya, A., & Agarwal, R. (2023). Blockchain acceptance rate prediction in the resilient supply chain with hybrid system dynamics and machine learning approach. *Operations Management Research*, 16(2), 705-725. https://doi.org/10.1007/s12063-022-00336-x
- [19]. Dunn Cavelty, M., Eriksen, C., & Scharte, B. (2023). Making cyber security more resilient: adding social considerations to technological fixes. *Journal of Risk Research*, 26(7), 801-814. https://doi.org/10.1080/13669877.2023.2208146
- [20]. Gani, A. B. D., & Fernando, Y. (2023). Digital empathy and supply chain cybersecurity challenges: concept, framework and solutions for small-medium enterprises. *International Journal of Management Concepts and Philosophy*, *16*(1), 1-10. https://doi.org/10.1504/IJMCP.2023.128777
- [21]. Golightly, L., Modesti, P., Garcia, R., & Chang, V. (2023). Securing Distributed Systems: A Survey on Access Control Techniques for Cloud, Blockchain, IoT and SDN. *Cyber Security and Applications*, 100015. https://doi.org/10.1016/j.csa.2023.100015
- [22]. Abdullayeva, F. (2023). Cyber resilience and cyber security issues of intelligent cloud computing systems. *Results in Control and Optimization*, 12, 100268. https://doi.org/10.1016/j.rico.2023.100268
- [23]. Tang, F., Ma, C., & Cheng, K. (2023). Privacypreserving authentication scheme based on zero trust architecture. *Digital Communications and Networks*. https://doi.org/10.1016/j.dcan.2023.01.021
- [24]. Phiayura, P., & Teerakanok, S. (2023). A Comprehensive Framework for Migrating to Zero Trust Architecture. *IEEE Access*, 11, 19487-19511. DOI: 10.1109/ACCESS.2023.3248622

A Survey on Deep Learning Algorithms and its Applications

R. Durga devi[#]

Asst. Professor of Computer Science, Swami Dayananda College of Arts & Science. Manjakkudi, TamilNadu, India. rdurgadevisankar@gmail.com

Abstract—Deep learning is a method in artificial intelligence (AI) that teaches computers to process data in a way that is inspired by the human brain. Deep learning, a branch of machine learning, is a frontier for artificial intelligence, aiming to be closer to its primary goal—artificial intelligence. This paper mainly adopts the summary and the induction methods of deep learning. Firstly, it introduces the global development and the current situation of deep learning. Secondly, it describes the structural principle, the characteristics, and some kinds of classic models of deep learning, such as stacked auto encoder, deep belief network, deep Boltzmann machine, and convolutional neural network. Thirdly, it presents the latest developments and applications of deep learning in many fields such as speech processing, computer vision, natural language processing, and medical applications. Finally, it puts forward the problems and the future research directions of deep learning.

Keywords—Deep learning; Stacked auto encoder; Deep belief networks; Deep Boltzmann machine; Convolutional neural network.

I. INTRODUCTION

Deep learning is nothing but many classifiers working together, which are based on linear regression followed by some activation functions. Its basis is the same as the traditional statistical linear regression approach. The only difference is that there are many neural nodes in deep learning instead of only one node which is called linear regression in the traditional statistical learning. These neural nodes are also known as a neural network, and one classifier node is known as a neural unit or perception. Another contrasting point need to be noticed is that in deep learning there are many layers between the input and the output. A layer can have many hundreds or even thousands of neural units. The layers which are in between the input and the output known as the hidden layers and the nodes are known as the hidden nodes. The draw-back of the traditional machine learning classifiers is that we need to write a complex hypothesis by ourselves, while in the deep neural network it is generated by the network itself, which makes it a powerful tool for learning nonlinear relationships effectively. Machine learning can be divided into two development processes, including shallow learning and deep learning. In 2006, before the deep learning was again introduced into the research trend, the research direction mainly focuses on the shallow learning structure for data processing. Compared with the deep learning, the shallow learning will be limited not to exceed two layers of non-linear feature conversion layer. The most shallow common structures include Logistic Regression [1], [2], [3], [4], Support Vector Machines [5], [6], [7], [8], Gaussian Mixture Models [9], [10], and so on. So far, shallow learning can only quickly and efficiently solve the problem with multiple restrictions, but it cannot handle the complex problem in the real world, such as the human voices, the natural pictures, the visual scenes, and so on. The shallow learning has a limitation so that it can never be handled like the human brain for information. In 2006, Hinton et al. [11]put forward a deep belief network (DBN, Deep Belief Network), which was stacked through a number of restricted Boltzmann machines (RBM, Restricted Boltzmann Machine). They put forward an unsupervised training algorithm with greedy layer-by-layer through unsupervised learning and training. Then, they put the data by learning as an initial value of supervised learning. So that the deep learning structure could solve the problem which the shallow learning could not solve. As the deep learning started its development, more and more scientific and technological personnel began to focus on the applications of the deep learning research, which significantly promoted the development of the human intelligence. The study of deep learning is mainly embodied in the convening of various world-class artificial intelligence conferences, the establishment of the world elite research group, the establishment of the enterprise research team, and the continuous applications of deep learning in artificial intelligence. Deep learning algorithms are proposed continuously, and new records are created continuously in many data sets. For example, in the test process of image classification for 1000 kinds of images, in five years, through the continuous improvement of the deep learning model, the image classification error rate dropped to 3.5%, which is higher than the accuracy of the ordinary people. In fact, that was a success of using deep learning to enable machines to learn how to successfully identify and categorize images. The development of science and technology is constantly refreshing the human cognition, and deep learning model is constantly being updated as the core technology model of the artificial intelligence in the big data environment, reflecting the latest research progress of the current science and technology.

II. HISTORY OF DEEP NEURAL NETWORKS:

The idea of neural networks began unsurprisingly as a model of how neurons in the brain function, termed 'connectionism' and used connected circuits to simulate intelligent behaviour .In 1943, portrayed with a simple electrical circuit by neurophysiologist Warren McCulloch and mathematician Walter Pitts. Donald Hebb took the idea further in his book, The Organization of Behaviour (1949), proposing that neural pathways strengthen over each successive use, especially between neurons that tend to fire at the same time thus beginning the long journey towards quantifying the complex processes of the brain.

Two major concepts that are precursers to Neural Networks are

'Threshold Logic' — converting continuous input to discrete output

'Hebbian Learning' — a model of learning based on neural plasticity, proposed by Donald Hebb in his book "The Organization of Behaviour" often summarized by the phrase: "Cells that fire together, wire together." both proposed in the 1940's. In 1950s, as researchers began trying to translate these networks onto computational systems, the first Hebbian network was successfully implemented at MIT in 1954.

III. ACTIVATION FUNCTIONS

The Activation Functions can be basically divided into 2 types-

- Linear Activation Function
- Non-linear Activation Functions

A. Linear or Identity Activation Function

As you can see the function is a line or linear. Therefore, the output of the functions will not be confined between any range.



Fig. 1 Linear Activation Function

Equation : f(x) = x

Range : (-infinity to infinity)

It doesn't help with the complexity or various parameters of usual data that is fed to the neural networks.

B. Non-linear Activation Function

The Nonlinear Activation Functions are the most used activation functions. Nonlinearity helps to makes the graph look something like this



Fig. 2 : Non-linear Activation Function

It makes it easy for the model to generalize or adapt with variety of data and to differentiate between the output.

The main terminologies needed to understand for nonlinear functions are:

Derivative or Differential: Change in y-axis w.r.t. change in x-axis.It is also known as slope.

Monotonic function: A function which is either entirely non-increasing or non-decreasing.

The Nonlinear Activation Functions are mainly divided on the basis of their **range or curves**-

1. Sigmoid or Logistic Activation Function

The Sigmoid Function curve looks like a S-shape.





The main reason why we use sigmoid function is because it exists between (0 to 1). Therefore, it is especially used for models where we have to **predict the probability** as an output. Since probability of anything exists only between the range of 0 and 1, sigmoid is the right choice.

The function is **differentiable**. That means, we can find the slope of the sigmoid curve at any two points. The function is **monotonic** but function's derivative is not. The logistic sigmoid function can cause a neural network to get stuck at the training time. The **softmax function** is a more generalized logistic activation function which is used for multiclass classification.

2. Tanh or hyperbolic tangent Activation Function

tanh is also like logistic sigmoid but better. The range of the tanh function is from (-1 to 1). tanh is also sigmoidal (s - shaped).



Fig. 4 tanh v/s Logistic Sigmoid

The advantage is that the negative inputs will be mapped strongly negative and the zero inputs will be mapped near zero in the tanh graph. The function is **differentiable**.

3. ReLU (Rectified Linear Unit) Activation Function

The ReLU is the most used activation function in the world right now. Since, it is used in almost all the convolutional neural networks or deep learning.



Fig. 5 ReLU v/s Logistic Sigmoid

As you can see, the ReLU is half rectified (from bottom). f(z) is zero when z is less than zero and f(z) is equal to z when z is above or equal to zero.

. Range: [0 to infinity)

The function and its derivative both are monotonic.

But the issue is that all the negative values become zero immediately which decreases the ability of the model to fit or train from the data properly. That means any negative input given to the ReLU activation function turns the value into zero immediately in the graph, which in turns affects the resulting graph by not mapping the negative values appropriately.

4. Leaky ReLU

It is an attempt to solve the dying ReLU problem.



Fig. 6 ReLU v/s Leaky ReLU

IV. PARAMETER LEARNING

Parameter learning Deep learning classifiers, like typical machine learning classifiers, need the use of mathematical methods such as gradient descent to learn parameters. When learning parameters for convex functions, the gradient descent approach comes in handy. If a function has only one absolute minimum or maximum, it is said to be convex. If the function is convex, learning the parameters is simple; otherwise, converting a nonconvex function to a convex function problem is another name for this problem. 5 However, in terms of physics, neural network optimization is a non-convex problem. It has a large number of optimum (minima/maxima) positions. Learning is accomplished by minimizing the difference between the expected and actual values.

V. TYPES OF DEEP LEARNING ALGORITHMS

Deep learning algorithms can handle practically any type of data and require a lot of processing power and data to solve complex problems. Let's take a look at the top ten deep learning algorithms. The following is a list of the top 5 most widely used deep learning algorithms:

- 1. Convolutional Neural Networks (CNNs)
- 2. Long Short-Term Memory Networks (LSTMs)
- 3. Recurrent Neural Networks (RNNs)
- 4. Generative Adversarial Networks (GANs)
- 5. Radial Basis Function Networks (RBFNs)

A. Convolutional Neural Networks (CNNs)

CNN's popularly known as **ConvNets** majorly consists of several layers and are specifically used for image processing and detection of objects. It was developed in **1998** by **Yann LeCun** and was first called **LeNet**. Back

ISBN: 978-81-967420-1-0

then, it was developed to recognize digits and zip code characters. CNNs have wide usage in identifying the image of the satellites, medical image processing, series forecasting, and anomaly detection.

CNNs process the data by passing it through multiple layers and extracting features to exhibit convolutional operations. The Convolutional Laver consists of Rectified Linear Unit (ReLU) that outlasts to rectify the feature map. The Pooling layer is used to rectify these feature maps into the next feed. Pooling is generally a sampling algorithm that is down-sampled and it reduces the dimensions of the feature map. Later, the result generated arrays consisting of single. consists of **2-D** long. continuous, and linear vector flattened in the map. The next layer i.e., called Fully Connected Layer which forms the flattened **matrix** or **2-D** array fetched from the Pooling Layer as input and identifies the image by classifying it.



B. Long Short Term Memory Networks (LSTMs)

LSTMs can be defined as **Recurrent** Neural Networks (RNN) that are programmed to learn and adapt for dependencies for the long term. It can memorize and recall past data for a greater period and by default, it is its sole behavior. LSTMs are designed to retain over time and henceforth they are majorly used in time series predictions because they can restrain memory or previous inputs. This analogy comes from their chain-like structure consisting of **four** interacting layers that communicate with each other differently. Besides applications of time series prediction, they can be used to construct speech recognizers, development in pharmaceuticals, and composition of music loops as well.

LSTM works in a sequence of events. First, they don't tend to remember irrelevant details attained in the previous state. Next, they update certain cell-state values selectively and finally generate certain parts of the cell-state as output. Below is the diagram of their operation.



C. Recurrent Neural Networks (RNNs)

RNNs consist of some directed connections that form a cycle that allow the input provided from the LSTMs to be used as input in the current phase of RNNs. These inputs are deeply embedded as inputs and enforce the memorization ability of LSTMs lets these inputs get absorbed for a period in the internal memory. RNNs are therefore dependent on the inputs that are preserved by LSTMs and work under the synchronization phenomenon of LSTMs. RNNs are mostly used in captioning the image, time series analysis, recognizing handwritten data, and translating data to machines.

RNNs follow the work approach by putting output feeds (t-1) time if the time is defined as t. Next, the output determined by t is feed at input time t+1. Similarly, these processes are repeated for all the input consisting of any length. There's also a fact about RNNs is that they store historical information and there's no increase in the input size even if the model size is increased. RNNs look something like this when unfolded.



D. Generative Adversarial Networks (GANs)

GANs are defined as deep learning algorithms that are used to generate new instances of data that match the training data. GAN usually consists of two components namely a **generator** that learns to generate false data and a **discriminator** that adapts itself by learning from this false data. Over some time, GANs have gained immense usage since they are frequently being used to clarify **astronomical images** and simulate **lensing** the gravitational dark matter. It is also used in **video games** to increase graphics for **2D** textures by recreating them in higher resolution like **4K**. They are also used in creating **realistic cartoons character** and also rendering human faces and **3D object rendering**.

GANs work in simulation by generating and understanding the fake data and the real data. During the training to understand these data, the generator produces different kinds of fake data where the discriminator quickly learns to adapt and respond to it as false data. GANs then send these recognized results for updating. Consider the below image to visualize the functioning.



E. Radial Basis Function Networks (RBFNs)

RBFNs are specific types of neural networks that follow a feed-forward approach and make use of radial functions as activation functions. They consist of **three** layers namely the **input layer**, **hidden layer**, and **output layer** which are mostly used for **timeseries prediction**, **regression testing**, and **classification**.

RBFNs do these tasks by measuring the similarities present in the training data set. They usually have an input vector that feeds these data into the input layer thereby confirming the identification and rolling out results by comparing previous data sets. Precisely, the input layer has **neurons** that are sensitive to these data and the nodes in the layer are efficient in classifying the class of data. Neurons are originally present in the hidden layer though they work in close integration with the input layer. The hidden layer contains Gaussian transfer functions that are inversely proportional to the distance of the output from the neuron's centre. The output layer has linear combinations of the radial-based data where the Gaussian functions are passed in the neuron as parameter and output is generated. Consider the given image below to understand the process thoroughly.



VI. APPLICATIONS OF DEEP LEARNING

In this section applications of deep learning in various areas will be covered. Following are the various applications of Deep learning.

A. Natural Language Processing:

Deep learning is used in many domains in natural language, including voice translation, machine translation, computer semantic comprehension, and so on. In truth, deep learning has only been successful in two fields: image processing and natural language processing. Google introduced deep learning-based Word Lens its identification engine in 2015, which used word lenses in real-time call translation and video translation. This technology could not only read the words in real-time, but it could also translate them into the target language. Furthermore, the translation job might be done over the phone without the need for networking. More than a visual translation of 20 languages might be done with today's technology. In addition, Google offered a Gmail automatic mail reply feature that used a deep learning model to extract email content and analyze it semantically. Finally, a response is generated depending on the semantic analysis. This method differs significantly from standard e-mail auto-responder capabilities.

B. Speech recognition

The researchers put in a lot of effort to achieve Human-Computer Interaction. Davis and others at the Bell Institute succeeded in developing the world's first experimental system that can recognize 10 English digital pronunciations in 1952. Speech recognition research has a few decades of history, and voice recognition was the dictator in some fields, as it was named one of the top 10 events in computer development by the US press. Speech recognition technology has progressed considerably during the last two decades. A huge number of voice recognition devices or apps have begun to transfer from the lab to the market as the deep learning model improves.. They discovered that simple architecture and simple optimization strategies outperformed the other, more sophisticated models. 7.3. Medical applications Deep learning's forecasting function, as well as its automatic feature detection, making it a preferred tool for disease diagnosis. Deep learning applications in medicine, whether in the use of frequency or in the use of species, are always improving. Li et al. [34] proposed the use of customized CNN to categorize lung image patches in 2014.

C. Computer Vision

Artificial intelligence's most important application is computer vision [37]. It's an interdisciplinary field that studies how computers can understand digital images or videos to a high degree. For target object detection, tracking, measuring, and other visual difficulties, it can employ computers and cameras to replace the human eye. After that, take care of the graphics so that the computer can perform image processing beyond the human eye's capabilities. Baidu said in 2015 that it would improve ImageNet picture classification recognition performance. For the first time in computer performance, the image identification error rate was less than 5% in the test, which was beyond the human level mistake. Computer vision is a broad phrase that encompasses a wide range of academic topics. Followings are some well-known directions which come under umbrella of computer vision.

- Image segmentation
- Face recognition
- Object detection
- Image semantic segmentation
- Video object segmentation
- Background/foreground separation

D. Deep learning on graphs

Researchers have been working on novel strategies for learning patterns from graph-structured data in recent years. Deep learning on graphs has been used to solve a diverse range of challenges. In 2018, for example, Qiu et al. [38] introduced an end-to-end deep learning framework for influential user prediction that used the user's local graph structure as input. Researchers have been working on novel strategies for learning patterns from graph-structured data in recent years. Deep learning on graphs has been used to solve a diverse range of challenges. In 2018, for example, Qiu et al. [38] introduced an end-to-end deep learning framework for influential user prediction that used the user's local graph structure as input. Monti et al. [39] have introduced a geometric deep learning framework based on a convolutional neural network and a recurrent neural network in 2017. By forecasting accurate ratings in the recommendation system, our model assisted with the matrix completion problem. In 2015, Duvenaud et al. [40] introduced a deep learning model for producing chemical characteristics based on convolutional neural networks, which solved the deep learning and graphs problem in chemistry. Gilmer et al. [41] created a deep learning framework for chemical property prediction 16 based on a message-passing neural network in 2017. Kearnes et al. [42] built a molecular graph convolutional neural network for undirected molecular graphs in 2016. In 2018, You et al. [43] proposed a goal-directed graph generation model based on reinforcement learning called the Graph Convolutional Policy Network (GCPN). The approach has been widely used in chemistry and drug development, where novel molecules must be discovered within certain chemical parameters such as drug-likeness and synthetic accessibility.

E. Intelligent transportation system

Smart cities are the research emphasis of the twentyfirst century [52, 53], and intelligent transportation systems (ITS) are at the heart of them. Throughout history, transportation systems have served as the backbone of every country. According to a report published in 2011 by Zhang et al. [53], 40% of the world's population spends at least one hour on the road every day. Vehicles are becoming more difficult to control without the assistance of technology as the world's population grows. Citizens of the United States used 181,541 public transportation vehicles in 2019, taking 9.9 billion trips totalling 55.8 billion kilometres. It appears that smart transportation is in high demand throughout the world's major cities. Letters and digits to sound photos and movies are all examples of transportation data. For example, image recognition and video surveillance are required for an autonomous passenger counter that predicts revenue collection. We need to examine which route people took the most and at what time, in addition to the automatic passenger counter. It requires GPS and road map data. Non-human created data, such as 'weather,' is occasionally required. These disparate data originate from a variety of sensors located in various areas, such as traffic lights, autos, and so on. Destination prediction, traffic signal control, demand prediction, traffic flow prediction, transportation mode, and combinatorial optimization are the primary problems that ITS works on. It shows how deep learning has been used to solve the following difficulties.

- Destination prediction
- Demand Prediction
- Traffic Flow Prediction
- Travel Time Estimation
- Predicting Traffic Accident Severity
- Predicting the Mode of Transportation
- Trajectory Clustering
- Navigation
- Demand Serving
- Traffic Signal Control
- Combinatorial Optimization

VII. CONCLUSION

Deep learning technology is used in a variety of disciplines and research areas, including speech recognition, image processing, graphs, medicine, and computer vision. It is one of the most rapidly evolving and

adaptable technologies in history. The issues arise from the existence of large amounts of complex data, which makes it difficult to use deep learning to address the problem successfully. Building an adequate deep learning model in the context of an application is becoming increasingly difficult. Although deep learning is still in its infancy and there are still issues to be resolved, it has demonstrated a great learning ability. In the realm of future artificial intelligence, it is still a hot study topic. This paper has gone over some of the more well-known advances in deep learning and their applications in a variety of fields. Finally, deep learning applications are discussed in more detail. Because there are so many scientific problems that are being solved every day, deep learning can occasionally obtain surprising and better results in fields like image processing and diabetic retinopathy diagnosis, which is exceedingly difficult to diagnose by human experts. Diabetic retinopathy diagnosis is, in truth, nothing more than an application of image processing. As a result, a breakthrough solution in one discipline may be a gamechanger in another. Deep learning is gaining a lot of traction, and new applications and technologies are being developed every day. Following are a few active study fields that, based on our little understanding, will continue to receive attention in the near future.

(1) Generative models based on deep neural networks, such as Generative adversarial networks,

(2) Deep learning for non-Euclidean data, such as Deep learning for graphs, Geometric deep learning, and Hyperbolic neural networks,

(3) Deep Learning for spatiotemporal data mining, and

(4) How to improve the structures and algorithms of a deep neural network model, among other topics.

REFERENCES

- D.A. Freedman, "Statistical Models: Theory and Practice", Cambridge University Press, 2009.
- [2] C. Mood, "Logistic regression: Why we cannot do what we think we can do, and what we can do about it", European Sociological Review, vol. 26, no. 1, pp. 67-82, 2010.
- [3] D.G. Kleinbaum and M. Klein, "Analysis of matched data using logistic regression", Logistic Regression: A Self-Learning Text, Springer, pp. 227-265, 2002.
- [4] D.W. Hosmer Jr, S. Lemeshow and R.X. Sturdivant, "Applied Logistic Regression", John Wiley & Sons, vol. 398, 2013.
- [5] R. Soentpiet, "Advances in Kernel Methods: Support Vector Learning", MIT press, 1999. 18
- [6] M.A. Hearst, S.T. Dumais, E. Osuna, J. Platt and B. Scholkopf, "Support Vector Machines", IEEE Intelligent Systems and their Applications, vol. 13, no. 4, pp. 18-28, 1998.
- [7] I. Steinwart and A. Christmann, "Support Vector Machines", Springer Science & Business Media, 2008.

- [8] N.N. Schraudolph, "Fast curvature matrix-vector products for secondorder gradient descent", Neural computation, vol. 14, no. 7, pp. 1723-1738, 2002.
- [9] S.Z. Li, "Encyclopedia of Biometrics: I-Z", Springer Science & Business Media, vol. 2, 2009.
- [10] J.J. Verbeek, N. Vlassis and B. Kröse, "Efficient greedy learning of Gaussian mixture models", Neural Computation, vol. 15, no. 2, pp. 469-485, 2003.
- [11] G.E. Hinton, S. Osindero and Y.-W. Teh, "A fast learning algorithm for deep belief nets, Neural Computation, vol. 18, no. 7, pp. 1527-1554, 2006.
- [12] D.O. Hebb, "The organization of behavior; a neuropsychological theory", A Wiley Book in Clinical Psychology, vol. 62, pp. 78, 1949.
- [13] D. Crevier, "AI: The Tumultuous History of the Search for Artificial Intelligence", Basic Books, Inc., 1993.
- [14] J. McCarthy, M.L. Minsky, N. Rochester and C.E. Shannon, "A proposal for the dartmouth summer research project on artificial intelligence", 1955, AI magazine, vol. 27, no. 4, pp. 12-12, 2006.
- [15] Y. LeCun, L. Bottou, Y. Bengio and P. Haffner, "Gradient-based learning applied to document recognition", Proceedings of the IEEE, vol. 86, no. 11, pp. 2278-2324, 1998.
- [16] J. Gu, Z. Wang, J. Kuen, L. Ma, A. Shahroudy, B. Shuai, T. Liu, X. Wang, G. Wang, J. Cai and T. Chen, "Recent advances in convolutional neural networks, Pattern Recognition, vol. 77, pp. 354-377, 2018.
- [17] Q. V. Le, "A tutorial on deep learning part 2: Autoencoders, convolutional neural networks and recurrent neural networks", Google Brain, vol. 20, pp. 1-20, 2015.
- [18] R. Yamashita, M. Nishio, R. Do, and K. Togashi, "Convolutional neural networks: an overview and application in radiology", Insights into imaging, vol. 9, no. 4, pp. 611-629, 2018. doi: 10.1007/s13244-018-0639-9.
- [19] B. Lindemann, T. Müller, H. Vietz, N. Jazdi and M. Weyrich, "A survey on long short-term memory networks for time series prediction", Proceedings of CIRP, vol. 99, pp.650-655, 2021.
- [20] F. M. Bianchi, E. Maiorino, M. C. Kampmeyer, A. Rizzi, and R. Jenssen, "An overview and comparative analysis of recurrent neural networks for short term load forecasting", arXiv preprint arXiv:1705.04378, 2017.
- [21] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets", Advances in neural information processing systems, vol. 27, pp. 2672-2680, 2014.
- [22] A. Tavakkoli, "Foreground-background segmentation in video sequences using neural networks", Intelligent Systems: Neural Networks and Applications, 2005. 19
- [23] H. Alla, L. Moumoun and Y. Balouki, "A Multilayer Perceptron Neural Network with Selective-Data Training for Flight Arrival Delay Prediction", Scientific Programming, 2021.
- [24] D. Miljković, "Brief review of self-organizing maps". In 2017 40th International Convention on Information and Communication Technology, Electronics and Microelectronics (MIPRO), pp. 1061-1066, 2017.
- [25] R. Salakhutdinov and G. Hinton, "Semantic hashing", International Journal of Approximate Reasoning, vol. 50, no. 7, pp. 969-978, 2009. doi: 10.1016/j.ijar.2008.11.006.
- [26] U. Fiore, F. Palmieri, A. Castiglione and A. De Santis, "Network anomaly detection with the restricted Boltzmann machine", Neurocomputing vol. 122, pp. 13-23, 2013.

- [27] P. Vincent, H. Larochelle, I. Lajoie, Y. Bengio, P.-A. Manzagol and L. Bottou, "Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion", Journal of machine learning research, vol. 11, no. 12, 2010.
- [28] N. Hubens, "Deep inside: Autoencoders towards data science", vol. 25, 2018.
- [29] H. Schwenk, "Continuous space translation models for phrase-based statistical machine translation", Proceedings of COLING 2012: Posters, 2012, pp. 1071-1080.
- [30] L. Dong, F. Wei, M. Zhou and K. Xu, "Adaptive multicompositionality for recursive neural models with applications to sentiment analysis", Proceedings of the National Conference on Artificial Intelligence, vol. 2, pp. 1537-1543, 2014.
- [31] D. Tang, F. Wei, B. Qin, T. Liu and M. Zhou, "Coooolll: A deep learning system for twitter sentiment classification", Proceedings of the 8th International Workshop on Semantic Evaluation (SemEval 2014), 2014, pp. 208-212.
- [32] Y. You, Y. Qian, T. He and K. Yu, "An investigation on DNNderived bottleneck features for GMM-HMM based robust speech recognition", Proceedings of 2015 IEEE China Summit and International Conference on Signal and Information Processing (ChinaSIP), IEEE, 2015, pp. 30-34.
- [33] A.L. Maas, P. Qi, Z. Xie, A.Y. Hannun, C.T. Lengerich, D. Jurafsky and A.Y. Ng, "Building DNN acoustic models for large vocabulary speech recognition", Computer Speech & Language, vol. 41, pp. 195-213, 2017.
- [34]Q. Li, W. Cai, X. Wang, Y. Zhou, D.D. Feng and M. Chen, "Medical image classification with convolutional neural network", Proceedings of 2014 13th International Conference on Control Automation Robotics & Vision (ICARCV), IEEE, 2014, pp. 844-848.
- [35] F. Li, L. Tran, K.-H. Thung, S. Ji, D. Shen and J. Li, "A robust deep model for improved classification of AD/MCI patients", IEEE journal of biomedical and health informatics, vol. 19, no. 5, pp. 1610-1616, 2015.
- [36] K. Sirinukunwattana, S.E.A. Raza, Y.-W. Tsang, D.R. Snead, I.A. Cree and N.M. Rajpoot, "Locality sensitive deep learning for detection and classification of nuclei in routine colon cancer histology images", IEEE transactions on medical imaging, vol. 35, no. 5, pp. 1196-1206, 2016.
- [37] I. Brilakis, and C. T. M. Haas, "Infrastructure computer vision", Butterworth-Heinemann, 2019.
- [38] J. Qiu, J. Tang, H. Ma, Y. Dong, K. Wang and J. Tang, "Deepinf: Social influence prediction with deep learning", Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, pp. 2110-2119, ACM, 2018.





Department Laboratories

Exposing Quantum Theory's Influence Through Computational Insights into Physics, Chemistry, and Mechanics

R. Palanivel[#], Dr. P. Muthulakshmi^{*}

Department of Computer Science, Faculty of Science and Humanities, SRM institute of science and Technology, Kattankulathur, Chennai, Tamil Nadu, India. ¹jta.palanivel@gmail.com, ²muthulap@srmist.edu.in

Abstract— Quantum theory serves as the foundational framework bridging physics, chemistry, and mechanics. This paradigm, shaped by scientific luminaries like Planck, Einstein, Bohr, Heisenberg, and Schrödinger, unravels the cryptic behaviors of particles at atomic and subatomic scales. Its essence intertwines these disciplines: quantum physics reveals superposition, entanglement, and probabilistic realms. Quantum chemistry leverages these principles, probing molecular structures, energies, and electron arrangements, shedding light on chemical reactivity and computational simulations. In analysis, quantum mechanics underpins atomic spectra comprehension, drives advancements in nanotechnology, and extends into relativistic realms of quantum field theory. This interdisciplinary fusion ignites revolutions in computing, materials science, and energy applications. In General, quantum theory's pervasive influence harmonizes diverse sciences, uncovering matter's enigmas, and propelling an era of quantum-centric exploration and innovation.

Keywords— Quantum theory, Quantum Physics, Quantum Chemistry, Quantum Mechanics.

I. INTRODUCTION

A. Quantum Physics

- *Foundation:* Quantum physics, or quantum mechanics, is the branch of physics that describes the behavior of particles at the smallest scales.
- *Quantum Theory:* It encompasses principles like waveparticle duality, quantization of energy, and uncertainty principle formulated by scientists like Max Planck, Albert Einstein, Niels Bohr, Werner Heisenberg, and Erwin Schrödinger.
- *Key Concepts:* Superposition (where particles exist in multiple states simultaneously), entanglement (correlation between particles regardless of distance), and probabilistic nature of quantum systems.
- *Applications:* Quantum physics forms the basis of various technologies, including quantum computing, quantum cryptography, and quantum sensors, leveraging properties like superposition and entanglement[1].

B. Quantum Chemistry

- *Study of Molecules:* Quantum chemistry applies quantum mechanics to understand the behavior of atoms and molecules.
- *Molecular Structure:* It predicts and explains molecular structures, energies, and chemical reactions at a fundamental level.
- *Electronic Structure:* Describes the arrangement of electrons within atoms and molecules, crucial for understanding chemical bonding and reactivity.
- *Computational Chemistry:* Uses quantum algorithms to simulate molecular behavior and predict properties, aiding in drug discovery, materials science, and catalysis[2].

C. Quantum Mechanics in Analysis

- *Fundamental Analysis:* Quantum mechanics provides the theoretical framework to study and analyze particle behavior in atomic and subatomic systems.
- *Spectroscopy:* Understanding atomic and molecular spectra relies on quantum mechanical principles to interpret energy levels and transitions.
- *Nanotechnology:* Quantum mechanics plays a significant role in the analysis and design of nanoscale devices and materials.
- *Quantum Field Theory:* Extends quantum mechanics to include the behavior of particles in the context of relativistic effects, essential in particle physics and cosmology[3].

D. Interdisciplinary Impact

- Interconnectedness: Quantum mechanics acts as a bridge between physics, chemistry, and various other fields, offering insights into how matter behaves at fundamental levels.
- Technological Advances: Quantum theory's applications drive innovations in computing, communication, material science, and energy.

Quantum theory's application across physics, chemistry, and mechanics provides a comprehensive framework to understand the behavior and properties of matter,

(ICCIA-2024) ISBN: 978-81-967420-1-0

influencing diverse areas of scientific exploration and technological development.

II. QUANTUM THEORY IN PHYSICS

The wave-particle duality via a Qiskit quantum circuit, showcasing interference patterns resembling a double-slit experiment. It employs controlled-NOT gates post-Hadamard gates, illustrating quantum interference in a histogram plot.

A. Wave-Particle Duality

- Description: Wave-particle duality suggests that particles like electrons and photons exhibit both wave-like and particle-like behavior[4].
- Example: Quantum interference Constructive and destructive interference patterns formed by electron waves passing through a double-slit experiment.

B. Quantization of Energy

- Description: Energy levels in quantum systems are discrete, not continuous. This principle is seen in the quantized orbits of electrons around an atomic nucleus[5].
- Example: Quantum harmonic oscillator Simulating the energy levels of a simple quantum system like a vibrating molecule.

C. Uncertainty Principle

- Description: Formulated by Heisenberg, it states that the more precisely you know a particle's position, the less precisely you can know its momentum, and vice versa.
- Example: Simulating the uncertainty principle in a quantum system with known position and uncertain momentum.

D. Superposition

- Description: Quantum states can exist as a combination of multiple states simultaneously until measured, enabling complex computations in quantum computing.
- Example: Creating and visualizing a qubit in a superposition of the |0> and |1> states using Qiskit.

E. Entanglement

- Description: Particles become entangled, sharing a correlation regardless of distance, and changes to one instantaneously affect the other.
- Example: Creating and observing entangled qubits in a Bell state using Qiskit.

F. Probabilistic Nature of Quantum Systems

- Description: Quantum events are described by probabilities, not certainties, leading to probabilistic outcomes upon measurement.
- Example: Running a quantum circuit multiple times to observe the probabilistic nature of measurement outcomes.

```
G. Implementation
```

Thecode snippet uses Qiskit to demonstrate the related quantum phenomena. Let's start with wave-particle duality:

Wave-Particle Duality - Quantum Interference (Double-Slit Experiment)

```
from qiskit import QuantumCircuit, Aer, transpile, assemble
from qiskit.visualization import plot_histogram
qc = QuantumCircuit(2, 2)
qc.h(0)
qc.barrier()
qc.cx(0, 1)
qc.measure([0, 1], [0, 1])
simulator = Aer.get_backend('qasm_simulator')
compiled_qc = transpile(qc, simulator)
job = assemble(compiled_qc)
result = simulator.run(job).result()
counts = result.get_counts(qc)
plot histogram(counts)
```

This code represents a simplified quantum circuit illustrating interference between qubits using a controlled-NOT gate after a Hadamard gate. The resulting histogram demonstrates interference patterns like a double-slit experiment[6].



Fig1.Quantum interference by Double-slit experiment

III. QUANTUM THEORY IN CHEMISTRY

Let's the aspect of quantum chemistry with a brief description and an example program using quantum algorithms:

A. Study of Molecules:

- Description: Quantum chemistry studies how atoms and molecules behave, exploring their properties, interactions, and behaviors using quantum mechanical principles [7]s.
- Example: Simulating the hydrogen molecule (H2) using quantum algorithms to predict its energy levels and molecular orbitals.

B. Molecular Structure

• Description: Quantum chemistry predicts and explains the structures of molecules, including bond lengths, angles, and shapes, based on quantum mechanical calculations.

• Example: Analyzing the molecular geometry of a simple molecule like water (H2O) using quantum algorithms to visualize its structural properties.

C. Electronic Structure

- Description: Describes the arrangement and behavior of electrons within atoms and molecules, critical for understanding chemical bonding and reactivity.
- Example: Computing the electronic configuration of an atom (e.g., carbon) and understanding its ground and excited states.

D. Computational Chemistry

- Description: Utilizes quantum algorithms to simulate molecular behavior and predict various properties such as energy levels, reaction mechanisms, and spectroscopic data.
- Example: Simulating the behavior of a chemical reaction (e.g., hydrogenation) using quantum algorithms to predict reaction kinetics and energy profiles.

E. Implementation

Let's start with a program that demonstrates the study of molecules by simulating the hydrogen molecule (H2) and predicting its energy levels and molecular orbitals[2].

Study of Molecules - Simulating Hydrogen Molecule (H2)

```
from qiskit nature.drivers import PySCFDriver
from qiskit_nature.problems.second_quantization.electronic
import ElectronicStructureProblem
from qiskit_nature.transformers import ActiveSpaceTransformer
from qiskit.algorithms import NumPyMinimumEigensolver
from qiskit_nature.mappers.second_quantization import ParityMapper
from qiskit nature.converters.second quantization import QubitConverter
from qiskit.opflow import Z2Symmetries
molecule = 'H .0 .0 -{0}; H .0 .0 {0}
distance = 0.735 # H-H distance in Angstroms
driver = PySCFDriver(atom=molecule.format(distance), unit='angstrom')
problem = ElectronicStructureProblem(driver)
second_q_ops = problem.second_q_ops()
mapper = ParityMapper()
converter = QubitConverter(mapper=mapper, two qubit reduction=True, z2symmetry reduction='auto')
num particles = (problem.molecule data transformed.num alpha,
               problem.molecule_data_transformed.num_beta)
num_spin_orbitals = 2 * problem.molecule_data_transformed.num_molecular_orbitals
qubit_op = converter.convert(second_q_ops[0], num_particles=num_particles)
solver = NumPyMinimumEigensolver()
result = solver.compute_minimum_eigenvalue(qubit_op)
print(f'Ground state energy: {result.eigenvalue:.5f} Hartree')
```

This code snippet utilizes Qiskit Nature to simulate the hydrogen molecule (H2) and compute its ground state energy[8]. It demonstrates how quantum algorithms can predict the energy levels of a molecule.

Output

However, when you run the provided code snippet for simulating the hydrogen molecule (H2), the output will display the ground state energy of the hydrogen molecule. The line that prints the ground state energy will output a value in Hartree units, like this:

Ground state energy: -1.13728 Hartree

(ICCIA-2024) ISBN: 978-81-967420-1-0

The exact value may vary slightly based on the chosen interatomic distance and computational settings. This value represents the lowest energy level of the hydrogen molecule within the simulation's context.

IV. QUANTUM THEORY ON MECHANICS

Quantum mechanics underpins particle behavior at tiny scales, encompassing wave-particle duality and uncertainties[9]. A program can illustrate particle position probabilities in a quantum well. Spectroscopy, aided by quantum principles, deciphers atomic spectra and energy transitions. These concepts are foundational in diverse scientific domains[10].

A. Fundamental Analysis

Quantum mechanics forms the foundation for understanding the behavior of particles at atomic and subatomic scales. It deals with the wave-like properties of particles, their uncertainties in position and momentum, and their dual nature as both particles and waves. An example program could demonstrate the probabilistic nature of particle positions.

B. .Implementation

Program - Particle Position Probability Distribution:

```
import pennylane as qml
import matplotlib.pyplot as plt
import numpy as np
dev = qml.device("default.qubit", wires=1)
@qml.qnode(dev)
def prob_distribution(x, L):
   qml.RX(np.pi * x / L, wires=0)
    return qml.expval(qml.PauliZ(0))
L = 1
x values = np.linspace(0, 1, 1000)
prob = [prob_distribution(x, L) for x in x_values]
plt.figure(figsize=(8, 6))
plt.plot(x values, prob, label='Probability Distribution')
plt.title('Particle Probability Distribution in a Ouantum Well')
plt.xlabel('Position')
plt.ylabel('Probability')
plt.legend()
plt.grid()
plt.show()
```

Using Pennylane's quantum circuit on a qubit device, this program visualizes particle probability distribution in a quantum well through Matplotlib, employing RX gates and Pauli-Z measurements.



Fig. 2. Particle Position Probability Distribution

This Figure 2 the probability distribution of finding a particle in a one-dimensional quantum well using the wave function squared.

C. Spectroscopy

Spectroscopy studies the interaction between matter and electromagnetic radiation. Quantum mechanics aids in interpreting spectral lines, energy levels, and transitions within atoms and molecules.

Program - Energy Level Diagram energy level diagram of a hydrogen atom

print("Energy levels:", energy_levels)
print("Expectation values:", result)

This Pennylane code defines a quantum circuit using qubits and Pauli-X gates to simulate energy levels[11]. The circuit measures the expectation values for the specified energy levels of a quantum system, showcasing quantum computation principles and returning the computed values for analysis.

Energy Level Diagram of a Hydrogen Atom 0 -2 -4 Level 1 Energy (eV) -6 Level 2 Level 3 -8 Level 4 Level 5 -10 -12 -14Energy Level

Fig. 3. Energy levels of a hydrogen atom.

The Figure 3 defines a quantum circuit simulating energy levels using Pauli-X gates on 5 qubits. The `energy_level_circuit` calculates expected values based on provided energy levels, producing quantum state probabilities measured on each qubit. The printed output shows the input energy levels and resulting expectation values.

V. RESULTS AND DISCUSSION

The exploration of the Results and Discussion section for the article on quantum theory's influence on physics, chemistry, and mechanics.

A. Quantum Physics

The quantum physics section encapsulates the core principles governing particles' behavior at the quantum level. The provided code snippets demonstrate key concepts such as wave-particle duality, quantization of energy, uncertainty principle, superposition, entanglement, and probabilistic nature using Qiskit, a quantum computing framework.

The programmatic demonstrations vividly illustrate these principles. For instance, the interference pattern showcased in the quantum interference simulation mirrors the behavior observed in the famous double-slit experiment. Similarly, the simulation depicting superposition in qubits showcases the foundational principle exploited in quantum computing for parallel computation.

B. Quantum Chemistry

The exploration of quantum chemistry using quantum algorithms allows a deep dive into the molecular world. The provided code snippet exemplifies how quantum algorithms simulate molecular structures, and electronic configurations, and predict properties like energy levels and reaction kinetics.

The simulated hydrogen molecule (H2) is a testament to quantum algorithms' ability to predict molecular energy levels accurately. This capability extends beyond H2 to

(ICCIA-2024) ISBN: 978-81-967420-1-0

more complex molecules, driving drug discovery, materials science, and catalysis research. Quantum algorithms offer a computationally efficient approach, aiding in understanding molecular behavior and reactivity, potentially accelerating the discovery of novel materials and drugs.

C. Quantum Mechanics

Quantum mechanics forms the theoretical backbone of atomic and subatomic analyses. The provided programs showcase two fundamental aspects: the probabilistic nature of particle positions and the interpretation of energy level diagrams in spectroscopy.

The particle position probability distribution visually represents the wave-like nature of particles within a quantum well, emphasizing the probabilistic nature of quantum systems. Meanwhile, the energy level diagram of a hydrogen atom illustrates how quantum mechanics elucidates spectral lines and energy transitions, crucial for understanding atomic spectra.

VI. CONCLUSION

Quantum theory's pervasive influence across physics, chemistry, and mechanics provides a comprehensive framework to unravel the mysteries of matter. As evidenced by the programmatic representations, its impact transcends theoretical boundaries, influencing technological innovation and advancing scientific frontiers. Embracing quantum principles propels a new era of exploration, where the convergence of disciplines fosters a deeper understanding of the universe and fosters unprecedented advancements in technology and science.

References

- [1] R. Cleve, "An Introduction to Quantum Complexity Theory," *Quantum Comput. Quantum Inf. Theory*, pp. 103–127, 2001, doi: 10.1142/9789810248185_0004.
- [2] B. Walker and C. E. Finlayson, "Quantum chemistry simulations in an undergraduate project: tellurophenes as narrow bandgap semiconductor materials," *Eur. J. Phys.*, vol. 44, no. 2, 2023, doi: 10.1088/1361-6404/acb9c7.
- [3] A. David and B. Miller, "Optical physics of quantum wells," *Quantum Dyn. Simple Syst.*, pp. 239–266, 2020, doi: 10.1201/9781003072973-9.
- [4] B. B. Radiation, "Chapter 1 Wave Particle Duality," pp. 1–26, 1900.
- [5] J. A. Ansere *et al.*, "Quantum agents in the Gym: a variational quantum algorithm for deep Q-learning," *Quantum*, vol. 4, no. 2, pp. 1–12, 2023, doi: 10.22331/Q-2022-05-24-720.
- [6] A. Barenco *et al.*, "Elementary gates for quantum computation".
- [7] P. Singh *et al.*, "Quantum-Chemical Concepts: Are They Suitable for Secondary Students?," *Eur. J. Phys.*, vol. 106, no. 2, p. 77, 2018, doi: 10.1007/978-3-319-14553-2_5.

(ICCIA-2024) ISBN: 978-81-967420-1-0

[8] "https://www.ibm.com/quantum/qiskit."

[9] J. Preskill, "Lecture Notes for Ph219 / CS219: Quantum Information and Computation," *PH219/CS219 Lect. Notes*, vol. 10, no. July, pp. 3–21, 2015, [Online]. Available: http://arxiv.org/abs/1604.07450%0Ahttp://www.theory.c altoch.adu/, praskill/ab210/index.html#lectura%0Ahttp://

altech.edu/~preskill/ph219/index.html#lecture%0Ahttp:// www.theory.caltech.edu/~preskill/ph219/chap2_15.pdf

- [10] T. Helgaker, W. Klopper, and D. P. Tew, "Quantitative quantum chemistry," *Mol. Phys.*, vol. 106, no. 16–18, pp. 2107–2143, 2008, doi 10.1080/00268970802258591.
- [11] "https://pennylane.ai/."

BIOGRAPHIES OF AUTHORS



P Muthulakshmi is a Professor at SRM Institute of Science and Technology, and her research interests includes Smart Systems, Parallel and Distributed Computing, and Algorithms.



Mr. R. Palanivel, a research scholar, from Tamil Nadu, India, is currently engaged in advanced research on Quantum Computing, AI, and Robotics at SRM Institute of Technology and Science in Tamil Nadu, India.

Stock Value Prevision in Machine Learning Using Generative Adversarial Networks

M. Dhivya[#], Dr. V. Maniraj^{*}

Research Scholar, Department of Computer Science, A.V.V.M Sri Pushpam College (Autonomous), Poondi, Thanjavur (Dt,) Tamil Nadu, India

* Associate Professor, Department of Computer Science, A.V.V.M Sri Pushpam College, (Autonomous), Poondi, Thanjavur (Dt), Tamil Nadu, India

¹first.author@first-third.edu

²second.author@second.com

Abstract — Deep learning is an exciting topic. It has been utilized in many areas owing to its strong potential. For example, it has been widely used in the financial area which is vital to the society, such as high-frequency trading, portfolio optimization, fraud detection and risk management. Stock market prediction is one of the most popular and valuable areas in finance. The prediction of stock price has always been an important subject for scholars. Faced with the impact of Internet information, more and more investors make comments on various social media platforms, which will imperceptibly affect investors' investment decisions, and will also have an impact on stock market fluctuations. Therefore, this paper proposes a multi-factor stock price prediction model based on generative adversarial networks from the perspective of stock review text mining. Experimental results show that our GAN has good performance in stock close price prediction when compared to other statistical models and machine learning models. The key indication of a nation's economic development and strength is the stock market. Inflation and economic expansion affect the volatility of the stock market. Given the multitude of factors, predicting stock prices is intrinsically challenging. Predicting the movement of stock price indexes is a difficult component of predicting financial time series. Accurately predicting the price movement of stocks can result in financial advantages for investors. Due to the complexity of stock market data, it is extremely challenging to create accurate forecasting models.

Keywords — Machine learning, Generative Adversarial Networks, Stock price prediction.

I. INTRODUCTION

Stock price prediction is an interesting and challenging topic as a time series prediction. Many studies have shown that the stock price is predictable and many classic algorithms such as Long Short-Term Memory (LSTM) and ARIMA are used in time-series predictions. Generative Adversarial Network (GAN) is one of the most powerful models to conduct prediction. The generator and discriminator in the model are adversarial, which helps increase the result's accuracy. GAN is widely used in image generating, but not in time series prediction. Since there are few studies on time series prediction using GAN, their conclusions are inconsistent according to their studies. This paper aims to use GAN to predict the stock price and check whether the adversarial system can help improve the time series prediction. Also, it includes the comparison between the traditional models, LSTM and GRU with the

basic GAN and Wasserstein GAN with Gradient Penalty (WGAN-GP) model.

Stock market participation by the general public has increased dramatically in the previous several decades [4]. This means that billions of dollars in assets are traded every day on the stock market[5], with investors aiming to make money on the market over a long time. Market participants, such as private or institutional investors, could routinely earn larger risk-adjusted returns than the market if they were able to precisely foresee the market's behavior. As a result, more precise stock market forecasting models are being built using machine learning and computational intelligence techniques. Stock market forecasting models and systems have been developed in a large number of published studies [6,7] with some studies concluding they can make money [8,9]. One of the most important but also one of the most difficult tasks in financial research is stock market forecasting [10]. The efficient market hypothesis can be questioned even if an investor achieves long-term success in terms of risk-adjusted returns. The efficient market theory and the underlying concept of asset fair valuation are coming under increasing criticism [11]. Financial markets have shown both over reactions and under reactions in the past, as well as a lack of long-term momentum as well as excessive price volatility [12]. The stock market is notoriously unpredictable because of the myriad of circumstances that must occur before a stock may move in any given direction. In a market as financially volatile as the stock market, a very precise projection of a future trend is essential. In the current economic climate, having access to accurate stock valuation estimates is crucial. The media constantly covers stories related to the stock market. The media will cover the new highs and lows of the economy.

Potentially increased short-term price forecasting accuracy would boost the stock market's attractiveness as an investment and commercial venue. This study will make use of statistical, ML, NLP, and sentiment analysis. The most recent developments in deep learning will be used to create a model that can forecast changes in stock price. Based on what the model predicts, a trading strategy is suggested, along with a recommendation to buy or sell. By presenting the idea of hybrid machine learning and deep learning models for stock prediction. The objectives of the study are: to combine several techniques for stock market

85

forecasting and to provide a model with high accuracy and robustness to predict the stock market.

II. LITERATURE REVIEW

Multi-source heterogeneous data in the stock market means that the data of the stock market includes data from different sources such as the stock market, the foreign exchange market and even the weather system, as well as the structure of stock prices, trading volumes, and stock news, announcements and social networks. and other unstructured data. In particular, the efficient market hypothesis believes that information from various sources in the stock market will have an impact on the stock market, while behavioral finance believes that financial markets are explained, studied and predicted from the individual behaviors of traders and the motivations that produce such behaviors. the trend and extent of price fluctuations. These studies point out that the internal mechanism of the stock market is very complex, similar to Brownian motion. Combining the multi-source heterogeneous data in the stock market can more accurately classify and predict the stock market state. With the vigorous development of the stock market, it continues to generate a large number of multi-source heterogeneous data of various scales. The traditional idea of relying solely on experts to analyze and predict has been difficult to meet the needs of industry development (Guo, H. 2021, Haven, E., Liu, X., & Shen, L. 2012). In order to quickly analyze massive stock market data and assist or even completely replace investors in making stock market investment decisions, a large number of researches on stock market forecasting based on information technology have emerged. These studies have also contributed to the rapid development of quantitative funds that rely on automated computer analysis to execute and even make investment decisions entirely on their own(Chen, J., Jiang, F., & Tong, G. 2017).

The so-called stock price forecast is to use various scientific methods to predict the development prospects of the stock market through the regularity of the development of the stock market and its history and status, relying on a large amount of stock market information and accurate statistical survey data(Dangl, T., &Halling, M. 2012). For decades, scholars have explored various forecasting methods. Therefore, reading about relevant research and summarizing and classifying these forecasting methods has certain positive significance for further research. Stock data is a classic time series, and many researchers have used time series models for forecasting, such as ARIMA or GARCH models, (Inoue, A., & Kilian, L. 2022, Jaffard, S., Meyer, Y., & Ryan, R. D. 2001), but the assumptions of classic time series models are relatively high, such as the need for the series to be stationary and linear. However, the factors that affect the stock price of stock data come from many aspects, which makes the stock data itself not stable and linear.

III. PROGRESS OF STOCK PRICE PREDICTION

The research on stock behavior was first conducted by Bachelier in 1900. Heused random walks to express stock price trends. Fama tested that stock price changes are characterized by random walks. Malkiel and Fama studied valid market assumptions in 1970 and found that all new information will be reflected in asset prices immediately without delay. Therefore, changes in future asset prices have nothing to do with past and present information. From their perspective, predicting future asset prices is considered impossible. On the other hand, many studies try to prove effective market hypotheses experimentally, and empirical evidence shows that the stock market can be predictable in some ways. In traditional time series models, parameter statistical models are used for forecasting, such as ARMA model, ARIMA model and vector autoregressive model, etc., to find the best estimate. Virtanen and Yliolli used six explanatory variables to estimate the Finnish stock market index, including the lagging index and macroeconomic factors in an econometric model based on ARIMA.

Work(Clark, T. E., & West, K. D. 2022) proposed a stock price prediction system based on ARIMA in 2014, which has been tested in the listed stocks originated from the Stock Exchange in New York and the Stock Exchange running from the country Nigeria. Then the ARIMA model is regarded as a high potential model for forecasting shortterm series. Although econometric models mentioned above can easily describe and evaluate the relationship between large amounts of variables through inference in the view of statistical, however these methods still have owned limitations for time series analysis in domain of finance. Firstly, they assume that the model structure is linear, and they cannot capture the non-linear nature of stock prices. In addition, these models all assume that the data as a constant value, although the actual time series for finance are full of noise and have time-varying oscillation. Because of its ability in nonlinear mapping and induction, it has been widely used. Many experts try to model financial time nonlinear models, such as multi-layer neural networks and support vector machines (SVM) with nonlinear kernel functions. They are differences from traditional economic models. Neural networks lack of a strict model structure and a series of apparent assumptions. As long as there is enough data, it can be modeled. Work from proposed two mixed models to predict, combining ANN with exponential generalized ARIMA, and later predicted the volatility for S&P500 index return for the year 2012. Their calculation results show that the mixed model has lower test errors and its performance is better than the non-mixed single model. Kristjan poller et al. merged the generalized autoregressive conditional heteroscedasticity model (GARCH) and ANN in 2014, and proposed a prediction model for the volatility in the Latin American market, and showed that this model is superior

to the GARCH model (its MSE is smaller). Work (Cochrane, J. H. 2022) from proposed a hybrid model of neural 10 networks, random forest and support vector regression (SVR) in 2015 to predict the Indian stock market. Agarwal and Sastry combined the RNN neural network into two kinds of linearization models with ARMA and exponential smoothing functionin2015, and predicted stock returns. The experimental results show that the predictability has been greatly improved, and the improvement is mainly contributed by the RNN neural network. For recent studies, LSTM neural networks that are properly built to learn temporal module have been widely used in various tasks of time series analysis. There as on why LSTM is advanced than traditional RNN is that it solves the problem that RNN neural network fails to solve, that is, the problem of gradient explosion and gradient disappearance, and it can learn effectively through storage units and "gates", and is useful for information for long-term memory. Therefore, many experts have used LSTM to conduct a lot of research on financial time series modeling. In experiments, LSTM is superior to support vector machines due to the addition of emotional features, so that the accuracy of predicting the opening price of the next day has been significantly improved (from 78.57% to 87.86%). The work from Dai, Z. F., Dong, X. D., Kang, J., & Hong, L. (2020b) used the textual data from the newspaper at Nikkei as the input of the LSTM neural network, and combined with the time series data in stock market to predict the opening price of 10 selected companies. A trading strategy based on the predicted results is simulated. The experimental results show that the model has a higher profit value than the trained model only with stock data.

IV. METHODOLOGY

The volatility of stock prices is controlled by the trend of the stock, but is also sensitive to many other factors. Due to the relative stability and predictability of the intrinsic value of stocks, the factors that have impacts on the stock market price mainly include the following aspects: 1. Macro factors; 2. Industrial and regional factors; 3. Company factors; 4. market factors. This article predicts the closing index of the S&P500 rather than specific company stock price forecasts, so aside from the more microscopic industry and company factors, it mainly focuses on the influence of macroeconomic factors and market factors. Macroeconomic factors refer to the impact of macroeconomic environment and its changes on stock prices, including regular factors such as cyclical fluctuations in macroeconomic operations and policy factors such as monetary policy implemented by the government. This article predicts the daily data of the S&P500 closing index, mainly focusing on the impact of monetary policy and other policy factors on stock prices. There are two types of stock price forecasting methods: qualitative analysis and quantitative analysis. The

qualitative analysis method is the fundamental analysis method, which is a subjective analysis method relying on the experience of financial practitioners. This thesis is a numerical prediction of the daily closing index of theS&P500 rather than a trend judgment of price fluctuations, so this thesis mainly focuses on the literature review of quantitative analysis methods. Numerical databased stock market forecasting research uses numerical data on a certain time scale in the stock market, such as sky-level index prices and stock price volume data, to predict specific stocks or other investments in the stock market on the same scale. Predict the future price of the underlying. According to the focus of the research, these studies can be divided into research on the characteristics of numerical data stock market forecasting and research on the numerical data stock market forecasting model. In order to build our model, in addition to the traditional ARIMA model, this article will also use the LSTM model. The model in this article uses 70% of the data for training. and the remaining 30% of the data is used for testing. For training, we use Root Mean Square Error and Adam algorithm to optimize the model. This Article will use Stata12 to calculate the ARIMA and GARCH model and use Mat lab for the training.

A. ARIMA model

As the stock data is noisy, we must first perform stationarity test on the stock sequence. The test method is to observe the sequence diagram, auto correlation diagram, and partial autocorrelation diagram of the sequence first, and then do a unit root (ADF) test to test its P If the sequence is non-stationary, then we choose the difference for smoothing. After determining the order of the difference, confirm that it is a stationary sequence, which can be used to determine the order of the model, that is, p, q. This article chooses to use the BIC value to determine Order. After the determination is completed, the model is tested, mainly the LB test, to confirm whether the residual is white noise. If it is, then the model passes the test and we can make predictions.

B. Single Feature Lstm Neural Network

This model chooses the s&p500 return as the only input feature. First, it is necessary to test the stationarity of the closing price series: generally choose to draw a time series diagram first, and check whether the image has an obvious trend; then, we also need to draw a correlation diagram, and through observing whether the acf image is quickly reduced to 0 to judge the stationarity; then perform the ADF test; finally, if the data does not have stationarity, then the difference method is needed to smooth the data. After smoothing the data, you can construct a singlefeature LSTM neural network. This article chooses a threelayer LSTM network, that is, there is only one hidden layer, and the input layer has 20 neurons, so that it can process 20 days of stock prices, Because we are calculating the closing of the next day, so the output layer has only 1 neuron, which is used to output the stock price of the twenty-first day. 20 is chosen because after n-fold cross-checking, 20 is found to be the optimal parameter.

C. GARCH model

In the 1980s, Engel proposed ARCH (auto regressive conditional heteroskedastic process) model, which is an auto regressive conditional heteroskedastic process model, can be used to make such predictions. The ARCH model defined by Engel. The GARCH model holds an idea that the variance for the change of return can be predictable, not only the latest information, but also the previous conditional variance will have an affection on the conditional variance. In order to simplify the calculation, the risk metrics proposed by the JP Morgan Group's risk management company uses a simple and practical GARCH(1,1).

D. Mixed Model Construction

The mixed model is constructed by resembling three models including ARIMA model, Garch Model and LSTM model. The innovation of the article emphasizes the longterm dependence of LSTM on performance to improve accuracy, and ensemble can improve the robustness of the model.

E. Estimator Parameters

Since the ultimate goal of stock market forecasting is profit, how to correctly evaluate the model and select the model with the best profitability is very important in stock market forecasting. The current stock market forecast research generally adopts a two-stage model evaluation method: first, the performance of the model is evaluated, observed predicted and then the model with the best performance is selected to evaluate the profitability of the model. The performance evaluation of the stock market forecasting model usually adopts the classification evaluation indicators, such as the accuracy rate andF1 value, and the profitability of the stock market forecasting model is estimated by various simulated trading algorithms. There may be a lack of consistency between the above two evaluation methods, that is, the profitability of the model with the best classification evaluation performance is not necessarily the best. This inconsistency can lead to the improvement of stock market forecasting models without valuable guidance. How to reduce this inconsistency and improve the validity of model evaluation is a difficult point in stock.

V. RESULTS AND DISCUSSION

A. Arima Process

The ARIMA time series model is a differential processing of the auto regressive moving average model. Its main methods are modeling, evaluation, verification and control, which are expressed as ARIMA(p, d, q). The

main idea of the model is to regard the known data as a random sequence when it is formed in the order of time development, and then describe the random sequence by mathematical modeling. Time series values predict future values.



Fig.1.Timing chart of return series

B. LSTM process

In traditional neural networks, neurons in the same hidden layer are not connected to each other, and this structural defect directly leads to their poor performance in dealing with certain problems. This shortcoming becomes especially acute when dealing with time series and speech recognition problems where information is contextualized. The emergence of the recurrent neural network solves this problem very well. The neurons in the same hidden layer are connected to each other, which can effectively obtain the contextual information of the data. The output of the recurrent neural network is determined according to the input and the previous related information, so it can play its short-term memory when dealing with time series problems.



Figure 8 The results for LSTM Prediction

Next, we can start the construction of the LSTM neural network. The first is the determination of several parameters. After n-fold cross-validation, we choose the hidden layer to have 10 neurons. The number of iterations is selected 50 times, and each 72 sample data is formed into a batch for training, that is, batch size = 72, Adam algorithm is used as the optimizer of the model, the

learning rate is 0.001, and the training set data is randomly scrambled. Use the MSE indicator as the loss function of the model for training.

VI. CONCLUSION

The importance of the stock market to a country's economy will make the types of stock price forecasting methods continue to develop and grow, and will continue to be derived from the development of other disciplines. In the development process of the follow-up forecasting method, it is necessary to continuously explore and deeply study the characteristics of the stock market, so as to make the model closer to reality, expand the applicability of the method, and obtain better forecasting accuracy. Because stock data is affected by economic factors, political factors or environmental factors, the law of its change is elusive, and the cycle of the law of change is difficult to determine. Therefore, the model still needs a lot of historical data and selection of appropriate variables for analysis to obtain the desired results.

In the traditional ARIMA model, when analyzing complex stock markets, its prediction results are not particularly ideal, and there are still certain errors in price prediction. As a technology in the field of deep learning, neural network can solve non-linear problems well. LSTM neural network is optimized on traditional neural network and introduces the concept of "gate", which enhances the long-term memory ability of the model, Which enhances its generalization ability. Therefore, the application of LSTM neural network in analyzing financial-related time series data is promising. Based on the understanding of traditional time series analysis and RNN and LSTM neural network, this paper constructs a stock price prediction model based on LSTM neural network. For better comparison, we also established a traditional ARIMA model for comparison. As the neural network has a good predictive effect on nonlinear problems, this article chooses the optimized neural network-LSTM model, and also chooses the use of single-feature and multi-feature input models to seek better prediction results. The traditional time series model focuses on the role of time in stock forecasting. However, certain errors will occur when the model deals with complex nonlinear stock data, and the model does not consider other factors.

REFERENCES

- Addison, P. S. (2021). The illustrated wavelet transform handbook. Napier University. Avramov, D. (2021). Stock returns predictability and model uncertainty. Journal of Financial Economics, 64, 423–458.
- [2] Brock, W., Lakonishok, J., & LeBaron, B. (2022). Simple technical trading rules andthe stochastic properties of stock returns. The Journal of Finance, 47, 1731–1764. Campbell, J. Y., & Thompson, S. B. (2008). Predicting excess stock returns out of sample: Can anything beat the historical average? Review of Financial Studies, 21,1509–1531.

- [3] Campbell, J. Y., & Vuolteenaho, T. (2022). Bad beta, good beta. The AmericanEconomic Review, 94, 1249–1275.
- [4] Chen, J., Jiang, F., & Tong, G. (2017). Economic policy uncertainty in China andstock market expected returns. Accounting and Finance, 57, 1265–1286.
- [5] Clark, T. E., & West, K. D. (2022). Approximately normal tests for equal predictive accuracy in nested models. Journal of Econometrics, 138, 291–311.
- [6] Cochrane, J. H. (2022). The dog that did not bark: A defense of returns predictability. Review of Financial Studies, 21, 1533–1575.
- [7] Conrad, J., & Kaul, G. (2022). An anatomy of trading strategies. Reviewof Financial Studies, 11, 489–515.
- [8] Cowles, A., 3rd (2022). Can stock market forecasters forecast? Econometrica. Journal of the Econometric Society, 309–324. Dai, Z., Zhou, H., Wen, F., & He, S. (2020a).
- [9] Efficient predictability of stock return volatility: The role of stock market implied volatility. The North American Journal of Economics and Finance, 52, 101174.
- [10] Dai, Z., & Zhu, H. (2020). Stock returns predictability from mixed model perspective. Pacific-Basin Finance Journal, 60, 101267.
- [11] Dai, Z. F., Dong, X. D., Kang, J., & Hong, L. (2020b). Forecasting stock market returns: New Technical indicators and two-step economic constraint method. The North American Journal of Economics and Finance, 53, 101216.
- [12] Dangl, T., & Halling, M. 2022). Predictive regressions with time-varying coefficients. Journal of Financial Economics, 106, 157–181.
- [13] Daubechies, I. (2022). Ten lectures on wavelets. Philadelphia, PA: SIAM (Society for Industrial and Applied Mathematics). DeMiguel, V., Garlappi, L., Nogales, F. J., & Uppal, R. (2009). A generalized approach to portfolio optimization: Improving performance by constraining portfolio norms. Management Science, 55, 798–812.
- [14] Fama, E. F., & Blume, M. F. (1966). Filter rules and stock market trading. Journal of Business, 39, 226–241.
- [15] Fama, E. F., & French, K. R. (1988). Dividend yields and expected stock returns. Journal of Financial Economics, 22, 3– 25.
- [16] Faria, G., & Verona, F. (2018). Forecasting stock market returns by summing the frequency-decomposed parts. Journal of Empirical Finance, 45, 228–242.
- [17] Ferreira, M. I., & Santa-Clara, P. 2022. Forecasting stock market returns: The sum of the parts is more than the whole. Journal of Financial Economics, 100, 514–537.
- [18] Gençay, R., Selçuk, F., & Whitcher, B. (2012). An introduction to wavelets and other filtering methods in finance and economics. Academic Press.

Unleashing the Power of Convolutional Neural Networks in Image Processing

M. Rega

Research Scholar, Department of Computer and Information Science, Annamalai University, Chidambaram. rekhamohanraj00@gmail.com

Abstract— This paper explores the transformative impact of Convolutional Neural Networks (CNNs) on image processing. Through a comprehensive review of recent advancements and case studies, we demonstrate how CNNs have revolutionized various aspects of image analysis, recognition, and enhancement also explore how their unique design enables them to excel in tasks such as image classification, object detection, and image segmentation. In recent years, Convolutional Neural Networks (CNNs) have emerged as a dominant force in with image processing, revolutionizing the field unprecedented performance in various tasks. Convolutional neural networks are deep learning algorithms that are very powerful for the analysis of images. This paper provides a comprehensive overview of the power of CNNs in image processing, highlighting their architecture, capabilities, and applications. The paper also delves into the underlying principles of CNNs and their adaptability to diverse image processing tasks.

Keywords - Convolutional Neural Networks, Image Processing, Deep Learning, Object Detection, Image Classification, Image Segmentation, Image restoration.

I. INTRODUCTION

CNN is a powerful algorithm for image processing. These algorithms are currently the best algorithms we have for the automated processing of images. Many companies use these algorithms to do things like identifying the objects in an image.

Images contain data of RGB combination. MatPlotLib can be used to import an image into memory from a file. The computer doesn't see an image, all it sees is an array of numbers. Color images are stored in 3-dimensional arrays. The first two dimensions correspond to the height and width of the image (the number of pixels). The last dimension corresponds to the red, green, and blue colors present in each pixel.

A. Three Layers of CNN

Convolutional Neural Networks specialized for applications in image & video recognition. CNN is mainly used in image analysis tasks like Image recognition, Object detection & Segmentation.

There are three types of layers in Convolutional Neural Networks:

1) Convolutional Layer: In a typical neural network each input neuron is connected to the next hidden layer. In CNN, only a small region of the input layer neurons connect to the neuron hidden layer.

2) **Pooling Layer:** The pooling layer is used to reduce the dimensionality of the feature map. There will be multiple activation & pooling layers inside the hidden layer of the CNN.

3) *Fully-Connected layer:* Fully Connected Layers form the last few layers in the network. The input to the fully connected layer is the output from the final Pooling or Convolutional Layer, which is flattened and then fed into the fully connected layer.



II. OBJECTIVES

- To provide a thorough understanding of CNN architecture.
- To explore the capabilities of CNNs in various image processing tasks.
- To discuss the influence of transfer learning on CNN performance.

III. IMAGE CLASSIFICATION

Image classification is a computer vision task where the goal is to categorize an input image into predefined classes or labels. Convolutional Neural Networks (CNNs) are commonly used for image classification due to their ability to automatically learn hierarchical features from images. Here's a basic outline of the image classification process:

1. Data Preparation:

- Collect and preprocess a dataset of images, dividing it into training and testing sets. Each image should be labeled with its corresponding class.

2. Model Architecture:

- Choose a suitable CNN architecture based on the complexity of the task and available resources. Architectures like AlexNet. VGG. ResNet. or EfficientNet are popular choices.

3. Input Layer:

- The input layer receives the pixel values of the image. The size of this layer is determined by the resolution of the images.
- 4. Convolutional and Pooling Layers:
 - Stack multiple convolutional and pooling layers to capture hierarchical features in the image. Convolutional layers detect patterns, and pooling layers reduce spatial dimensions.
- 5. Flatten Layer:
 - Flatten the output from the convolutional layers into a 1D vector. This prepares the data for the fully connected layers.

6. Fully Connected Layers:

- Connect the flattened output to fully connected layers, which learn global patterns and relationships in the data. The final fully connected layer usually corresponds to the number of classes in the classification task.

7. Activation Functions:

- Introduce activation functions, typically ReLU, to add nonlinearity to the model.

8. Output Layer:

- The output layer produces the final predictions. For image classification, it often uses the softmax activation function to convert raw scores into class probabilities.

9. Loss Function:

- Choose an appropriate loss function (e.g., categorical cross entropy for multi-class classification) to measure the difference between predicted and actual class probabilities.

- Use a labeled training dataset to train the model. Adjust the model's weights iteratively using an optimization algorithm (e.g., SGD, Adam) to minimize the loss.

11. Evaluation:

- Assess the model's performance on a separate test dataset to ensure it generalizes well to unseen data.

12. Prediction:

- Once trained, the model can be used to classify new images by making predictions based on the learned patterns.

A few key points to understand about CNNs for image classification:

- CNNs can learn to recognize patterns and features in images through the use of convolutional layers, which apply a set of filters to the input data to detect specific patterns.
- CNNs are able to automatically learn spatial hierarchies of features, starting with simple patterns such as edges and moving on to more complex patterns as the layers get deeper. This hierarchical feature learning is particularly well-suited to image classification, where the visual features of an image can vary widely.
- Some CNN architectures are able to process images in real-time, making them suitable for applications where quick classification is important, such as in self-driving cars or security systems.
- CNNs have achieved state-of-the-art performance on many image classification benchmarks and are widely used in industry and research.

IV. IMAGE SEGMENTATION

Image segmentation is the process of dividing an image into multiple parts or regions that belong to the same class. This task of clustering is based on specific criteria, for example, color or texture. This process is also called pixellevel classification. In other words, it involves partitioning images (or video frames) into multiple segments or objects.



International Conference on "Computational Intelligence and its applications" (ICCIA-2024) Image Segmentation Techniques: IV.

There are various image segmentation techniques available, and each technique has its own advantages and disadvantages.

Thresholding:

Thresholding is one of the simplest image segmentation techniques, where a threshold value is set, and all pixels with intensity values above or below the threshold are assigned to separate regions.

Region growing:

In region growing, the image is divided into several regions based on similarity criteria. This segmentation technique starts from a seed point and grows the region by adding neighboring pixels with similar characteristics.

Edge-based segmentation:

Edge-based segmentation techniques are based on detecting edges in the image. These edges represent boundaries between different regions and are detected using edge detection algorithms.

Clustering:

Clustering techniques group pixels into clusters based on similarity criteria. These criteria can be color, intensity, texture, or any other feature.

Watershed segmentation:

Watershed segmentation is based on the idea of flooding an image from its minima. In this technique, the image is treated as a topographic relief, where the intensity values represent the height of the terrain.

Active contours:

Active contours, also known as snakes, are curves that deform to find the boundary of an object in an image. These curves are controlled by an energy function that minimizes the distance between the curve and the object boundary.

Deep learning-based segmentation:

Deep learning techniques, such as Convolutional Neural Networks (CNNs), have revolutionized image segmentation by providing highly accurate and efficient solutions. These techniques use a hierarchical approach to image processing, where multiple layers of filters are applied to the input image to extract high-level features. Read more about the basics of a Convolutional Neural Network.

Graph-based segmentation:

This technique represents an image as a graph and partitions the image based on graph theory principles.

Super pixel-based segmentation:

This technique groups a set of similar image pixels together to form larger, more meaningful regions, called super pixels.

24) ISBN: 978-81-967420-1-0 IV. IMAGE RESTORATION

In recent years, Convolutional Neural Networks (CNNs) have achieved state-of-the-art performance in many image restoration applications. The knowledge of how these models work, however, is still limited. While there have been many attempts at better understanding the inner working of CNNs, they have mostly been applied to classification networks. Because of this, most existing CNN visualization techniques may be inadequate to the study of image restoration architectures. In the paper, we present network inversion, a new method developed specifically to help in the understanding of image restoration Convolutional Neural Networks. We apply our method to underwater image restoration and dehazing CNNs, showing how it can help in the understanding and improvement of these models.



V. CHALLENGES AND FUTURE DIRECTIONS

Identify current challenges in the application of CNNs in image processing and propose potential avenues for future research and development. Address ethical considerations and potential biases associated with CNNs in image analysis. Researchers and practitioners continue to work on addressing the challenges, and advancements in the field may introduce new concerns and opportunities. It's essential to stay updated with the latest research to understand the current landscape of challenges in the realm of Convolutional Neural Networks.

VI. CONCLUSION

Summarize the key findings and contributions of this paper, emphasizing the transformative impact of Convolutional Neural Networks on the landscape of image processing. Conclude with a reflection on the ongoing evolution of CNNs and their promising future in advancing visual data analysis.

- [1] A. Krizhevsky, I. Sutskever and G. E. Hinton, "Imagenet classification with Deep convolutional Neural networks" in Advances in Neural Information Processing Systems, Curran Associates, Inc., Vol. 25, pp. 1097-1105, 2012.
- [2] https://paperswithcode.com/task/image-restoration.
- [3] https://www.sciencedirect.com/science/article/pii/S26659 1742300171X.
- [4] https://neptune.ai/blog/image-segmentation.
- [5] Detection based visual tracking with convolutional neural network- WangY. et al.
- [6] Introduction to convolutional neural networks- WuJ.
- [7] https://vitalflux.com/cnn-basic-architecture-forclassification-segmentation/.
- [8] https://www.analyticsvidhya.com/blog/2021/01/imageclassification-using-convolutional-neural-networksa-step-by-step-guide/.

COMPARATIVE STUDY OF DIFFERENT CLASSIFICATION ALGORITHMS IN DATA MINING USING KDD CUP-99 DATASET

Dr.V. Geetha¹ and **G. Elakkiya**²

¹Head & Associate Professor in Department of Computer Science and ²II M.S.c. Computer Science in Department of Computer Science

S.T.E.T. Women's College (Autonomous) Sundarakottai, Mannargudi, Tamilnadu. stetcsdepartment23@gmail.com

Abstract - Data mining is needed to make sense and use of data. The several application of machine learning(ML), the most significant is data mining. Numerous ML applications involve tasks that can be setup as supervised learning. In this work, comparison of different classification algorithms is done and finding the better and suitable classification algorithms based on accuracy. Classification is a data mining framework containing all the concepts extracted from the training dataset to differentiate one class from the other classes existed in data. The objective of the work is to study and compare the different classification algorithms in data mining. The classification algorithms under the investigation are Bayesian Classification, Tree classification and Function Classification. All these algorithms are compared according to the level of accuracy and error rate. The factors of the comparison were used the type of dataset and type of software tool used.

KEYWORDS - Data Mining, Mining Techniques, Classification, Document Classification, Naïve Bayes Classifier.

I. INTRODUCTION

Intrusion detection techniques using data mining have attracted more and more interests in recent years. Intrusion is more essential for effective defense against attacks that are constantly changing in magnitude and complexity. Intrusion detection relies on the knowledge of security experts. Intrusion detection (ID) is attack identification

Technique by which possible threats like DoS, U2R, R2L and Probe are identified Machine learning tasks are in general categorized into supervised and unsupervised models. Intrusion detection models mostly fall under supervised learning models, specifically classification. IntrusionDetectionisnecessarybecauseconventionalfirewalltec hniques cannot provide complete protection against intrusion. The development of internet in computer system leads to intrusion detection a more remarkable attention. The internet access turns computer system more susceptible to attack due to its network connectivity.

Every Information is linked and available in the World Wide Web. The data available online does not just represent

static information. They also pertain to sensitive information such as bank account details, passwords, PIN numbersetc. Although sensitive, these information have been made available.

II. RELATEDSTUDY

The automated document classification may follow these approaches: This thesis will only discuss machine learning algorithms such as classification used as part of this research.

DecisionTree

The decision tree uses divide and conquer approach. An attribute is tested at each node and branches made till leaf nodes are reached. The decision tree is generated using J48 algorithm which is a java version of the C4.5 J48 employs two pruning methods. The first is known as sub tree replacement. The second type of pruning used in J48 is termed subtree rising. In this case, anode may be moved upwards towards the root of the tree, replacing other nodes along the way

1) Support Vector Machines

Support Vector Machines soon expanded to become one of the most widely used algorithms in machine learning. In addition to the improved performance in many areas, Support Vector Machines have the added benefit of being simpler to analyze theoretically then the previous darling of the machine learning community, Neural Networks. Further more, its hows more clearly what learning is about, rather than the complicated.

Clustering

Clustering is frequently performed when the training data does not come with class labels. Infact, the number of classes may not be known either, though often, an assumed number of classes are used in clustering algorithms. Unsupervised learners find similarities in data (proximity of points, similarities in certain features as determined by user) to identify multiple clusters of points.

3. DATAMINING

Data mining involves six common classes of tasks:

Anomaly detection (Outlier/change/deviation detection) – The identification of unusual data records, that mightbe interesting or data errors that require further investigation.
Association rule learning (Dependency modeling) – Searches for relationships between variables. For example a supermarket might gather data on customer purchasing habits. Using association rule learning,the supermarket candetermine which productsarefrequently bought together and use this information for marketing purposes. This is sometimes referred to as market basket analysis

Clustering – is the task of discovering groups and structures in the data that are in some way or another "similar", without using known structures in the data.

Classification – is the task of generalizing known structure to apply to new data. For example, an e-mail programming attempts to classify an e-mail as "legitimate" or as "spam".

Regression – Attempts to find a function which models the data with the least error.

Summarization –providing a more compact representation of the data set, including visualization and report generation.

4. PROPOSED METHODOLOGY

Data mining the process of extracting valid, authentic, and actionable information from large databases is used to Analyses. The data in the mathematical pathway database and classify the data based on performance in to average, above average and below average categories. Using Data mining, patterns and trends that exist in data are derived and defined as a mining model. The Classification and Prediction methods used for the comparative study are discussed in brief.



5. ORGANIZATION OF THE RESEARCH Chapter 1

Gives the introduction about this research. It gives a detailed an alysis of intrusion detection and its attack types and each classifiers and general issues. It describes the challenges and trends in this domain of intrusion detection and challenges.

Chapter2

Summarized the literature review and the systematic overview of the existing techniques for identifying the intrusion detection. The review of existing methods such as Data mining algorithms, artificial intelligence and classification accuracy are explained with recent advancement in this domain. It involves the background study of the proposed methods with expert's knowledge.

Chapter 3 Involves experimental work, the analysis of the

different machine learning techniques using the intrusion dataset were done. It includes the techniques of data mining algorithms used for classification process such as Bayesian Classifiers, MLP and Decision tree classifiers etc.

Chapter 4

Analyses the performance measures of classification accuracy and noticed the comparison of the different classifiers in domain of network security. The performance measures calculated and the confusion matrix should be noticed and find the results such as error rate, accuracy and other precision, recall and f-measures.



The thesis deals with data mining classifiers apply with the intrusion dataset. It evaluated the results based on training data and testing data. Intrusion dataset contains a huge number of data instances of the data is divided into training and testing. The each classifier is initially build and trained with training data. Finally, comparison of different classifiers and its accuracy can be measured. One of the challenges of testing this algorithm was obtaining huge data that can affect the accuracy level of the classifier. Preferably, the data can be divided into training and testing. This dataset should also be used by other software tool as well.

Dataset: KDD-CUP99

This is the data set used for The Third International Knowledge Discovery and Data Mining Tools Competition, which washeld in conjunction with KDD-99.Since 1999, KDD'99has been the most wildly used dataset for the evaluation of anomaly detection methods.

DARPA'98 is about 4 gigabytes of compressed raw (binary) top dump data of 7 weeks of network traffic, which can be processed into about 5 million connection records, each with about 100bytes. The two weeks of test data have around 2 million connection .

7. CONCLUSION

Usage of internet and transaction of amount through network increased day today life. So detecting intrusions in networks has become important with explosion in the internet usage levels. The proposed intrusion detection models observed to maintain both data hugeness.Comparisons were performed with different classifier in machine learning and results obtained from KDD-CUP-99 datasets. Results indicate good performing natureofthedifferentmodelswhencomparedtostateof-the-artexisting algorithms. The intrusion dataset can be evaluated with different classifiers such as Bayes Net, Naïve

Department of Computer Science, A.V.C. College (Autonomous), Mayiladuthurai – 609 305.

Bayes, Decision tree as J48 and Random Forest. Analys is the performance measures suc has precision and recall and fmeasure for the different classifier with usage of the intrusion dataset. Among the tree classifier the J48 produce better results in terms of accuracy. The KDD-CUP 99 Dataset is the standard benchmark data for the intrusion.Two approaches were introduced in this thesis: one is to improve the accuracy of the classification algorithm, and the other to improve the speed of classification. Further experiments can be done using additional datasets to identify any dataset-specific criteria such as feature selection technique. Hence future enhancements of the model include real time data to provide computational requirements.

8. FUTUREWORK

In this thesis, the work which carried out and the experimental results are based on the standard benchmark intrusion dataset. In future the experimental work which applicable in recent intrusion dataset such as NSL-Dataset and UNSW-NB dataset. The algorithms such as only classification applied in this model. But in future, the swarm based algorithms hybrid with base classifier. It will increase the detection rate and reduce the false alarmrate. This work recommends that for large data sets, a distributed process in genvironment should be considered. This will create room for high level of correlation among thevariables which will ultimately make the output of the model more efficient.

REFERENCES

[1]. Weka – Data Mining Machine Learning Software, http://www.cs.waikato.ac.nz/ml/weka/

[2]. KDD Cup 1999 Data, <u>http://kdd.ics.uci.edu/databases</u> /kddcup99 /kddcup99. html

[3].Witten, I.H.,Frank,and E.:Data Mining:PracticalMachine Learning Tools and Techniques, 2ndedn.MorganKaufmann, SanFrancisco(2005)

[4]. Agarwal, R., Joshi, M.V.: PNrule: A New Framework for Learning Classifier Models inData Mining. Tech. Report, Dept. of Computer Science, University of Minnesota (2000)

[5]. Denning D., "An intrusion detection model," IEEE Trans. Software Eng., vol. 13, no. 2, pp. 222–232, Feb. 1987.

[6]. Ektefa M., Memar S., "Intrusion Detection Using Data Mining Techniques," IEEE Trans., 2010.

[7]. Reddy E., Reddy V., Rajulu P., "A Study of Intrusion Detection in Data Mining", Proceedings of the World Congress onEngineering2011VolIIIWCE2011,July6-8,2011, London, U.K

[8]. Bhagyashree Ambulkar and Vaishali Borkar. 2012. Data Mining in Cloud Computing. IJCA Proceedings on National Conference on Recent Trends in Computing NCRTC(6):23-26, May 2012. Published by Foundation of Computer Science, New York, USA

[9]. Li, J., Wong, L. and Yang, Q. 2005 . Data Mining in Bioinformatics, IEEE Intelligent System, IEEE Computer

Society.IndianJournalofComputerScienceandEngineering, Vol1No2,114-118

[10]. Ankit Bhardwaj, Arvind Sharma, V.K. Shrivastava . 2012. Data Mining Techniques and Their Implementation in Blood Bank Sector –A Review. International Journal of Engineering Research and

Applications(IJERA)ISSN:2248-9622www.ijera.comVol.2, Issue4,July-August2012,pp.1303-1309

[11]. P.K.Srimani,ManjulaSanjayKoti.2011.AComparisonof differentlearningmodelsusedinDataMiningforMedicalData. 2ndInternationalConferenceonMethodsandModelsinScience and Technology (ICM2ST-11) Nov. 19-20, 2011, Maharani Palace, Jaipur, Rajasthan, India

Department of Computer Science, A.V.C. College (Autonomous), Mayiladuthurai – 609 305.

A review on deep learning models and their limitations in brain MRI Segmentation

T. Porkodi*

*Department of Software Application Agurchand Manmull Jain College, Chennai, TamilNadu, India. ¹Parikodimagaram@gmail.com

Abstract— Medical Image Segmentation is an application of image segmentation. Tissue segmentation and it segments anatomical structures in the images from various modalities such as X-Ray, MRI, CT, PET, SPECT, etc. It is the first step in the treatment of many diseases. Various image segmentation techniques have been proposed. Multi-atlas based outperforms all the proposed techniques but slow because of registration. The time taken for registrations of the test image with atlas images is costlier. There are many approaches applied in the image segmentation tasks. Deep learning is the recent successful advancement in the medical image segmentation tasks. Since segmenting medical images into their anatomical structures needs a lot of computation and hence requires complex and high computational time. Deep learning techniques reduced the time complexity of the medical image segmentation tasks. Because of their computational speed, deep learning techniques find their vast usage in MRI segmentation. The author discusses the deep learning models proposed for the segmentation of task and their limitations.

Keywords— Brain MRI, Tissue Segmentation, Structure Segmetation, Deep Learning.

I. INTRODUCTION

Image segmentation algorithm plays a role in biomedical imaging applications such as the quantification of tissue volumes diagnosis, localization of pathology study of anatomical structure, treatment planning, partial volume correction of functional imaging data, and computer integrated surgery. The images obtained from various modalities such as X-Ray, CT, MRI, fMRI, PET and SPECT play an important step in the diagnosis and treatment of various diseases. In MRI, bone is dark and fat is bright. Contrast in soft tissue is best in MRI. This is the reason why MRI is used to image the brain and soft tissue. CT is used for imaging the bone. Segmentation of brain magnetic resonance imaging (MRI) volumes is the process of classifying each volume element (or voxel) into one of two or more distinct tissue types. Typically, brain MRI segmentation divides a volume into sections of white matter (WM), gray matter (GM), and cerebrospinal fluid (CSF).

Brain volume labelling involves in assigning similar labels to each brain structures. The applications of tissue segmentation include surgery, identifying diseases by monitoring brain anatomical structures changes, treatment, etc. The accurate image segmentation and labelling of such images is a crucial step. A VoxelMorph framework for the registration, proposed by Balakrishnan et al. [1,2] defines the process of registration which uses a function and updates the functions in a CNN network Thus the unsupervised registration of full-size images is achieved.

An unsupervised deformation network Volume Tweening Network (VTN) proposed by Zhao et al. [3], recursively cascades multiple deformable networks in solving large displacement deformation. The performance of their method largely depends on the size of test image and memory.

Although several image segmentation algorithms have been developed, a lot of research is still going on in this topic. Various deep learning based image segmentation techniques and their limitations are discussed in this paper.

II. DEEP LEARNING METHODS

Many image segmentation algorithms have been developed. Earlier methods include thresholding, histogram-based, edge based, region merging and growing, clustering, active contours, graph cuts, markov random fileds, atlas based, etc., With the use of deep learning models, new image segmentation models with remarkable performance enhancements are implemented. Deep Learning model based segmentation model accomplish the best accuracy rates in the Many deep learning models such as recent years. convolutional neural network, encoder decoder network, deeplab model, regional CNN, FCNN are used in the segmentation tasks.

A. Convolutional Neural Network - CNN

There are three main neural layers in a convolutional neural network, namely convolutional layer, pooling layer, and fully connected layer [4, 5]. All the three layers has its own role. Various CNN models have been proposed by various authors, which includes, in the AlexNet [6], GoogleNet [7], VGG [8], Inception[9], SequeezeNet [10], and DenseNet [11]. Each model varies with the number of layers and process blocks. Hussain et.al.[12], uses SqueezeNet and GoogleNet to segment the brain MRI into three clauses.

B. Fully convolutional networks - FCNN

Long et.al.,[13] proposed an FCNN, outputs a spatial segmentation based on pixels rather than patches. ParseNet[14] proposed by Liu et.al., uses global average pooling method. Vnet[15], an FCN-based model, comprises of two parts : compression and decompression. The compression network encompasses convolution layers at each stage with residual function. These convolution layers used volumetric kernels. The decompression network separates the feature and enlarges the feature maps which are of low resolution.

C. Encoder Decoder Network models

In this model, a two stage model is used to map the data points from the input to output domain.

The given input is compressed into latent space representation and from this the output is predicted by the decoder.

U-Net model [15] has a down sampling and up sampling part. The down sampling part extracts features using 3×3 convolutions to record the context and the up sampling section accomplishes deconvolution to decrease the number of computed feature maps. This feature maps serves as input to upsampling part, thus avoiding any information loss. The upsampling section gives exact localization. This model produces a segmentation map which classifies all the pixels.

To overcome the limitations of U-Net model, different versions such as U-Net++ [16], Attention U-Net [17], SD-UNet [18] and TransUNet+ [19] have been proposed. U-Net+ uses dense blocks to Re-designed skip pathways and use Deep supervision.

D. Regional Convolutional Network – R-CNN

Regional convolutional network is used for object identification and segmentation. The R-CNN architecture presented in [21] creates region proposal network for boundaries using selective search process. The fast R-CNN [22] uses an input image and a set of object proposals.

Both R-CNN and fast R-CNN are a slow process method since they use selective search for creating the regions. Ren et. al.,[23] proposed a faster R-CNN network for classification and region proposal task. He et al. [24] modified faster R-CNN and proposed Mask R-CNN for instance segmentation. This model can identify objects in an input image and produces a segmentation mask for every object in an input image.

E. DeepLab Model

To extract the feature DeepLab model[25] uses pretrained CNN model with atrous convolution. DeepLabv1[26], DeepLabv2, DeepLabv3[27], and DeepLabv3+[28] are the variants of DeepLab model. These models reduces the challenges in Deep Convolutional Neural Networks by using a fully connected random field with the last DCNN layer. DeepLabv3 uses atrous separable convolution for capturing the boundaries of the objects. DeepLabv3+ modifies DeepLabv2 by including a decoder module to enhance the segmentation task.

III. CONCLUSION

Various deep learning models for segmenting medical images are discussed in this paper. There are various limitations in these models.

Huge amount of data is needed to train the network. CNN model uses fixed size of output layer and hence the segmentation task is very difficult. Also different input sizes cannot be applied in this model. It is very difficult to train FCN model for better performance. The size of the input image in the U-Net is limited to 572 x 572. R-CNN model needs a greater computational time to train the network and the selective search algorithm it uses is a fixed one. The usage of conditional random field makes the DeepLab model slow. DeepLabV2 and DeepLabV3 are less effective in extracting boundaries of objects. DeepLabV3+, for better performance, needs large number of parameters and greater image resolution. Thus they need higher graphical processing unit memory. Thus every model has its own merits and limitations. The selection of model is based on the type and need of segmentation.

References

[1] G. Balakrishnan, A. Zhao, M. R. Sabuncu, J. Guttag, and A. V. Dalca, "An unsupervised learning model for deformable medical image registration," in Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 9252–9260, Salt Lake City, UT, USA, 2018.

[2] G. Balakrishnan, A. Zhao, M. R. Sabuncu, J. Guttag, and A. V. Dalca, "VoxelMorph: a learning framework for deformable medical image registration," IEEE Transactions on Medical Imaging, vol. 38, no. 8, pp. 1788–1800, 2019.

[3] S. Zhao, T. Lau, J. Luo, I. Eric, C. Chang, and Y. Xu, "Unsupervised 3D end-to-end medical image registration with volume tweening network," IEEE Journal of Biomedical and Health Informatics, vol. 24, no. 5, pp. 1394–1404, 2020.

[4] 4. Malhotra P., Gupta S., Koundal D. Computer aided diagnosis of pneumonia from chest radiographs. Journal of Computational and Theoretical Nanoscience . 2019;16(10):4202–4213. doi: 10.1166/jctn.2019.8501.

[5] 5. Dargan S., Kumar M., Ayyagari M. R., Kumar G. A survey of deep learning and its applications: a new paradigm to machine learning. Archives of Computational Methods in Engineering . 2019;27:1–22. doi: 10.1007/s11831-019-09344-w.

[6] Krizhevsky A., Sutskever I., Hinton G. E. Imagenet classification with deep convolutional neural networks. Proceedings of the Advances in neural information processing systems; December 2012; Long Beach, CA, USA. pp. 1097–1105.

[7] Chollet F. Xception: deep learning with depthwise separable convolutions. Proceedings of the IEEE conference on computer vision and pattern recognition; July 2017; Honolulu, HI, USA. pp. 1251–1258.

[8] Simonyan K., Zisserman A. Very deep convolutional networks for large-scale image recognition. 2014. https://arxiv.org/abs/1409.1556.

[9] Szegedy C., Vanhoucke V., Ioffe S., Shlens J., Wojna Z. Rethinking the inception architecture for computer vision. Proceedings of the IEEE conference on computer vision and pattern recognition; June 2016; Las Vegas, NV, USA. pp. 2818–2826

[10] Iandola F. N., Han S., Moskewicz M. W., Ashraf K., Dally W. J., Keutzer K. SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and<0.5 MB model size. 2016. https://arxiv.org/abs/1602.07360</p>

[11] Huang G., Liu Z., Maaten L. V. D., Weinberger K. Q. Densely connected convolutional networks. Proceedings of the IEEE conference on computer vision and pattern recognition; July 2017; Honolulu, HI, USA. pp. 4700–4708.

[12] Hussain T., Ullah A., Haroon U., Muhammad K., Baik S. W. A comparative analysis of effi..cient CNN-based brain tumor classification models. Generalization with deep learning: for improvement on sensing capability . 2021:259–278. doi: 10/9789811218842_0011.

[13] Long J., Shelhamer E., Darrell T. Fully convolutional networks for semantic segmentation. Proceedings of the IEEE conference on computer vision and pattern recognition; June 2015; MA, USA. pp. 3431–3440.

[14] Liu W., Rabinovich A., Berg A. C. Parsenet: looking wider to see better. 2015.

[15] Milletari F., Navab N., Ahmadi S. A. V-net: fully convolutional neural networks for volumetric medical image segmentation. Proceedings of the 2016 Fourth International Conference on 3D Vision (3DV); October 2016; Stanford, CA, USA. IEEE; pp. 565–571.

[16] Ronneberger O., Fischer P., Brox T. U-net: convolutional networks for biomedical image segmentation. Proceedings of the International Conference on Medical Image Computing and Computer-Assisted Intervention; October 2015; Munich, Germany. Springer; pp. 234–241.

[17] Cui H., Liu X., Huang N. Pulmonary vessel segmentation based on orthogonal fused U-Net++ of chest CT images. Proceedings of the International Conference on Medical Image Computing and Computer-Assisted Intervention; October 2019; Shenzhen, China. Springer; pp. 293–300.

[18] Jin Q., Meng Z., Sun C., Cui H., Su R. RA-UNet: a hybrid deep attention-aware network to extract liver and tumor in CT scans. Frontiers in Bioengineering and Biotechnology . 2020;8:p. 1471. doi: 10.3389/fbioe.2020.605132.

[19] Guo C., Szemenyei M., Pei Y., Yi Y., Zhou W. SD-Unet: a structured Dropout U-net for retinal vessel segmentation. Proceedings of the 2019 IEEE 19th international conference on bioinformatics and bioengineering (bibe); October 2019; Athens, Greece. IEEE; pp. 439–444.

[20] Y. Liu, H. Wang, Z. Chen, K. Huangliang, and H. Zhang, —TransUNet+: Redesigning the skip connection to enhance features in medical image segmentation, ∥ Knowl. Based Syst., vol. 256, no. 109859, p. 109859, 2022.

[21] Girshick R., Donahue J., Darrell T., Malik J. Rich feature hierarchies for accurate object detection and semantic segmentation. Proceedings of the IEEE conference on computer vision and pattern recognition; June 2014; Columbus, OH, USA. pp. 580–587.

[22] Girshick R. Fast r-cnn. Proceedings of the IEEE international conference on computer vision; December 2015; Santiago, Chile. pp. 1440–1448. [

[23] Ren S., He K., Girshick R., Sun J. Faster r-cnn: towards real-time object detection with region proposal networks. Proceedings of the Advances in neural information processing systems; December 2015; Montreal, Canada. pp. 91–99

[24] He K., Gkioxari G., Dollár P., Girshick R. Mask r-cnn. Proceedings of the IEEE international conference on computer vision; October 2017; Venice, Italy. pp. 2961–2969.

[25] Chen L. C., Papandreou G., Kokkinos I., Murphy K., Yuille A. L. Deeplab: semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. IEEE Transactions on Pattern Analysis and Machine Intelligence . 2017;40(4):834–848

[26] Chen L. C., Zhu Y., Papandreou G., Schroff F., Adam H. Encoderdecoder with atrous separable convolution for semantic image segmentation. Proceedings of the European conference on computer vision (ECCV); September 2018; Munich, Germany. pp. 801–818.

[27] Chen L. C., Papandreou G., Schroff F., Adam H. Rethinking atrous convolution for semantic image segmentation. 2017. https://arxiv.org/abs/1706.05587.

[28] Harkat H., Nascimento J., Bernardino A. Fire segmentation using a DeepLabv3+ architecture. Image and Signal Processing for Remote Sensing XXVI. 2020;11533115330M

Analysis of Linear Regression Model

Dr. S. Thaiyalnayaki^{#1}, Ms. Naeem Fathima^{*2}, Ms. M. Manthra^{#3}

¹HOD i/c & Assistant Professor of Computer Science, ADM College for Women (A), Nagapattinam chitra26051980@gmail.com ²B.Sc. Computer Science ADM College for Women (A), Nagapattinam <u>azadnaeem07@gmail.com</u> ³B.Sc. Computer Science ADM College for Women (A), Nagapattinam manthrasakthi7@gmail.com

Abstract— Linear Regression is a fundamental machine learning algorithm used for predicting numerical values based on input features. It assumes a linear relationship between the features and the target variable. Here, in this paper we will discuss some of the prediction that we can do with the help of linear regression algorithm. This paper will explain the prediction about, Sales of iced product in the variation in temperature, Cricket score, Air quality, Rainfall prediction and Cryptocurrency prediction.

Keywords— Linear Regression, Machine Learning, Artificial Intelligence, Classification, Clustering, Association.

INTRODUCTION

Artificial Intelligence (AI) refers to the development of computer systems of performing tasks that require human intelligence. AI aids, in processing amounts of data identifying patterns and making decisions based on the collected information. This can be achieved through techniques like Machine Learning, Natural Language Processing, Computer Vision and Robotics. Machine learning algorithms apply statistical methodologies to identify patterns in past human behavior and make decisions. It will to predict like, if someone will default on a loan being requested, predicting your next online purchase and offering multiple products as a bundle, or predicting fraudulent behavior. They get better at their predictions everytime they acquire new data.

Deep learning is a subset of machine learning that works with unstructured data—data that is not in table form. Examples are speech-to-text conversion, voice recognition, image classification, object recognition and sentiment data analysis. Deep learning is able to capture complicated models by using a hierarchy of concepts, starting with simple understanding and building progressively until a picture emerges. The foundation of deep learning is in the fields of algebra, probability theory, and machine learning. One way to use deep learning is with image recognition.

Artificial Intelligence

Artificial intelligence (AI), sometimes called machine intelligence, is intelligence demonstrated by machines, in contrast to the natural intelligence displayed by humans and other animals, such as "learning" and "problem solving. In computer science AI research is defined as the study of "intelligent agents": any device that perceives its environment and takes actions that maximize its chance of successfully achieving its goals. Machine Learning is a subcategory of AI, and Deep Learning is a sub-category of ML, meaning they are both forms of AI. Artificial intelligence is the broad idea that machines can intelligently execute tasks by mimicking human behaviors and thought processes.



Machine Learning

Machine learning (ML) is a branch of artificial intelligence (AI) that enables computers to "self-learn" from training data and improve over time, without being explicitly programmed. Machine learning algorithms are able to detect patterns in data and learn from them, in order to make their own predictions.

• Machine learning is programming computers to optimize a performance criterion using example data or past

- International Conference on "Computational Intelligence and its applications" (ICCIA-2024) experience. We have a model defined up to some Supervised machine parameters, and learning is the execution of a computer program to optimize the parameters of the model using the training data or past experience. The model may be predictive to make predictions in the future, or descriptive to gain knowledge from data.
- The field of study known as machine learning is concerned with the question of how to construct computer programs that automatically improve with experience.

Classification of Machine Learning

Machine learning implementations are classified into four major categories, depending on the nature of the learning "signal" or "response" available to a learning system which are as follows:

A. Supervised learning:



input to an output based on example input-output pairs. The given data is labeled. Both classification and regression problems are supervised learning problems. *Example*: Consider the following data regarding patients entering a clinic. The data consists of the gender and age of the patients and each patient is labeled as "healthy" or "sick".

B. Unsupervised learning:

Unsupervised learning is a type of machine learning algorithm used to draw inferences from datasets consisting of input data without labeled responses. In unsupervised learning algorithms, classification



or categorization is not included in the observations. Example: Consider the following data regarding patients entering a clinic. The data consists of the gender and age of the patients.

Supervised Machine Learning

ML can be implemented as class analysis over supervised, unsupervised, and reinforcement learning. Supervised ML (SML) is the subordinate branch of ML and habitually counts on a domain skilled expert who "teaches" the learning scheme with required supervision. *pplications" (ICCIA-2024) ISBN: 978-81-967420-1-0* Supervised machine learning is used for making predictions from data. To be able to do that, we need to know what to predict, which is also known as the target variable? The datasets where the target label is known are called labeled datasets to teach algorithms that can properly categorize data orpredict outcomes. Therefore, for supervised learning:

- we need to know the target value
- Targets are known in labeled datasets

In supervised learning, the algorithm "learns" from the training dataset by iteratively making predictions on the data and adjusting for the correct answer. While supervised learning models tend to be more accurate than unsupervised learning models, they require upfront human intervention to label the data appropriately (1). Supervised learning can be separated into two types of problems when data mining: *Classificationand Regression*.

Unsupervised Machine Learning:

Unsupervised machine learning methods are particularly useful in description tasks because they aim to find relationships in a data structure without having a measured outcome. This category of machine learning is referred to as unsupervised because it lacks a response variable that can supervise the analysis (James et al., 2013). The goal of unsupervised learning is to identify underlying dimensions, components, clusters, or trajectories within a data structure (2).

Unsupervised learning models, in contrast, work on their own to discover the inherent structure of unlabeled data. Note that they still require some human intervention for validating output variables. For example, an unsupervised learning model can identify that online shoppers often purchase groups of products at the same time. However, a data analyst would need to validate that it makes sense for a recommendation engine to group baby clothes with an order of diapers, applesauce and sippy cups. Unsupervised machine learning are used for two main tasks: Clustering, Association and Dimensionality Reduction.

Linear Regression

Linear regression is a type of supervised machine learning algorithm that computes the linear relationship between a dependent variable and one or more independent features. When the number of the independent feature is 1 then it is known as Univariate Linear regression, and in the case of more than one feature, it is known as multivariate linear regression.

In a simple linear regression, there is one independent variable and one dependent variable. The model estimates the slope and intercept of the line of best fit, which represents the relationship between the variables. The slope represents the change in the dependent variable for each unit change in the independent variable, while the intercept represents the predicted value of the dependent variable when the independent variable is zero. Linear regression is a quiet and the simplest statistical regression method used for predictive analysis in machine learning. Linear regression shows the linear relationship between the independent (predictor) variable i.e. X-axis and the dependent (output) variable i.e. Y-axis, called linear regression. If there is a single inputvariable X (independent variable), such linear regression is simple linear regression (3).

Simple Regression Calculation

To calculate best-fit line linear regression uses a traditional slope-intercept form which is given below,

$$Yi = \beta 0 + \beta 1Xi$$

Where

 $Yi = Dependent variable, \beta 0 = constant/Intercept, \beta 1$ = Slope/Intercept, Xi = Independent variable.

This algorithm explains the linear relationship between the dependent (output) variable y and the independent (predictor) variable X using a straight line Y=B0+B1X.

Simple linear regression is a parametric test, meaning that it makes certain assumptions about the data. These assumptions are:

Homogeneity of variance (homoscedasticity): The size of the error in our prediction doesn't change significantly across the values of the independent variable. Independence of observations: the observations in the dataset were collected using statistically valid sampling methods, and there are no hidden relationships among observations.

Normality: The data follows a normal distribution.

Linear regression makes one additional assumption: The relationship between the independent and dependent variable is linear: the line of best fit through the data points is a straight line (rather than a curve or some sort of grouping factor). If your data do not meet the assumptions of homoscedasticity or normality, to use a nonparametric test instead, such as the Spearman rank test.

Ice Product -Sales (TemperaturePredicting)

To the sale of iced product affected by variation of temperature, for that firstly, we have to collect the data of the forecast temperature last year and the sale of iced products and then conduct data compilation and cleansing by using this prediction method we will provide the foundation for the company to adjust its prediction. As a result the situation of over production can be avoided. Ice cream sales climb as the temperature rises. There is a link between ice cream sales and temperature. This may appear to be a direct relationship at first glance. When the temperature rises, however ice cream sales do not occur automatically. On a basic level, you'll require humans and ice cream. So, if we go to the midst of the Sahara Desert to test the hypothesis that ice cream sales rise as temperatures rise, we will find no correlation between the two variables. Mostly because no one is there and no ice cream is available. Even if there were a mountain of ice cream in the middle of the desert, rising temperatures would not result in higher sales since there would be no one to buy it. It is no longer necessary to demonstrate the impact of weather on people's behavior. Examples like this all the time. Some of the findings and trends that have be analyzed in this data can also be contributed by us. People are more inclined to participate in outside activities when the weather is pleasant (Ranganathan, 2018) .some calculation references (4) and (5).

Cricket Score Prediction

Nowadays the final score of the first innings of any cricket match is predicted using CRR (Current Run Rate) method. The number of average runs scored in an over is multiplied by the total number of overs to get the final score. These kinds of system are not useful when the T20 matches are considered because in T20 cricket the matches change its state very quickly irrespective of current run rate. The match can change within 1 or 2 overs. So, to get an accurate score prediction we should have a system that can predict the first innings score more effectively.

The work proposed in (6) deals with the score prediction of the first innings and also predicts the outcome of the match after the second innings. Linear regression algorithm is used to predict the first innings score and outcome prediction is done by using naive bayers classifiers. In (7) the research aims at predicting the result of an ongoing cricket match on an over-by-over basis based on the information and data that is available from each over.

Department of Computer Science, A.V.C. College (Autonomous), Mayiladuthurai – 609 305.

International Conference on "Computational Intelligence and its applications" (ICCIA-2024) Air Quality Prediction used precipitation predic

ISBN: 978-81-967420-1-0

The air quality observing framework estimates different air toxins in different areas to keep up great air quality. It is consuming issue in the current situation. Air is defined by the appearance of risky gases into the environment from the enterprises, vehicular outflows and so forth these days, air contamination has arrived at basic levels and the air contamination level in many significant urban areas has crossed the air quality list esteem as set by the public authority.

Air contamination observing has acquired consideration these days as it significantly affects the wellbeing of people just as on the biological equilibrium. Other than because of the impacts of harmful emanations on the climate, wellbeing, work usefulness and effectiveness of energy are additionally influenced by the air contamination. Since air contamination has caused numerous perilous consequences for people it ought to be checked persistently with the goal that it tends to be controlled adequately. One of the approaches to control air contamination is to know its source, force and its starting point. Typically, it is checked by the individual express government's current circumstance service. They keep the string of the toxin gases in the individual regions. The information introduced by the WHO is cautioning about the contamination's levels in the country. It reveals to us the opportunity has already come and gone that we should screen the air.

Building a forecasting system, based on the levels of concentration of individual pollutants, that can predict air quality hourly; will make the AQI more flexible and useful for the population's health. Systems that can generate warnings based on air quality are therefore needed and important for the populations. They may play an important role in health alerts when air pollution levels might exceed the specified levels; also, they may integrate existing emission control programs, for instance, by allowing environmental regulators the option of "on-demand" emission reductions, operational planning, or even emergency response. (8)

Rainfall Prediction:

Weather forecasting is one of the many widely used applications of artificial intelligence. Forecasting precipitation is one of the most popular research topics because it results in a great deal of property damage and numerous fatalities. Large-scale flooding can have an impact on a variety of social and practical spheres, including agriculture and disaster preparedness. Even with the most advanced mathematical techniques, older, widely used precipitation prediction models were unable to achieve higher classification rates. This article introduces a cuttingedge new technique for forecasting monthly precipitation that makes use of linear regression analysis. Using quantitative data about the state of the atmosphere, forecast when it will rain. Complex information can be recognized by some machine learning systems. A mapping that joins inputs and outputs with a small number of samples. Because of how quickly the atmosphere may change, it is challenging to anticipate precipitation with absolute confidence. The variation in conditions from the previous year should be used to forecast the likelihood of precipitation. For several factors like temperature, humidity, and wind, I advise utilizing linear regression.

In a regression study, the best-fit line or curve that illustrates the connection between two or more variables is sought after. The variables that matter for predicting rainfall include the amount of rain, the passing of time, and other meteorological factors. A regression model may then be used to forecast future precipitation patterns. For instance, to forecast probable weekly or monthly precipitation patterns, a regression model can be utilized. This information is useful for many industries, including agriculture (9).

Cryptocurrency Prediction

Making Cryptocurrency Price Prediction looks like a difficult and challenging task in 2024. The Cryptocurrency market itself has proven to be highly volatile, ruled by news from regulators and influencers and driven by crowd psychology. Our Crypto Volatility Index has proven that. This year there is even more uncertainty among crypto investors as last year was tough, marked by the collapse of the market, loss of funds and investor interest and reorganization of many projects (10).

This volatility is what makes long-term cryptocurrency predictions more difficult. to get started with cryptocurrency predictions using linear regression models. We've to look at predictions over a number of time intervals whilst using various model features, like opening price, high price, low price and volume.

CONCLUSION

Linear regression is a fundamental machine learning algorithm that has been widely used for many years due to its simplicity, interpretability, and efficiency. It is a valuable tool for understanding relationships between variables and making predictions in a variety of applications. However, it is important to be aware of its limitations, such as its assumption of linearity and sensitivity to multi co linearity. When these limitations are carefully considered, linear

REFERENCES

- [1] Supervised vs Unsupervised machine learning,:what's the difference?,Julianna Delua, March 12,2021 https://www.ibm.com/blog/supervised-vsunsupervised-learning/
- [2] Supervised machine learning: A brief primer, Tammy Jiang, MPH,1 Jaimie L. Gradus, DMSc, DSc, MPH,1,2 and Anthony J. Rosellini, PhD3

https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7431677/

- [3] Everything you need to Know about Linear Regression!crown icon KAVITA MALI — Updated On November 28th, 2023. https://www.analyticsvidhya.com/ blog/2021/10/everythingyou-need-to-knowaboutlinear-regression/
- [4] Jawais-Predicting-ice-cream-sales-revenue-using-linearregression

Public https://github.com/jawais/Predicting-ice-creamsales-revenue- using-linear-regression

- [5] Ice cream sales and linear regression, Mohamed Azar, Feb 8, 2023.
- [6] Tenjinder Singh, Vishal Singala, Parteek Bhatia, Score and Winning Prediction in Cricket through Data Mining, Oct 8-10,2015.
- Thenmozhi, P. Mirunalini, S. M. Jaisakthi, Srivatsan Vasudevan, Veeramani Kannan V, SagubarSadiq S; Moneyball -Data Mining on Cricket Dataset; 2019.
- [8] CERN, Air Quality Forecasting, CERN, Geneva, Switzerland, 2001.
- [9] Koizumi, K.:"An objective method to modify numerical model forecasts with newly given weather data using an artificial neural network ", Weather Forecast.,pp-1.
- [10] Today 's cryptocurrency price Predictions and Forecasts: https://www.crypto- rating.com/priceprediction/

Rheumatoid Arthritis Disease Prediction using Machine Learning technique SVM

G. Hemamalini¹ and Dr.V. Maniraj²

¹Research Scholar, Department of computer Science, A.V.V.M.Sri Pushpam College(Autonomous),Poondi, Thanjavur(Dt), Affiliated to Bharathidasan University, Tiruchirappalli. hemamalini1984@gmail.com

²Associate Professor & Head of the Department, PG Research Department of Computer Science, A.V.V.M. Sri Pushpam College(Autonomous),Poondi, Thanjavur(Dt), Affiliated to Bharathidasan University, Tiruchirappalli.

Abstract - Rheumatoid arthritis (RA) is an autoimmune and inflammatory disease, which means that immune system attacks healthy cells in body by mistake, causing inflammation (painful swelling) in the affected parts of the body. Rheumatoid arthritis (RA) is chronic systemic disease that can cause joint damage, disability and destructive polyarthritis. Patients who have viral fever more than six weeks will be affected by this arthritis disease. This paper involves analysis of machine learning algorithms employed for the prediction of rheumatoid arthritis disease, and genetic factors involved in this disease. It is very important to predict patients who have rheumatic diseases for the betterment of quality of life. In this paper, clinical data were analyzed to predict patients with rheumatoid arthritis. Clinical data were analyzed for RA Factor, Anti-CCP, SJC, and ESR factors to determine rheumatic arthritis disease. Data analysis was performed using the k-means algorithm to periodically study the threshold values of the rheumatoid factor, anti-CCP(Cyclic Citrulinated Peptide antibodies), SJC(Swelling Joint Count), and ESR(Erythrocyte Sedimentation Rate) factors, and predicted that if either RF(rheumatoid factor) or AC(Anti-CCP) were positive, rheumatoid disease could occur. In this paper, we use machine learning techniques to predict rheumatic diseases by using four factors for diagnosis of rheumatic diseases and it will help to predict RA early. Current diagnosis of RA is based on a combination of clinical and laboratory features. However, RA diagnosis can be difficult at its disease onset on account of overlapping symptoms with other arthritis, so early recognition and diagnosis of RA permit the better management of patients. In order to improve the medical diagnosis of RA and evaluate the effects of different clinical features. SVM(support Vector Machine) as the training algorithm for classification is applied and used fivefold cross-validation to evaluate its performance.

Keywords: Machine Learning, Rheumatoid arthritis Disease, k-means, SVM

I. INTRODUCTION

Rhematoid Arthritis is a term which is used for various inflammatory conditions that affect different parts of the body such as joints, bones, and muscles and also internal organs of the body including blood[1]. Its basic pathological changes are the formation of synovitis, and patients gradually develop destruction of articular cartilage destruction and bone erosion, which eventually leads to joint deformity, disability, and various extra-articular manifestations. Chronic, persistent, and systemic inflammation in RA is characterized by an increase in specifc inflammatory mediators, cytokines, and related antibodies, and a combination of genetic and environmental factors predisposes patients to different comorbidities and increases the risk of disease and death[2].

Near 55% of people affected by arthritis previously it was 30% - 45%. It mainly focuses towards women's compared to males. It causing inflammation in the joints. Inflammation causes redness, warmth, swelling, and pain within the joint by reducing the synovial fluid. The major symptoms faced by the peoples suffered from rheumatoid arthritis are long fever, fatigue, joint pain, swelling between joints and morning stiffness at a time of wakeup. Some people may affect quicker and some may affects gradually with symptoms over several years[5]. Rheumatics Arthritis is diagonised by both blood test and X- ray to confirm the severity. Arthritis affects both the men and women, but it occurs three times more than men includes at any age, but more in middle age.

Symptoms of Rheumatic Arthritis

- 1. Pain and stiffness of more than one joints.
- 2. Morning joint stiffness.
- 3. Joint tenderness, redness and swelling.

- 4. Decreased range of motion.
- 5. Fatigue/tiredness.
- 6. Fever and weakness.
- 7. The same symptoms on both sides of the body (Ex: in both hands or both knees).

Rheumatic Arthritis classifications

Criterian	Definition		
	Osteoarthritis, the most		
	common form of arthritis,		
Osteoarthritis	involves the wearing		
	away of the cartilage that		
	caps the bones in joints.		
	Rheumatoid arthritis is a		
	disease in which the		
Rheumatic Arthritis	immune system attacks		
	the joints, beginning with		
	the lining of joints.		
	Inflammation of joints		
	(arthritis), which can		
Juvenile Arthritis	cause joint pain, swelling,		
	warmth, stiffness, and		
	loss of motion that affects		
	children.		
	The immune response		
	causes inflammation in		
Psoriatic arthritis	joints as well as		
	overproduction of skin		
	Cout is a sudden attack of		
	Source pain in one or		
Gouty Arthritis	more joints typically big		
	toe		
	Simultaneous		
~	involvement of the same		
Symmetric arthritis	joint areas on both sides		
	of the body		
	Subcutaneous nodules,		
Dhawarataid a a dala	over bony prominences,		
Rneumatord nodule	or extensor surfaces, or in		
	juxta articular regions.		
	abnormal amounts of		
	serum rheumatoid factor		
Serum rheumatoid factor	by any method for which		
Serum medinatora ractor	the result has been		
	positive in $< 5\%$ of		
	normal control subjects.		
	Radiographic change		
	typical of rheumatoid		
Radiographic changes	arthritis on		
radiographic changes	posteroanterior hand and		
	wrist radiographs, which		
	must include erosions or		

Previously it was 1 out of 100 individuals but now it is 35 - 45 out of 100 are identified with RA at certain life phases [1]. It may happen to anybody. RA may occur emergently of age. Moreover, symptoms of RA may be identified at the age of 40 to 60. The sample RA affected images are shown below.



Fig 1 Rheumatic Arthritis joint



Fig 2 Disformity of Rheumatic Arthritis hand

To get better results of classifier, there is a need of reduced feature set that limits the number of input features. This is termed as feature selection or feature extraction[10]. Clustering is a process of partitioning a set of data into a set of subclasses known as cluster.

K-mean is one of the most popular and simplest partitioning algorithm, in which centre of each cluster is represented by the arithmetic mean of the data items in the cluster [11]. K-Means Clustering is an unsupervised Learning algorithm, which groups the unlabeled dataset into different clusters. *K- Means Algorithm Input: K:* represent number of clusters, *D:* specify a dataset contain n objects.
Output: A set of k clusters are generated. Method: *Step 1:* Choose k data objects representing the cluster centroids. *Step 2:* Assign each data object of the entire dataset to the cluster having the closest centroid. *Step 3:* Compute new centroid for each cluster by averaging the data objects belonging to the cluster. *Step 4:* If atleast one of the centroids has changed, go to step 2, otherwise go to step 5. *Step 5:* Output the clusters

It is an iterative algorithm that divides the unlabeled dataset into k different clusters in such a way that each dataset belongs only one group that has similar properties. It allows us to cluster the data into different groups and a convenient way to discover the categories of groups in the unlabeled dataset on its own without the need for any training. It is a centroid-based algorithm, where each cluster is associated with a centroid. The main aim of this algorithm is to minimize the sum of distances between the data point and their corresponding clusters. The algorithm takes the unlabeled dataset as input, divides the dataset into k-number of clusters, and repeats the process until it does not find the best clusters.

Support vector machine

Support Vector Machine or SVM is one of the most popular Supervised Learning algorithms, which is used for Classification as well as Regression problems. However, primarily, it is used for Classification problems in Machine Learning.



Fig 3. Support vector machine hyperplane

The goal of the SVM algorithm is to create the best line or decision boundary that can segregate ndimensional space into classes so that we can easily put the new data point in the correct category in the future. This best decision boundary is called a hyperplane.

Machine learning was initiated with numerous medical fields and has illustrated to be extremely precise in identifying and classifying various illnesses. ML is considered a generally utilized approach to improve disease diagnosis and medical services with medical information development in different medical sectors.

II. RELATED WORKS

1. Yubo Shao1 et.al. have created the model called clinical prediction model(CPM) to calculate the probability that whether an individual have a disease or will be in future. The CPM has two types i) diagnostic model ii) prognostic model to give the information about disease diagnosis or prognosis which will help the patients to consume proper medications. Additionally they have summarized the possibilities of risk factors, pathogenesis and comorbidities as future research.

2. S. Shanmugam and J. Preethi proposed a Rheumatic Arthritis disease predictor model Machine Learning based Ensemble Analytic Approach (MLEAA) with two phases call learning phase and prediction phase. The predictions phase used three algorithms called Ababoost, SVM, ANN. This model worked under the big data environment using hadoop and mapreduce. The dataset for prediction were collected from hospitals in Coimbatore through blood serums and general investigations. Finally the Rheumatic Arthritis disease is predicted by using Machine Learning based Ensemble Analytic Approach (MLEAA) model with the calculated value of voting system which leads to predict the RA disease earlier.

3. Shanmugam Sundaramurthy et al. have created the ensemble classifiers to predict the Rhematoid Arthritis disease. The dataset were collected from Sakthi Rheumatology center that holds 1000 patient profiles (750-RA affected and 250 non-affected profiles). SVM, Ada-boosting, and random subspace, were used in this work. Data classification is done with 10-fold cross-validation and evaluation is done with metrics like Accuracy, Precision, and ROC. The values of these metrics were compared with baseline algorithms and various ensemble classifiers. Optimality of these algorithms and ensemble provides improvement in RA disease Prediction.

4. Sandeep Kaur and Dr. Sheetal Kalra introduced the algorithm called Hybrid K-Means and also proposed the Support Vector Machine algorithm to predict any

Department of Computer Science, A.V.C. College (Autonomous), Mayiladuthurai – 609 305. PROCEEDINGS 107

disease. Hybird K-Means is used to select the initial centroids, number of clusters and to enhance the K-Means algorithm. It also used for dataset dimension reduction. The output of the Hybrid K-Means is given as input to Support vector machine for classification. These process were done under MATLAB and the authors analysed that the accuracy is increased and execution time is reduced.

5. Jihyung Yoo et al. predicted the disease call Rheumatoid Arthritis using Machine Learning. rheumatoid factor, anti-CCP, SJC, and ESR all these four factors from blood test were taken to predict whether the RA disease is positive or negative in patients. K-Means algorithm is used to find the threshold values of four factors periodically. Authors stated the earlier diagnosis and prediction of RA disease will improve the patients quality of life.

III. PROPOSED MODEL

In this paper, enhanced K-means algorithm and Support Vector Machine model is proposed for disease prediction to improve the efficiency and accuracy for prediction. The enhanced K-means algorithm is applied for dimensionality reduction to remove outliers and noisy data. The initial centroids are randomly selected in case of simple K-means algorithm but it is not so in proposed algorithm. The enhanced K-means algorithm is follows.





Algorithm:

Input:

D (d1, d2, d3.....dn): specify a data set containing n objects.

Output: A set of k clusters are generated. Method: *Step1: Read the authenticated dataset.*

Step2: Plot the silhouette values of data points in the dataset on 2 D plane and identify the number of clusters based on average silhouette value. The number of clusters showing highest silhouette value is chosen as a value of k for that data.

Step3: Partition the dataset into k equal parts.

Step 4: The arithmetic mean of each part is taken as the centroid point.

Step 5: Compute the Euclidean distance of each datapoint di to all the centroids as edist (di, cj)

Step 6: For each di, examine the closest centroid and assign di to that centroid.

Step 7: Set Near_edist[i] = edist(di, cj) //cj: nearest centroid.

Step 8: For each cluster j, recalculate the centroids.

Step 9: Repeat 10. For each data-point di

Step 10.1: Compute its distance from the new centroid of the present nearest cluster.

Step 10.2: If this distance is less than or equal to the previous distance, the data-point stays in that cluster, *Else*

10.2.1: Compute edist (di, cj) from all cluster centroids; End for.

10.2.2: Assign the data-point di to the cluster with the nearest Centroid.

10.2.3: Set Near_edist[i] = edist(di, cj); End for loop.

Step 11: Take best average sum of all Euclidean distances and obtain the final output.

Step 12: Train the SVM classifier using reduced dataset.

Step 13: Classify the new data using SVM classifier.

The proposed work is to select the initial centroids by partitioning the data into k equal parts and then the arithmetic mean of each part is taken as the centroid point. The efficiency and accuracy of enhanced Kmeans algorithm is more than simple K-means. It is unsupervised learning algorithm that is used to solve the sound known clustering problems by partitioning data points into k clusters where each data item belongs to the cluster with its nearest mean. It is relatively efficient and provides best results for distinct datasets. In Kmeans clustering algorithm, there are always k cluster and each cluster always contains at least one item.

	SJ	Rheumatoid	Anti-	ES
ID	С	Factor	CCP	R
KB-R-11-0	2	6.2	31.2	52

Department of Computer Science, A.V.C. College (Autonomous), Mayiladuthurai – 609 305. PROCEEDINGS

108

KB-R-11-1	2	159.7	-1	65
KB P 11 2	7	129.7	20.7	26
KD-K-11-2	/	129.7	29.1	20
KB-R-11-3	3	88.6	150.9	40
KB-R-11-4	0	8	225.7	16
KB-R-11-5	0	72.7	50.9	13
KB-R-11-6	12	3	17.9	26
KB-R-11-7	0	30.4	6.5	64
KB-R-11-8	0	70.6	54.3	39
KB-R-11-9	5	8.4	135.4	120
KB-R-11-				
10	6	48.3	72.5	5
KB-R-11-	4	49.9	15 1	50
	4	40.0	13.1	39
кв-к-11- 12	3	30.4	123.7	8
KB-R-11-	_			-
13	1	6.6	1	30
KB-R-11-				
14	5	6.6	121.7	91
KB-R-11-				
15	4	4.1	81.9	28





Fig 5. Data Cluster



Fig 6. Classification using SVM



Fig 7 Accuracy Graph



Fig 8 Time Graph

IV. CONCLUSION

An enhanced K-means and Support Vector Machine algorithm for disease prediction is proposed in this paper to improve the efficiency and accuracy for prediction. For the prediction process initially dataset is clustered and after finding the centroids using enhanced K-means and given as input to the support vector machine for the classification. The initial centroids are randomly selected in case of simple Kmeans algorithm but it is not so in proposed algorithm. The proposed work is to select the initial centroids by partitioning the data into k equal parts. The final result shows that the efficiency achieved by proposed algorithm is better than simple K-means algorithm. The K-means achieved the accuracy of 82% and the hybrid algorithm achieved the accuracy of 92% on the same dataset. RA diagnosis can be difficult at its disease onset on account of overlapping symptoms with other arthritis, so early recognition and diagnosis of RA permit the better management of patients.

REFERENCES

1. Maleeha Imtiaz a, Syed Afaq Ali Shah b,* , Zia ur Rehman, May(2022). A review of arthritis diagnosis techniques in artificial intelligence era: Current trends and research challenges.

2. Yubo Shao1,2,3,4, Hong Zhang1,2,4, Qi Shi1,2,4, Yongjun Wang1,2,4* and Qianqian Liang1,2,4*(2023). Clinical prediction

models of rheumatoid arthritis and its complications: focus on cardiovascular disease and interstitial lung disease, 25:159

3. Linlu Bai1 , Yuan Zhang2 , PanWang2 , Xiaojun Zhu1 , Jing-Wei Xiong1* & Liyan Cui2 (2022). Improved diagnosis of rheumatoid arthritis using an artificial neural network, 12:9810.

4. S. Shanmugam and J. Preethi Design of Rheumatoid Arthritis Predictor Model Using Machine Learning Algorithms,(January 2018). Chapter in SpringerBriefs in Applied Sciences and Technology.

5. Shanmugam Sundaramurthy, Dr. Saravanabhavan C, Dr. Pravin Kshirsagar (November 2022).

Prediction and Classification of Rheumatoid Arthritis using Ensemble Machine Learning Approaches.

6.Sandeep Kaur, Dr. Sheetal Kalra(2016)IEEE. Disease Prediction using Hybrid K-means and Support Vector Machine.

7. Jihyung Yoo *, Mi Kyoung Lim *, Chunhwa Ihm**, Eun Soo Choi***,Min Soo Kang(2017).A Study on Prediction of Rheumatoid Arthritis Using Machine Learning, Volume 12, pp. 9858-9862.

8. Jenny Ann Verghese1, D.Pamela2*, Prawin Angel Michael3, R.Meenal* (2021). Rheumatoid arthritis detection using image processing, conference series 1937 (2021) 012037.

9.Huang, C. L., Liao, H. C., & Chen, M. C. (2008). Prediction model building and feature selection with Support Vector Machines in breast cancer diagnosis. Expert Systems with Applications, 34(1), 578-587.

10.Shah, S., & Singh, M. (2012, May). Comparison of a time efficient modified K-mean algorithm with K-mean and K-medoid algorithm. In Communication Systems and Network Technologies (CSNT), 2012 International Conference on (pp. 435-437). IEEE.

11. Kim, K.J.; Tagkopoulos, I. Application of machine learning in rheumatic disease research. Korean J. Intern. Med. 2019, 34, 708–722.

A Comprehensive Analysis of Novel Intrusion Detection Systems and Intrusion Prevention Systems for Blockchain Technology

C. Ananth^{1,} S. Sathiyarani², and Dr.N. Mohananthini³

¹AssistantProfessor / Programmer, ²Research Scholar, ³AssociateProfessor ^{1, 2} Department of Computer and Information Science, Annamalai University, Annamalainagar, India. ¹ananth.prog@gmail.com, ²sathiyaranilect@gmail.com

³ Department of Electrical and Electronics Engineering, MuthayammalEngineering College, Rasipuram, India. ³mohananthini@yahoo.co.in

ABSTRACT -Nowadays, the Blockchain Technology (BC) environments are Growing and Rising in popularity The number of devices connected to the Decentralized network continues to rise. Blockchain is an interconnected network of numerous devices in which data is gathered from the environment by transferring over the internet without human support and intervention. Blockchain technology enhances internet-based interaction with real-world applications, incorporating P2P networks, Decentralized, Smart contracts, Digital ledgers, and communities to create a smarter environment. Blockchain appliances serve various purposes in various environments, including healthcare, education, military, agriculture, and commerce. They hold great promise for enhancing social and corporate life, but are vulnerable to attacks due to resource limitations and network nature, making them a soft target. Blockchain security is ensured through various technologies like the Intrusion Detection System (IDS) and Intrusion Prevention System (IPS), which protect it from various attacks, ensuring privacy and reliability. This paper examines various **IDS/IPS** proposals for Blockchain technology between 2019 and 2024, identifying their strengths, shortcomings, and challenges. It identifies areas for improvement and outlines the research direction for future researchers, paving the way for new avenues of study.

Keywords: Blockchain, Intrusion Detection System (IDS), Intrusion Prevention System (IPS).

I. INTRODUCTION

Blockchain, a technology developed by Nakamoto in 2008, is the foundation of Bitcoin's infrastructure and has been increasingly used in various fields, particularly security. It is present in traditional networks, IoT, and cloud computing. Blockchain's unique feature is its ability to function in decentralized and distributed environments. eliminating the need for a trusted third party to manage the network. Cryptocurrency networks have also adopted blockchain technology, using it as the foundation for secure financial transactions within networks. A blockchain is a digital data structure that contains a continuously expanding log of transactions and their chronological order. It is a shared and distributed database that contains digital transactions, data records, and executables. Transactions are aggregated into blocks, which are time-stamped and cryptographically linked to previous blocks. This chain of records determines the sequencing order of events, or the 'blockchain'. The term is also used to describe digital consensus architectures, algorithms, domains of applications built on such or architectures. [1]

Blockchain facilitates transactions by encrypting information with public and private keys and verifying them through peer-to-peer networks, ensuring sufficient balance for successful completion.

Blockchain technology offers decentralization, autonomy, integrity, immutability, verification, faulttolerance, anonymity, auditability, andtransparency in a trustless environment, attracting significant academic and industrial attention in recent years due to its advanced features [8].

In a blockchain network, multiple transactions are verified and added to a mempool, forming a block. To prevent disruption, nodes use a consensus algorithm to ensure that each new block is the only version of the truth agreed upon by all nodes. Miners, selected to add a block, receive rewards and a hash code for the block, ensuring secure attachment to the blockchain. The process of adding a new block to the

Department of Computer Science, A.V.C. College (Autonomous), Mayiladuthurai – 609 305. PROCEEDINGS 111

blockchain involves obtaining its hash value and authentication. This process ensures that blocks are cryptographically linked, and the transaction is completed, with the transaction details permanently stored in the blockchain for easy retrieval and confirmation [2].

Digital currencies like Bitcoin, Litecoin, Ethereum, and Ripple are integrated into resilient ecosystems, but they face challenges in detecting attacks due to their intricate infrastructure. Conventional intrusion detection systems (IDSs) cannot detect blockchain- related attacks, prompting researchers to use blockchain technology to improve IDSs and enhance attack detection [5].

The field of Intrusion Detection Systems (IDS) and Intrusion Prevention Systems (IPS) began in 1986 with an academic paper. Over the past 20 years, significant improvements have been made in these products. In the 1990s, firewalls were effective but lacked "deep packet inspection" capabilities. The emergence of new threats like SQL injections and cross-site scripting (XSS) attacks in the early 2000s highlighted the limitations of firewalls, leading to the adoption of Intrusion Detection Systems (IDS) as a security best practice. The adoption of IPS began to grow in 2005, leading to increased vendor support. The competitive landscape changed as more vendors entered the IPS market, focusing on performance concerns. The timeline presented in the narrative shifts to 2006, indicating a transition to discussing

Table 1 , Complementative analysis of 1D5s and 1

the history and evolution of IDS/IPS beyond 2005. From 2006 to 2010, faster-combined intrusion detection and prevention systems were adopted, followed by next-generation intrusion prevention systems from 2011 to 2015, and now, next-generation firewalls. [6] In 2020 and beyond, Next-Gen SIEM leverages cloud and big data platforms to enhance scalability, productivity, and accuracy.[7] The literature presents two main solutions for detecting or preventing attacks: an Intrusion Detection System (IDS) and an Intrusion Prevention System (IPS). IDS is a precautionary measure, allowing the system to raise an alarm in case of intrusion, while IPS is a punitive measure, requiring action in case of intrusion. However, false positives can block legitimate users [9]. Table 1 represents a detailed comparative analysis of IDS and IPS systems. [10] The remaining parts of this paper are organized as follows; Section I introduces the background of blockchain technology and its features. Section II introduces the background of intrusion detection, collaborative IDSs, and the challenges regarding data and trust management. Section III, this text provides an overview of intrusion prevention strategies for data and trust management challenges. Section IV discusses how blockchain technology can be applied to solve the challenges in intrusion detection. We then discuss some open issues point out future directions in Section V, and conclude our work in Section.

Key Contribution	Intrusion Detection System (IDS)	Intrusion Prevention System (IPS)
	IDS is a monitoring tool that compares	IPS is a control-based solution that accepts or
	network packets against a threat signature	rejects network packets based on
Scope	database or machine learning baseline,	predetermined rulesets, enabling it to function

	designed for detection and surveillance,	as an IDS but not vice versa.	
	taking minimal action when a threat is		
	detected.		
Location and Range	time across an enterprise network, scanning packets for compromise indicators and flagging detected threats or anomalies. If a security policy violation, like a port scanner, ransomware, or malware, is detected, IDS alerts human security personnel for further action.	intersection of internal networks and the internet. It detects threats and stops malicious traffic flow, alerting security personnel. Unlike IDS, IPS has a limited range but can rely on IDS to expand its surveillance range. However, its range can be limited.	
Types	Host-based IDS (HIDS) is deployed at the endpoint level to protect individual devices from cyber threats by monitoring network traffic, running processes, and system logs. It only protects the host machine, providing granular visibility into its workings. Network-based IDS (NIDS) monitors the entire enterprise network, tracking all traffic to and from every device and making decisions based on packet metadata and content. NIDS has a wider viewpoint than HIDS, providing more contextual information and detecting widespread threats, but may not offer granular visibility into the devices they secure.	Host-based IPS (HIPS) is a cybersecurity software that monitors events and thwarts attacks at the device level. At the same time, Network-based IPS (NIPS) is deployed within the enterprise network infrastructure, ensuring all data is monitored and threats are thwarted before they reach their targets. Wireless IPS (WIPS) monitors radio waves for unauthorized access points and automatically takes countermeasures.	
Intervention Level Required	IDS relies on IT teams or security systems to prevent threats and can scan networks for known and unknown threats. However, it cannot independently address identified threats. If IPS isn't implemented, IDS would require a dedicated human resource to handle malicious traffic, requiring additional resources for effective threat management.	Unlike IDS, IPS is a proactive cybersecurity solution that uses a database of threat signatures or an ML-powered behavior model to detect and prevent cybersecurity violations, autonomously stopping threats before they cause damage.	
Configuration	IDS operates in inline mode and can be configured by security teams to perform specific actions upon detecting a threat. These actions include creating a log, transmitting a notification, or communicating a command to a router or firewall. Logging provides forensic information for analysis and can be used to update router, firewall, and server policies to prevent recurring events. Enterprises typically set up IDS to handle logs and alerts while fighting threats.	IPS is a network security tool placed behind the firewall, configured to operate as an end host or inline mode. It can sometimes raise false alarms due to harmless anomalies caught in its filter. It can recognize normal network traffic and detect threats without disrupting daily network operations by fine-tuning its configuration.	

II. INTRUSION DETECTION SYSTEMS FOR BLOCKCHAIN

This section reviews seven recent research studies that propose novel IDSs for Blockchain, detailing their methodology, mechanisms, and datasets. It also discusses the strengths and weaknesses of these works, highlighting the strengths of each method and the challenges they pose in detecting attacks targeting blockchain networks as shown in Table 3.



Figure 1: Architecture of IDS

IDS is a device or software that uses various detection methods to detect system attacks and notify the system's administrator. It can be a standalone device or a network system that performs local analysis. IDSs offer three crucial security services: data confidentiality, data availability, and data integrity. Data confidentiality ensures data is stored securely, availability ensures data is accessible for authorized users, and integrity ensures data consistency with other system data.[15]

Blockchain-based Intrusion Detection System (BIDS): Blockchain-based intrusion detection systems use blockchain's transparency and immutability to detect and prevent unauthorized network activities. Research has explored the integration of blockchain and intrusion detection systems, with relevant information found in research papers and articles. [16] IDS can be categorized into Network-based Intrusion Detection Systems (NIDS) and Host- based Intrusion Detection Systems (HIDS), each with its unique features and capabilities.

A. Network-based Intrusion Detection Systems (*NIDS*): NIDS is a real-time network traffic monitoring tool that can detect suspicious patterns or anomalies across multiple systems and be classified into signature-based and anomaly-based detection.

Signature-Based NIDS: This type of detection uses predefined patterns or signatures of known threats to identify malicious activity,

but may struggle with new or unknown threats.

Anomaly-Based NIDS: This type of alert system sets a baseline for network behavior, detecting deviations and alerting when detected, effective in detecting new or unseen attacks but may also generate false positives.

B. Host-based Intrusion Detection Systems (*HIDS*): HIDS monitors system logs, file changes, and user activities within a host or computer system to detect unauthorized access or malicious behavior atthe host level.

Behavior-Based HIDS: This type of analysis focuses on observing and analyzing program and process behavior on a host to detect deviations from normal behavior, potentially detecting zero- day attacks or previously unknown threats.

Log-Based HIDS: This type of security method uses predefined rules or signatures to analyze system and application logs for suspicious activity signs. [32][28]

The challenges of Intrusion Detection Systems (IDS) in blockchain technology include adaptingtraditional intrusion detection systems to the unique features of blockchain technology, such as its decentralized structure, consensus mechanisms, privacy considerations, smart contract complexities, and dynamic network topology. These challenges require innovative solutions and specialized approaches to intrusion detection within blockchain environments. Ongoing research aims to develop effective IDS mechanisms that can accommodate the distinct characteristics and security considerations of blockchain technology, enhancing the security posture of blockchain networks by effectively integrating intrusion prevention mechanisms.

A. PERFORMANCE METRICS OF INTRUSIONDEDUCTION SYSTEM

a). Accuracy (AC): Measures the overall correctness of the IDS in detecting both intrusionsand non-intrusions. [33],[34]

Accuracy =
$$\frac{TP+TN}{TP+TN+FP+FN}$$

b). Precision (P): Measures the accuracy of the predicted attacks, indicating how many predicted

attacks were attacks. [36]

Precision =
$$\frac{TP}{TP+FP}$$

c). True Negative Rate (TNR) or Specificity: Measures the proportion of correctly identified normal values. [37]

d). False Positive Rate (FPR): Indicates the ratio of normal points incorrectly identified as attacks. A high FPR suggests low IDS performance. [38]

e). True Positive Rate (TPR) or Recall: Measures the ratio of correctly predicted attacks to the actualnumber of attacks. [36]

f). False Negative Rate (FNR): Represents the ratio of actual attacks not detected by the IDS. [38]

$$FNR = FN / (TP + FN)$$

g). Confusion Matrix: A table (not provided) that summarizes the classification results, including true positives (TP), false positives (FP), true negatives (TN), and false negatives (FN). [18],[19]



Figure 2: Architecture of IPS

III. INTRUSION PREVENTION SYSTEMS FORBLOCKCHAIN

Intrusion Prevention Systems (IPS) are network security applications that monitor network or system activities for malicious activity. They identify, collect information, report, and attempt to block or stop such activity. IPS and Intrusion Detection Systems (IDS) are considered an augmentation of each other, as they operate network traffic and system activities for malicious activity. IPS records events, notifies security administrators, and produces reports. They can also respond to detected threats by stopping the attack, changing the security environment, or altering the attack's content.

An IPS analyzes real-time network traffic, compares it against known attack patterns, and blocks suspicious traffic when detected, ensuring network security. [29]

Intrusion Prevention Systems (IPS) can be classified into two main types based on their detection mechanisms.

A. Host-Based Intrusion Prevention System (HIPS): HIPS is a security solution that focuses on protecting individual computing devices from security threats. It operates directly on the host system, monitoring and analyzing activities at the operating system and application levels. The primary goal is to prevent, detect, and respond to security incidents that may compromise the host's integrity and confidentiality.

Ref.	Method	Detection Mechanism	Dataset	Attacks	Characteristics / Strengths	Limitations / Challenges
[20]	machine learning, federated learning	blockchain- based federated forest for SDN-enabled intrusion detection	CAN BUS	fuzzy, DoS, and Impersonatio n	our model outperformed other detectors such as Logistic regression (LR), and the K- nearest neighbor TL approach. The LR, DT, K-NN, and LR, recorded lower accuracies of 0.342, 0.969, and 0.568, respectively. The highest model attack detection rate	We used blockchain technology to reduce the risk of poisoning the models. The testbed shows efficient use of memory and CPU resources for the proposed system.
[21]	machine learning	HybridChain IDS integration of TEE and blockchain. Cheetah Optimization Algorithm	UNSWN B15, DAS- CIDS, NSL- KDD, CIDDS- 001	Brute Force Attack, Syn Flood Attack, Phishing Attack	is about 0.981. The proposed HybridChain-IDS framework achieves 14% better accuracy compared to existing works. 2) The proposed work reached a 15% high detection rate compared to existing works. The proposed work achieves a 15% low false alarm rate when compared to existing works. accuracy is 98.4	proposed work HybridChain-IDS achieves a low false alarm rate with a compare to existing works. The proposed HybridChain-IDS has achieved better performance in In terms of precision, detection rate, false alarm rate, and accuracy, there are several factors to consider.
[22]	machine learning, and deep learning	anomaly detection techniques.	KDD-99	Distributed denial of services (DDoS) attacks	some of the datasets like KDD-99 consist of many attacks records, i.e., 80.14% in the training set and 80.52% in the test set. Whereas some of the datasets like CDX 2009 consist of a comparatively smaller number of attack records, i.e., 0.76%.	RobustmachinelearningModel,generalizabilityofmodel,real-timeanalysis,resourceconstraints of IoTDevices, the longertraining time ofIntrusiondetectionmodel.we plan toimplementthesesolutions anddevelop a robust and

Table 2. Analysis of IDS and IPSs for Blockchain

						generalized intrusion detection model.
[23]	federated learning	MetaCIDS	CIC- IDS2017	zero-day attack	The detection model in MetaCIDS is It is effective for detecting both multi- class and zero-day attacks. with 95– 99% accuracy and detection rate in most test cases, which outperforms various ML models such as LightGBM and Tabnet.	the attacks to different types of devices might be different, which potentially leads to a lower detection accuracy of the model. Therefore, our future work will investigate this impact to improve the performance of the IDS models accordingly.
[24]	Deep learning,	Blockchain- Assisted IoT Healthcare A system using Ant Lion Optimizer with Hybrid Deep Learning (BHS- ALOHDL)	BHS- ALOHD L technique is tested on ToN- IoT and CICIDS- 2017 datasets	All CICIDS Attacks	the BHS-ALOHDL method obtains maximal outcomes with an <i>accuracy</i> of 99.55%, <i>precn</i> of 99.55%, <i>recal</i> of 99.55%, and <i>Fscore</i> of 99.55%.	The proposed work is a simulation analysis of the The BHS-ALOHDL method is tested on two benchmark datasets and the outcomes indicate the remarkable performance of the BHS-ALOHDL technique over other models
[25]	Deep learning,	Secure Federated Intrusion Detection Model. The model uses Bidirectional Long Short- Term Memory (BiLSTM)	CIDDS	Cyber - attacks	The mentioned study achieved an accuracy of 97% and a False Positive Rate (FPR) of 0.0021, while the BiLSTM in the SecFedIDM-V1 outperforms these metrics.	The authors express the intention to extend both the architecture and the web application in the future, with a focus on covering a broader range of novel network attack classes
[26]	federated learning	Blockchain- Based Federated Forest for SDN- Enabled In- Vehicle	CAN BUS DATASE T	Evasion attack, Jacobian- based Saliency Map Attack	FGSM with perturbation higher than 0.08, SVM attack, and DT attack have the most significant impact on	The proposed integration of a statistical test as an adversarial detector and subsequent augmentation of the

		Network Intrusion Detection System (BFF- IDS)	(JSMA), SVM-attack, and DT- attack	the model, degrading accuracy from over 97.5% to below 34%. JSMA non-targeted attack is found not to impact the model's confidence.	BFF-IDS with detected adversarial samples is effective against adversarial examples.
[27]	federated learning	Blockchain- based Intrusion - Detection and Prevention System (BIDPS)	FILELESS ATTACK	Without the aid of Sysmon, BIDPS achieves a success rate of 76.7% for fileless attacks. With the integration of Sysmon, the success rate for fileless attacks improves to 95.2%	the integration of a memory detection toolkit with BIDPS is proposed. This integration aims to improve the system's capacity to detect and prevent APT attacks, particularly those involving memory-based techniques.

HIPS uses techniques like signature-based detection, behavior analysis, and real-time response mechanisms to identify and mitigate potential threats. It provides an additional layer of defense by monitoring and controlling activities within the host environment. HIPS solutions protect the host system from various security threats, including malware and unauthorized access. [30]

B. A Network-Based Intrusion Prevention System (NIPS): NIPS is a security solution deployed at the network layer of the OSI model, analyzing traffic to detect and prevent malicious activities in real-time. It uses signature-based and anomaly- based detection techniques to identify known attack patterns and take immediate preventive actions, such as blocking or alerting, to mitigate the impact of attacks. NIPS aims to provide centralized defense against external threats, ensuring network security and integrity. [31]

Intrusion Prevention Systems (IPS) in blockchain face challenges in adapting to the unique features of blockchain technology, such as decentralization, consensus algorithms, privacy considerations, secure smart contract execution, encrypted transactions, and dynamic network topology. These challenges require innovative solutions and specialized approaches to intrusion prevention within blockchain technology. Continuous research and development aim to enhance the security posture of blockchain networks by effectively integrating intrusion prevention mechanisms. An IPS is crucial in preventing sophisticated cyber-attacks, so it's essential to choose one with essential features. The system offers comprehensive security measures, including IPS vulnerability protection, antimalware protection, command-and-control protection, automated security measures, uniform policy management, and automated threat intelligence.[17]

B. PERFORMANCE METRICS OF INTRUSION PREVENTION SYSTEM

a) Blocking Rate: This metric measures the percentage of malicious activities or attacks that were successfully blocked by the IPS. A higher blocking rate indicates better prevention capabilities.

> Blocking Rate= Number of Blocked Attacks / Total Number of Attacks

b) False Positive Rate (*FPR*): The False Positive Rate represents the percentage of normal or legitimate activities that were incorrectly flagged asmalicious. A lower FPR is desirable to minimize false alarms and avoid unnecessary disruptions.

> FPR= Number of False Positives / (Number of Legitimate Actions + Number of False Positives)

c) Throughput: Throughput measures the volume of network traffic that the IPS can effectively process within a specified time frame. It is typically measured in bits per second (bps) or packets per second (pps). Higher throughput values indicate better processing capacity.

Total Data Processed Throughput= Time

IV. SECURITY CHALLENGES AND ISSUES IN **BLOCKCHAIN**

Previous studies have focused on using blockchain technology to improve Intrusion Detection Systems (IDSs) efficiency in various network settings. However, most models face challenges related to blockchain technique, IDS approach, or network environment. This paper emphasizes the need for further research to enhance IDS performance using blockchain technology.

CHALLENGE/ISSUES	DESCRIPTION
Signature verification	A complex cryptographic calculation is required for the signature transaction in a
	blockchain.
Energy consumption	The miner adds new blocks to a blockchain, consuming power to validate the
	transaction's expanding value.
Slow	Blockchain blocks must be encrypted and verified before being broadcasted over
	networks.
High cost	Blockchain's high initial capital cost and energy consumption in large transaction
	volumes increase the overall cost of maintaining each transaction.
Scalability	The immutability of the blockchain prevents nodes from deleting any blocks,
	resulting in incremental growth of the blockchain size over time.
Massive false alerts:	IDS generates accurate alerts for security administrators, but false alarms pose a
	challenge due to immature signals and inaccurate profiles, affecting detection
	performance and increasing workload.
Inaccurate profile	Anomaly-based detection faces challenges due to traffic dynamics and limited
establishment:	training data, resulting in inaccurate machine learning classifications, especially
	when using labelled attack data.
Overhead traffic with	The existence of overhead packets is frequently observed during periods of high
limited handling capability:	detection systems, resulting in the significant discording of network packets and
	an increased computational load.
Security and privacy:	Blockchain applications require smart transactions with identifiable individuals,
	raising privacy and security concerns. It also attracts cyber-criminals, making it
	vulnerable to DDoS attacks.
Regulations And	Advanced technology often surpasses regulations, leading to lagging. Bitcoin's
Management	lack of standardized protocols has improved efficiency, but blockchain
	applications are expected to eventually conform to regulatory frameworks.

Table 3, Challenges/Issues of Blockchain and IDS/IPS. [3,4, 5]

V. FUTURE DIRECTION

Blockchains, an emerging technology with disruptive potential across various industries, is expected to gain credibility through proof-of-concept implementations. While blockchain can significantly impact intrusion detection, it's primary applications focus on balancing advantages and costs, demonstrating its potential for significant advancements. Blockchains are ideal for data sharing because they can record events, medical records, and transaction processing. They can improve performance in large distributed detection systems or networks by promoting trust and ensuring data privacy amongall parties involved.

Alert Exchange discusses using blockchains to secure and exchange genuine alerts from different nodes, a topic of significant interest and importance for future research due to the limited practical system implementations.

Blockchains, initially designed for cryptocurrencies, should not be seen as a solution in search of a problem. While traditional solutions are crucial, it's essential to stay updated with emerging technologies and maintain a balance on a case-bycase basis. [3]

VI. DISCUSSION

This study reviews recent studies published between 2019 and 2023 on innovative Intrusion Detection Systems (IDS) and Intrusion Prevention Systems (IPS) for Blockchain networks. The research focuses on the proposed methodology for building an IDS/IPS for Blockchain networks, the mechanisms used, and the datasets used for model assessment. It also highlights the strengths and limitations of the works, identifying remaining problems, establishing the research direction, and paving the way for future studies. Some researchers may focus on optimizing the analyzed works or combining two techniques to enhance the performance of IDS/IPS systems for Blockchain. The study provides guidelines and summarizes the main conclusions.

VII. CONCLUSION

Blockchain technology has revolutionized security measures in various industries. This study explores blockchain structure, IDSs, and IPSs, and compares existing models. However, there is limited research and a lack of standardized approaches. Enhancing blockchain security and privacy is crucial incomputer security, making it essential to implement Intrusion Detection Systems (IDS) and Intrusion Prevention Systems (IPS) to ensure network security and privacy. The authors suggest that the CIDS architecture is the most suitable for developing a comprehensive framework for IDSs and IPSs based on blockchain due to its ability to facilitate data sharing through a peer-to-peer network, a crucial feature within a blockchain structure. This paper analyzes eleven IDS/IPS systems introduced between 2019 and 2024 in IoT networks, identifying their features, strengths, weaknesses, and challenges. It identifies unresolved issues and outlines research direction, paving the way for future researchers to explore new research avenues.

REFERENCES

[1] Mattila, Juri, et al. *Industrial blockchain platforms: An exercise in use case development in the energy industry*. No. 43. ETLA WorkingPapers, 2016.

[3] Meng, Weizhi, et al. "When intrusion detection meets blockchain technology: a review." *Ieee Access* 6 (2018): 10179-10188.

[4] Chiba, Z., et al. "A deep study of novel intrusion detection systems and intrusion prevention systems for Internet of Things Networks." *Procedia Computer Science* 210 (2022): 94-103.

[5] Al-E'mari, Salam, et al. "Intrusion Detection Systems Using Blockchain Technology: A Review, Issues and Challenges." *Computer Systems Science & Engineering* 40.1 (2022).

[6] https://www.secureworks.com/blog/the- evolutionof-intrusion-detection-prevention

[7] https://seniordba.files.wordpress.com/2021/04/sie mhistory.png.

[8] Guo, Huaqun, and Xingjie Yu. "A survey on blockchain technology and its security." *Blockchain: research andapplications* 3.2 (2022): 100067.

[9] Mishra, Nivedita, and Sharnil Pandya. "Internet of Things applications, security challenges, attacks, intrusion detection, and future visions: A systematic review." *IEEE Access* 9 (2021): 59353-59377.

[10] https://www.spiceworks.com/it security/networksecurity/articles/ids-vs-ips

[11] Agrawal, Gaurav, Shivank Kumar Soni, and Chetan Agrawal. "A SURVEY ON ATTACKS AND APPROACHES OF INTRUSIONDETECTION SYSTEMS." *International Journal of Advanced Research in Computer Science* 8.8 (2017).

[12] Inedjaren, Youssef, et al. "Blockchain-based distributed management system for trust in VANET." *Vehicular Communications* 30 (2021):100350.

[13] https://www.paloaltonetworks.com/cyberpedia/ what-is-anintrusion-prevention-system-ips

[14] Ghorbani, Ali A., Wei Lu, and Mahbod Tavallaee. *Network intrusion detection and prevention: concepts and techniques.* Vol. 47. Springer Science & Business Media, 2009.

[15] Sultana, Nasrin, et al. "Survey on SDN based network intrusion detection system using machine learning approaches." *Peer-to-Peer Networkingand Applications* 12 (2019): 493-501.

[16] I. Aliyu, M. C. Feliciano, S. Van Engelenburg, D. O. Kim and C. G. Lim, "A Blockchain-Based Federated Forest for

SDN-Enabled In-Vehicle Network Intrusion Detection System," in IEEE Access, vol. 9, pp. 102593-102608, 2021, doi: 10.1109/ACCESS.2021.3094365.

[17] A. A. M. Sharadqh, H. A. M. Hatamleh, A. M. A. Alnaser, S. S. Saloum and T. A. Alawneh, "Hybrid Chain: Blockchain Enabled Framework for Bi-Level Intrusion Detection and Graph-Based Mitigation for Security Provisioning in Edge Assisted IoT Environment," in *IEEE Access*, vol. 11, pp. 27433-27449, 2023, doi: 10.1109/ACCESS.2023.3256277.

[18] N. Mishra and S. Pandya, "Internet of Things Applications, Security Challenges, Attacks, Intrusion Detection, and Future Visions: A Systematic Review," in *IEEE Access*, vol. 9, pp.59353-59377, 2021, doi: 10.1109/ACCESS.2021.3073408.

[19] V. T. Truong and L. B. Le, "MetaCIDS: Privacy-Preserving Collaborative Intrusion Detection for Metaverse based on Blockchain and Online Federated Learning," in *IEEE Open Journal of the Computer Society*, vol. 4, pp. 253-266, 2023, doi: 10.1109/OJCS.2023.3312299.

[20] H. Alamro *et al.*, "Modeling of Blockchain Assisted Intrusion Detection on IoT Healthcare System Using Ant Lion Optimizer With Hybrid Deep Learning," in *IEEE Access*, vol. 11, pp. 82199-82207, 2023, doi: 10.1109/ACCESS.2023.3299589.

[21] E. B. Mbaya *et al.*, "SecFedIDM-V1: A Secure Federated Intrusion Detection Model With Blockchain and Deep Bidirectional Long Short- Term Memory Network," in *IEEE Access*, vol. 11, pp. 116011-116025, 2023,doi:10.1109/ACCESS.2023.3325992.

[22] I. Aliyu, S. Van Engelenburg, M. B. Mu'Azu, J. Kim and C. G. Lim, "Statistical Detection of Adversarial Examples in Blockchain-Based Federated Forest In-Vehicle Network Intrusion Detection Systems," in *IEEE Access*, vol. 10, pp. 109366-109384, 2022, doi:10.1109/ACCESS.2022.3212412.

[23] Alevizos, Lampis, et al. "Blockchain-enabled intrusion detection and prevention system of APTs within zero trust architecture." *IEEE Access* 10 (2022): 89270-89288. Stallings, W. (2017). "Network Security Essentials: Applications and Standards." Pearson.

[24] https://www.geeksforgeeks.org/intrusion- preventionsystem-ips

[25] Northcutt, S., & Novak, J. (2002). "Network Intrusion Detection: An Analyst's Handbook." NewRiders.

[26] Cheswick, W. R., Bellovin, S. M., & Rubin, A.

D. (2003). "Firewalls and Internet Security: Repelling the Wily Hacker." Addison-Wesley.

[27] Bejtlich, R. (2007). "The Tao of Network Security Monitoring: Beyond Intrusion Detection." Addison-Wesley.

[28] Farnaaz, Nabila, and M. A. Jabbar. "Random forest modeling for network intrusion detection system." *Procedia Computer Science* 89 (2016): 213-217.

[29] Iervolino, Iunio, et al. "Quantitative risk analysis for the Amerigo Vespucci (Florence, Italy) airport including domino effects." *Safetyscience* 113 (2019): 472-489.

[30] Anbar, Mohammed, et al. "A machine learning approach to detect router advertisement floodingattacks in

next-generation IPv6 networks." Cognitive Computation 10 (2018): 201- 214.

[31] Elhamahmy, M. E., Hesham N. Elmahdy, and Imane A. Saroit. "A new approach for evaluating intrusion detection system." *CiiT International Journal of Artificial Intelligent Systems and Machine Learning* 2.11 (2010): 290-298.

[32] Abdullah, B., et al. "Performance evaluation of a genetic algorithm-based approach to network intrusion detection system." *International Conference on Aerospace Sciences and Aviation Technology*. Vol. 13. No. AEROSPACE SCIENCES & AVIATION TECHNOLOGY, ASAT-13, May 26–28, 2009. The Military Technical College, 2009.

[33] Gupta, Neha, Komal Srivastava, and Ashish Sharma. "Reducing false positive in intrusion detection system: a survey." *International Journal of Computer Science and Information Technologies* 7.3 (2016): 1600-1603.



A Comprehensive Review of Literature on Current and Future Research Fields **Related to Robotics**

Dr.S. Jayaprakash¹ and J.P. Keerthana²

¹Research Supervisor in Department of Computer Science and ²Research Scholar in Department of Computer Science, Edayathangudy .G.S.Pillay Arts & Science College Nagapattinam, Tamilnadu. ¹jayaprakashsoundar@gmail.com and ²jayakeerthana2094@gmail.com

Abstract— Throughout antiquity and the middle ages, the primary function of ROBOTs was amusement. The development of industrial robots saw a boom in the 20th century. Automation has reduced the need for humans to carry out dangerous and routine tasks, increasing productivity. Beyond the state of the art emerging research fields are also of particular interest. In order to accomplish this, this paper reviews recent technical scientific bibliographies using a systematic literature review methodology. It also identifies current and upcoming research fields related to robotic roles, types, and applications such as automated harvesters have proven beneficial to farmers, while advancements in assistive surgical robotics have benefited the medical field, under extreme pressure the aquatic robots are managing complex underwater environments, disaster response robots and engaged learners work together to prevent the spread of damage during emergencies and educational robots have the potential to improve learning experiences by promoting student collaboration, problem-solving and active engagement.

Keywords—agriculture robotics, healthcare robotics, underwater robotics, search & rescue robotics, educational robotics.

I. INTRODUCTION

Robots and machines with learning capabilities may have a wider range of uses in the future. Future robots would be better suited to more difficult and dynamic activities if they could learn new processes, adapt to their environment and change their behavior. Technical vision and computer learning technologies for recognition and selective spraying of weeds, the use of spot treatment systems with a micro dose of chemicals, the development of robots with a certain level of versatility, which makes it possible to perform other field works in addition to spraying [1]. Digital agriculture is seen as a key to mastering this challenge. By deploying sensors and mapping fields, farmers can begin to understand their crops at a micro scale, conserve resources, and reduce impacts on the environment [2].Indeed, significant refinements have been achieved by integrating perception, decision making, control, execution techniques. However, most agricultural robots continue to require intelligence solutions, limiting them to small-scale applications without quantity production because of their lack of integration with artificial intelligence.

[3]. When autonomous robots collaborate with humans, social skills are necessary for adequate communication and cooperation. Considering these facts, endowing autonomous and social robots with decision-making and control models is critical for appropriately fulfilling their initial goals [4]. The development of AI and robotic technologies also takes place in the context of neoliberal modes of governance adopted across different national contexts, including privatization and cuts to public funding combined with political rhetoric's that promote citizens' responsibility and control over their own lives [5].Radiation protection aims to lessen the unfortunate consequences of ionizing radiation by lessen redundant radiation exposure [6].Autonomous and social robots with decision-making and control models is critical for appropriately fulfilling their initial goals. This manuscript presents a systematic review of the evolution of decisionmaking systems and control architectures for autonomous and social robots in the last three decades [7]. A team of marine robotic agents for the purpose of cooperative marine litter detection and mapping, while also including interested citizens in the loop and allowing them to serve as operators. Two Autonomous Surface Vehicles (ASVs), a Remotely Operated Vehicle (ROV), and a Smart Buoy were deployed in a real marine environment to demonstrate the cooperative abilities of this system [8]. By contrast, underwater biomimetic robots show better stability and maneuverability in harsh marine environments [9]. search and rescue robot which can work in semi- autonomous and wireless modes and can be used in harsh physical environments of disaster regions to hold out the given tasks more effectively by the utilization of advanced and economic sensors [10]. The current state of the art in ground and aerial robots, marine and amphibious systems, and humanrobot control interfaces and assess the readiness of these technologies with respect to the needs of first responders and disaster recovery efforts [11]. Rescue robots can replace rescue workers in search and rescue missions in unknown environments and hence play an important role in disaster relief. To realize the advantages of high carrying capacity, easy control and simple structure to adapt to a complex disaster scene, a novel wheel-legged rescue robot, integrated with rescue function modules, was designed for the present work [12]. STEM learning and transferable skills were the most popular educational goals when applying robotics [13].Although the Virtual and Physical Robots (VPR) strategy is not always better than the Physical Robots(PR), it has unique advantages on facilitating students' higher-order

Department of Computer Science, A.V.C. College (Autonomous), Mayiladuthurai – 609 305.

thinking in solving complex problem, as well as reducing their cognitive load [14].Robots and automation systems that rely on data or code from a network to support their operation i.e., where not all sensing, computation and memory is integrated into a standalone system [15].

II. ROBOTICS TECHNOLOGY

Machine vision technology is used in much agricultural robotic advancement to identify crops, avoid hazards and even determine when they are ready to be harvested. However, doctors can now conduct various testing procedures and provide healthcare remotely thanks to teleoperated robots that are outfitted with cutting-edge technologies like edge deployment, haptics, 5G, and virtual reality.

Because SONAR sensors operate using sound waves and can see in low light, they are frequently used by underwater robots, where light-based imaging presents difficulties. Typically, sensors and actuators are installed in firefighting robots to enable them to recognize and react to fire-related events. Temperature, air quality, and infrared cameras are among the sensors that fire robots use. Robot on wheels made to teach kids about electronics, programming, and robotics. Thanks to kid-friendly Scratch-based software, it is simple to assemble and operate. More skilled users can build more intricate robots thanks to its Makeblock platform compatibility and electronic components built on the opensource Arduino ecosystem.



Fig. 1. Types of Current and Upcoming Robotics.

By 2030, deep learning algorithms and neural networks will be widely used in robots, improving their ability to see, make decisions, and adapt to their surroundings. One fascinating application of this integration is in self-driving automobiles.

III. CURRENT AND UPCOMING ROBOTICS

Field robots are machines designed to operate in uncontrolled environments such as air, forests, mines, underwater, and farms. These applications require close attention to engineering details as well as cutting edge robotics concepts. It also identifies current and upcoming research fields as shown in Fig.1. related to robotic roles, types, and applications.

A. Agricultural Robotics

The global population explosion has made it challenging for agricultural companies. In order to meet the dietary requirements of billions of people worldwide, more food production is required. The advancement of robotics and automation technology is not exclusive to farmers' production capabilities. In various ways, agricultural robots are increasing production yields for farmers. Creative and innovative way are being used to deploy technology, including drones, autonomous tractors, and robotic arms. With the help of robots and these machines, farmers can track weather changes, precipitation, pest infestations, soil moisture, pollution rates, fertilization, pests, disease, and many more. Farmers and food manufacturers are turning to robots and autonomous machines as a response to this need in the coming future.

a) Ecorobotix : This fully autonomous drone is powered by the sun and also has a lightweight GPS tracker. The <u>robot</u> uses its complex camera system to target and spray weeds. Seeding drones are primarily used for cover crops. Adopting the practice of planting cover crops is a soil health practice. Cover crops help to limit soil erosion and reduce pollution in water runoff. For spraying and seeding drones, the payload capacity is attention to the size of the tank and the spray width, powerful battery and lift capability. The machine can cover three hectares of land per day. The robot's upper part is covered with photovoltaic solar panels that provide a steady energy supply and other sensors that enable them to avoid obstacles.

b) *FendtXaver:* It is a technology that enables farmers to deploy a swarm of small <u>robots</u> in a field, which are then directed to fulfill a particular task. The robotic system consists of several parts. It uses also satellite-based navigation to relay their exact position, helping operators to optimize <u>plating</u> operations. Fendt's field robotic system is energy efficient due to its low weight and low-maintenance origin.

c) SmartCore : It is able to navigate fields and gather samples from specific locations autonomously. Detection algorithms and GPS assist the machine in collecting samples from specific locations. After collecting the sample, SmartCore takes it to the edge of the field to be shipped to a lab. The self cleaning hydraulic auger that this robot uses is a significant advantage in ensuring accurate samples and demonstrating ground composition.

B. Healthcare Robotics

Robots are transforming surgery methods in the medical field, accelerating the delivery of supplies and cleaning, and freeing up medical staff members to focus on patient care and interaction. Intel offers a wide range of technologies for the development of medical robots, including modular robots, autonomous mobile robots, and surgical assistance robots. Medical robots can assist with intelligent therapeutics, frequent and personalized monitoring for patients with chronic diseases, minimally invasive procedures, and social engagement for senior citizens. Nurse robots are capable of doing duties that human nurses do not tired out. Advances in robotics, robots may soon be able to remove plaque from arteries, collect tissue samples, conduct lab tests autonomously, and fight cancerous tumors. In the future, robots may also provide targeted medication, engage in conversation with patients regarding their symptoms, and handle minor medical issues.

a) Robotics in radiotherapy : Robotics was introduced to the radiotherapy field in the 1990s. The first system treated tumors

Department of Computer Science, A.V.C. College (Autonomous), Mayiladuthurai – 609 305.

precisely in a variety of locations by mounting a linear accelerator on a robotic arm that could move around the body. Since then, robotics has been further incorporated into radio surgery and radiotherapy. For example, before treatment starts, the patient is precisely positioned on robotic treatment couches. They also enable remote patient repositioning by medical professionals without requiring them to visit the patient's room.

b) Prosthetic Robots : The goal of this relatively new use of medical robots is to give their wearers functional limbs that resemble real ones. Even though there are already robotic prostheses on the market, patients still have to pay a high price for them as this technology advances. Neuromusculoskeletal prostheses, which are implanted into the patient's muscles and nerves and operate via bidirectional interfaces connected to their neuromuscular system, are an example of advancements in this field.

c) Robots with Social Skills : Social robots in a hospital setting provide cognitive support to patients, particularly the elderly and children, by encouraging social interaction and showing patients how to perform specific motor tasks. These increasingly autonomous, human-like robots can carry out their duties and engage in natural interactions with patients and medical personnel. These social robots could fill the gap in patient social interaction that hospitals across the globe are experiencing due to a scarcity of nurses.

C. Underwater Robotics

ROVs and other underwater robots are connected to a surface ship through a long power and communication cable and controlled by pilots aboard the ship. Untethered autonomous underwater vehicles (AUVs) operate via an onboard computer that has been preprogrammed. Underwater robots must be watertight, and unmanned vehicles that are utilized for underwater research must contain computers and other electronic equipment. This indicates that since the equipment is enclosed in a covering that keeps water out, it cannot be harmed by water.

Unmanned underwater vehicles (UUV) to take clear pictures and films of marine life and the ocean floor equipped with cameras, sonar, and other sensors. These robots are capable of exploring parts of the ocean that are too risky or challenging for people to visit. Diverse forms and sizes are available for underwater robots, which can be equipped with a wide range of sensors and instruments to gather copious amounts of data from deep-sea habitats. Underwater robots are being utilized more and more in deep-sea resource exploration (including cleaning up oil spills), pipeline maintenance, shipwreck investigations, surveys, oceanographic sampling, underwater archaeology, under-ice surveys, and the study of natural phenomena like hurricanes and tsunamis. Technological advancements, high-speed internet controls and autonomous navigation are now possible, providing nearly continuous service.

Underwater drones are being trained by artificial intelligence to function with little assistance from humans.

AI enables robots to use preprogrammed logic to solve issues. Restarting a device or avoiding an obstruction that wasn't in its intended route could be the easiest ways to solve the problem. Underwater vehicles that remain submerged, accomplish jobs, recharge, and emerge again without any hindrance are now undergoing testing.

a).*Eelume:* It is a Self-propelled and shaped like asnake robot that lives its entire life underwater. An ROV operating for underwater maintenance is typically operated by a crew on a vessel that must remain on site for the duration of the project. Without a surface vessel present, Eelume remains stationary on the ocean floor, prepared for subsea inspection, maintenance, and repair, or IMR. They can operate on both ends of the unit for tasks like illumination, inspection, and manipulation because they are flexible and self-propelled, allowing them to squeeze into tight spaces. Additionally, the system is modular, with parts to support long-range travel as well as a range of underwater tasks.

b) Blueye: It is a compact, lightweight underwater drone designed for quick deployment for sub aquatic asset inspection. The drone can be deployed from nearly anyplace, weighs roughly 20 pounds, and has a 5-hour battery life. Ship inspections can be finished without waiting for divers or specialized ROV operators, which can be an expensive and time-consuming task, by sending a Blueye into the water. For quick outcomes, the drone provides operators and decision makers with live video streaming. The seawater intakes of the Kristin semi-submersible platform are inspected by Equinor using Blueve to look for vegetation growth such as mussels and algae that could obstruct the intakes for the platform's fire fighting system. The crew can instantly know the condition of the intakes underwater thanks to Blueye's images. It can also be applied to marine research, aquaculture, and all kinds of underwater inspections, according to Blueye Robotics. Three guest ports on the most recent model allow for the connection of additional external devices, such as sonar, camera and sensors.

c) Undersea Shuttle: The subsea shuttle concept is an underwater shuttle measuring 135 meters that can carry out various tasks, including returning CO2 to reservoirs. It uses a hybrid of drone and submarine technology to move cargo underwater without interference from other ships or the weather. It could transport CO2 to reservoirs that aren't served by pipelines because it has a 300 km range and produces no emissions. In a similar vein, the PowerX electricity transport vehicle that we discussed in our article on offshore power generation could deliver renewable energy to places without the cable infrastructure to receive it. Like a shuttle tanker but underwater and without a crew, it could also carry materials or oil to places without pipelines.

D. Search & Rescue Robotics

Rescue operations through reconnaissance, mapping, debris removal, supply delivery, medical care, or casualty evacuation. Search and rescue teams frequently face extreme danger when entering affected areas following a natural disaster, such as an earthquake or flood, in order to locate victims and survivors. Enabling durable robots to enter areas too dangerous for people and rescue dogs is the aim of developing these machines for use in rescue operations In order to focus all efforts on areas where victims are known to be, robots can be used to assess the situation, locate people who may be trapped, and communicate the location back to the rescue teams. Additionally, food and medical supplies are being carried by robots, which will concentrate resources where they are most needed Durability and usability of robots, how to build robots that are easily transportable, that perform well in all weather conditions and have long-lasting power, and that can navigate themselves and have sufficient sensors to identify victims are the primary research questions in the field of mobile robotics for search and rescue missions.

a) *Smart Remote Controlled Emergency:* These standard rescue apparatus used by private rescue organizations and government emergency response teams in locations such as lakes, rivers, reservoirs and the ocean is the smart remote-controlled emergency water rescue robot. Additionally, it's the gear that satisfies the military's requirement for systematic function customization during particular rescue missions thanks to its marine standard.

b) Firefighting Robot: The Firefighting Robot is a small, lightweight emergency response robot that helps firefighter's battle tall buildings, particularly in hazardous situations where people shouldn't be inside. With the development of technology and the growing demand for safer and more effective firefighting techniques, the use of firefighting robots is growing in popularity. These robots can navigate through smoke-filled environments and identify hotspots that might be invisible to the human eye thanks to their advanced features, which include sensors, cameras, and thermal imaging technology. They are perfect for use in firefighting scenarios because they are made to withstand harsh conditions and extreme temperatures. This robot's ability to recognize various forms of fire and avoid obstructions is quite advantageous. Due to their struggles, many firefighters are unable to carry out their duties, which increase the number of fatalities on missions and in incident-related situations.

c) Land Rescue Robots: At the moment robots aren't nearly as nimble as humans when it comes to traversing uneven and unpredictable ground (like that you would expect to find after an earthquake or flood). While there are numerous companies trying to develop biped robots that can walk, run, and scramble over uneven ground as well as humans, the technology is not yet there. As the video from the 2015 DARPA world challenge (which focused on developing rescue robots) makes it clear, while some big strives have been made, for many tasks robots are nowhere near as fast or dextrous as humans. These robots might have some limited uses, (such as in a place too radioactive for a human to survive), but when it comes to search and rescue on the ground, humans are likely to be used as the first search and rescue efforts. On land, one major advantage robots currently bring to the table is their size. Small robots can be built to fit into places humans can't. Robots can travel through small tunnels underground, pass through small gaps, or fit into tiny pockets of air beneath fallen buildings. Currently, the small normally remote controlled robots used in this capacity tend to be the same robots used for inspections in small spaces at industrial facilities, or a slightly modified version of them. The company Inuktun is a major producer of these robots and their equipment was used after 9/11 at the twin towers site and after Hurricane Katrina. These small robots with their tank-like treads can

easily travel into small and dangerous spaces. This video demonstrates how one was used after Katrina.

E. Educational Robotics

Educational robots can enhance learning outcomes in a number of ways. They can promote student participation, problem-solving and teamwork as instruments for active learning. Including robotics in the classroom can help kids develop their critical thinking and creative skills. Additionally, robots can serve as a scaffold for the social skill-building of children, particularly shy or special needs children.

STEAM (Science, Technology, Engineering, Arts, and Mathematics) resulting in a multidisciplinary learning process through the creation of actual projects that are grounded in actual circumstances. In this sense, STEAM education promotes practice, experimentation, and working with these subjects with kids and teenagers. This approach to learning involves educational centers developing real- world projects that address real-world issues in order to help kids and teens to integrate various cross-disciplinary curriculum ideas.

In the context of early childhood education, educational robotics provides students with all the resources they require to build and program a robot that can easily perform a range of tasks. More expensive and advanced robots are also available for use in secondary and tertiary education. In any case, the students' age is always taken into account when determining the level of discipline complexity.

a) OWI 535: It is a robotic arm suitable for young people aged 13 or over. It can lift objects weighing up to 100 grams and has a wide variety of movements on which students can program customizations. This robot is also recommended for vocational training cycles.

b) NAO : NAO is one of the world's most popular educational robots. It is a 58-cm high humanoid robot that is constantly evolving. As well as two cameras and four microphones, it has a great many sensors that allow it to interact with the environment in a similar way to humans. NAO can observe, listen, have conversations and teach any subject. Its facilities and several programming levels enable its integration into the learning process of students from age 5 up to university level.

c) Robo Wunderkind: It consists in a set of blocks that the children can connect as they wish to build their own robot. Each block has a function identified with a color (camera microphone, motion sensors) and after building their robot, the children can use an app to program it to react to certain noises, avoid obstacles or play music when someone approaches, among other functions.

IV. ANALYSIS OF ROBOTICS

A.Agricultural Robotics

• Lower labor costs: Because of the high cost of crops, managing agricultural fields is very expensive. One major benefit that robots offer is a decrease in labor costs. The costs associated with using human laborers can be reduced by using agribots, as robots are more productive at what they

PROCEEDINGS

do. When an agribot is used to plant seeds on a oneacre farm, it can finish the task three times faster than ten laborers, increasing profit margins.

- Elevated yields: Buyers are assured better and higher crop yields by agricultural robots. Every seed that is planted needs to develop into a fully developed plant in the right environment. As a result of the farming process being optimized for maximum productivity, the harvest is increased. Here, robots increase crop yield while reducing losses.
- **Reduced mistakes:** Agribots operate at faster speeds and with tighter tolerances despite not having a vacation period. They complete tasks at higher speeds and with greater quality while making fewer mistakes. Additionally, they can easily maneuver around ponds, trees, rocks, and other obstacles without causing harm to crops. When sowing seeds, the robots are able to measure distance. They are able to calculate how much water and pesticides need to be sprayed on the designated application areas. All of these with very few mistakes result in higher yields.
- **Reduced use of pesticides**: Metered doses of pesticides are applied to the impacted crops by robots. Since the pesticides are only applied to the farm's affected areas, less waste is produced overall. They cut the amount of pesticides used on farms by about 80%. When the chemicals are administered manually, the agribots shield people from the damaging effects as well.

B. Healthcare Robotics

- Superior Medical Care: Intelligent therapeutics, frequent and tailored monitoring for patients with chronic diseases, minimally invasive procedures, and social engagement for senior patients are all made possible by medical robots. Furthermore, nurses and other caregivers can provide patients with more empathy and human interaction as robots reduce workloads, which can improve patients' long-term wellbeing.
- Ideal Medical Procedures: Autonomous mobile robots, also known as AMRs, reduce the physical strain of repetitive tasks on human labor and ensure more dependable processes. These robots can help guarantee that supplies, equipment, and medication are available when needed, thereby mitigating staffing shortages and related issues. They do this by keeping inventory and placing orders on time. washing and sanitizing AMRs enable hospitals to quickly clean and ready rooms for incoming patients, freeing up staff members to focus on patient-centered and value-driven work.
- Safety in the Workplace: Healthcare workers no longer have to worry about certain workplace hazards because nurse robots can perform heavy

lifting tasks like moving hospital beds and lifting patients. In hospitals where there is a risk of pathogen exposure, AMRs are used to transport supplies and linens in order to help protect healthcare personnel. Robotic cleaning and disinfection reduces pathogen exposure and hospital-acquired infections.

• Streamlined Clinical Procedures: Autonomous mobile robots, or AMRs, ensure more consistent processes while alleviating the physical demands on human workers. By monitoring inventory and placing prompt orders to make sure supplies, equipment, and medication are available where they are needed, these robots can assist in addressing staffing shortages and other challenges. Cleaning and disinfection AMRs free up staff time to concentrate on patient-centered, value-driven work by enabling hospital rooms to be swiftly cleaned and prepared for new patients.

C. Underwater Robotics

- **Investigating Marine Life:** Underwater robots are used by biologists and oceanographers to gather information on salinity, temperature, and ocean currents as well as to investigate animals and their habitats.
- Finding Lost Vehicles: Underwater robots can also be used to search for shipwrecks and plane crashes. Researchers can study wrecked planes, shipwrecks, and other debris without disturbing the sites by using robots to explore them. Additionally, they have the ability to map and survey the region, giving important information for future research projects or recovery missions.
- Upgrading Offshore Infrastructure: The oil and gas sector is one of the main applications for underwater robots. They are able to maintain and inspect offshore rigs and pipelines, making these facilities safer and more effective to operate.
- Looking for Trash: AUVs with artificial intelligence (AI) image recognition software can search the ocean floor for rubber, metal and plastic waste. Scientists can organize cleanup efforts to concentrate on the area where the trash is located and note its location once the robot relays the information.

D. Search & Rescue Robotics

• Enhanced protection: Rescue robots are capable of assessing hazardous situations and providing medical assistance without putting themselves in danger. By putting robots in danger instead of people, rescue operations can be conducted more safely and effectively. Some of these features

Department of Computer Science, A.V.C. College (Autonomous), Mayiladuthurai – 609 305.

International Conference on "Computational Intelligence and its applications" (ICCIA-2024)

E. Educational Robotics

ISBN: 978-81-967420-1-0

include sophisticated sensors, cameras, and other technologies that can recognize and report hazardous materials in the environment. Robots can now operate in difficult and even remote locations with the help of this technology, providing rescuers with information that will help them make more informed decisions about what to do next.

- **Enhanced precision :** People in need of assistance can be precisely located by robots, even in challenging environments like collapsed buildings. Because of their accuracy and capacity to function in hazardous environments, they have grown in significance during emergencies, such as natural disasters. Additionally, because they can be controlled remotely and work for extended periods of time without requiring a break, robots offer the advantage of not endangering human lives. In addition, these robots are built to react fast to circumstances, choosing the best course of action and giving disaster relief teams real-time updates. Rescue teams may be able to save more lives and lessen the damage caused by natural disasters with this improved accuracy and efficiency.
- Quickly Reach Distant Areas: The role of rescue robots in supporting emergency response operations is growing. These robots can swiftly and effectively deploy resources in life-threatening situations by traveling to remote locations, including crime scenes, war zones, and damaged areas. Rescue robots are an essential tool in the fight to save lives and property because they make it possible to swiftly enter dangerous areas. In addition, rescue robots can conduct survivor searches and reconnaissance missions. They can also be used to provide those in need with necessities. Rescue robots can mean the difference between life and death in a disaster scenario and have the potential to completely change how emergency response teams function.
- Simpler mapping: The mapping and assessment of hazardous environments by first responders has been transformed by rescue robots. Responders can swiftly and safely map out an unpredictable area with the use of autonomous robots, which enables them to quickly assess the situation and plan a course of action. Rescue robots can also be programmed to carry out a wide range of functions, including creating intricate threedimensional maps of the surroundings, recognizing people and objects, and even supporting the actual rescue effort. This lets responders act more swiftly and effectively to save lives by cutting down on the time and resources needed to map out a challenging environment.

- Fosters Logical Thinking and Problem Solving: Educational robotics involves young learners in creating, programming, setting up, and controlling robots. This activity requires the use of logical thinking and problem-solving skills, in other words, the development of computational thinking. Children need to figure out how to program a robot to perform specific tasks. As they face different challenges and overcome obstacles, children learn to analyze the situations they encounter, identify the best solution in each case, and adapt to the needs of the moment. This adaptability and reasoning ability are part of a set of fundamental skills that can be applied not only in school but also in various aspects of life.
- Promotes Innovation and **Creativity:** Programming and robotics are frequently perceived as "closed" fields that restrict kids' creativity. They actually have the opposite effect, though, encouraging inventiveness and creativity from the very beginning of schooling. Children can experiment with different ideas and explore their creativity by designing and customizing their robots. Educational robotics encourages students to explore all of their ideas and approach problems from various angles in order to accomplish a particular goal by teaching them that there is more than one way to accomplish a task and that there is never a single right answer.
- Facilitates Practical STEM Learning: Educational robotics is an effective way to introduce scientific and technological concepts in a practical and tangible way. Students can acquire knowledge of mechanics, electronics, coding and other STEM disciplines while interacting and having fun with robots. It is a hands-on experience that helps improve their understanding of theoretical concepts, presents them in an engaging and motivating manner for young learners, and sparks an interest in STEM fields, steering them away from the notion that these are overly complicated subjects.
- Enhances Collaboration Skills: Robotics projects can foster skills such as teamwork. Children learn to communicate, collaborate, and share ideas as they work together to build and program their robots. This collaboration can be applied in the workplace, where the ability to work as a team is essential. Additionally, this way of working also teaches them to appreciate and harness all their individual strengths to achieve common goals.
- Boosts Confidence in Digital Skills: In a digitized world, having technology-related skills is essential. Educational robotics helps young learners develop confidence in their digital skills and their ability to interact with technology practically, turning them into creators rather than just consumers of products. As children gain technological experience in creating and programming robots,

International Conference on "Computational Intelligence and its applications" (ICCIA-2024)

they become more comfortable working with electronic devices and software, making them increasingly capable of facing any technological challenge.

V. FUTURE ENHANCEMENT

Robotics will boost output and economic expansion while giving many people new work opportunities around the world. Nonetheless, there are worries that by 2030, thirty percent of all jobs could be automated, resulting in massive job losses. Robots will, however, take over more difficult, repetitive manual labor jobs as a result of their everincreasing precision, improving transportation and healthcare while freeing up people to pursue personal growth.

ACKNOWLEDGMENT

The authors would also like to thank the anonymous referees for their valuable comments and helpful suggestions. The work is supported by the Edayathangudy G.S.Pillay College Arts and Science College, Nagapattinam and really sincere thanks to my guide Dr.S.Jayaprakash M.Sc., M.Phil., Ph.D., research supervisor and associate professor in department of computer science for his timely suggestion, valuable guidance and sustained interest at every stage of this article.

REFERENCES

- [1] World trends in the creation of robots for spraying crops , S Bykov E3S Web of Conferences, 2023.
- [2] Digitalisation in agriculture–From the perspective of a global agricultural machinery producer, J Horváth, B Schmitz -Hungarian Agricultural Engineering, 2019.
- [3] Recent advancements in agriculture robots: Benefits and challenges, C Cheng, J Fu, H Su, L Ren Machines, 2023.
- [4] A Systematic Review Of Research Into How Robotic Technology Can Help Older People, M Shishehgar, D Kerr, J Blake – Smart Health, 2018 Elsevier.
- [5] Beyond the hype: acceceptable future for AI and robotic tecnologies in healthcare G De Togni, S Erikainen, S Chan. - Ai & Society, 2023 – Springer.
- [6] A Short Review on the Imaging Technology in Radiation Therapy SMN Raja, SA Othman, RM Roslan - e-Jurnal Penyelidikan dan, 2023.
- [7] A systematic literature review of decision-making and control systems for autonomous and social robots M Maroto-Gomez, F Alonso-Martín, M Malfaz of Social Robotics, 2023 – Springer.
- [8] Cooperative marine litter detection and environmental monitoring using heterogeneous robotic agents A Babić, F Ferreira, N Kapetanović OCEANS 2023 - ieeexplore.ieee.org.
- [9] Underwater Undulating Propulsion Biomimetic Robots: A Review, G Li, G Liu, D Leng, X Fang, G Li, W Wang, Biomimetics, 2023.
- [10] Internet Rescue Robots For Disaster Management, KMB Punith, S Sumanth, MA Savadatti , International Journal of Wireless and Microwave Technologies, 2021.
- [11] The Current State and Future Outlook of Rescue Robotics, J Delmerico, S Mintchev, A Giusti, B Gromov, K Melo, T Horvat, C Cadena, M Hutter, Journal of Field Robotics, 2019, Wiley Online Library.
- [12] Design And Analysis For A Multifunctional Rescue Robot With Four-Bar Wheel-Legged Structure M Ning, Z Ma, H Chen, J Cao, C Zhu, Y Liu, Y Wang, Advances in Mechanical Engineering, 2018, journals sagepub.com.
- [13] Trends and research foci of robotics-based STEM education: a systematic review from diverse angles based on the technology-based learning model D Darmawansah, GJ Hwang, Education, 2023 stemeducationjournal.springer.

- [14] An exploration of combining virtual and physical robots in robotics education, B Zhong, J Zheng, Z Zhan - Interactive Learning Environments, 2023 - Taylor & Francis.
- [15] A Systematic literature Review of Decision-Making and control System for Autonomous and Social Robots, Macros Maroto-Gemoze , Fernando Alonso-Martin ,Maria Malfazl ,Alvaro Castro-Gnozale, Miguel Angel Salichs, International journal of Social Robots,2023,Springer.

An Secure Commercial Enterprise Transaction System Using Alk (Alpha Keys Authentication) Technique

N.RUBA

Research Scholar, Department of Computer Science, KhadirMohideen CollegeAdirampattinam, Thanjavur Affiliated to Bharathidasan University Thiruchirappalli, Tamilnadu, India rubaanand17@gmail.com

Dr.A.SHAIK ABDUL KHADIR

Research Supervisor, Head & Associate Professor, Deptof Computer Science, KhadirMohideen College, Adirampattinam, Thanjavur Affiliated to Bharathidasan University, Thiruchirappalli, Tamilnadu, India hiqmath4u@gmail.com

Abstract-In today's business world, it is imperative to develop a new approach to secure transactions between customers and entrepreneurs. Therefore, the current use of a specific identification number (Leg) to verify the validity of the customer's identity in all online sales systems, which is vulnerable to unauthorized access and illegal money extortion from online sales, requires other reliable transmission paths of the authentication stoner. Introducing the tri-league authentication model with three layers of word, point and nonce authentication (OTP). The smoker's identity is verified by word, period and OTP. The result is a three-factor authentication model for online sales.Initial keys and some special character keys were entered on the numeric keypad for authentication. The sale involved an advanced collection of security points during the paperless sale.

Key words-Online Transaction, Authentication, Password, Fingerprint, One-Time Password, Security.

I. INTRODUCTION

The advent of the internet has now created a new dimension to service provisioning in the context of buying and selling. In this regard, online shopping, e-commerce, e-banking, and the newest era of cloud computing has brought about a dramatic change in business transactions. However, owing to current market competition by organizations and private enterprises, customers are becoming ever more demanding both in terms of the quality of goods and services that they receive while seeking for flexibility in business transactions. In particular, their interaction with these organization/administration in ensuring a good transaction platform will be highly dependent on how well designed such a platform is integrated. Since the internet plays a major role in e-commerce platforms, it is important to solve security vulnerability problems in such online platforms. According to [1], online retailing is a form of electronic commerce which allows consumers to directly buy goods or services from a seller over the Internet using a web browser. Alternative names are: e-shop, e-store, Internet shop, web shop, web-store, online store, and virtual store [1]. An online shop evokes the physical analogy of buying products or services from a retailer or shopping center.


Fig 1- Secured Transaction using Bio-metrics

This process is referred to as business-toconsumer (B2C) online shopping. In the case where a business entity buys from another business entity, the process is called business-tobusiness (B2B) online shopping. Statistics have shown that eBay and Amazon are the largest online retailing corporations in the world (both based in the United States). Generally, in Nigeria today, the traditional way of offline business transactions presents a lot of limitations to prospective customers and business owners. Some of these include security vulnerabilities, improper account auditing, poor inventory documentation, inflexibility and poor service delivery, etc. A new paradigm has emerged with the advent of E-commerce and its associated technologies to address these issues. Since the issue of security vulnerability will always constitute enormous constraint for numerous users on the E-commerce domain, the need for a re-engineered e-commerce platform with a robust encryption algorithm that will allow access and protect the users and the web owners will gain wide acceptance. This follows the fact that secure communication is an intrinsic requirement of today's world of on-line transactions. Whether exchanging financial, business or personal.

II. LITERATURE SURVEY

Patil, Chandrekar, Chavan and Chaudhri[10] proposed an Online transaction system built on the technology of embedded system. It uses an 8bit AT Mega 16 micro controller developed by Microchip technology and the original verifying method (the use of PIN) to authenticate users. Reference [10] aimed at improving Online transaction security through the use of biometrics technology (fingerprint). This Online transaction system is related to the model presented (three-tier authentication model) in the use of biometric authentication. However, the system by [10] employed PIN authentication as a second factor authentication. The use of PIN is susceptible to replay attack and illegal access to customers' credentials. If a false acceptance rate error occurs with the fingerprint device, a criminal with the correct PIN of an account holder can easily access the customer's account illegally.

Jaynthi and Sarala [11] developed an Online transaction system that uses PIN for user authentication. The system sends an approval SMS alert to the corresponding mobile phone number of an account holder upon a successful authentication. An acceptance message received from the account holder grants access to the user else access is denied. Simultaneously, the image of the person who made the transaction is sent to the e-mail account of the bank and that of the account holder. If any misuse of card or Online transaction hijack occurs, the system automatically alerts both the bank manager and the police by switching on a buzzer. The system achieves this through the use of GSM technology and Internet communication network.

Iwasokun and Akinyokun [12] developed a fingerprint based authentication framework for Online transaction. This Online transaction system is based on fingerprint authentication, eliminating the use of PIN and Online transaction card for authentication. The Internet serves as the operational environment and platform for the system. The thumbprint database of customers is available on the Internet. User verification involves enrollment, enhancement, feature extraction and matching. This work is related to the current study in the use fingerprint authentication. However, it has some defects such as the use of PIN as a means of authentication. The use of PIN is considered to be unreliable, in that if false acceptance rate occurs, the security of the system will be greatly compromised. Again, hosting sensitive information as customers' fingerprint database on the Internet could be risky as well since cyber crime has been prevalent in recent times.

Shimal and Jhunu [13] presented an enhanced Online transaction security system using twolevel authentication where PIN and OTP were both used for user authentication. This second level authentication (the use of OTP) was employed if a customer wishes to exceed a specified withdrawal limit otherwise the customer is authenticated using only PIN. This Online transaction system operates in two modes. The first mode operates like the traditional Online transaction system when a customerspecified withdrawal limit is yet to be attained. The second mode is an enhancement on the traditional Online transaction system. It is only used when a customer wishes to exceed the withdrawal limit.

Malviya [14] developed an Online transaction authentication model which uses face recognition technique to authenticate users for improved security. The Online transaction system consists of embedded camera that recognizes the face standing about 2 feet far in front of the system and performs matches against the facial database. The findings of Mwaikali [15] identified insecurity as one of the major challenges facing Online transaction users in Tanzania.

III. METHODOLOGY A. E-COMMERCE

Electronic commerce is called e-commerce. This refers to electronic media and the Internet for trading in goods and services. E-commerce includes a business with access to the Internet and computers, such as a E.g. Electronic Data Interchange (EDI). E-commerce refers to the website of the provider who exchanges goods or services with the user directly from the platform. The gateway uses a wireless shopping cart or shopping cart to pay by credit card, debit card, or Electronic Funds Transfer (EFT).Payment Transmission, which enables e-commerce, financial transactions, bricks and clicks, and traditional credit card payments, is an ecommerce application service for online transaction services. The most important variables of online transactions are payment channels, which include credit cards, debit cards, bank purchases and electronic funds transfers.

There is a need for payment gateways for the sustainable e-commerce of the future and the environment is shifting from cash to digital currency.

Analytics is an experimental method of converting data into information to make decisions. Analytics helps organizations collect, organize, review, and comment on their customers. The massive increase in data has led companies to rely on research to understand customer behavior. Retailers need access to realtime information to calculate the ROI of online and channel mergers. There are basic analytics for ecommerce players; average order size, cart size measurement, conversion rates and a deeper analytical approach are required to better understand customers. Companies constantly use social networks to promote their products. Social media includes blogs and computer applications that allow you to use your computer or mobile phone to connect to the internet and share information.Social networks are more important in product development and remind customers of different occasions. Product or service reviews are also helpful. It offers a branding tool to build a trusted consumer group, posts, word of mouth, etc.

B. FINANCIAL TRANSACTION

Financial transaction processing, refers to a class of systems that facilitate and manage transaction-oriented applications, typically for data entry and retrieval transaction processing. [13] The term "transaction" has different meanings in different scenarios. It typically refers to data entry and retrieval transactions in several industries, including banking, airlines, mail-order, supermarkets, and manufacturers. In the context of business or commercial transactions. OLTP(Online transaction processing) refers to processing in which the system responds immediately to user requests. An Automatic Teller Machine (Online transaction) for a bank is an example of a commercial transaction processing application. In computer science, transaction processing is information processing that is divided into individual, indivisible operations, called transactions.

Financial transaction processing increasingly requires support for transactions that span a network and may include more than one company. For this reason, new OLTP (Online Transaction Processing) software uses client/server processing and brokering software that allows transactions to run on different computer platforms in a network. In large applications, efficient OLTP (Online Transaction Processing) may depend on sophisticated transaction management software (such as CICS) and/or database optimization tactics to facilitate the processing of large numbers of concurrent updates to an OLTP (Online Transaction Processing) -oriented database. For even more demanding Decentralized database systems, OLTP brokering programs can distribute transaction processing among multiple computers on a network. OLTP (Online Transaction Processing) is often integrated into Service-Oriented Architecture (SOA) and Web services.

C. PROPOSED TECHNIQUE

The proposed system developed and implemented a an Online Transaction model using three-tier The investigative phase of the OOADM was deployed as the paradigm for systematic study in order to obtain information on the current trends in the research area of Online Transaction. The information obtained necessitated the definition of a high-level model. The system uses three different layers of authentication to validate Online Transaction users' identity to foster improved security. The three authentication mechanisms used are: password, biometric identifier (fingerprint) and OTP. The system is made up of alpha keys, numeric keys and some special character keys for authentication. The Online Transaction was interfaced with a fingerprint reader for improved security. In addition, the system also has a card reader, cash dispenser, screen, fingerprint scanner, and bank database. When the system is idle, a greeting message is displayed, the keys on the keypad remain.

ALGORITHM

<u>Algorithm: Client Access Validation</u> while true: get UI and session_key get_all UI list for id in UI_list ifid=UI generateaccess_key sendaccess_key else Commit client endfor DESCRIPTION The port of the process i

The next of the process is to validate access of the client. The client must connect to the access server next. Then the client should bind its user id and session key and sent it to the Access Server. The Access Server validates the client's user id and the session id and then generates the access id for the valid client. And sends that access id to the client. Using this access id the client can request for services from the service Online Transaction.

Algorithm: Secured Key Factor Authentication while true: get request parse request $parameters[] \leftarrow extract \ parameter(request)$ for value in parameters: getaccess key end for ifaccess key is valid generateaccess_id generatekey_factor key factor \leftarrow user id+access id+key id encryptkey_factor with access_key sendkey_factor else reject access

Algorithm: Alpha Key Factor Authentication

thread_ while true: getaccess_request parseaccess_request extractaccess_id validateaccess_id ifaccess_id=server_access_id provideaccess_id else discard_user_id

DESCRIPTION

From the previous step the client will be having the access key to access the Transaction. The Online Transaction has to verify the access of the client to the server. The client must send his/her access key to the Online Transaction. Using this access key the Online Transaction generate and encrypt the key factor for the client. Then encrypts the key factor using the access key of the client. Then the encrypted key factor is sent to the client. Only the valid client can decrypt the key factor and extract the access id from the key factor. Then using this access id the client can request for the web service. The service Online Transaction validates the access id, user id and Alpha keys foe the secured Transaction.



Fig.2. Work Flow

IV. RESULT AND DISCUSSIONS

The system provides strong security with the use of biometric identifier and alphanumeric characters for password. This password becomes very difficult if not impossible to be guessed correctly by fraudsters Online Transaction theft will be reduced since a person's biometric is not transferable. This is required before a successful authentication process. The level of security provided by the system will make it impossible for would-be perpetrators. This will discourage Online Transaction fraud. The problem of replay attack is completely eliminated with the use of OTP. Customers' confidence will be restored on the use of Online Transaction to meet their banking needs.

Three-Tier Authentication Model seeks to design an Online Transaction system with three layers of authentications – the use of password, fingerprint and OTP. Therefore, the system is interfaced with a fingerprint scanner for biometric authentication and it is capable of generating token as OTP. In addition, the system introduced alphabets and special characters to the existing numeric keypad of an Online Transaction system. The design is aimed at providing robust security to the existing card-based Online Transaction system by eliminating the problem of identity theft through the introduction of password as a substitute for PIN, and the use of fingerprint and OTP for second and third tierauthentication respectively.

The result of the proposed Three-Tier Authentication Model for Online Transaction a system with improved security, interfaced with fingerprint scanner for biometric а authentication and an Online Transaction keypad with a modified form factor. The incorporation of alphabets and special character keys to the existing numeric keys changed the form factor of the keypad. The system is also capable of generating OTP for third-tier authentication to eliminate any possibility of replay attack. The system was evaluated alongside the existing system in terms of speed and the level of security each provides.

A. SECURITY

The existing systems employ only one means of verifying customers' identity. In the case of identity theft, where a successful guess is made on a customer's PIN by fraudsters or where customers' debit cards and PINs are stolen or forcefully taken from them, cash are withdrawn from the Online Transaction through illegal means. This undoubtedly leads to huge financial loss to both the customer and the bank. However, the new system provides improved security on Online Transaction system by employing the use of three different authentication mechanisms. The essence is to cover up every loophole which could lead to identity theft. It is obvious that the three security protocols can never fail at the same time, hence, eliminating the problem of identity theft completely. Again, to prove the level of security the new system provides, different wrong passwords, OTPs and fingerprint templates were tried on the system but access was denied in all cases. This is an indication that the new system provided robust security and cannot be hacked by criminals whose aim is to withdraw customers' cash illegally in a short time.

B. TIME CONSUMPTION

Time taken to complete user authentication was collated for two categories of users who underwent authentication three times both in the existing system and the proposed system. The result was represented using a line chart in Figure.3.



V. CONCLUSION

The problem of identity theft, unauthorized access to customers' account details and illegal withdrawal of cash from the Online Transaction will be completely eliminated with the adoption of the proposed Three-Tier Authentication Model as the current use of PIN for Online Transaction user's verification and identification is marred with some level of insecurity. This Three-Tier Authentication Model uses password, biometric identifier and OTP to verify the validity of user's identity at three different layers of authentication. These three authentication mechanisms must be in the affirmative before access is granted to the user. The adoption of the new system by financial institutions will strengthen the security of Online Transaction systems and restore the confidence of customers. The study will no doubt foist a sense of futility on wouldbe perpetrators. This will discourage Online Transaction fraud. Bank customers are reassured that their account details and cash cannot be tampered with, hence, better service delivery which will attract many customers to use Online Transaction.

REFERENCES

 B. Komolafe (2017, Jan.). Nigerians withdraw N4.7 trillion through online Transaction in 2016 [Online]. Available: Http://www. vanguardngr.com/2017/01/nigerianswithdrawn4-7-trillion-online Transactions-2016.

- [2] N.Y. Asabere, R.O. Baah and A.A. Odefiya, "Measuring standards and service quality of Automated Teller Machines (online Transactions) in the banking industry of Ghana," International Journal of Information and Communication Technology Research, vol 2, issue 3, pp 216–226, 2020.
- [3] P.S. Rose and S.C. Hudgins, Management and Financial Services, 9th ed. New York: McGraw-Hill, 2013.
- [4] J. Hota. (2020) "Window-based and webenabled online Transaction: issues and scopes," The IUP Journal of Information Technology, vol. 3, issue 4, pp 52-59. Available: http://www.academia.edu/5043734/windowsbased-andweb-enabled-online Transactionsissues-and-scopes
- [5] H.A. Hayder (2011) "Implementing additional security measure on online Transaction through biometric [Online]. Available: http://www.etd.uum.edu.my/2576
- [6] Diebold Incorporated (2020). online Transaction fraud and security. [Online]. Available: http://securens.in/pdfs/KnowledgeCenter/5_onl ineTransaction %20Fraud% 20and%20Security.pdf
- [7] S.A. Adelewo. (2010, August). "Challenges of automated teller machine (online Transaction) usage and fraud occurrences in Nigeria. A case study of selected banks in Minna metropolis," Journal of Internet Banking and Commerce, vol. 5, issue 2, pp 10-20. S.T. Bhosale and B.S.

Sawant, "Security in e-banking via cardless biometric online Transactions," International Journal of Advanced Technology and Engineering Research (ITATER), vol. 2 issue 4, pp 9-12, July 2020.

- [8] Gemalto (2011, Feb). One-Time-Password Solution for Secure Network Access. [Online]. Available:http://www.gemalto.com/brochuressi te/downloadsite/Documents/ent_otp_secure_ac cess.pdf
- [9] B. Patil, B.S. Chandrekar, M.P. Chavan and B.S. Chaudhri, "RBI 3X – fingerprint based online Transaction," International Journal of Advanced Research in Computer and Communication Engineering, vol. 5, issue 3, pp 577 – 581, March 2016.
- [10] P. Jaynthi and S. Sarala, "Enhanced online Transaction security using differentiated passwords with GSM technology," International Journal of Innovative Research in Engineering & Science, vol. 5, issue 4, pp 28 – 35, May 2015.
- [11] G.B. Iwasokun and O.C. Akinyokun. (2013) "A fingerprint-based authentication framework for online Transaction," Journal of Computer Engineering and Information Technology, vol. 2, issue 3. Available: http/dx.doi.org/ 10.4172/2324-9307.1000112.
- [12] Shimal and D. Jhunu. (2011). "Designing a biometric strategy (fingerprint) measure for enhancing online Transaction security in Indian e-banking system," International Journal of Information and Communication Technology Research, vol. 1, issue 5. http://www.esjournals.org.
- [13] Malviya. (2014, Dec.). "Face recognition technique: Enhanced safety approach for online Transaction," International Journal of Scientific and Research Publications, vol 4, issue 12. Available: http://www.ijsrp.org
- [14] E.J. Mwaikali, "Assessment of challenges facing customers in Automated Teller Machine in the banking industry in Tanzania: A case of some selected banks in Tanzania," International Journal of Research in Business and Technology, vol. 4, issue 3, pp 480-488, 2014.
- [15] S. Subasree And N. K. Sakthivel, "Design Of A New Security Protocol Using Hybrid Cryptography Algorithms", IJRRAS 2 (2), February 2010 Subasree&Sakthivel, Design of a New Security Protocol.
- [16] Udo G.J., "Privacy and Security Concerns As Major Barriers for E-commerce: A Survey Study," Information Management & Computer Security, vol. 9, no.4, pp.165-174, 2020.
- [17] 84. Roca J.C., Garcia JJ., de la Vega JJ., "The Importance of Perceived Trust, Security and Privacy in Online Trading Systems," Information Management & Computer Security, vol. 17, no. 2, pp. 96-113, 2009.
- [18] Chen Y-H., Barnes S., "Initial Trust and Online buyer behavior," Industrial Management & Data Systems, vol. 107, no. 1, pp. 21-36, 2007.
- [19] Abdulghader.A.Ahmed.Moftah, SitiNorul Huda Sheikh Abdullah, Hadya.S.Hawedi, "Challenges Of Security, Protection And Trust On E-Commerce: A Case Of Online Purchasing In Libya" International Journal Of

Advanced Research In Computer And Communication Engineering, Vol. 1, Issue 3, May 2020.

A Systematic Review on Privacy Preservation of Big Data in Cloud

R. Sathya

I M.Sc. (CS) Student, Department of Computer Science, A.V.C.College (Autonomous), Mannampandal, Mayiladuthurai, Tamil nadu, India. sathyaravi1201@gmail.com

Abstract - With the recent advancement in the field of information technology and internet, the amount of data being generated, processed and stored is increasing very rapidly to the scale of Big Data. This data contains sensitive private information of individuals; their privacy needs to be protected from attackers. Many organizations demand efficient solutions to store and analyze huge amount of information. Cloud computing as an enabler provides scalable resources and significant economic benefits in the form of reduced operational costs. The key challenge in cloud computing environments is privacy preservation. This paper reviews on a systematic privacy preservation of big data in cloud.

Keywords – Big Data, Privacy Preservation, Cloud, Anonymization.

I. INTRODUCTION

During recent years, data production rate has been growing exponentially [1, 2]. Many organizations demand efficient solutions to store and analyze these big amounts of data that are preliminary generated from various sources such as high throughput instruments, sensors or connected devices. For this purpose, big data technologies can utilize cloud computing to provide significant benefits, such as the availability of automated tools to assemble, connect, configure and reconfigure virtualized resources on demand. These make it much easier to meet organizational goals as organizations can easily deploy cloud services.

This shift in paradigm that accompanies the adoption of cloud computing is increasingly giving rise to security and privacy considerations relating to facets of cloud computing such as multi-tenancy, trust, loss of control and accountability [3]. Consequently, cloud platforms that handle big data that contain sensitive information are required to deploy technical measures and organizational safeguards to avoid data protection breakdowns that might result in enormous and costlydamages.

Sensitive information in the context of cloud computing encompasses data from a wide range of different areas and disciplines. Data concerning health is a typical example of the type of sensitive information handled in cloud computing environments, and it is obvious that most individuals will want information related to their health to be secure. Hence, with the proliferation of these new cloud technologies in recent times, privacy and data protection requirements have been evolving to protect individuals against surveillance and database disclosure.

This paper presents an overview big data privacy issues and challenges, privacy techniques of big data, privacy preservation big data solutions in the cloud.

II. BIG DATA PRIVACY ISSUES AND CHALLENGES

The big data could be a current innovate technology, which is rapidly adopted by numerous industry to anticipates how case patterns and client behavior. The security and privacy challenges cover the entire spectrum of the Big Data life cycle (Figure 1) i.e. data production source or devices, the data, data processor, storage of data, and data transport and data usage on different devices[13].



Figure 1: Security and Privacy challenges in Big Data system

A. Big Data Privacy Issues

Following are the application that shows how privacy issues are created in big data:

1) Privacy issue in Big Mobile Data

Everything is available on mobile nowadays. People are sharing a lot of information on mobilephones. Mobile sends data to the service provider without user's knowledge. Identifying the person using his mobile data and the details provided by the service provider is very easy. Therefore, privacy in mobile data is very important. Text message analysis is an example of unstructured big data analytics in mobile. This method is used in WhatsApp to identify the mobile number.

2) Privacy issue in Health Care Data

Big data analytics and genome research having real-time access to the patient record to take adecision. Electronic Health Record (EHR) helped a lot to digitize the health care system make accurate and complete EHR. EHR have personal information. Therefore, the privacy-preserving analysis is required and data need to be anonymized before data analysis. Ex. Pathology report

3) Privacy issue in Social Media Data

Social media is the biggest revaluations in past decade. Lots of information is being shared by people on social media. Sometimes, people close to you share some information about you, which you don't want disclose on social media. This may lead to privacy violation of an individual. Ex. Facebook, Twitter etc.

4) Privacy issue in Web Usage Data

Intel wants to make its internal website dynamic based on web usage data of all the users of the website. With browser information and IP address from web usage data, any user be identified and whatever activities he is performing on line may be detected. Thus, privacy is required.

B. Big Data Privacy Challenges

A gathering of privacy and security issue must be considered before building a big data situation. We have the following most significant challenges, when managing with big data.

1) Random Distribution

Big data analytics in view of parallelism, the extensive data is stored and processed in a different cluster, which is gathering of disseminated servers around the world.

2) Privacy

Current big data analytics care for all data with the similar need and do not combine unique activities, similar encryption or impaired processing [10]. Thus, programmer and malicious node gain accesstothecluster.

3) Computations

The main objective of big data is to separate helpful bits of knowledge performing the particular calculation. It is significant to secure and protect this computation.

4) Integrity

Big data contain a substantial volume of a substance. For creating decision based on big data, it's necessary to make sure the validity and trust level of that data in order to avoid on the suspect or compromised records.

5) Communication

Big data is put away data in a few hubs, which are distributed around in the world. All communication linking nodes and cluster is completed by the ordinary public and private network [4].

6) Access control

In big data context, any modification in cluster's states such expansion or erasure of nodes ought to checked by confirmation mechanism to shield the system. Sometimes, database items fall under restrictions and practically no users can see the secret information, like personal info in medical records, etc.

III. PRIVACY-PRESERVATION TECHNIQUES FOR BIG DATA

Privacy in big data has raised significant issues transfer into seeing the need for proficient privacy preservation methods. In this segment, we need to examine three privacypreserving techniques: **Data Anonymization**, **Differential Privacy**, **Notice**, and **Consent**.

There are some common terms used in privacy field of this method:

Identifier attributes include information that uniquely and directly distinguishes individual such as full name driver license, social security numbers.

- Quasi-identifier attributes a set of information such as birth date, gender, age, zipcode. That can be combined with other external data in order to reidentifyindividuals.
- Sensitive attributes are private and personal information.
- Intensive attributes are the general and the innocuous information.
- Equivalence classes are sets of all records that consist of the same value on the quasi-identifiers.

A. Anonymization

Data Anonymization is that the method of fixing information which will be used or revealed during an approach that stops the identification of key information. It is typically referred as data de-identification. These are method utilized in this: Key items of confidential data are until obscured during approach that maintains data privacy and unleashes data publically by Anonymization. Example, a hidden attribute like full name, license number, voter id etc. The main drawback with data anonymization is that data might look anonymous however re-identification is often done simply by linking by linking it to other different external data [4]. It is shown that re-identification of anonymous medical records is done exploitation external constituent list data. Example attribute as if genders, date of birth, the zip code that can join without side data to the re-identify individual call quasi-identifier attributes. There are privacypreserving methods of anonymization: k-anonymization, T-closeness, L-diversity.



Figure 2: Differential privacy

B. Differential Privacy

Differential privacy is a mathematical framework for ensuring the privacy of individuals in datasets. It can provide a strong guarantee of privacy by allowing data to be analyzed without revealing sensitive information about any individual in the dataset.

Differential Privacy could be technique enable analysts to remove helpful answer as of database containing individual information, donating solid individual privacy protection[5][6]. The aim of this method is to limit the chances of individual distinguishing proof while querying data. The phases of differential privacy are personating in the Figure 2.

C. Notice and Consent

The foremost common privacy preservation method for

web services is notice and consent [12]. Whenever individual access a new application or services, a notice stating privacy issues is displayed. Theend user needs to consent the notice before exploitation service. This technique empowers a private to makesure his privacy rights. It puts the burden of privacy preservation on the individual [7]. Once applied to bigdata, this technique posesvarious challenges.

IV. PRIVACY CONSIDERATIONS OF PROCESSING SENSITIVE DATA

The security issues in cloud computing lead to a number of privacy concerns. Privacy is a complex topic that has different interpretations depending on contexts, cultures and communities.

"Privacy is the claim of individuals, groups, or institutions to determine for themselves when, how, and to what extent information about them is communicated to others."

The International Association of Privacy Professionals (IAPP) refers to privacy as the appropriate use of information under the circumstances.

The security and privacy responsibilities of cloud providers include integrating solutions to ensure legitimate delivery of cloud services to the cloud consumers. The security and privacy features that are necessary for the activities of cloud providers are described in Table1.

SecurityContext	Description
Authentication	Authentication and
and Authorization	authorization of cloud
	consumers using pre-defined
	identification schemes
Identity and	Cloud consumer provisioning
Access Management	and deprovisioning via
	heterogeneous cloud service
	providers
Confidentiality, Integrity,	Assuring the confidentiality
Availability (CIA)	of the data objects,
	authorizing data
	modifications and ensuring
	that resources are available
	when needed
Monitoring and	Continuous monitoring of the
Incident Response	cloud infrastructure to ensure
	compliance with consumer
	security policies and auditing
	requirements
Policy	Defining and enforcing rules
Management	to enforce certain actions such
	as auditing and proof of
	compliance

Privacy	Protect personally identifiable
	information (PII) within the
	cloud from adversarial attacks
	that aim to find out the
	identity of the person that the
	PII relates to

Table 1: Securityand PrivacyFactors of theCloud Providers

V. PRIVACY-PRESERVING BIG DATA SOLUTIONS IN THE CLOUD

Over the time, organizations have collected valuable information about the individuals in our societies that contain sensitive information, e.g. medical data. Researchers need to access and analyze such data using big data technologies [8-10] in cloud computing, while organizations are required to enforce data protection compliance.

There has been considerable progress on privacy preservation for sensitive data in both industry and academia, e.g., solutions that develop protocols and tools for anonymization or encryption of data for confidentiality purposes. This section categorizes work related to this area according to different privacy protection requirements.

"Outsourcing privacy" is the concept where a database owner updates the database over time on untrusted servers. This definition assumes that database clients and the untrusted servers are not able to learn anything about the contents of the databases without authorized access. The authors [11] implement a server-side indexing structure to produce a system that allows a single database owner to privately and efficiently write data to, and multiple database clients to privately read data from, an outsourced database.

Homomorphic encryption is another privacypreserving solution that is based on the idea of computing over encrypted data without knowing the keys belonging to different parties. To ensure confidentiality, the data owner may encrypt data with a public key and store data in the cloud. When the process engine reads the data, there is no need to have the DP's private key to decrypt the data. In private computation on encrypted genomic data [12], the authors proposed a privacy-preserving model for genomic data processing using homomorphic encryption on genome-wide association studies.

Anonymization is another approach to ensure the privacy of sensitive data. SAIL provides individual-level information on the availability of data types within a collection. Researchers are not able to cross-link (which is similar to an equality join in SQL) data from different outside studies, as the identities of the samples are anonymized.

VI. CONCLUSION

Big data privacy is a critical component in today's digital world where data is generated, accessed and shared widely with each other. It is now mandatory to promise privacy in big data analytics. Privacy measures should now give emphasis on the uses of data instead of data collection. This paper gives a good insight on

138

overview of privacy issues and challenges, privacy preserving techniques for big data. Also this paper focused on privacy preserving big data solutions in the cloud.

REFERENCES

- A.Szalay and J. Gray, "2020 Computing: Science in an exponential world," *Nature*, vol.440, pp. 413-414, Mar. 2006.
- [2] Lynch, "Big data: How do your data grow?," *Nature*, vol. 455, pp. 28-29, Sept. 2008.
- [3] S.Pearson, "Privacy, security and trust in cloud computing," in *Privacy* and Security for Cloud Computing, Computer Communications and Networks, pp. 3-42, Springer London, 2013.
- [4] L.Sweeney, "K-anonymity: A model for protecting privacy", International Journal on Uncertainty, Fuzziness, and Knowledge Based system, pp. 557-570, 2002.
- [5] J.Salido, "Differential Privacy for everyone", White paper, *Microsoft Corporation*, 2012.
- [6] O.Heffets and K.Ligett, "Privacy and data-based research", NBER Working paper, September 2013.
- [7] F.H.Cate, V.M.Schonberger, "Notice and Concent in a World of Big Data", *Microsoft Global Privacy Summit Summary Report and Outcomes*, November 2012.
- [8] S.Sharma, U.S.Tim, J.Wong, S.Gadia, S.Sharma, "A Brief Review on Leading Big Data Models," *Data Science Journal*, 13(0), pp. 138-157. 2014.
- [9] S.Sharma, U.S.Tim, J.Wong, S.Gadia, R.Shandilya, S.K.Peddoju, "Classification and comparison of NoSQL big data models," *International Journal of Big Data Intelligence (IJBDI)*, Vol. 2, No. 3, 2015.
- [10] S.Sharma, R.Shandilya, S.Patnaik, A.Mahapatra, "Leading NoSQL models for handling Big Data: a brief review," *International Journal* of Business Information Systems, Inderscience, 2015.
- [11] Y.Huang and I.Goldberg, "Outsourced private information retrieval," in Proceedings of the 12th ACM Workshop on Workshop on Privacy in the Electronic Society, WPES '13, (New York, NY, USA), pp. 119-130, ACM, 2013.
- [12] K.Lauter, A.Lopez-Alt, and M.Naehrig, "Private computation on encrypted genomic data," Tech. Rep. MSR-TR-2014-93, June 2014.
- [13] M.Gostev, J.Fernandez-Banet, J.Rung, J.Dietrich, I.Prokopenko, S.Ripatti, M.I.McCarthy, A.Brazma, and M.Krestyaninova, "SAIL - a software system for sample and phenotype availability across biobanks and cohorts," *Bioinformatics*, vol. 27, no. 4, pp. 589-591, 2011.

An Outlook on Role and Challenges of Big Data Analytics in Health Care

J. Rakkesh Kumar¹ and Dr. M. Hemamalini²

¹ I M.Sc., Department of Computer Science,

A.V.C. College (Autonomous), Mannampandal, Email-ID: u202530rakkeshkumar@gmail.com ² Assistant Professor, Department of Computer Science, A.V.C. College (Autonomous), Mannampandal, Email-ID: maliniavcce@gmail.com

Abstract-

Big data analytics is the process of collection, examining and analyzing large amount of data to discover market trends, insights and patterns that can helps to make better business decision. It is a growing area with the potential to provide useful insight in healthcare. Data is being produced at 2.5 quintillion bytes a day. While many dimensions of big data still present issues in its use and adoption, such as managing the volume, variety, velocity, veracity, and value, the accuracy, integrity, and semantic interpretation are of greater concern in clinical application. It has provided tools to accumulate, manage, analyse, and assimilate large volumes of disparate, structured, and unstructured data produced by current healthcare systems. Potential areas of research within this field which have the ability to provide meaningful impact on healthcare delivery are also examined. The utilization of large volumes of medical data while merging multimodal data from different sources is discussed. This paper provides an overview of Big Data and applicability of it in healthcare system.

Keywords: Big data, Healthcare and Big data analytics.

I. INTRODUCTION

It has been considered that the production of data will be 44 times greater in 2020 than it was in 2009. "Big data refers to the tools, processes and procedures allowing an organization to create, manipulate, and manage very large data sets and storage facilities". A collection of large and complex data sets which are difficult to process using common database management tools or traditional data processing applications. To make effective use of the big data in healthcare system, a perceptive of what the 2.5 quintillion bytes of data consists of, where they exist in, are they raw, processed or derived artifacts. There are five dimensions of big data [1].

 Volume: This is the management of the terabytes or peta bytes of data data.(Feldman, 2012).

- Variety :The data in many forms such as structured, semi structured and unstructured (Feldman, 2012)
- Velocity: Data in motion such as the frequency of data that is produced, processed, and analyzed (Feldman, 2012).
- Veracity: The data in doubt (ie. Data inconsistencies, Incompleteness. (Clifford, 2008).



Figure 1. Five dimensions of Big Data

Big data is not just about size. It finds insights from complex, noisy, heterogeneous, longitudinal, and voluminous data and it answers to unexplored areas. The Challenges are capturing, storing, searching, sharing and analyzing. The big data challenges in health care are

- Inferring information from diverse patient sources.
- Understanding unstructured clinical observations in the exact perspective.
- Efficiently managing the large volumes of medical imaging data.
- Examine genomic data is a composite task and Capturing the patient's behavioral data.

II. BIG DATA FRAMEWORK

A. Information Extraction

Real world clinical data is noisy and varied in nature and sometimes correlated attributes. This data resides in multiple databases. The sources and technique for big data in Health care system are

- Electronic Health Records (EHR) data
- Healthcare Analytic Platform
- Resources

The collection, integration, and analysis of such big, complex, and noisy data in healthcare are a challenging task. For this reason, healthcare information systems can be considered as a form of big data not only for its sheer volume, but also for its complexity and diversity which makes traditional data warehousing solutions prohibitively cumbersome and ill suited for large scale data exploration and modelling. We examine how a big data framework can be leveraged to extract and pre-process data. The Hadoop is our big data framework to archive performance, scalability and fault tolerance for our task. The data can be obtained from various sources such as International Classification of Diseases, Current Procedural Terminology, Lab Results, Medication and Clinical Results. ICD stands for International Classification of Diseases. ICD is a hierarchical terminology of diseases, signs, symptoms, and procedure codes maintained by the World Health Organization

(WHO). CPT stands for Current Procedural Terminology created by the American Medical Association. CPT is used for billing purposes for clinical services [1].

The data can be obtained from lab results. The standard code for lab is Logical Observation Identifiers Names and Codes (LOINC®). The Challenges for lab are many lab systems still use local dictionaries to encode labs and diverse numeric scales on different labs. Some data can be missed. The order of a lab test can be predictive, for example, BNP (B-type Natriuretic Peptide -Blood Test) indicates high likelihood of heart failure. Next the data can be obtained from medication. Standard code is National Drug Code (NDC) defined by Food and Drug Administration (FDA), which gives a unique identifier for each drug. But it is not used universally by EHR systems. There are too specific, drugs with the same ingredients but different brands have different NDC. Clinical notes contain rich and diverse source of information. The challenges for handling clinical notes are some are Ungrammatical, short phrases, Abbreviations, Misspellings, Semistructured information i.e. Copy-paste from other structure source , Lab results, vital signs and SOAP notes(Subjective, Objective, Assessment, Plan) [1]. Table 1 shows the summary of common EHR data [2] and Figure 2 shows the healthcare analytic platform.

	ICD	СРТ	Lab	Medication	Clinical Notes
Availability	High	High	High	Medium	Medium
Recall	Medium	Poor	Medium	Inpatient: High Outpatient: variable	Medium
Precision	Medium	High	High	Inpatient: High Outpatient: variable	Medium
Format	Structured	Structured	Almost Structured	Structured and Un Structured	Un Structured
Pros	Easy to work	Easy to work , high precision	Data validity is high	Data validity is high	Doctors suggestions
Cons	Disease code often used for screening, so disease might not be there	Missing data	Data normalization and ranges	Prescribed not necessary taken	Difficult to process



Figure 2.Healthcare Analytic Platform

*1.Text Mining in Healthcare:*Text Mining helps with information overload and overlook and discovers unsuspected links from the huge amount of literature and supports medical research. [3] It integrates knowledge from many sources and enhances clinical decision support systems and supports translational medicine. It reduces costs and errors in handling information [4].

B. Feature Selection:

Feature Selection is the process of selecting a set of relevant features for construction the model.

The basic principle behind the usage of feature selection technique is that the data contains many features that are either redundant or immaterial and can thus the irrelevant information can be removed without loss of information [5]. The selected risk factors are predictive of the target condition [6]. There is a minimal correlation occurs between data driven risk factors and knowledge driven risk factors [7]. Figure 3 shows the combining knowledge and data driven risk factor.





C. Predictive Modeling.

Two types of predictive models are used. Regression techniques are used for continuous outcome and classification techniques are used for categorical outcome [8]. We took case study as Heart failure on set prediction. Patient similarity learns a customized distance metric for a specific clinical context [9].



1. Early detection of Heart failure: Heart failure (HF) is a complex disease. It reduces the cost for payers and improves the existing clinical guidance of Heart failure prevention. It is slow or potentially reverse disease progress for providers. It also improves the quality of life and reduces mortality for provider. It gives huge societal burden. There are 5 millions HF patients in US, 0.5 million new cases each year and 20% life time risk after 40 years old. The main aim of predictive modelling design is to classify HF cases against control patients. For example consider 50,625 Patients (Geisinger Clinic PCPs). This study shows 4,644 HF case patients and 45,981 controlled patients

matched on age and gender. Big data could save the health care industry up to \$450 billion. But additional things are essential too [1].

- Patients should take more energetic steps to develop their health and promoting a coordinated approach to care in which all caregivers have right to use the same information.
- Any professionals who treat patients must have correct performance record for achieving the best outcomes.
- Improving value, quality and identifying new approaches to health-care delivery.

III. TOOLS USED FOR BIG DATA FRAMEWORK

New tools and programming paradigms for such data intensive applications influence the distributed computation model. Apache Hadoop [10] is one such distributed framework that implements a computational paradigm. In MapReduce [11], the application is divided into many fragments of work. Each application may be executed on a number of compute nodes in a cluster of data intensive applications. In combination with a distributed system such as Hadoop, the Apache Mahout [12] framework gives a useful set of machine learning libraries. These libraries used for executing modelling tasks such as classification and clustering though there is need to develop advanced domain specific applications of these algorithms. Mahout was intended to work in combination with Hadoop to scale for compute clusters and large datasets. The Hive [13] and Cassandra [14] are now accessible for distributed query processing and exploratory analyses; although few case studies are available that shows their use in the healthcare setting.

IV. CONCLUSION

Big data analytics is a promising right direction which is in its infancy for the healthcare domain. Healthcare is a data-rich domain. As more and more data is being collected, there will be increasing demand for big data analytics. Unraveling the "Big Data" related complexities can provide many insights about making the right decisions at the right time for the patients. Efficiently utilizing the colossal healthcare data repositories can yield some immediate returns in terms of patient outcomes and lowering care costs. Data with more complexities keep evolving in healthcare thus leading to more opportunities for big data analytics.

REFERENCES

[1] Jimeng Sun, Chandran K.Reddy," Big Data Analytics for Health care", Tutorial Presentation at the SIAM International Conference on Data Mining, Austin, TX, 2013.

[2] Joshua C. Denny Chapter 13: Mining Electronic Health Records in the Genomics Era. PLoS Comput Biol. 2012 December; 8(12):

[3] Meystre et al. Extracting Information from Textual Documents in the Electronic Health Record: A Review of Recent Research. IMIA 2008

[4]

http://www.nactem.ac.uk/event_slides/Ananiadou2 21009.pdf

[5] Bermingham, Mairead L.; Pong-Wong, Ricardo; Spiliopoulou, Athina; Hayward, Caroline; Rudan, Igor; Campbell, Harry; Wright, Alan F.; Wilson, James F.; Agakov, Felix; Navarro, Pau; Haley, Chris S. (2015). "Applications of highdimensional feature Selection: evaluation for genomic prediction in man", *Sci. Rep.* **5**.

[6] Jimeng Sun, Jianying Hu, Dijun Luo, Marianthi Markatou, Fei Wang, Shahram Ebadollahi, Steven E. Steinhubl, Zahra Daar, Walter F. Stewart. Combining Knowledge and Data Driven Insights for Identifying Risk Factors using Electronic Health Records. AMIA (2012).

[7] Dijun Luo, Fei Wang, Jimeng Sun, Marianthi Markatou, Jianying Hu, Shahram Ebadollahi, SOR: Scalable Orthogonal Regression for Low-Redundancy Feature Selection and it's Healthcare Applications. SDM'12

[8]http://www.ijarcsse.com/docs/papers/Volume_5/ 6_June2015/V5I6-0570.pdf

[9] Jimeng Sun, Fei Wang, Jianying Hu, Shahram Edabollahi: Supervised patient similarity measure of heterogeneous patient records. SIGKDD Explorations 14(1): 16-24 (2012).

[10] The Apache Software Foundation., http://hadoop.apache.org/common/credits.html.

[11] Ghemawat D.J .MapReduce: simplified data processing on large clusters. In: Proc of OSDI, 2004.

[12] Owen S. and Anil R. Mahout in Action. Manning Publications Co., Greenwich, Connecticut, 2010.

[13] Thusoo, A., Sarma, J. S., Jain, N., Shao, Z., Chakka, P., Anthony, S., Liu, H., Wyckoff, P., and Murthy, R. Hive—a warehousing solution over a Map-Reduce framework. In VLDB, 2009.

[14] Wikipedia, http://en.wikipedia.org/ wiki/ Apache_Cassandra. International Conference on "Computational Intelligence and its applications" (ICCIA-2024) ISBN: 978-81-967420-1-0

Nanorobotics Using Artificial Intelligence

Ragavi. R¹ & Subastri. P²

1st year MCA Department of Computer Applications, A.V.C.College of Engineering, Mannampandhal, Mayiladuthurai-609 305. ¹ragavirahul753@gmail.com ²subasripari31@gmail.com

Abstract- Nan robotics are the emerging field in science and technology which help to develop the mini robots and operates nano scale. There is a task between those technology which plays a realistic in medical and electronics field oftechnology. These robots are in the scale of from one to 1000 nano meters, which it is equal to the billion to trillion of nano robots in the covenants. It has capability to cure the cancel cells and biological function to reproduce the substance in the human body.

Keywords-technology, nanometer, nanoscale, styling, insert

I. INTRODUCTION (NANOROBOTICS)

Nano robots are short and emerging technology in the field of machine. It has a small compound which nearly it works on the scale of (10-9) meters. Now a days, in the engineering world, building Nanoid robots are that so easy, which people need a knowledge about a machine and robots. Nano robotsare opposed to microbotics. Building a nano robot which preferred size are approximately ranging in size from 0.1 to 10 micrometers The terms nanobot, nanoid, nanite, nanomachine and nanomite, such devices currently under research and development.

II. WHAT IS NANOTECHNOLOGY

A. Definition

It is an innovative hidden sensor which technology is using Nani scale and nano meters. In this technology we can assume that it can also hack DNA substance and plays an important role in medical fields. It also toxic and independent which it either comes to the chemical field of positions. Another definition, it defines the interaction between the human and the nano biotics. It is mainly used in the secondary and higher medical fields, used to cure cancers and pre precursor of biotics form of substance. It allows a resource with para medical center.

B. Application of Nanorobotics

It was mainly used for medical fields. And Biological machine which is also used to destroy and cure cancer cell. It also found that it may use for monitoring the diabetic stages of a person. It can perform in higher accuracy. It increases quality of manufacturing the products.



Fig (I) nano robots' interaction with biologicalcells.

III. HOW IT WORKS?

It supports many other technical supports to generate those functions like nano scale sensors, nano scale accurator, which used to detect the presence of molecular cells. In this sensor and accurator Or detest the signals of molecules which need an aid to recover the cells. In here, nano robots interaction plays an role to cure and help to recover the aid cells. It also has the potential and positive significant to benefit the society.



Fig (ii) Nano molecular with bio-cells interaction.

The science behind this mechanism is quite complex. Passage of cells across the blood endothelium, a process known as transmigration, is a mechanism involving engagement of cell surfacereceptors to adhesion molecules, active force exertion and dilation of the vessel walls and physical deformation of the migratingcells. The science behind this mechanism is quite complex. Passage of cells across the blood endothelium, a process known as transmigration, is a mechanism involving .

engagement of cell surface receptors to adhesion molecules, active force exertion and dilation of the vessel walls and physical deformation of the migrating cells. By attaching themselves to migrating inflammatory cells, the robots can in effect "hitch a ride" across the blood vessels, bypassing the need for a complex transmigration mechanism of their own.[1]

IV. HISTORY

In 1959, A physicist Richard Feynman is American Nobel prize, He was the first who speak through nano technology and found application to produce a nano subscription at the california instituteof technology. In the 21st century I t denotes that a specific area is consolidated and accepted the perseverance of nano technology and micro manufactures, organic chemistry behind its working techniques. It includes other areas such as micro- Organisms, around billions and trillions of dollars are invested to produce these technologies in the market place through NNI (National Nanotechnology initiative). This may turn the environment sector to drive a economic growth.

V. Types of nanotechnology

There are Different types of Nano technologybeing classified in below:

□ GRAPHITE BASED

In this Biological machine, such it used for the medical fields and enhances the development of resources. There is mechanism which need a carbon molecule in the body to represent the nano meters to detect the data and information in the cells and needa scale range between one to 100 nano meters in size.

□ SURFACE BOUND PHRASE

A large number of molecules are enough to detect the machine in the biological function which is either to re produce the molecular and a aid the defective cells and repent it's functions.

Dry nanotechnology

Dry nanotechnology uses inorganic materials, including metals and semiconductors, to create items used by electrical and mechanical engineers to promote development in manufacturingtechniques (Madon).[2]

□ Wet nanotechnology

Wet nanotechnology is an upcoming new subdiscipline of nanotechnology that will be dominated by different types of wet engineering. The procedures will take place in aqueous solutions and are quite similar to those employed in biotechnology

/ bio-molecular manufacturing, which is primarily concerned with the creation of biomolecules such as proteins and DNA/RN. Working up to large masses from small ones is the goal of wet nanotechnology (also known as wet nanotech).[3]

VI . EXAMPLES AND APPLICATIONS OF NANOTECHNOLOGY

They are Mainly found in these areas: R&D

Society Nanotechnology and nanomaterial's can be applied in all kinds of industrial sectors. They are usually found in these areas:

Electronics

Nano technologies in Electronics works faster, and need quantum computing to research more.it also allows to produce new display technologies in presence of the nano materials.

Biomedicine

It Improving early diagnosis and also help to re produce the ideal treatment of neurogenic cancer and helps to relieve the cancer cells to aid negative cells and repent the positive significant. It can't be harmful to human because now those technology is improved; human has to certain live in these zones now a days.

VII Nanotechnology in the future

In the future, nanotechnology could also enable objects to harvest energy from their environment. New nanomaterials and concepts are currently being developed that show potential for producing energy from movement, light, variations in temperature, glucose and other sources with high conversion efficiency.[4]



A. Figures and Tables

VIII. ACKNOWLEDGMENT

I would like to express my special thanks of gratitude, who gave me the golden opportunity to hiswonderful project of "Nano robotics using AI" who also helped me in completing my project. I came to know about so many things, im really thankful. Secondly, I would like to thank my parents and friends who helped me in a lot in finalizing this project within the limited frame.

REFERENCES

- [1] https://en.m.wikipedia.org/wiki/Nanorobotics.
- [2] https://builtin.com/robotics/nanorobotics
- [3] https://worldnanotechnologyconference.com/pro gram/scientificsessions/nanochemistry-and- wet-nanotechnology.
- [4] https://www.google.com/amp/s/phys.org/news/2 016-03-waysnanotechnology-future.amp.
- [5] https://www.cureus.com/articles/108503-the-use-of-nanorobotics-inthe-treatment-therapy-of-cancer-and-its-future-aspects-a-review#!/.
- [6] R. Nicole, "Title of paper with only first word capitalized," J. Name Stand. Abbrev., in press.

Big Data Analytics: Review & Recommendation

Raja Thangavelu,

Department of Computer Science and Applications, ARC Viswanathan College, Sitharkadu, Mayiladuthurai raja.thangavelu.chennai@gmail.com

Abstract — In the past few decades, the explosive growth of internet and related applications has resulted in generation of data of multiple petabyte scale levels, - the Big Data. Given the volume, velocity and variety of data, in Big Data, the ability to handle such data and bring out business value out of them has been a challenge and an area of exploration, analysis and research. This paper aims at reviewing a limited set of technology, available, including Cloud Computing and Artificial Intelligence (AI). The gamut of database (DB) systems to store such humongous raw data, the methods, techniques and tools to extract, transform and load data, analytical tools, machine learning models etc., provides multiple options depending on the business objective to be achieved. An approach of raw data handling, treating, and aggregating data as needed, whereas staying with the basic tenets of software application architecture definition and delivering the business objective, is discussed in this paper. This also includes analytics, selection of appropriate analytical models, and data visualization of the analytical results. As a conclusion, this paper presents a possible recommended generic architecture covering all these aspects, such that the immediate, varied business objectives are achieved and still stay relevant to deliver the future analytical / business needs too.

Keywords—analytics, big data, data science, data visualization

I. INTRODUCTION

In today's connected world, we are generating more data in two days than in decades of history. Approximately, 90% of world's data is generated in the last few years. Businesses are vying to get the best of these data, to further grow their businesses. In the process, Big Data has come into greater prominence. Data is generated from various sources – legacy to new devices.

Raja Thangavelu is with ARC Viswanathan College, Sitharkadu, Mayiladuthurai (phone: 0091-8903060470; e-mail: raja.thangavelu.chennai@gmail.com). This paper will brief about what is Big Data, the technology available for the various aspects of Big Data Analytics solution. Finally, would recommend a generic architecture for such a solution.

As indicated above, more data is generated continuously almost in all our everyday transactions. This not only results in huge volume of data but also in variety of data, data types in structured, semi structured, unstructured data, the velocity with which such data flows in, the veracity of such data, and the value such data mean, in different contexts. These are the characteristics of Big Data. Thus storing Big Data and analysis of it was technically challenging with normal databases designed for structured data.

The development of Big Data technologies presented a whole lot of information & opportunity to businesses, which was earlier restricted to structured data. Businesses seek to take advantage of all the other types of data too, to arrive at better business strategies and decisions.

In a survey conducted in late 2021, by a consulting firm, across executives from 94 large companies, 91.7% said they will be increasing their investments in Big Data projects and other data & AI initiatives and about 92.1 % indicated that they are getting measurable improved outcomes thru these initiatives.[5]

II. BIG DATA TECHNOLOGIES AND TOOLS

Hadoop distributed processing framework, kicked off, in real earnest the Big data practice by being the first open source platform that could handle such varied sets of data. A whole lot of other technologies like NoSQL were developed. Some have eclipsed, many others are still being used by many organizations. The technologies and tools that are almost common options for now are:

- · Raw Data Storage.
- Apache Hadoop Distributed File System built on a cluster system that let BigData data types to run parallel. It can process from one server to multiple computers. 2) Amazon cloud Simple Storage Service and 3) Google Cloud Storage.

• NoSQL databases. Examples include Dynamo DB, Cassandra, Couchbase, CouchDB, HBase, MarkLogic Data Hub, MongoDB, Redis and Neo4j.

• Processing engines. Examples include Spark, Hadoop MapReduce and

• Stream processing platforms such as Flink, Kafka, Samza, Storm and Spark's Structured Streaming module.

• Structured Query Language (SQL) query engines. Examples include Drill, Hive, Presto and Trino.

• Data Lake, Data Fabrics, Data Mesh and Data Warehouse (DW) platforms. Examples include Amazon Redshift, Delta Lake, K2view, Atlan, Denodo, Talend Data Catalog, Google BigQuery, Kylin and Snowflake.

• Platforms and managed services. Examples include Amazon Elastic MapReduce (Amazon EMR), Azure HDInsight, Cloudera Data Platform and Google Cloud Dataproc.

A. Some of the Google Big Data Offerings: Google BigQuery

BigQuery is a data warehouse that processes and analyzes large data sets using SQL queries. These can also store streaming data for real-time analytics. BigQuery is a serverless offering. The model demands Enterprises to pay only for the storage and compute they consume.

Google Cloud Dataproc

Cloud Dataproc is a managed Apache Hadoop and Spark service for batch processing, querying, streaming and machine learning. It can be fully integrated with other Google big data services

Google Cloud Data Catalog

Cloud Data Catalog is a data discovery service that enables enterprises to create a catalog by capturing technical and business metadata.

Google Cloud Data Transfer

Cloud Data Transfer moves small and large amounts data -- physically and virtually -- to Cloud Storage, BigQuery and Cloud Dataproc.

Google Cloud Bigtable

Bigtable is a managed NoSQL database service that can handle massive loads of data. Cloud Bigtable uses a low-latency storage stack and is available globally. It supports open source HBase Application Programming Interface (API), which makes applications more portable. It is used for marketing, financial, and Internet of Things (IoT) data.

B. Some of the Amazon Big Data Offerings:

Amazon Web Services (AWS) also provides for a set of data analytics tools. At present, there are few primary AWS products for analytics: Elastic MapReduce (EMR), Kinesis, Redshift, Data Pipeline and Machine Learning.

Big Data storage and access

Amazon's is well equipped for storing big data, including the Amazon's Simple Storage Service (S3). In order to enable faster access and processing, Amazon's no SQL database, DynamoDB provides for it. Elastic File System can be another useful tool in big data projects, scaling up to handle large flows of data.

Big Data processing and visualizing

Amazon Kinesis provides for building visual dashboard or application to monitor data as soon as it comes in. Third-party products like Tableau provide connectivity to EMR and other AWS products. Amazon Machine Learning provides visualization tools and helps create models to react to real-time data.

C. Others

APACHE Cassandra

APACHE Cassandra is an open-source NoSQL distributed database. It is highly scalable and highly available without compromising on speed and performance. It was created by Facebook

Qubole

It's an open-source tool can fetch data using ad-hoc analysis in machine learning. Qubole is a data lake platform with reduced time and effort in moving data pipelines. It is capable of configuring multi-cloud services such as AWS, Azure, and Google Cloud.

APACHE Spark

This is another framework that is used to process data and perform numerous tasks on a large scale. It offers easy-to-use APIs that provide easy data pulling methods and is highly suitable for Machine Learning (ML) and AI today.

Mongo DB

It's an open-source NoSQL platform and a documentoriented database that is used to store a high volume of data. It uses collections and documents for storage and its document consists of key-value pairs

Apache Storm

It's a robust, user-friendly tool used for data analytics and has no programming language barrier. It can handle data in fault-tolerance and horizontally scalable methods.

Statistical Analytical System (SAS)

It is one of the best tools for creating statistical modeling used by data analysts. Statistical Analytical System or SAS allows a user to access the data in any format (SAS tables or Excel worksheets). Besides that it also offers a cloud platform for business analytics called SAS Viya and also to get a strong grip on AI & ML.

Rapid Miner

It's a fully automated visual workflow design tool used for data analytics. It's a no-code platform and users aren't required to code for segregating data. Rapid Miner can be used to easily deploy their ML models to the web or mobile.

III. POSSIBLE GENERIC ARCHITECTURE

The Fig.1 below presents a possible typical generic logical architecture right from the data sources to information delivery thru various means / devices. This architecture is defined considering the basic tenets of *"Functional Decomposition"* and *"Separation of Concerns"*.

As could be seen the data sources can be myriad depending on the type of industry / business being served by the Big Data Analytics Architecture.

This Architecture recommends edge computing capability for the data sources/ devices to enable a minimal data cleansing, data correction etc. to be done at the source itself. For example, for customer data, the address may not be filled in completely or correctly. In such case the data can be processed thru software like Informatica's AddressDoctor and shall be corrected/ completed at the source itself. Such an approach will ensure that the computing power at different levels of the architecture is utilized correctly and optimally.



Fig.1.A generic typical Big Data Analytics Architecture. Note: Some elements shown may not be applicable for certain real time implementation.

Considering the fact that the businesses have been running lot of legacy systems, applications, data and other systems till the internet/ social media based information explosion, All those data, including those data which were never analyzed – called as "Dark Data" – isin a way a gold mine to be explored, and should be considered as we move forward and adapt new technologies and new architecture. This is ensured in this architecture.

As could be seen, the data sources include datasets, Enterprise Data Warehouse, Relational Database Management System (RDBMS) databases etc. and other structures, unstructured, semi-structured and streaming video, audio and textual data.

The architecture considers a Data Lake as an essential element, as Data Lake can store diverse data and naturally Big Data requirements. Therefore support any architecture going forward which would like to use the enterprise wide data would need to necessarily consider a Data Lake. Organizations can use Data Lake information in many ways, and the data sources do not need a predefined purpose to qualify for ingestion into a Data Lake. Analysts explore, experiment and evaluate Data Lake information to identify its benefits and use cases. The strategy for a Data Lake implementation is to ingest and analyze data from virtually any system that generates information. In a Data Lake, analysts apply schemas after the ingestion process is complete.

Inspite of this, in memory database would be needed for quick responding requirements / analytics. This can also include databases that support in-database analytics. Hence, the architecture includes these too. The architecture includes a "Data Fabric" for accessing and using the data, in the Data Lake. A data fabric abstracts away the technological complexities engaged for data movement, transformation and integration, making all the data available all across the enterprise. Data Fabric being semantic in nature, it enables technical users & business users to see value of the enterprise data pooled in the Data Lakeand extract, transform and use for their specific requirements.

A Data Mesh is also considered in the architecture considering that there will be domain specific applications with specific requirements calling for custom syntactic code, and query structures to serve their needs. Though, we expect this to becoming more and more less in development, as the concept of No Code Low Code gains currency, but might still be needed/ continue to exist, for maintenance and operation of legacy systems.

Data Lake, Data Fabric and Data Mesh are evolving continuously, and there are overlaps in their definition at this point of time.

However, the author of this paper believes and for the purpose of this paper, Data Lake is considered as a repository of raw/ near-raw data, Data Fabric & Data Mesh is considered as a data architecture framework enabling data access, security, governance, integration, quality, replication, virtualization, caching, federation, querying, and analysis. An enhanced version may also include ML/AI models that automate building of data pipelines for other data consuming applications / analytics applications.

The author is of the view that considering the growing nature of the generated data, which is further enlarged thru Generative AI, the above consideration of Data Lakes, Data Fabric and Data Mesh will hold good for long, for almost all the future application development and enterprise application architecture.

As in Fig.1, and as briefed above, the Big Data is accessed by the analytical, reporting and other applications / toplayers thru the Data Access Layer.

As shown in the architecture, a layer of data analysis and modeling layer consisting of Online Analytical Processing (OLAP) cubes, analytical, AI/ML etc. engines based on the business needs shall be implemented on distributed computing machines. These shall also include edge computing to ensure the data is prepared and made available as required by these engines / functional models. Legacy / new OLAP cubes might continue to provide some multi-dimensional reports as may be required by the business, till they are completely replaced by other cost effective means.

With the advancement and clarity in the usage of the varied analytical and ML algorithms, the engines can be typically segregated/ grouped as shown in the architecture based on business needs and computing power requirements.

Traditional Analytics Engine

This would serve some of the pre-defined metrics, business measures on predefined dashboards, as required by the business. These might include statistical measures and models, if any. However, the data pipeline for this might be built and improved with time using AI/ML at the Data Fabric level.

AI/ML Analytics Engine

This employs advanced analytics algorithms to deliver business insights, Predictive and Prescriptive analytics from the underlying data. Some of the algorithms as applied for typical problem scenarios are as under in TABLEI:

Deep Learning (DL) Engine

Considering that Big Data includes streaming data, social media, audio, video etc., Deep Learning (DL) algorithms are used for analysis and decision making. It uses artificial neural networks to work on the data and do tasks. This mimics the functioning of the human brain. Ideally this shall be implemented such that higher computing power is made available for this.

Natural Language Processing (NLP) Engine

As businesses are becoming more global, more languages are being supported by the enterprise software applications. Also, business dealings becoming more multi-modal thru chatbots, messaging platforms, automated speech recognition etc., the need for language processing and to respond as required, become an essential component of Big Data Architecture. For example, ChatGPT is a DL model using NLP. With Large Language Models (LLM), NLPs are becoming more capable. It plays a role in sentiment analysis.

AI, ML, DL, NLP and ASR are quite intermingled, the below –Fig. 2 gives a visual representation of the same:

Table I : Some Algorithms and Applicable problem

 Scenarios



Fig. 2. Visual Representation of the overlapping relationship of AI, ML, DL, ASR, and NLP

ASR /

Real-time Analytics Engine

In today's data driven world, there are instances that require instant analysis of streaming data and quick decision making. For example, credit card fraudulent transaction detection and remedial action. Real-time database and In-memory database may be used to get the required business outcome.

Apart from the above, a control layer is considered to ensure security, access control, load balancing and session management for the entire Big Data Analytics system.

A logical API layer is considered to separate the presentation layer from the underlying analysis, modeling, analytical, control and other functional models.

The presentation layer can comprise of any User Interface (UI) client like websites, Dashboards, Handhelds, IOT devices, reports, alarms, etc.

Data Governance and Data ops are also need to be considered appropriately along with the enterprise data architecture definition.

As could be seen the architecture is following the Model-View-Controller(MVC) architecture pattern, however the recommended implementation is detailed as above. Some aspects, data source, models etc. may vary from this based on the implementation / business scenario.

IV. CONCLUSION

As globally, more and more data is being generated; Big Data will not be big anymore (or) may need to be redefined. However, the above generic architecture may become more main stream & reference for all application development as more and more data, AI is expected in all applications.

Going forward with more data to be crunched and with development of affordable Quantum computing, the infrastructure might move to Quantum computing. Cloud based solutions have begun to be more affordable and secure too. Therefore, Big Data Analytics may move more to the cloud and less to in-premise infrastructure. Newer business models like DaaS (Data as a Service) will become more preferred. Data Lake, Data Fabric and Data Mesh might get more matured and more main stream as unstructured data becomes more dominant type to be stored. The architecture will also make better use of the hitherto un-analyzed "Dark data", which might unlock benefits to the business. Edge computing and NLP will continue to deliver value as shown in the architecture. Data governance will become more important and critical. The traditional software and IT roles will be redefined to suit this growth in data, data driven applications, AI, ML and analytics. Data analysis might become more common and may become the basic skill needed in every industry. Altogether, in the next few years, Businesses along with AI growth will demand more Data analysis, management and Big data analytics.

REFERENCES

1) David Dietrich, Barry Heller, Belbel Yang, *Data Science* and *Big Data Analytics: Discovering, Analyzing, visualizing and presenting data,* John Wiley& Sons, 2015

2) Christopher Long, Shaurya Rana, *Reference Architecture* to enable self-service Analytics, Gartner Research, April 2022

3) Paul C. Zikopoulos, Chris Eaton, Dirk DeRoos, Thomas Deutsch, George Lapis, *Understanding Big Data Analytics for Enterprise classHadoopand Streaming Data*, McGrawHill, 2012

4) Enterprise Strategy Group Research Report, 2023 Technology Spending Intentions Survey, TechTarget, Nov 2022

5) TechTarget – *The ultimate guide to Big Data for Businesses*,https://www.techtarget.com/searchdatamanage ment/pro/The-Ultimate-Guide-to-Big-Data-for-Businesses?offer=Content_OTHR-PillarPage_Theultimateguidetobigdataforbusinesses

6) Navadeep Singh Gill, 10 Latest Trends in Big Data Analytics for 2023 / Ultimate Guide, Xenonstack.com/ttps://www.xenonstack.com/blog/latesttrends-in-big-data-analytics 7) Yuvaraj, *10 most Popular Big Data Analytics Tools* geeksforgeeks.org https://www.geeksforgeeks.org/10-most-popular-big-data-analytics-tools/

8) AI Analytics Vs Traditional Analytics: 3 Essential Differences, aberdeen.com https://www.aberdeen.com/blog-posts/blog-ai-analytics-vstraditional-essential-differences/

9) Gulbahar Karatas, *Data Fabric* 2024: Modern Data Integration Components Guide, research.aimultiple.comhttps://research.aimultiple.co m/data-fabric/

10) Alex Woodie, Data Mesh Vs Data Fabric: Understanding the Differences, datanami.com https://www.datanami.com/2021/10/25/data-mesh-vs-datafabric-understanding-the-differences/

11) Artificial Intelligence : AI vs ML vs NLP, sonix.ai https://sonix.ai/articles/difference-between-artificialintelligence-machine-learning-and-natural-languageprocessing

12) Kathleen Walch, *What a big data strategy includes and how to build one*, techtarget.com https://www.techtarget.com/searchdatamanagement/feature /How-to-build-an-enterprise-big-data-strategy-in-4-steps

13) Ronald Schmelzer, *Top trends in Big Data for 2023 andBeyond*,techtarget.comhttps://www.techtarget.com/sear chdatamanagement/feature/*Top trends in Big Data for 2023 and Beyond*

14) Donald Farmer, 6 essential big data practices for businesses, techtarget.com https://www.techtarget.com/searchbusinessanalytics/tip/6essential-big-data-best-practices-for-businesses

15) Mary K.Pratt, *Building a big data architecture – Core Components, best practices*,techtarget.comhttps://www.techtarget.com/search datamanagement/feature/Building-a-big-data-architecture-Core-components-best-practices

16) Ron Karjian, *Enterprise data lakes hold the key to actionable insights*, techtarget.com https://www.techtarget.com/searchdatamanagement/ehandb ook/Enterprise-data-lakes-hold-the-key-to-actionable-insights

17) George Lawton, *Data quality for big data: Why it's a must and how to improve it,* techtarget.com https://www.techtarget.com/searchdatamanagement/feature /Data-quality-for-big-data-Why-its-a-must-and-how-to-improve-it

18) Deep Learning Vs Machine Learning: The ultimate battle, turing.com/kb/ultimate-battle-between-deep-learning-and-machine-learning

ROBOTIC INNOVATIONS IN HEALTHCARE: ENHANCING PATIENT CARE AND EFFICIENCY

Dr.R. Senthamarai¹ and Dr.A. Senthil Kumar²

¹Ph.D. Part time Research Scholar, Department of Computer Science and ²Assistant Professor, Department of Computer Science Tamil University Thanjavur, Tamilnadu.

senthukrishna17@qmail.com

Abstract - Robots are increasingly being adopted in healthcare to carry out various tasks that enhance patient care. This research explores the transformative role of robotics in the healthcare sector, focusing on the integration of robotic technologies to improve patient care and operational efficiency. By examining the deployment of robots in tasks ranging from surgery assistance to patients monitoring, we aim to assess the impact on healthcare outcomes and resource utilization. Our findings highlight the potential of these advancements to enhance medical services, reduce human error, and contribute to the overall improvement of healthcare deliverv. Through a comprehensive analysis of case studies and emerging technologies, the study provides insights into the future landscape of robotics in healthcare. Emphasizing the potential benefits and challenges associated with this evolving integration.

Keywords - Robotics, Healthcare, Surgery, Rehabilitation, Socially Assistive, Imaging Assistance.

I. INTRODUCTION

In recent years, the intersection of robotics and healthcare has given rise to a new era of innovation, revolutionizing the way medical services are delivered and patient care is administered with technological advancements reaching unprecedented heights, robots are increasingly finding applications in diverse healthcare settings from surgical theatres to patient care facilities. This integration promises not only to alleviate the burdens on healthcare professionals but also to enhance the overall quality and efficiency of healthcare services.

The deployment of robots in healthcare is driven by a myriad of factors, including the quest for precision in medical procedures, the need to minimize human error, and the optimization of resource utilization. This convergence of robotics and healthcare presents a compelling narrative of a future where machines and humans collaborate synergistically to address the complexities of modern medical challenges.

This introduction sets the stage for an in-depth exploration of the various facets of robotic innovations in healthcare. From surgical robotics and rehabilitation assistance to Artificial Intelligence driven diagnostics and patient monitoring, the scope of these innovations is vast and holds significant promise for transforming the healthcare landscape.

Since the advent of the COVID-19 pandemic, the healthcare industry has been flooded with novel technologies to assist the delivery of care in unprecedented circumstances. Staff vacancy levels increased, social restrictions curtailed many traditional means of care delivery [1,2] and stringent infection control measures brought new challenges to human-delivered care [3]. Although many of the challenges that the pandemic brought onto healthcare have subsided, staff burnout, an increasingly elderly population, and backlog strains caused by the pandemic have meant that staff shortages persist across healthcare systems across the world [4,5].

This paper aims to establish the types of robots being used in healthcare and identify where they are deployed by way of qualitative analysis of the literature. Through this, predictions can be made for the future of robotics.

II. IDENTIFIED ROBOTS AND THEIR ROLES IN HEALTHCARE

Surgical:

Surgical robots can be used to assist in performing surgical procedures. Their specific roles within surgery are varied, ranging from instrument control to automated surgical table movement.

Some of the identified robots can also assist with biopsy. For example, the iSR'obot Mona Lisa can assist with visualisation and robotic needle guidance in prostate biopsy. This robot prospectively in a group of 86 men undergoing prostate biopsy with the researchers primarily evaluating detection of clinically significant prostate cancer [6].

Rehabilitation and Mobility:

Rehabilitation and mobility robots are those that can physically assist or assess patients to aid in achieving goals. They can function to improve dexterity, achieve rehabilitation targets or aid in mobilisation. These robots may be used in the inpatient setting as well as in community rehabilitation centres.

Most are used for their ability to provide physical support to patients, assisting with rehabilitation. This can include single-joint or whole-body support. Others may be used for posture training through robotic tilt tables or for mobilisation through robotic wheelchairs.

The most common robot, Lokomat, is a gait orthosis robot that can be used for rehabilitation in disorders such as stroke. Its primary role is to increase lower limb strength and range of motion. One study that evaluated this robot came from Husemann et al. [7] who carried out a randomised controlled trial with 30 acute stroke patients and compared those receiving conventional physiotherapy alone to those receiving conventional plus Lokomat therapy and evaluated outcomes such as ambulation ability. The second most studied robot, HAL, is a powered exoskeleton with multiple variants including a lower limb and single-joint version. Studies predominantly explore its use in neurological rehabilitation, but research is also present in areas such as post-operative rehabilitation.

Radiotherapy:

Radiotherapy robots can be used to assist with delivery of radiotherapy. This review identified one robot in this group: Cyber knife. This robot can assist with application of radiotherapy and image guidance to manage conditions such as liver and orbital metastases.

Telepresence:

A core feature of the telepresence robotic group is the ability to allow individuals to have a remote presence through means of the robot. The robot may be used for activities such as remote ward rounds, remote surgical mentoring or remote assessment of histology slides. Croghan et al. [8] used this robot for surgical ward rounds with a remote consultant surgeon and compared the experience to conventional ward rounds.

Interventional:

Separate from their surgical counterpart, robots from this group are used to assist with interventional procedures. Their function can range from catheter guidance to stent positioning. Arya et al. [9] carried out a case–control study comparing the Niobe system with conventional manual catheter navigation and evaluated effectiveness and safety in managing atrial fibrillation.

Socially Assistive:

Socially assistive robots can take multiple forms, such as humanoid or animal-like, and work to provide support in areas traditionally done by humans such as companionship and service provision. Pepper is a humanoid robot with a touch screen, capable of interacting with people through conversation. Boumans et al. [10] explored the use of Pepper in outpatient clinics with a randomised clinical trial.

Pharmacy:

There are a group of robots with the specific role of assisting with the management and delivery of pharmacy services. This includes drug storage, dispensing and compounding. For example, a robot may assist in preparation of cytotoxic drugs with the goal of reducing errors and minimising operator risk. The BD Rowa Vmax is an automated system that allows for storage of medication and dispensing at the request of a user. Berdot et al. [11] used this system in a teaching hospital pharmacy and evaluated the return on investment including the rate of dispensing errors. The APOTECA Chemo system can be used to automate the production of chemotherapeutic treatment. Buning et al. [12] explored the environmental contamination of APOTECA Chemo compared to conventional drug compounding.

Imaging Assistance:

Robots in this group have been specifically used for their ability to assist in carrying out imaging in different areas of medicine. They predominantly include robotic camera holders in theatre but can also include robotic microscopes in neurosurgery and transcranial magnetic stimulation robots. Robotic camera holders may be controlled by various inputs such as voice and a joystick.

Disinfection:

Robots may be used to disinfect clinical areas such as the ward or outpatient clinic. Systems use ultraviolet (UV) light for disinfection of rooms, with the UVD-R being able to move autonomously. UVD-R was explored by Astrid et al. [13] who analysed its ability to disinfect waiting rooms in hospital outpatient clinics and compared this to conventional manual disinfection.

Delivery and Transport:

There exists a role for robots in the transfer of items between areas. The TUG Automated Delivery System is a robot that after being loaded by an operator was used to autonomously deliver drugs from the pharmacy department to the ICU.

III. CHALLENGES OF ROBOTICS IN HEALTHCARE

Integrating robotics into healthcare, while promising, poses several challenges.

Cost implications:

The upfront costs associated with acquiring and implementing robotic systems can be substantial, making it challenging for many healthcare facilities especially smaller ones, to invest in such technologies.

Training and skill gaps:

Healthcare professionals may require specialized training to operate and interact with robotic systems effectively. Addressing this skill gap is crucial for the successful integration of robotics into routine medical practices.

Ethical considerations:

The use of robots in sensitive healthcare tasks, such as patient care and surgery, raises ethical questions. Issues related to privacy, consent, and the potential dehumanization of healthcare interactions need careful consideration.

Interoperability:

Ensuring seamless integration and interoperability of robotic systems with existing healthcare infrastructure and electronic health records poses a significant challenge standardization effort are crucial to overcome compatibility issues

Regulatory compliance:

The regulatory landscape for healthcare robotics is evolving, and ensuring compliance with existing regulations and standards is a complex task. Striking a balance between innovation and regulatory adherence is essential.

Patient acceptance:

Patient acceptance and trust in robotic technologies may vary. Gaining patient confidence and addressing concerns about the use of robots in healthcare settings are vital for successful adoption.

Maintenance and reliability:

Robotic systems require regular maintenance to ensure optimal performance. Unforeseen technical issues or malfunctions could disrupt healthcare services, emphasizing the need for robust maintenance protocols.

Limited adaptability:

Some robotic systems may have limitations in adapting to dynamic and unpredictable healthcare environments. Flexibility and adaptability are crucial for robots to effectively navigate diverse patient conditions and healthcare scenarios. **Data security:**

With the integration of robotics and Artificial intelligence, there are concerns about the security of patient data. Ensuring robust cyber security measures to protect

Department of Computer Science, A.V.C. College (Autonomous), Mayiladuthurai – 609 305.

sensitive health information becomes imperative.

Rural and underserved areas:

The deployment of robotic technologies may be challenging in rural or underserved areas due to infrastructure limitations, lack of resources, and limited access to specialized training for healthcare professionals.

Understanding and addressing these challenges is pivotal for the successful implementation of robotic innovations in healthcare, ensuring that the benefits of these technologies are maximized while mitigating potential risks and obstacles.

IV. CONCLUSION

In conclusion, the integration of robotics into healthcare represents a paradigm shift with the potential to revolutionize patient care and the overall healthcare landscape. The journey so far has showcased remarkable advancements in surgical precision, diagnostics, rehabilitation, and patient assistance.

The prospect of more widespread robotic-assisted surgeries, coupled with remote surgical capabilities, promises to bridge geographical gaps in access to specialized medical expertise telemedicine, empowered by robotics, could redefine the delivery of healthcare services, making them more accessible and responsive, especially in time of crises.

Rehabilitation robotics and the development of intelligent prosthetics hold the promise of restoring mobility and independence for individuals with physical impairments. The synergy between robotics and Artificial intelligence in diagnostics not only enhances the accuracy of medical assessments but also lays the groundwork for personalized treatment plan and tailored to individual patient profiles.

Companion robots, designed to provide emotional support and monitor well-being, could become integral in addressing the social and mental health aspects of patient care, particularly in aging populations. Additionally, the automation of routing tasks in laboratories and pharmacies through robotics contributes to the efficiency and precision of healthcare processes.

However, as we embrace the potential benefits of robotics in healthcare, it is imperative to acknowledge and address the associated challenges, cost considerations, ethical implications, and the need for robust regulatory frameworks must be navigated carefully toe ensure responsible and equitable deployment of these technologies.

In essence, the future of robotics in healthcare holds the promise of a more interconnected, precise, and patient-centric healthcare system. The ongoing collaboration between human expertise and robotic capabilities is poised to redefine medical practices, enhance patient outcomes, and contribute to a healthcare landscape that is not only efficient but also compassionate and accessible to all.

REFERRENCES

[1] Schmitt N, Mattern E, Cignacco E, Seliger G, König-Bachmann M, Striebich S, et al. Effects of the COVID-19 pandemic on maternity staff in 2020 – a scoping review. BMC Health Serv Res. 2021;27(21):1364. doi: 10.1186/s12913-021-07377-1. [PMC free article] [PubMed] [CrossRef] [Google Scholar]

[2] White EM, Wetle TF, Reddy A, Baier RR. Front-line nursing home staff experiences during the COVID-19 pandemic. J Am Med Dir Assoc. 2021;22(1):199–203.

doi: 10.1016/j.jamda.2020.11.022. [PMC free article] [PubMed] [CrossRef] [Google Scholar]

[3] Luciani LG, Mattevi D, Cai T, Giusti G, Proietti S, Malossini G. Teleurology in the time of COVID-19 pandemic: here to stay? Urology. 2020;140:4–6. doi: 10.1016/ j.urology.2020.04.004. [PMC free article] [PubMed] [CrossRef] [Google Scholar]

[4] Duffy SW, Seedat F, Kearins O, Press M, Walton J, Myles J, et al. The projected impact of the COVID-19 lockdown on breast cancer deaths in England due to the cessation of population screening: a national estimation. Br J Cancer. 2022;126(9):1355–1361. doi: 10.1038/s41416-022-01714-9. [PMC free article] [PubMed] [CrossRef] [Google Scholar]

[5] Ya-Ping Jin, P, Canizares M, El-Defrawy S, Buys YM (2022) Predicted backlog in ophthalmic surgeries associated with the COVID-19 pandemic in Ontario in 2020: a time-series modelling analysis. Can J Ophthalmol. 0(0). https://www.un.org/development/desa/pd/es/news/world-

population-ageing-2020-highlights. [PMC free article] [PubMed]

[6] Miah S, Servian P, Patel A, Lovegrove C, Skelton L, Shah TT, et al. A prospective analysis of robotic targeted MRI-US fusion prostate biopsy using the centroid targeting approach. J Robotic Surg. 2020;14(1):69–74. doi: 10.1007/s11701-019-00929-y. [PMC free article] [PubMed] [CrossRef] [Google Scholar]

[7] Husemann B, Müller F, Krewer C, Heller S, Koenig E. Effects of locomotion training with assistance of a robotdriven gait orthosis in hemiparetic patients after stroke: a randomized controlled pilot study. Stroke. 2007;38(2):349– 354. doi: 10.1161/01.STR.0000254607.48765.cb. [PubMed] [CrossRef] [Google Scholar]

[8] Croghan SM, Carroll P, Ridgway PF, Gillis AE, Reade S. Robot-assisted surgical ward rounds: virtually always there. BMJ Health Care Inform. 2018;25(1):41–56. doi: 10.14236/jhi.v25i1.982. [PubMed] [CrossRef] [Google Scholar]

[9] Arya A, Zaker-Shahrak R, Sommer P, Bollmann A, Wetzel U, Gaspar T, et al. Catheter ablation of atrial fibrillation using remote magnetic catheter navigation: a casecontrol study. Europace. 2011;13(1):45–50. doi: 10.1093/ europace/euq344. [PubMed] [CrossRef] [Google Scholar]

[10] Boumans R, van Meulen F, van Aalst W, Albers J, Janssen M, Peters-Kop M, et al. Quality of care perceived by older patients and caregivers in integrated care pathways with interviewing assistance from a social robot: noninferiority randomized controlled trial. J Med Internet Res. 2020;22(9):e18787. doi: 10.2196/18787. [PMC free article] [PubMed] [CrossRef] [Google Scholar]

[11] 1. Berdot S, Korb-Savoldelli V, Jaccoulet E, Zaugg V, Prognon P, Lê LMM, et al. A centralized automateddispensing system in a French teaching hospital: return on investment and quality improvement. Int J Qual Health Care. 2019;31(3):219–224. doi: 10.1093/ intqhc/mzy152. [PubMed] [CrossRef] [Google Scholar]

[12] Werumeus Buning A, Geersing TH, Crul M. The assessment of environmental and external cross-contamination aduthurai – 609 305. PROCEEDINGS 154

Department of Computer Science, A.V.C. College (Autonomous), Mayiladuthurai – 609 305.

in preparing ready-to-administer cytotoxic drugs: a comparison between a robotic system and conventional manual production. Int J Pharm Pract. 2020;28(1):66–74. doi: 10.1111/ijpp.12575. [PMC free article] [PubMed] [CrossRef] [Google Scholar]

[13] Astrid F, Beata Z, Van den Nest M, Julia E, Elisabeth P, Magda DE. The use of a UV-C disinfection robot in the routine cleaning process: a field study in an Academic hospital. Antimicrob Resist Infect Control. 2021;10(1):84. doi: 10.1186/s13756-021-00945-4. [PMC free article] [PubMed] [CrossRef] [Google Scholar]

[14] Ohmura Y, Suzuki H, Kotani K, Teramoto A. Comparative effectiveness of human scope assistant versus robotic scope holder in laparoscopic resection for colorectal cancer. Surg Endosc. 2019;33(7):2206–2216. doi: 10.1007/s00464-018-6506-4. [PubMed] [CrossRef] [Google Scholar]

[15] Summerfield MR, Seagull FJ, Vaidya N, Xiao Y. Use of pharmacy delivery robots in intensive care units. Am J Health Syst Pharm. 2011;68(1):77–83. doi: 10.2146/ ajhp100012. [PubMed] [CrossRef] [Google Scholar]

[16] Turagam MK, Petru J, Neuzil P, Kakita K, Kralovec S, Harari D, et al. Automated noncontact ultrasound imaging and ablation system for the treatment of atrial fibrillation: outcomes of the first-in-human VALUE trial. Circ: Arrhythm Electrophysiol. 2020;13(3):e007917. [PubMed] [Google Scholar]

ARTIFICIAL NEURON NETWORK: REVIEW

¹K. Kavitha–Assistant Professor of Computer Science Department -A.D.M. College for Women(Autonomous), , Nagapattinam. ²S. Muthulakshmi – III B.Sc. Computer Science - A.D.M. College for Women (Autonomous), Nagapattinam. computeradmc2023@gmail.com ³ M. Nithisha - III B.Sc. Computer Science, A.D.M College for Women (Autonomous), Nagapattinam.

ABSTRACT

In recent years, deep artificial neural networks (including recurrent ones) have won numerous contests in pattern recognition and machine learning. This historical survey compactly summarises relevant work, much of it from the previous millennium. The grid-connected system was evaluated in terms of power, energy, specific yield, capacity factor, and cost of energy, and payback period. This paper compares TWO different artificial neural network approaches for computer network traffic forecast, such as: 1) multilayer perceptron (MLP) using the back propagation as training algorithm; (2) recurrent neural network (RNN). Internet traffic prediction is an important task for many applications, such as adaptive applications, congestion control, admission control, anomalydetection and bandwidth allocation. To Streamline the diagnostic process in daily routine and avoid misdiagnosis, artificial intelligence methods(especially computer aided diagnosis and artificial neural networks) can be employed.

KEYWORDS

Building Integrated PhotoVoltaic(BIPV), Computer Network Traffic Prediction, ANN Medical Diagnosis, ANN Psychology.

INTRODUCTION

Artificial intelligence, or AI, refers to the simulation of human intelligence by software-coded heuristics. Nowadays this code is prevalent in everything from cloud-based, enterprise applications to consumer apps and even embedded firmware. When most people hear the term artificial intelligence, the first thing they usually think of is robots. That's because big-budget films and novels weave stories about human-like machines that wreak havoc on Earth. But nothing could be further from the truth. The applications for artificial intelligence are endless. The technology can be applied to many different sectors and industries. AI is being tested and used in the healthcare industry for suggesting drug dosages, identifying treatments, and for aiding in surgical procedures in the operating room.

Deep learning is a subset on Machine Learning, which is essentially a neural network with three or more layers. These neural networks attempt to simulate the behavior of the human brain—albeit far from matching its ability—allowing it to "learn" from large amounts of data. While a neural network with a single layer can still make approximate predictions, additional hidden layers can help to optimize and refine for accuracy.

Deep learning drives many Artificial Intelligence (AI) applications and services that improve automation, performing analytical and physical tasks without human intervention. Deep learning technology lies behind everyday products and services (such as digital assistants, voice-enabled TV remotes, and credit card fraud detection) as well as emerging technologies (such as self-driving cars).

Machine learning (ML) is defined as a discipline of artificial intelligence (AI) that provides machine the ability to automatically learn from data and past experiences to identify patterns and make predictions with minimal human intervention. This article explains the fundamentals of machine learning, its types, and the top five applications. It also shares the top 10 machine learning trends in 2022.

ARTIFICIAL NEURAL NETWORK:

ANNs are simple processing structures, which are separated into strongly connected units called artificial neurons (nodes). Neurons are organised into layers, one layer has multiple neurons and any one neural network can have one or more layers, which are defined by the network topology and vary among different network models (Haykin, 1998). Neurons are capable of working in parallel to process data, store experimental knowledge and use this knowledge to infer new data. Each neuron has a synaptic weight, which is responsible for storing the acquired knowledge. Network knowledge is acquired through learning processes (learning algorithm or network training) (Haykin, 1998). In the learning process, the neural network will be trained to recognise and differentiate the data from afinite set.

BUILDING INTEGRATEDPHOTOVOLTAIC(BIPV):

One of the most attractive applications of PV systems is building integrated photovoltaic (BIPV), which has undergone rapid developments in recent years. However, BIPV has experienced quick advancements as of late; regardless, BIPV systems have been progressively created.1 The performance of the PV unit is considered the primary and vital issue as the performance of these units

156

deteriorates over time as they are affected by weather conditions (temperature, relative humidity, dust, etc.).2 The high temperature of the photovoltaic (PV) unit by 10°C results in a 5% decrease in electrical efficiency.3 Also, dust accumulation is strongly depending on the location of the installed PV system, wind speed, etc..4 High temperature and dust accumulation cause losses in module performance and are attributed to external conditions.

TRAFFIC PREDICTION

Several types of ANN have been studied for network traffic prediction. There are several studies in feed-forward neural networks, such as MLP (Oliveira et al., 2014; Cortez et al., 2012; Ding et al., 1995), but many studies aim RNN (Hallas and Dorffner, 1998) because of the internal memory cycles that it has, facilitating learning temporal and sequential dynamical behaviour and making it a good model for time series learning.

MLP

One of commonest architectures for neural networks is the MLP. This kind of ANN has one input layer, one or more hidden layers, and an output layer. Best practice suggests one or two hidden layers (de Villiers and Barnard, 1993). This is due to the fact that the same result can be obtained by raising the number of neurons in the hidden layer, rather than increase the number of hidden layers (Hornik et al., 1989). MLPs are feed-forward networks, where all neurons in the same layer are connected to all neurons of the next layer, yet the neurons in the same layer are not connected to each other. It is called feed-forward because the flow of information goes from the input layer to the output layer.

The training algorithm used for MLP is the back propagation, which is a supervised learning algorithm, where the MLP learns a desired output from various entry data.

RNN

RNNs are neural networks that has one or more connections between neurons that forms cycles. These cycles are responsible for storing and passing the feedback of one neuron to another one, creating an internal memory that facilitate learning of sequential data (Hallas and Dorffner, 1998;Haykin, 1998). The cycles can be used anywhere in the neural network and in any direction, e.g., it can have a delayed feedback from the output to the input layer; a feedback loop from one hidden layer to another layer or to the same layer and any combination of it (Haykin, 1998).

MEDICAL DIAGNO

Artificial neural network provides a powerful too help doctors analyze, model and make sense of complex clinical data across a broad range of medical applications. Most application of AI neuron network medicine are classification problems. An artificial neural network (ann) is a computational model that attempts to account for the parallel nature of the human brain. It is a network of highly interconnected processing elements(neurons) operating in parallel. These elements are inspired by biological nervous systems.

PSYCHOLOGY (ANN)

Psychological, and emotional wellbeing are all considered to be components of one's mental health. It affects how someone thinks, feels, and responds to circumstances. When one has good mental health, it is easier to perform efficiently and reach their full potential (1). Preschool, adolescence, and adulthood are all included in the definition of mental health. Anxiety, social phobia, depression, panic disorder, substance dependence, and specific illnesses are factors that contribute to mental health problems that result in mental illness. The mental health status of adolescents in India is a topic of great concern and importance. Adolescents are children within the age group of 10-19. According to the National Mental Health Survey of India (2015-2016) (2), the prevalence of psychiatric disorders among adolescents of ages 13-17 is 7.3% and in the US it is 27.9% (3). This problem is further aggravated by stigma and lack of awareness surrounding mental health and a treatment gap of 95% in common mental disorders which is greater than the treatment gap for severe disorders (76%).

CONCLUSION

Neural Networks are suitable for predicting time series mainly because of learning only from examples, without any need to add additional information that can bring more confusion than prediction effect. Neuron networks areable to generalize and are resistant to noise. On the other hand, it is generally not possible to determine exactly what a neural network learned and it is also hard to estimate possible prediction error. However, neural networks were often successfully used for predicting time series. They are ideal especially when we do not have any other description of the observed series.

REFERENCES:

- [1] Deep Learning in Neural Networks, Jurgen Schmidhuber The Swiss AI Lab IDSIA Istituto Dalle Molle di Studi sull'Intelligenza Artificiale University of Lugano & SUPSI Galleria 2, 6928 Manno-Lugano Switzerland 8 October 2014. Goole Scholoar cited by 49.
- [2] Experimental and deep learning artificial neural network approach for evaluating grid - connected photovoltaic systems Hussein A. Kazem1,2 | Jabar Yousif1 | Miqdam T. Chaichan3 | Ali H.A. Al - Waeli,Google Scholar cited by 43.
- [3] Computer network traffic prediction, Tiago Prado Oliveira, Jamil Salem Barbar and Alexsandro Santos soares, Federal

International Conference on "Computational Intelligence and its applications" (ICCIA-2024)

University of Uberl
ndia (UFU),(2016) , Google Scholar cited by 144. $\ensuremath{\mathsf{}}$

[4] Filippo Amato ,December 2013, Google Scholar.

[5] B.A Garro, R. Vazquez, 19 June 2015 ,Goole Scholar.

[6] https://www.coursera.org - deep learning-ai.

[7] https://techttarget.com - ai.

[8] Dongky Lee, Jinhwa Jeong, sungHoon Yon, young tae chae, (2019), Google Scholar.

Artificial Intelligent with IOT

V. Parvathi¹ and P.Pragadeeswari²

Ist year MCA Department of Computer Applications, A.V.C.College of Engineering, Mannampandhal, Mayiladuthurai-609305 ¹parvathiguna56@gmai.com , ²pragadeeswari13@gmail.com

Abstract--In recent years, the use of the Internet of Things (IOT) has explosive, and cyber security concern have a increased as well as. Progressive of cyber security is Artificial Intelligence (AI), which is used for the development of complex algorithms to protect networks and systems, including IOT systems. The research in artificial intelligence (AI) has been attempt to give an answer to this problem. Thus, IOT with AI can become a huge development. This is not just about saving money, bright think, cut down human effort, or any trending promotion. This is much more than that easing human life. However, some serious problem like the security concerns and ethical issues which will go on provoking IOT. The big picture is not how fascinating IOT with AI seems, but how the common people perceive it -a benefit, a burden, or a hazarad. However, This review paper compiles information from several other surveys and research papers regarding IOT, AI, and attacks with and against AI and explores the relationship between these three topics with the purpose of comprehensively presenting and abstract relevant literature in these fields.

Keywords - Artificial Intelligence, Internet of Things, Data Security, Smart Homes

I. INTRODUCTION

Internet of Things (IOT) enables the interconnection between tons of devices, industrial machines, processes, and users to exchange data without any central coordination. However, lots of data is ever multiplex in the cache, processing, and inferencing processes. Therefore, Artificial Intelligence (AI) has become the most promising combination with IOT for better use, storage and avoid the uncertainty management .In decision making. AI in IOT is playing a significant role and can improve the value of diverse types of data sensed and collected by IOT devices. For proper utilization of this diverse type of data will offer an efficient solution for the development of products and services to achieve the user's expectation from different sectors. Despite the various advantages of the integration of AI with different intelligent systems for various industrial applications, the appropriate application of AI poses several challenges with respect to data quality, data volume, integration, and accuracy of the inferences drawn from the collected data.

II. ARTIFICIAL INTELLIGENCE

Artificial Intelligent (AI) is a set of technologies that enable computers to perform a variety of advanced functions, including the ability to see, AI technology is widely used throughout industry, government, and science. In general AI system work by ingesting large amount of labled training data , analyse the data for correlation and patterns and using this patterns to make predictions about future state.

III. INTERNET OF THINGS (IOT)

What we had since 1991 was "Internet of Computers (IOC)" and Itgradually grew in size as more and more people started using it. With the advent of pocket phones and connected devices, the Internet of Devices started and eventually grew larger as mobile phones, computers, laptops and tablets became cheaper and more accessible to the common man.IOT is the collective network of connected devices and the technology that facilitates communication between the devicesthemselves.



IV. AIENABLED IOT

AI enabled IOT represent a powerful synergy between artificial intelligence and interconnected device. By integrating AI into IOT system, device gain the ability to analyze and interpret data, enhancing overall efficiency and functionality. This convergence allows for smarter decision –making, predictive analystics, and automation based on realtime data. The combination of AI and IOT has transformative implications across various sectors, from optimizing smart homes and industrial processes to enabling more responsive and intelligent system in healthcare, transportation, and beyond.

V. HOW IOT AND AI ARE CHANGING OUR LIVES?

Society has advanced exceedingly well thanks to AI and IOT. We now have much well again technologies like autonomous vehicles, robotic cutting tools, home apps, Siri, Alexa, and other everyday items. When AI and IOT are combined to create AIOT (Artificial Intelligence of Things), IOT devices may assess data and make proactive, wise, and clear-cut judgments without the need for human intrusion. These gadgets will eventually build up into intelligent, communicative, and strong machines that can analyze data and reach decisions more quickly and specifically than before. It makes sense that people can be trusted more, and it is not easy to become more reliant on artificial intelligence. But, AIOT can close this gap and assist a smooth shift from trusting people to trusting science.

VI. BENEFITS OF AI WITH IOT

Efficiency: Optimizing processes, reducing downtime, and enhancing resource utilization.

Predictive Analytics: Anticipating issues before they occur through data analysis.

Cost Savings: Streamlining operations and minimizing resource wastage.

Scalability:The low-end sensors that make up the complete IOT system might not always be 100% accurate. However, the chances of scalability can be greatly improved with an AI filter that evaluates the data and provides a better version for additional processing.

Saving money and time:Using predictive maintenance, you can predict potential damages and failures far in advance. Lowering operational expenses is another aspect of AI and IOT's benefit for your operations. The IOT system's incorporation of AI eliminates the need for an explicit method to address data problems, saving money and time.



VII. APPLICATIONS OF AI WITH IOT

In Industry-One of the trending sectors using solutions like IoT, AI, face recognition, deep learning, robots, and many more is industry 4.0. Factory robots are gaining intelligence thanks to sensors that can communicate data via implanted sensors. Also, since the robots come with AI structures, they may gain insight from new data. This approach not only reduces costs and time but also enhances production over time.

In Agriculture: As a high-tech solution, smart farming plays a key role in the clean and sustainable production of food. Some of the IoT gadgets used in intelligent farming include agricultural drones, animal tracking systems, and smart greenhouses. Plus, keeping an eye on soil, humidity, and growth conditions, mainly promises better harvests. Continuous monitoring ensures minimized crop loss with real-time updates.

In Smart Homes: The AIOT can learn about each family member's living patterns using the data from all linked devices inside your home. Depending on your choices, it will automatically set the lighting, temperature, and curtains when it recognizes that you are in your room. Plus, it can identify every family member through facial recognition, making sure that no stranger can unlock the doors.

In Security: AIOT technology uses Real-time data analytics to determine whether something odd has happened in an area. This is a trend that exists primarily in the retail industry, where consumer profiling allows surveillance cameras to secure public shopping by detecting previous offenders who pose a security risk to the organization.

In Transportation: AIOT transmits traffic data across the network to the control center using systems for recognition, communication, and surveillance. Information from many sources is brought into the traffic control center and integrated.

In Healthcare: IOT and AI work well together to support healthcare in different ways, such as ultraaccurate diagnostics, telehealth and remote patient treatment, and reducing office work in healthcare centers. Plus, such technology can prevent many sudden critical conditions by continuously monitoring patient records.

VIII. CHALLENGES OF AI WITH IOT

Data Security: Data security is essential because IOT and AI-enabled devices gather sensitive and vital data about users and clients. However, data privacy in IOT and AI can become tricky since there have been many cases of security breaches in using IOT devices, and adding AI may devices, and adding AI may complicate things more.

Impact on Human Live: AIOT devices and software are capable of making decisions that have a big impact on people's lives. Due to the potential for biased or inaccurate judgments, this might raise some ethical issues. As a result, it is essential to make sure people make ethical decisions when developing and using such tools.

IX. FUTURE OF AI WITH IOT

AI and IOT will dominate the near and distant future. AI and IOT grow quickly in our evergrowing environment. So, it is now essential for all businesses to integrate these technologies into everyday activities to boost productivity, minimize human errors, and achieve maximum profitability.

As a leading international hardware & software development company, we are dedicated to helping businesses to optimize their operational efficiency with IOT, AI, and AIOT solutions. We have the knowledge and expertise to help you navigate your way in a fiercely competitive landscape.

X. CONCLUSION

The integration of AI and IoT is reshaping industries, providing unprecedented insights, and paying the way for a more connected and intelligent future. In future, people will be wearing intelligent gadgets, eating intelligent capsules that judge the impact of the medicine on the body, living inside intelligent homes, and so on. This sounds like a science fiction, but this is what all the present research is about. Everything will be smart and will be connected to the Internet. All branches of science will collaborate to create something of a big value. We will have a 'smart cyber revolution'. However, there is still a debate on whether we are heading towards a creative destruction or not. For instance, machines are now able to take on less-routine tasks, and this transition is occurring during an era in which many workers are already struggling.

Nonetheless, with the right policies we can get the best of both worlds: automation without rampant unemployment. Eventually, human ingenuity changes the role of productive work. Educational opportunities will be promoted and there will be more skilled labor with re-skilling and up-skilling. As we will continuously deploy AI models in the wild we will be forced to re-examine the effects of such automation on the conditions of human life. Although these systems bring myriad benefits, they also contain inherent risks, such as privacy breach, codifving and entrenching biases, reducing accountability and hindering due process and increasing the information asymmetry between data producers and data holders. The IOT-CPS is a diverse and complex network. Keeping track of every unethical or security breach incident will be difficult. Any failure or bugs in the software or hardware will have serious consequences. Even power failure can cause a lot of inconvenience. So, we may need another AI system on top of such AI enabled IoT to monitor its whereabouts each instant.

REFERENCES

[1] .R. S. Michalski, J. G. Carbonell, and T. M. Mitchell, Machine Learning: An Artificial Intelligence Approach. Springer Science & Business Media, 2013.

[2]. I. H. Witten and E. Frank, Data Mining: Practical Machine Learning Tools and Techniques. Morgan Kaufmann, 2016.

[3]. L. Monostori, B. Kádár, T. Bauernhansl, S. Kondoh, S. Kumara, G. Reinhart, O. Sauer, G. Schuh, W. Sihn, and K. Ueda, "Cyber-Physical Systems in Manufacturing," CIRP Annals, vol. 65, no. Manufacturing," CIRP Annals, vol. 65, no. 2, pp. 621–641, 2016.

[4]. E. A. Lee and S. A. Seshia, Introduction to Embedded Systems: A Cyber-Physical Systems Approach. MIT Press, 2016.

[5]. Q. F. Hassan, A. R. Khan, and S. A. Madani, Internet of Things: Challenges, Advances, and Applications. Chapman & Hall/CRC Computer and Information Science Series, CRC Press, 2017.

[6]. G. Fortino and P. Trunfio, Internet of Things based on Smart Objects: Technology, Middleware and Applications. Springer, 2014.

[7]. L. T. Yang, B. Di Martino, and Q. Zhang, "Internet of Everything," Mobile Information Systems, vol. 2017, 2017.

Prediction of DiabeticKidney Disease Using Deep learning Techniques

Mr.S. SenthilKumar¹ and Dr. T.S.Baskaran²

¹Research Scholar and²Associate Professor & Research Supervisor,

PG & Research Department of Computer Science, A.V.V.M Sri Pushpam College (Autonomous), Poondi - 613503, Thanjavur, (Affiliated to Bharathidasan University, Tiruchirappalli-620024) TamilNadu, India. ¹sksen88@gmail.com and ²t s baskaran@yahoo.com

Abstract-Diabetic Kidney Disease is one of the most critical illnesses nowadays and proper diagnosis is required as soon as possible. Deep learning technique has become reliable for medical treatment. With the help of a machine learning classifier algorithms, the doctor can detect the disease on time. For thisperspective, Diabetic Kidney Disease prediction has been discussed in this article. Diabetic Kidney Disease dataset has been taken from the UCI repository.Our proposed method is based on deep neural network which predicts the presence or absence of Diabetic kidney disease with a high accuracy. Compared to other available algorithms, the model we built shows better results which is implemented using the crossvalidation technique to keep the model safe from over fitting. This automatic Diabetic kidney disease treatment helps reduce the kidney damage progression, but for this Diabetic kidney disease detection at initial stage is necessary.. Diabetic kidney disease (DKD) is among the significant contributor to morbidity and mortality from non-communicable diseases that can affected 10-15% of the global population.Different deep-learning techniques have been used for effective classification of diabetic kidney disease from patients' data.

Keywords: Diabetic Kidney Disease, Deep Learning, Prediction.

I.INTRODUCTION

Diabetic Kidney Disease (DKD) is a condition resulting in insufficient kidney function, where patients have to live with a compromised quality of life. DKD is a substantial financial burden on patients, healthcare services, and the government. Treatments of the ESRD with Renal Replacement Therapy are either expensive (hemodialysis and peritoneal dialysis) or complex (transplantation)with the availability of biomedical data, the use of machine-learning techniques in healthcare for developing disease prediction models has become common. Further, methods such as deep learning and techniques like ensemble learning have greatly improved the predictive power of machine learning models.At the patient level, a physician can assess the onset of CKD using laboratory tests by looking at standard parameters such as the glomerular filtration rate (eGFR) and the albumin-creatinine ratio [7]. On the other hand, from the public health perspective, laboratory data is typically not available on a large scale. However, two types of data can generally be extracted from the insurance companies' databases: diagnoses and medications for each patient's visit at the hospital Common approaches for developing disease prediction models clinical and laboratory data from sources developed a predictive model for kidney disease among patients with hypertension They proposed a neural network framework based on Bidirectional long short-term memory and auto-encoders to encode the textual and numerical information, respectively. Data using an ensemble feature selection method to predict the risk of kidney disease among diabetes patients predicted kidney-related outcomes among diabetes patients In this paper, we aimed to develop machine-learning models that predict the onset of DKD within the next 6 and 12 months. The model is based on the insurance claims data (age, sex, comorbidities, and medication) over an observation period of 24 months. Further, we aim to assess the reliability of the models by identifying the comorbidities and medications that impact the development of DKD the Most of previous researches focused on two classes, which make treatment recommendations difficult because the type of treatment to be given is based on the severity of DKD.

II. RELATED WORKS

A remarkable number of researches have been conducted in the area of machine learning for building models that assist in predicting different types of different types of diseases and health related problems, using different machine learning algorithms. This section presents a review of some of conducted research in the area of machine learning in Diabetic Kidney Disease. The disease detection model uses real-world machines and deep learning classification algorithms to ascertain the components that lead to the onset of disease. Data collected from the real world must be routinely updated and organized before computing can take care of it. While it is understandable that sometimes real-world data would have errors in the actual metrics used, this does not mean that these mistakes can be ignored or ignored. Thus, information preparation takes the raw data and cleans it up by cycling it, removing the errors and sparing it from a second round of examination.Data is being preprocessed, it undergoes a series of operations. During data cleaning, many methods are used to get rid of old or irrelevant information. For example, increases the quantity of information available, while omitting commas and other cryptic symbols, decreases it. Information refers to the process of integrating data from several sources. Any discrepancies are then promptly addressed by updating the data.

III. METHODOLOGY

This methodology depicts the design of the method to be exploited to carryout the experiment. It incorporates data collection, data pre-processing, andtarget variable selection.



Patient's kidney disease record is selected as the source of data for thiswork.This dataset is collected from General Hospital Common diagnostics lab. It contains 400 patients record with 11attributes/parameters: Age, Gender, Sodium, Potassium, Chloride,Bicarbonate, Urea, Creatinine, Urea Acid, Albumin and Classificationincluding a target variable classified into a binary classification of DKDand non DKD disease.

Table1: Dataset Attributes

Attribute	Attribute Description		
Sex	Gender		
Age	Age		
Sod	Sodium		
Pot	Potassium		
Chl	Chloride		
Bica	Bicarbonate		
Urea	Urea		
Cre	Creatinine		
UA	Urea Acid Alb Albumin		
Class	{Kidney Disease, No Kidney Disease}		

V. DATA PRE-PROCESSING

Data Pre-processing represent the most important task in data mining techniques, it involves cleaning, extraction and transformation of data into a suitable format for machine execution, Raw data contains information, bad formats, missing invalid information and it leads to disaster in prediction with machine learning. The dataset used had some missing cells which was replaced using simple imputation with the mean value of the attribute and the attribute Sex was converted to numeric values as '1'and '0' for Male and Female respectively to make it possible for the machine to process since the machine will not understand string values.

IV. DATA COLLECTION

VI. TARGET VARIABLE

DKD dataset has got a lot of useful variables which are key and necessary for the identification of the disease in patients. we used our variables based on the test and method used by the Hospital in determining the occurrence of kidney disease which is Ten (10) blood test and that are the once we used as our input variables. Names too were among the test, but we decided to remove it since it has no impact on our test and the privacy of the patients must be kept.

VII. PREDICTION USING DNN

DNN is a subset of Artificial Neural network which simulate the structureand functionalities of biological neural network consisting of an input,weights and activation function, the structure of DNN has an input, hiddenlayer and an output.In DNN, a is referred to the output, where Wi and Xiare the weight and input respectively.



Fig 1: Graph for Validation Accuracy

VIII. EVALUATION OF THE MODEL

In this work, the performance is measured by Accuracy, specificity, sensitivity, kappa statistic, precision, F1 score, ROC Score and recalldescribed as follows.

IX. CONFUSION MATRIX

Confusion matrix indicates statistical suitability of the model and its compatibility with the dataset; itcan also be defined as a table layout that is specifically used forvisualization of algorithm performanceTable 2 shows the summary of confusion matrix.

Classification		Observation		
		Negative	Positive	
Observatio ns	Negative	True Negative(T N)	False Positive(FP)	
	Positive	False Negative(FN)	True Positive(TP)	

- Accuracy: It is used to identify the number of correctly predicted data points out of all data points. It is defined as the number of all correct predictions made divided by the total number of predictions made, it is expressed as;
- Sensitivity: (Recall or True Positive Rate):it is defined as the proportion of actual positive cases that got predicted as positive. It is a ratio of true positive to the sum of true positive and false negative. In medical diagnosis, test sensitivity (Recall) is the ability of a test to correctly identify those with the disease, it is expressed as;
- Specificity: (True Negative Rate): it is defined as the proportion of actual negative that got predicted as the negative. it is calculated as the number of correct negative predictions divided by the total number of negatives. It is also called true negative rate, it is expressed as;
- Cohen Kappa: This is a classifier performance measure between two sets of classified data. Kappa result values are between 0 to 1. The results become meaningful with increasing values of kappa. it measures how closely instances classified by the machine learning classifier matched the data labeled as the truth, it is expressed as;



Fig 2: Feature Importance

Confusion Matrix

- Precision: It is defined as the fraction of relevant instances among the retrieved instances. This is given as the correlation number between the correctly classified modules to entire classified fault prone modules, it is expressed as;
- Recall/ Sensitivity:Recall is a metric that quantifies the number of correct positive predictions made out of all positive predictions that could have been made, it is expressed as;
- F1 Score: This is the harmonic mean between precision and recall. Range for f1-score is from 0 to 1. It describes the preciseness (how many records can be correctly classified by the model) and robustness (it avoids missing any significant number of record) of a model. The expression of F1-score is as follows;

X. CONCLUSION:

The main goal of this research is to use DNN model for the prediction ofkidney disease to high degree of accuracy. We succeeded in classifyingkidney disease dataset into DKD and non-DKD with 98% overall accuracywhen the model was tested with a set of data that were not used during thetraining process. The adopted DNN model proved to be efficient andsuitable for the prediction of kidney disease. The study also highlighted the importance of the features used in the prediction of kidney disease. This revealed that from the 10 attributes, Creatinine and Bicarbonate are the attributes with highest influence on DKD.

X1. FURTHER RESEARCH

Incorporating more data to have large dataset will provide more accuracyand efficiency; hence more dataset is needed to accommodate enoughsamples. Having enough samples will make the prediction wider to captureand identify regions and areas with CKD vulnerability. This paper usedDNN model only, having other models to compare techniques mayprovide a better performance.

REFERENCES

- WHO. (2006). World Health Organization: Definition and Diagnosis of Diabetes Mellitus and Intermediate Hyperglycemia. Geneva, World Health Org. WHO2.
- [2] World Kidney Day. (2017). Chronic Kidney Disease -World Kidney Day. ISN – Global Operations Center.
- [3] Ulasi, I. I., &Ijoma, C. K. (2010). The enormity of chronic kidney disease in Nigeria: The situation in a teaching hospital in south-east Nigeria. Journal of Tropical Medicine, 2010.
- [4] Shafi, N., Bukhari, F., Iqbal, W., Almustafa, K. M., Asif, M., & Nawaz, Z. (2020). Cleft prediction before birth using deep neural network. Health Informatics Journal,1(18) 1460458220911789.
- [5] Sharma, S., &Parmar, M. (2020). Heart Diseases Prediction using Deep Learning Neural Network Model. International Journal of Innovative Technology and Exploring Engineering (IJITEE), 9(3).
- [6] Scholar, P. G. (2018). Chronic Kidney Disease Prediction Using Machine Learning. International Journal of Computer Science and Information Security (IJCSIS), 16(4).
- [7] National Kidney Foundation (NKF). (2015). Global Facts: About Kidney Disease. In National Kidney Foundation
- [8] Lote, C. J. (2013). Principles of renal physiology. In Principles of Renal Physiology
- [9] Kriplani, H., Patel, B., & Roy, S. (2019). Prediction of Chronic Kidney Diseases Using Deep Artificial Neural Network Technique. In Computer Aided Intervention and Diagnostics in Clinical and Medical Images (pp. 179-187). Springer, Cham. Chicago
- [10] Kumar, S. (2018). Chronic Kidney Disease Prediction Using Machine Learning. International Journal of Computer Science and Information Security (IJCSIS), 16(4).
- [11] Geri, G., Stengel, B., Jacquelinet, C., Aegerter, P., Massy, Z. A., &Vieillard-Baron, A. (2018). Prediction of chronic kidney disease after acute kidney injury in ICU patients: study protocol for the PREDICT multicenter prospective observational study. Annals of Intensive Care, 8(1), 77.
- [12] Ge, Y., Wang, Q., Wang, L., Wu, H., Peng, C., Wang, J. & Yi, Y. (2019). Predicting post-stroke pneumonia using deep neural network approaches. International Journal of Medical Informatics, 132, 103986.
- [13] Chimwayi, K. B., Haris, N., Caytiles, R. D., &Iyengar, N. C. S. N. (2017). Risk Level Prediction of Chronic Kidney Disease Using Neuro- Fuzzy and Hierarchical Clustering Algorithm (s). International Journal of Multimedia and Ubiquitous Engineering, 12(8), 23-36
- [14] Başar, M. D., & Akan, A. (2018). Chronic Kidney Disease Prediction with Reduced Individual Classifiers. Journal of Electrical and Electronics Engineering, 18(2), 249-255.
- [15] Arafat, F., Fatema, K., & Islam, S. (2018). Classification of chronic kidney disease (ckd) using data mining techniques (Doctoral dissertation, Daffodil International University.
Agriculture Data Analysis using Machine Learning Techniques

M. Menaha,

ResearchScholar, Department of Computer Science, Adaikala Matha College, Vallam, Thanjavur. menaha.m1989@gmail.com

Abstract-India is generally an agricultural country. Now a days the most important emerging field in the real world is agriculture and it is main occupation and backbone to our country. The most important domain is Agriculture inbroadly cultivating countries like India. The modern technologies can change the situation of farmers and decision making in agricultural field in a better way. Many machine learning algorithms are available that provide food for specifically to crop yield prediction research purposes as well as provide predictive insights based on critical variables like rainfall or temperature. These results are highly valuable inputs for farmers when determining which crops would be most suitable given local conditions and their yields thus improving agricultural productivity. Predicting the crop yield prior to its harvest would help farmers to take appropriate steps. We attempt to resolve the issue by building a user-friendly prediction system. The results of the prediction are suggested to the farmer such that suitable changes can be made in order to improve the produce. This paper focuses on the analysis of the agriculture data and finding optimal yield to provide an insight before the actual crop production using Machine Learning techniques.

Keywords: Agriculture, Machine learning, SVM, KNN

I. INTRODUCTION

Today, India is one of the leading producers across the world in the agriculture sector[1]. Agriculture is the backbone of Indian Economy. In India, majority of the farmers are not getting the expected crop yield due to several reasons. The agricultural yield is primarily depends on weather conditions. Rainfall conditions also influences the rice cultivation. In this context, the farmers necessarily requires a timely advice to predict the future crop productivity and an analysis is to be made in order to help the farmers to maximize the crop production in their crops. Yield prediction is an important agricultural problem. Every farmer is interested in knowing, how much yield he is about expect. In the past, yield prediction was performed by considering farmer's previous experience on a particular crop[2]. The volume of data is enormous in Indian agriculture. The data when become information is highly useful for many

purposes. One of the best ways of predicting unknown values is by use of machine learning algorithms. This work intends to develop crop prediction model using machine learning. The application intends to predict crop yield so it could help farmer to choose best seeds for plantation. There are plenty of ML algorithms which could be used, algorithms like Regression analysis, Support Vector Machine, Neural Networks, K-Nearest Neighbor (K-NN) can be utilized. In this paper, machine learning algorithms, K-NearestNeighbours and Support Vector Machines are used topredict the most suitable crop. The crop production depends on various factors which changes with everysquare meters and mainly depends on the geography of theregion, weather conditions, soil type, and humidity and so on. Huge data sets can be used for predicting their influence on he major crops of that particular district or state. Machinelearning techniques have advanced considerably over thepast several decades.

II. RELATED WORKS

Yield forecasting is an important service in the field of agriculture. In this section we highlight a few works done in the field of agriculture by using machine learning.

In [3] Predicting yield of the crop using machine learningalgorithm in International Journal of Engineering Science Research Technology focuses on predicting the yield of the crop based on the existing data by using Random Forest algorithm. Real data of Tamil Nadu were used for building the models and the models were tested with samples. Random Forest Algorithm can be used for accurate crop yield prediction.

In [4] Random forests for global and regional crop yield prediction. PLoS ONE Journal. Our generated outputs show that RF is an effective and adaptable machine-learning method for crop yield predictions at regional and global scales for its high accuracy and precision, ease of use, and utility in data analysis. Random Forest is the most efficient strategy and it outperforms multiple linear regression (MLR). In[5]. Prediction On Crop Cultivation. International Journal of Advanced Research in Computer Science and Electronics Engineering (IJARCSEE) Volume 5, Issue 10, October 2016. Presently, soil analysis and interpretation of soil test results is paper based. This in one way or another has contributed to poor interpretation of soil test results which has resulted into poor recommendation of crops, soil amendments and fertilizers to farmers thus leading to poor crop yields, micro-nutrient deficiencies in soil and excessive or less application of fertilizers. Formulae Match Crops with Soil. Fertilizer to Recommendation.

In [6].A Study to Determine Yield for Crop Insurance using Precision Agriculture on an Aerial Platform. Symbiosis Institute of Geoinformatics Symbiosis International University 5th & 6th Floor, Artur Centre, Gokhale Cross Road, Model Colony, Pune – 411016. Precision agriculture (PA) is the application of geospatial methodologies and remote sensors to identify variations in the field and to deal with them using different strategies. The causes of variability of crop growth in an agricultural field might be due to crop stress, irrigation practices, incidence of pest and disease etc. The Paper is implemented using Ensemble Learning (EL).

In [7]. Random Forests for Global and Regional Crop Yield Predictions. Institute on the Environment, University of Minnesota, St. Paul, MN 55108, United States of America. The generated outputs show that RF is an effective and different machine-learning method for crop yield predictions at regional and global scales for its high accuracy. The Paper is Implemented using k-nearest neighbour, Support Vector Machine (SVM).

III. PROPOSED SYSTEM

In the proposed system, supervised learning algorithmsare used to form a model which will help us in providing choices of the most feasible crops that can be cultivated inthat region along with its estimated yield. Two of thealgorithms used here is K-Nearest Neighbor and SupportVector Machine. The main stages involved in the processare dataset collection, pre-processing the data, featureextraction and classification. The dataset used for this project is collected from variousonline sources like Kaggle.com and data.govt.in. Someimportant features or the parameters which has the highestimpact on the agricultural vield considered that is rainfall,humidity,temperature,area,yield,soiltype,locat ion and price. After the selection of the dataset, it has to be pre-processedinto a form. The final step is transforming the selected data. Thepreprocessed data here is then transformed into data thatis ready for

machine learning algorithms by using variousengineering features like scaling, feature aggregation and so on.

IV. METHODS

A)KNN Algorithm

KNN is a supervised machine learning algorithm. It learns by analogy. It is a simple but a powerful approach for makingpredictions. In the project, according to the input given, the

dataset is preprocessed to obtain the extracted dataset whichis our training set. Test data is selected randomly from this training set. K-most similar records to the test record iscalculated. Euclidian distance is calculated for finding thesimilarity between the records. Once k neighbours arediscovered, a summarized predictions are made by returningthe most common outcome.

B) SVM Algorithm

SVM algorithm is also a supervised machine learningalgorithm which is used for both classification and regression problems. Hyperplane is the criteria that SVM uses to segregate the two classes. Finding the support vectors ie. the nearest data points to the hyperplane helps in giving the most optimal hyperplane. A training model isbuilt by importing the SVC classifier from sklearn SVM module.predict the values using the SVM algorithm. It has a higher acuuracy. It works well with all limited datasets.Kernel SVM contains a non-linear transformation function convert the complicated non-linearly seperable data into linearly seperable data.

V. RESULTS AND DISCUSSION

The implementation of KNN and SVM algorithm is done and both the algorithms are compared in terms oftheir accuracy, Execution time in seconds from epoche, and n terms of their Precision and Recall scores. ClassificationReport is a report that is used to measure the quality ofpredictions from а classification algorithm. Precision, Recall, F1score are some of the metrics given in theClassification report. Precision is the ability of a classifiernot to label an instance positive that is actually negative. Recall is the ability of a classifier to find all positive instances. It can be analyzed the accuracy of SVM, precision and recall values of SVM is higher compared to that of the KNN classifier. That SVM takes less time to execute compared to that of KNN.

VI. CONCLUSION

Farmersare still not connected with the modern technologies. Itefficiently bridges the gap between the rural farmers andthe modern technologies. Machine learning algorithms haveproved very effective in predicting the crops and its yield.From the comparison analysis of KNN and SVM, SVM worksbetter than KNN for the dataset chosen.

REFERENCES

- [1] https://www.investopedia.com/articles/investing/100615/4co untries-produce-most-food.asp.
- [2] http://www.fao.org/fileadmin/templates/ess/documents/meet ings_and_workshops/GS_SAC2013/Improving_methods_fo r_crops_estimates/Crop_Yield_Forecasting_Methods_and_E arly_Warning_Systems_Lit_review.pdf
- [3] P. Priya, U. Muthaiah M. Balamurugan. Predicting yield of the crop using machine learning algorithm. International Journal of Engineering Science Research Technology.
- [4] J. Jeong, J. Resop, N. Mueller and team. Random forests for global and regional crop yield prediction. PLoS ONE Journal.
- [5] Shweta K Shahane, Prajakta V Tawale. Prediction On Crop Cultivation. IInternational Journal of Advanced Research in Computer Science and Electronics Engineering (IJARCSEE) Volume 5, Issue 10, October 2016.
- [6] BaisaliGhosh. A Study to Determine Yield for Crop Insurance using Precision Agriculture on an Aerial Platform. Symbiosis Institute of Geoinformatics Symbiosis International University 5th & 6th Floor, Atur Centre, Gokhale Cross Road, Model Colony, Pune – 411016.
- [7] Jig Han Jeong, Jonathan P. Resop, Nathaniel D. Mueller, David H. Fleisher, Kyungdahm Yun, Ethan E. Butler, Soo-Hyung Kim. Random Forests for Global and Regional Crop Yield Predictions. Institute on the Environment, University of Minnesota, St. Paul, MN 55108 United States of America.
- [8] R. Nagini, Dr. T.V. Rajnikanth, B.V. iranmayee, "Agriculture Yield Prediction Using Predictive AnalyticTechniques, 2nd International Conference on Contemporary Computing and Informatics (ic3i),2016.
- [9] Dr. Rakesh Poonial, Sonia Bhargava "Prediction of Crops Methodology using Data Mining Techniques", International Journal of Innovative Research in Science, Engineering and Technology, Vol. 6, Issue 10, October 2017.
- [10] Mehta D R, Kalola A D, Saradava D A, Yusufzai A S, "Rainfall Variability Analysis and Its Impact on CropProductivity - A Case Study", Indian Journal of Agricultural Research, Volume 36, Issue 1, 2002, pages:29-33.

A Survey Concerning the use of Electronic Medical Record Search Engines in the Healthcare Industry

S. Narmatha¹ and **Dr. V. Maniraj²** ¹*Research Scholar* and ² *Research Advisor*

¹Department of Computer Science, A.V.V.M Sri Pushpam College (Autonomous) Poondi, Thanjavur. Affiliated to Bharathidasan University, Tiruchirapalli, Tamil Nadu, India. imnarmatha2792@gmail.com
²PG and Research Department of Computer Science, A.V.V.M Sri Pushpam College (Autonomous) Poondi, Thanjavur. Affiliated

PG and Research Department of Computer Science, A.V.V.M Sri Pushpam College (Autonomous) Poonal, Inanjavur. Affiliated to Bharathidasan University, Tiruchirapalli, Tamil Nadu, India. maniraj_vee@Yahoo.co.in

Abstract - Healthcare organisations generate data at a rapid pace. This has several positive outcomes, but it also has some negative ones. A tremendous quantity of electronic health record data is being produced by the fast expansion of free-text clinical papers in the healthcare industry. The patient's medical record is accessible to anyone. A comprehensive record known as a medical chart is kept of a patient's essential clinical data and medical information, such as vital signs, prescriptions, demographics, exams, treatment plans, improvement notes, risks, vaccination dates, allergies, radiological pictures, and laboratory and test results. Information is gathered from a variety of sources, including as administrative databases used for bill payment or care management, patient surveys, and medical records. A medical database is an electronic collection of patient records.

Healthcare practitioners used terms like "electronic health record," "computer-stored patient record," and "computerised medical record" to describe this kind of data. It is not an easy chore to get medical and healthcare domain information, what with the time required and the need for accurate results. Access to some patient information is subject to authorization since health records are considered confidential. Accordingly, we have compiled a large number of papers that deal with EHR data retrieval. And thus, we have covered the many search engines, information retrieval strategies, and ways that may be used to obtain health data stored in medical databases. Next, we'll go over some of the restrictions and potential problems.

Keywords - Natural Language Processing, EHR, MESH.

I. INTRODUCTION

One online-based tool that helps consumers find information on the web is a search engine. The three most common reasons individuals use the search engine are to shop, do research, and find enjoyment. Several research initiatives have uncovered the three basic categories of search engines. Informational, transactional, and navigational are some of them. To better understand people's search demands, all the main search engines analyse queries and find each client's intent to determine which one the user has selected.

II. RELATED WORK

A large number of publications devoted to patient Discover how to get reliable health data from a medical database. The needs and constraints of paperbased data gathering gave rise to the current standards for managing drug data in clinical trials. However, a plethora of electronic tools have emerged to facilitate the gathering and analysis of data pertaining to medications. The first area of AI that allows computers to comprehend spoken language is known as natural language processing (NLP). Second, syntactic matching is the process of determining which keywords to use in response to a search by analysing the words the user actually types into the engine. This is a phrase match and would be precise. Thirdly, semantic matching involves determining the purpose of the searcher and then matching terms to those inquiries.

III. METHODOLOGY

Natural Language Processing: Using "notational language" input by ophthalmologists during patient encounters, this article describes how to get structured data pertaining to glaucoma diagnosis and progression [1]. Using natural language processing methods, such as the one used in this research, it is possible to filter and split raw query input into useful analytical types. Users' level of understanding and the complexity of the material they seek out on the website may be better understood with the use of this data, which can be obtained and analysed [2]. One example is a health information site that would benefit from natural language search engines. This would provide users, particularly those with less background knowledge or language skills, easier access to the information they need. Additional thorough investigation will be conducted on future requests. Research has shown that electronic health records (EHRs) may help clinicians improve patient outcomes, speed up treatment, and decrease medication errors. Several publications that were reviewed highlighted the potential advantages of

Department of Computer Science, A.V.C. College (Autonomous), Mayiladuthurai – 609 305. PROCEEDINGS

NGS 169

electronic health records (EHRs) by showing how clinicians may be helped in monitoring complicated events after surgery by extending natural language processing to electronic data [3].

Researchers and doctors are unable to quickly and effectively investigate large numbers of clinical interactions because clinical notes have to be manually scrutinised to get information. The ability to manually comprehend EHRs is a product of natural language processing, which looks at the context of medical record phrases and words before granting them access to high speed computers. From automated quality assessment to comparative effectiveness studies, natural language processing has many potential applications. Their hypothesis that natural language processing (NLP) methods, especially POS tagging, can improve IR models and thereby raise the identification rate in the biomedical industry is borne out by the results. The suggested machine learning and the point-of-sale category selection are further tested to ensure their efficacy. The large improvements show that POS tagging is useful for biological IR applications. They want to look at the effects of different POS taggers, which is the first limitation of the proposed approach.

There may be previously undiscovered insights in EHR data that may be uncovered using NLP-based approaches, as they explain. They used only one natural language processing technology for their analyses. Reviewing and comparing different natural language processing (NLP) techniques for different use scenarios should be a goal of future research. The generalizability of their results is questionable as they just used data from a single healthcare facility.

Syntactic Search: Medical or clinical data refers to health-related information that is associated with standard patient treatment or is part of a clinical trial program. Patient and illness registries are useful tools for collecting and tracking clinical information of specific patient groups. An electronic health record (EHR) is a digitally recorded, standardized compilation of a patient's medical history.. Utilizing previously collected data for purposes unrelated to the original study is an example of data reuse. For free-text documents in EMRs, EMERSE is a strong and spontaneous search engine to use. We also spoke about few other publications that were on the same subject.

Search Engines: Over 90% of the patients in the University of Michigan Health System's registry are found via the manual screening process. An automated computer system may do this labor-intensive and time-consuming task accurately [6]. In order to manually search for cancer-related terms in free-text medical data, they created a technology. They compiled 800 SNOMED codes and more than 2,500 words and phrases into their own lists. When building the registration database, the Case Finding Engine (CaFE) scans the text for relevant

phrases and flags them for the registrar to review. Their registration team has given the caFE high marks for accuracy and efficiency. Significant enhancements have already been accomplished as a consequence of the registrars' suggestions. If more research focused on certain areas, its dependability and acceptability might be further enhanced. The writers discuss the procedure and how they used the Star Tracker medical database. The end aim is a search feature that integrates demographic, clinical, and diagnostic data to help users categorise patients [7].

The optimal option, however time-consuming and resource-intensive, is to build an enterprise-level data warehouse. A different approach that is compatible with several older systems and doesn't need the integration of those databases beforehand might easily and cheaply bring substantial value. With the Star Tracker search engine, you can use your current hardware and database systems to do comprehensive population-level searches. To quickly find and gather information for this research, they plan to employ technical advancements to manage massive amounts of unstructured textual data. Having access to information quickly is of the utmost importance in the medical industry [8]. The future generation of healthcare for patients will need better ways to retrieve sensitive information that is hidden in millions of patient records. Using IR approaches to extract relevant phrases from record explanations instead of relying just on item titles is one potential improvement to the current method of keyword selection from the PHR.

Data Warehousing and Reuse of Clinical Data : In order to help with the creation of PHRs and to establish the worth of PHR features, they provide a framework. Although this paradigm is exclusive to PHRs, the authors have extended it to evaluate other forms of healthcare IT [9]. The offered PHR framework, together with the related technique, should provide a comprehensive evaluation of PHR value. Coding prescription data in multi-site research scenarios is reviewed by the author. Medication data classification, reporting, and analysis was suggested by the author in their study [10]. Additional research is required to evaluate the practicality, accuracy, repeatability, and extensibility of various classification subsets, as well as to improve the incorporation of drug categorization into data management and analysis processes. In order to permit their reuse whenever required, methods are required to standardize, aggregate, and query data covered by EHRs [11].

The authors of this piece suggest an EHR-based DW (Data Warehouse) setting that would facilitate data reuse, interoperability, and rapid data aggregation across many contexts. Data coming from the EHR may be modeled, transformed, integrated, standardized, and aggregated for reuse using the methods and technologies mentioned in this article. Acknowledging the significance of clinical information modeling processes

Department of Computer Science, A.V.C. College (Autonomous), Mayiladuthurai – 609 305. PROCEEDINGS 170

for reusing clinical data as well as for standard health care delivery. Patient information retrieval is a timeconsuming and labor-intensive part of healthcare data management. They proposed a medical information system data warehouse design [12] based on the novel methods.

The proposed design keeps clinical data up-to-date and makes it easier for data analysts and clinical managers to mine and analyse warehoused data. One big drawback of the data warehouse is the higher upkeep needed because of the massive volume of data.

Using Emerse: Aneurysms in people who have had abdominal transplants are the focus of this study, which aims to describe their frequency and typical symptoms. Additionally, it is centred on locating arterial aneurysms in the medical records of individuals who have had more than eleven years of liver or kidney transplants [13]. They were limited to seeing EMERSE users' electronic medical records at their facility. Some patients' records may have been kept on paper or in systems that did not use computers, or they may have been transferred to other institutions in the state. According to this study's findings, people with HNSCC may benefit greatly from taking antacids on a daily basis. Investigating the root reasons of this correlation might pave the way for the creation of innovative, low-toxicity treatment and prevention plans [14]. Finding patients' medication consumption required reviewing their charts and extracting data from CareWeb using the EMERSE application. Utilizing this tailor-made technology, they crafted intricate but precise search queries to identify medications used and when (before or after therapy), together with clinical, baseline demographic, and histological data from this cohort.

Additionally, treatment regimens tailored to patients in the contemporary era are also necessary. Without cost-cutting measures for recently authorised DAAs (Direct-acting antiviral drugs), screening programmes and access to effective treatment regimens would not significantly reduce disease burden. The EMERSE electronic medical record search engine, created by the University of Michigan, was used to examine the medical records. The study's protocol was finally greenlit by their IRB. The authors detail the nine years that the University of Michigan was involved in the Electronic Medical Record Search Engine, a full-text search engine designed to aid in the retrieval of information from documents stored in electronic health records (EMERSE). This approach improves the efficiency of medicinal IR. Enhanced sensitivity and specificity have been shown in several evaluation studies. The effectiveness of chart review is also significantly enhanced. Integrating IR capabilities into EHRs is challenging, however. The different requirements of those seeking to extract information from archived health records may be better understood as a result.

Others: The CER Hub is a platform for the systematic and scalable extraction, consolidation, and analysis ofmassive amounts of clinical data from several institutions.

The challenges of doing comparative effectiveness research across several institutions may be alleviated with the help of CER Hub. This tool allows doctors to quickly find and analyse groups of patients who are similar to themselves in order to get clinical insights. Stratified survival analysis and physician-driven cohort selection are both made possible by the MRLU. This method has revealed some well-documented clinical patterns in melanoma, such as the incidence of BRAF mutations, the survival rate of patients with BRAF mutant tumours after BRAF inhibitor treatment, and trends depending on sex. To determine the optimal timing and method for using this feature inside the EMR clinical workflow to direct clinical decision-making, more study is required.

Open-Source Search Engines: Although there is a wealth of information available via web search engines, there is a lack of functionality to automatically connect the data collected from these engines to particular biological uses. CDAPubMed is a platform-agnostic application that makes it easy to search for literature based on keywords in certain EHRs. The conventional CDAPubMed interface seamlessly incorporates CDAPubMed, making it seem like an extension of a web browser. Thanks to CDAPubMed's open-source nature and modular design, the biomedical informatics community will be able to contribute to the tool in the future and integrate it with other systems. It is compatible with various systems and may be used for free for nonprofit reasons. At the French cancer treatment centre Leon Berard Center, a full-text search engine is now standard operating procedure.

The multi-level modelling method is used by this application in open EHR. Performing fluctuations in GastrOS used to take half as long. Software maintainability may be enhanced with the help of openEHR model-driven development. The fact that just one developer made updates to each programme and that neither of them knew about the other's modifications is a drawback of the study's design. Having said that, prior to the research, the second author—who was responsible for developing GastrOS and implementing the CR lacked both domain knowledge and expertise with open EHR deployment.

Clinical Correlation: Integrating relatively unstructured data from sources such as research databases, discharge letters, clinical notes, and diagnostic reports with organised medical data from sources such as vital statistics, medicines, and test results was the primary emphasis of the authors' work. We will be merging textual patient records with additional data points like as diagnoses, meds, and test results in order to analyse data

Department of Computer Science, A.V.C. College (Autonomous), Mayiladuthurai – 609 305. PRO

in a time-based framework in future work. A combination of structured and unstructured (free-text) data from electronic health records were analyzed using the algorithm that was used to diagnose exfoliation syndrome (XFS).

This allowed us to extract an XFS likelihood score from the EHRs of all patients. An improved method for detecting XFS has been created, tested, and validated here; it seems to be superior to the standard method of studying individuals with ocular problems using EHR data. In addition, the illness would go undetected if the practitioner made a mistake in diagnosing the patient or failed to document the XFS-specific symptoms. Last but not least, it should be mentioned that more EHRs can enhance the algorithm's ability to detect this disease with more sensitivity and specificity.

Emerse: DEPARTMENT OF MEDICAID'S PROGRAM The EMR's EMERSE search engine is a potent tool for retrieving free-text materials. The demand for searching medical records for research and data abstraction led to the development of EMERSE, an electronic medical record search engine. It features a user-friendly layout that even someone with little computer knowledge can use, and it gives you a lot of options for making complicated search queries. If a medicine recall affects certain individuals, this might help identify them. Further concerns around privacy would have to be handled. With doctors under growing time pressure, EMERSE is also useful in direct patient care. We would much appreciate a streamlined process for examining a patient's medical history for noteworthy events.

Others: Their ability to be adaptable and flexible in the face of changing needs is being hindered by the centralised system design concept. These requirements are a direct outcome of the methods' application environment and evolving user requirements. Make a model and see how it holds up in future situations that include real-world problems. In order to make any meaningful use of EHR data, individual patient groups need accurate and comprehensive data. Their research reveals untapped potential for better healthcare data collection and secondary use (e.g., research and surveillance) in existing systems.

There must be a heavy emphasis on data quality throughout the lengthy clinical trials. Both patient recruitment and trial data documentation are greatly aided by medical informatics, because to the growing amount of electronically accessible patient data. They demonstrated that the intensive care unit's electronic medical record could be automatically used to reuse data, which would reduce time while gathering research data. Due to the increased likelihood of typos, data quality might be compromised. They hope that researchers will always have the option to input data in bulk via webbased data entry interfaces and existing electronic data sources when designing multi-center data collection systems for next studies. A framework is used to represent the textual and temporal characteristics concurrently of the clinical narrative text, treating it as a sequence of documents. The flexibility of the framework is shown by its generalizability and structural flexibility. Eventually, they want to find a more expressive and effective way to portray the time-related aspects of EHRs by delving more into their temporal structure. They also want to continue working on a more sophisticated method of combining textual and temporal similarities.

According on the kind of the research question, electronic health records (EHRs) can provide the necessary up-to-date data for positive outcomes and the therapeutic impact of the modification. There are benefits and drawbacks to every current method for using EHR data to aid research. The ideal way to integrate researchoriented data collection into the institution's standard paediatric urologic clinical practise is to collaborate on this matter. One common argument in favour of visual retrieval in the healthcare sector is the relative smallness of data sets, which is driven by privacy concerns. Nonetheless, due to scientific challenges and the increasing availability of massive data sets, medical visual information retrieval has made significant progress. This study has a number of caveats, one of which is that it only covers the six years from 2011 to 2017 in terms of the number of texts that were analysed. It is difficult to be thorough and methodical in this arena since the particles are scattered throughout several research fields. Since it seems difficult to assess visual data without all the data that affects the visuals, it is vital to combine multimodal data whenever possible. Therefore, strategies for merging data from different sources, known as fusion techniques, will be important moving forward. The low level of clinical use is another problem that needs fixing.

Semantic Search Clinical Data: Using several different algorithms, Bio Patent Miner can locate data about patents in the field of biomedicine. Physiologically relevant words are located and analyzed by the algorithm, which then develops links between them. They want to enhance the Connection Annotator's memory by adding more templates for identifying other related patterns in phrases and conduct user testing with domain experts to confirm the effectiveness of their approaches By establishing secure and dependable wireless connections between healthcare professionals, patients, and hospital records, on-demand mobile healthcare services might be delivered directly to patients' homes. Several semantic similarity measures that are relevant to the web have been used in the research. On the other hand, it's not easy to tell how close the two statements are semantically. Conventional methods based on ontologies state that the two concepts should have resided on the same branch of

Department of Computer Science, A.V.C. College (Autonomous), Mayiladuthurai – 609 305. PROCEEDINGS 172

the tree (s). There is a gigantic growing machine that is the internet, and it never stops.

SEMANTIC ANALYSIS: Every sector, including healthcare, has seen a meteoric rise in the amount of digital information during the last 20 years. Faster patient diagnoses are possible thanks to the retrieval of important information about treatments and the development of clinical issues from medical records. Information extraction using semantic analysis was the major focus of this work, which aimed to make a good contribution in this area [39]. In addition, any language may make use of this paradigm provided that the database dictionaries are provided with the necessary information. In the future, they want to delve further into complicated terms, such as correcting spelling mistakes, identifying negative phrases, and analysing assertions based on probability and speculation. Although it is not a simple task, integrating medical knowledge sources may aid in information retrieval. A novel medical information retrieval system was proposed, using a two-stage query expansion approach, to enhance performance in this research.

Medical Records Search Engine: Standard search engines can't understand the user's intent behind their query. People worry a lot about this when they are seeking health-related information. In it, they detail a search engine that tailors its results to each individual's medical history in order to help people discover the health information they need online. They put this tool through its paces with 18 volunteers, and the results are detailed here. Using information retrieval (IR) techniques to extract relevant terms from record descriptions, as opposed to relying just on item titles, might enhance the PHR's context keyword selection process. Many different types of information retrieval procedures use well-known online search engines to get their initial data. These engines include Google, Yahoo!, and Live Search, among others. Users of search engines may see these returned results as having no relevance to their original query. In order to find better, more relevant results in the field of electronic medical records (EMRs), this effort intends to study and build a question-and-answer search engine.

Despite the potential limitations, one strength is the use of genuine patient situations inside a basic system to achieve a realistic aim. In order to better understand how a semantic search may assist a physician in their information collecting work, this research set out to investigate this question. In the semantic search task, players were to find comparable data in another patient record using the semantic search interface. With this solution, doctors are able to make better judgments by removing the performance-limiting factors associated with information overload. Results may change if subjects used a different electronic health record system, as the research was conducted in a single EHR environment. Instead of doing the search tasks at random, everyone followed the same sequence. The semantic search capabilities were shown to consumers via a short video sample. In order to determine how effective the semantic search tool is, they suggest using a larger sample size.

Others : In the future, they want to build and verify the model to incorporate other health information systems. Two realistic situations requiring the retrieval of patient information were presented to the participants (N=10). The first time around, participants got the patient-specific information question right. The second scenario included sharing a semantic search tool with participants, which could identify terms inside a patient's electronic health record. In a subsequent semi-structured interview, participants were questioned about their present EHR use. In therapeutic contexts, semantic search skills might be a helpful method to reduce cognitive load for similar patient-specific information needs, according to this research [46]. One possible way to improve trust and, by extension, the sense of correctness, is to provide a description or list of the searchable objects in the EHR. It would be an interesting study to find out how the level of perceived and actual accuracy develops with time.

They provide a system that takes free-text clinical reports, analyses them for patterns and outcomes, correlates those patterns with ideas in other knowledge sources, and then tailors the display of the user's record to their specific information requirements. Two main problems with current search engines are the lack of an explanation for why the returned resources are relevant to the query and the restricted user engagement with the list of resources. Nobody knows how to make their searches more specific so they get the intended results. In order to facilitate user interactions, this research [48] suggests an information retrieval system that uses domain ontology to locate a set of relevant resources and uses graphical explanations of query results. In addition to a number of optimization and mathematical issues, reformulation brings up serious problems with user input in terms of maintaining comprehension and effectively engaging with the IRS.

Collaborative Search: The joint search was categorised using UMLS and electronic health data. To facilitate communication between computer systems and EHRs, the Unified Medical Language System (UMLS) compiles a number of files and programmes that integrate several biological and health vocabularies and standards.

UMLS Concept: Discover Common Data Elements (CDEs) in the eligibility criteria of several clinical trials studying the same disease using a human-computer collaborative technique. Clinical trial eligibility requirements for two representative diseases, cardiovascular problems and breast cancer, were compiled into a collection of free-text criteria in order to identify disease-specific eligibility

Department of Computer Science, A.V.C. College (Autonomous), Mayiladuthurai – 609 305. PROCEEDINGS 173

criteria CDEs. Within the eligibility criteria text, a semantic annotator [49] is used to identify terms from the Unified Medical Language Systems (UMLS). Verbs produced a small number of false positives (for example, "arrange" and "repair"). In order to increase the speed of retrieval, their method [50] effectively utilises medical domain knowledge. They provide a collaborative search architecture for building an electronic medical record search system that can adaptably use medical domain knowledge. Their innovative method for expanding the external concept space and efficient method for enriching domain knowledge utilising MetaMap and UMLS serve to illustrate this characteristic.

Electronic Health Record: Potentially limiting users' ability to make effective use of full-text search engine solutions is their familiarity with search engines and/or medical domain expertise. They created and evaluated a "collaborative search" tool to facilitate user participation and teamwork to document, improve, and share EHR search knowledge across different fields and persons. Therefore, they suggest that researchers and practitioners look into this and maybe other social informationforaging mechanisms that are widely used online to enhance the efficiency and accuracy of healthcare information retrieval. A timestamp, or creation date, is associated with every entry in clinical data. To enhance the performance of EHR searches, the authors of this paperpropose a method to include the correlation between EHR time distributions in the IR system. Using more sophisticated combination tactics will be a part of their future endeavour. In the second half, we'll go further with temporal information analysis, for example by using it in conjunction with medical ontology analysis to boost the efficiency of EHR searches. Information retrieval is a critical process that must be carried out to guarantee the correct utilisation and understanding of available medical data. To organise search results according to various aspects and types of user intent, this research developed a ranking algorithm. Annotating textual content in EHR repositories with meanings and intention elements was their notion.

Query-Based Search: Further, in query-based search, journals are classified into the following groups. First query log analysis refers to the study of user searches that aims to delve into information-seeking behaviour, system functioning, and search topic trends. By converting an information demand into a question via an interactive procedure called query formulation (QF), a user may get access to an information access system, such as a search engine (e.g., Google).

Query Log Analysis: The researchers use domain mapping to determine how UMLS knowledge exploration has affected medical domain information retrieval. Automated document and search expansion using UMLS hierarchical semantic relations and the massive image collection CLEFMcd's text to UMLS concepts [54]. As an external knowledge base, they looked at how the UMLS Metathesaurus may be used for medical information retrieval.. The search process requires knowledge about the user's history, which log analysis does not provide. Nevertheless, it is a fast way to learn about a user's behaviours. If we want to know how people use search engines and what they're searching for, we need to supplement log analysis with survey and observational studies. The research also has limitations related to the semantic categorization of questions. Using a semi-manual method is the only option for resolving this issue. A pre-trained classifier algorithm sorts a subset of queries into predetermined categories, which are then reviewed by hand. Consequently, providing intelligent query suggestions to assist information retrieval from electronic health data is a key problem and topic of study. Findings from this research may help in developing a robust method for retrieving data from EMRs.

Query Formulation: Their results show that one way to increase user autonomy when querying EHR data is to let biomedical researchers do reference interviews. Based on pre-determined topics, they generated search criteria. Topics relevant to this review that aren't searchable using the query derived from the pre-selected topics may be excluded by this approach, which may prefer selfselected themes. They may have neglected important works from the past since they only looked at literature from the last four years. But they think that their thorough citation search should have fixed the problem. In this work, they provide a detailed account of the problems encountered in the field of communication. They come up with a mixed-initiative system that may help clinical researchers and clinical data interact more effectively by providing a structured framework for the formulation of queries. There is a notable lack of study on the topic of communication in clinical research query mediation. Thorough and methodical research is necessary to provide useful descriptions of dialogue behaviour in HCI and UHC conversations.

Duplicate Document Detection: This framework assists in identifying duplicate documents by comparing their content similarities. The term "document transformation" refers to the steps used to convert documents into a more analytically-friendly format. This study laid the groundwork for document transformation with its Automatic and Appropriate File Name Generation methodology. It is possible for documents to have the same size, contents, or name.

If two documents have the same exact text, we call them duplicates. Search quality is diminished when such duplicate documents are included in the search results. During the process of generating search results, it is possible to identify and prohibit numerous documents with similar data using duplicate document detection. During global analysis, the indexing process scans the content of every document for duplication. When data is collected from many sources, algorithms are crucial for identifying instances of duplication. Deleting duplicate documents improves search accuracy and reduces runtime. Duplicate document identification speeds up indexing and searching. The importance of detecting duplicate papers and avoiding having the same identification document recorded in the database numerous times was suggested.

Mapping the character pictures according to their position and shape related to the text lines was done using the character Shape Coding (CSC) approach. It was a sturdy and cost-efficient approach. An array of linked and bound boxes was created from the page picture depending on the size of the page's components. The removal of large components occurred because they were considered to be non-textual information. When making comparisons, only textual information was considered. With datasets containing up to 100,000 records, this strategy proved to be efficient. A method for detecting near-duplicates using a sentence-level detection algorithm was suggested. To find instances of plagiarism and potential citations, this approach proved quite useful. At first, phrase word counts were analysed using a sliding window. We also compared the efficiency and effectiveness of document fingerprinting and shingling approaches using word-level characteristics. By providing a framework for explanation and formalization, this study offered a solution to the problem of duplicate document identification.

Along with a supplementary approach for string matching, the author has introduced four separate models. It was a case of full-content, partial-layout, full-content, and full-layout duplication for the models. The system's resilience was verified via testing it on a set of sample data. The limitations of the old-fashioned Optical Character Recognition (OCR) Model were overcome in this study. Identified identical or very similar documents using the fingerprint approach and signature stability assessment. Level and duplicate variations were the primary foci of this effort. To find duplicates, they have tested with 50 million papers. The authors have previously investigated the practicality of using digital signatures to deal with document duplication found duplicate documents using I-Match collection statistics with dataset samples ranging from 30 MB to 2 GB. In

comparison to current state-of-the-art technologies such as syntactic filtering, shingles, digital syntactic clustering, and its super shingle, the solutions offered correct findings in a fraction of the time (DSC-SS). Detecting duplicates in document image databases was accomplished using uncorrelated OCR results..

Additionally, it identified duplicates in databases with several false duplicates and established an empirical link between duplication models, comparison measures, and duplicate detection. Presented a strategy for assessing fingerprinting and information retrieval approaches in [65]. An analysis of the efficacy for accurately identifying co-derived documents has been conducted by the authors. Created a new approach to detecting duplicate documents by using the trie-tree structure in [66]. Documents downloaded from the internet had their 64-bit fingerprints saved in this structure. To improve the accuracy of spam email detection, the suggested strategy was used. An effective methodology for content-based automated file renaming has been suggested in paper [67]. The three main components of this methodology are the following: content analysis, semantic proof of file names, and file name categorization. They have suggested new algorithms for every step. All things considered, the files' contents have been effectively renamed using the suggested approach.

IV. RESULTS AND DISCUSSION

Everything from a patient's test findings and diagnosis to their treatments and follow-up appointments are part of their electronic health records, or EHRs. Because there are several ways to define the same issue, researchers may find it difficult to use EHR data searches. The statistics may refer to cancer as Carcinoma, for instance. It is possible for the EHR to have misspellings or acronyms. People may find a plethora of useful information with search engines. Therefore, it is important for researchers to develop a powerful search engine that can efficiently sift through electronic health record (EHR) data for research and health organisations. The clinician's frequent use of synonymous and antonymous words and phrases, as well as opposite and avoided expressions, contributes to this complexity. Uncertainties are also caused by a lack of standard language and punctuation use. Additionally, it makes dealing with context-sensitive interpretations fundamentally more difficult for computer systems. Since extra privacy issues need to be handled, EMERSE has a constraint.

It is intended that a rapid method would be used to review the patient's history for significant relevant acts. A free-text EHR search engine is EMERSE. For complicated activities like synthesising patient cases and abstract research data, it is a useful tool that helps researchers and practitioners retrieve information from

Department of Computer Science, A.V.C. College (Autonomous), Mayiladuthurai – 609 305. PROCEEDINGS 175

electronic health records more rapidly. Due to the impossibility of visually representing the data, multimodal data will be used for retrieval in the future. There will be a need in the future for a variety of approaches and methodologies that can integrate data from diverse sources. Another issue that needs fixing is the underutilization of clinical data. To prevent unauthorised parties from obtaining patients' confidential health information, a medical search engine that prioritises data security is essential. Natural language processing and other automated data extraction techniques need to standardise and arrange clinical notes so that they may be easily understood and used. To further increase the value of disorganised clinical data, search engines or IR systems may provide an efficient, adaptable, and scalable solution. The difficulties in deciphering the medical text should be minimal. In order to make any meaningful use of EHR data, it is necessary to get complete and accurate information on individual patient groups.

V. CONCLUSION

Several search strategies for retrieving information from EHRs have been identified, along with their advantages and disadvantages, according to this literature review. Using discrete developments in numerous approaches to collect the necessary Data for a user query, this research gives an abrupt examination of syntactic and semantic search engines. These days, the medical industry generates massive amounts of clinical data, making it difficult to get the relevant information for user inquiries. In order to tackle these obstacles quickly and effectively, EHR search engine technology was considered along with a number of other concerns and potential improvements. However, there is still a lack of knowledge about the various user needs for information mining from narrative clinical documents, and building information recovery features into EHRs is very dangerous.

REFERENCES

[1] Ethun, Cecilia G., et al. "Frailty and cancer: implications for oncology surgery, medical oncology, and radiation oncology." *CA: a cancer journal for clinicians* 67.5 (2017): 362-377.

[2] Calabresi, Paul, Philip S. Schein, and Saul A. Rosenberg. "Medical oncology: basic principles and clinical management of cancer." (1985).

[3] Schrag, Deborah, and Morgan Hanger. "Medical oncologists' views on communicating with patients about chemotherapy costs: a pilot survey." *Journal of Clinical Oncology* 25.2 (2007): 233-237.

[4] McDermott, Ultan, and Jeff Settleman. "Personalized cancer therapy with selective kinase inhibitors: an emerging paradigm in medical oncology." Journal of Clinical Oncology 27.33 (2009): 5650-5659.

[5] Muscaritoli, Maurizio, et al. "Prevalence of malnutrition in patients at first medical oncology visit: the PreMiO study." Oncotarget 8.45 (2017): 79884.

[6] Tong, Ho, Elisabeth Isenring, and Patsy Yates. "The prevalence of nutrition impact symptoms and their relationship to quality of life and clinical outcomes in medical oncology patients." Supportive care in Cancer 17 (2009): 83-90.

[7] Newell, Sallie, et al. "How well do medical oncologists' perceptions reflect their patients' reported physical and psychosocial problems? Data from a survey of five oncologists." Cancer: Interdisciplinary International Journal of the American Cancer Society 83.8 (1998): 1640-1651.

[8] Phelps, Ruby M., et al. "NCI-navy medical oncology branch cell line data base." Journal of cellular biochemistry 63.S24 (1996): 32-91.

[9] Chang, Victor T., et al. "Symptom and quality of life survey of medical oncology patients at a veterans affairs medical center: a role for symptom assessment." Cancer: Interdisciplinary International Journal of the American Cancer Society 88.5 (2000): 1175-1183.

[10] Shanafelt, Tait D., et al. "The well-being and personal wellness promotion strategies of medical oncologists in the North Central Cancer Treatment Group." Oncology 68.1 (2005): 23-32.

[11] Cherny, Nathan I., and Raphael Catane. "Attitudes of medical oncologists toward palliative care for patients with advanced and incurable cancer: report on a survey by the European Society of Medical Oncology Taskforce on Palliative and Supportive Care." Cancer 98.11 (2003): 2502-2510.

[12] Newell, Girgis, and Ackland. "The physical and psycho-social experiences of patients attending an outpatient medical oncology department: a cross-sectional study." European journal of cancer care 8.2 (1999): 73-82.

[13] Fadul, Nada, et al. "Supportive versus palliative care: What's in a name? A survey of medical oncologists and midlevel providers at a comprehensive cancer center." Cancer 115.9 (2009): 2013-2021.

[14] Braun, Ilana M., et al. "Medical oncologists' beliefs, practices, and knowledge regarding marijuana used therapeutically: a nationally representative survey study." Journal of Clinical Oncology 36.19 (2018): 1957.

[15] Stoffel, Elena M., et al. "Hereditary colorectal cancer syndromes: American Society of Clinical Oncology clinical practice guideline endorsement of the familial risk–colorectal cancer: European Society for Medical Oncology clinical practice guidelines." Journal of clinical oncology 33.2 (2015): 209.

A Comprehensive Review on Improving Plant Leaf Disease Detection Accuracy through Computer Vision Techniques

S. Ubaidulla¹ and **Dr. S. Mary Vennila**² ${}^{1}Guest \ Lecturer \ and {}^{2}Associate \ Professor \ and$

¹Guest Lecturer and ²Associate Professor and Research Supervisor,

¹Department of Computer Science, Governemtn Arts & Science College Perumbakkam, Chenna. ubaidulla1989@gmail.com

²PG and Research Dep[artment of Computer Science, Presidency College, Chennai, India. maryvennila13@gmail.com

Abstract - Agriculture serves as the primary source of sustenance for the growing global population, even in the face of rapid demographic expansion. The early prediction of plant diseases is crucial in agriculture to ensure a continuous and abundant food supply. Unfortunately, accurately foreseeing diseases during the initial stages of crop development remains a challenge. This review paper offers a comprehensive examination of recent advancements and methodologies aimed at enhancing the accuracy of plant leaf disease detection through sophisticated image preprocessing techniques in computer vision applications. The increasing demand for efficient and automated plant disease diagnosis necessitates a thorough understanding of preprocessing steps to optimize the quality of input data. The review systematically categorizes and analyzes various image preprocessing stages, including acquisition, color space conversion, contrast enhancement, noise reduction, segmentation, and data augmentation. Each stage's contribution to refining input images, extracting relevant features, and ultimately improving the precision of disease detection models is critically assessed. A key strength of this review lies in its survey of recent experimental studies that have implemented diverse preprocessing pipelines. These studies, encompassing a wide range of plant species and diseases, provide valuable insights into the effectiveness of preprocessing techniques in different contexts.

Keywords—plant leafdiseases, computer vision, image preprocessing, image segmentation, feature extraction, accuracy.

I.INTRODUCTION

In today's farming world, using technology in combination with traditional practices is crucial. One standout area is using computer vision to spot plant leaf diseases early, a key factor in ensuring there's enough food globally. With the world's growing population, it's vital to predict and control crop diseases from the beginning to keep the global food supply stable.

This review looks closely at how improving the accuracy of plant leaf disease detection is impacted by image preprocessing techniques using computer vision [21]. We're trying to understand recent advancements and methodologies to see how they help in agriculture. By focusing on these preprocessing methods, we're recognizing their essential role before we can identify diseases accurately. This exploration sets the groundwork for understanding how computer vision models tailored for agriculture become more effective.

We'll explore the diverse world of plant disease detection, unraveling the complexities of image preprocessing methods, evaluating how they affect data quality, and understanding how they contribute to precision agriculture goals [22].

From capturing detailed images to extracting features and training machine learning models, each part of this review aims to provide a complete and nuanced perspective.

This exploration is urgent due to the challenges in modern agriculture – balancing the need for more food with efficient resource use. Image preprocessing, acting as a gateway to accurate disease detection, can help farmers and stakeholders proactively protect crops.

II. RELATED WORK

Extensive efforts have been dedicated to the detection of leaf diseases through image processing over the years, and this field continues to captivate researchers for further exploration. In recent times, there has been a notable surge in the significance of automatic crop disease detection employing image processing and machine learning techniques.

In the study by P. Krithika et al. [6], the preprocessing involved image resizing, contrast enhancement, and color-space conversion. Subsequently, K-Means clustering was applied for segmentation, and feature extraction using GLCM was carried out. Classification was performed using a multiclass SVM.

In the work of R. Meena et al. [7], color space conversion and enhancement processes were

PROCEEDINGS

177

implemented, converting primary leaf colors into L*A*B*. Segmentation was accomplished through the K-Mean clustering algorithm, and feature extraction and classification were executed using GLCM and SVM, respectively.

In the study by Bharat et al. [9], images captured with a digital camera underwent image enhancement using a median filter, followed by segmentation with K-Mean clustering and classification using SVM. Pooja et al. [8] employed k-Mean clustering and Otsu's detection for segmentation, transitioning from RGB to HSI. Subsequently, boundary and spot detection algorithms were applied for further segmentation.

Rukaiyya et al. [10] performed pre-processing through contrast adjustment and normalization, with color transformation into YCBCR and Bi-level thresholding for segmentation. Features were extracted using GLCM and HMM for classification [11].

Chaitali et al., [12], applied image segmentation for background subtraction and adopted a classification approach utilizing KNN, ANN, and SVM methods. In KNN, the classification is based on the nearest distance between trained and testing subjects [13].

Varun et al., [14], developed a model incorporating thresholding techniques and morphological operations. A multiclass SVM serves as the classifier, and for segmentation, a set of marks generated from the analysis of color and luminosity components in the LAB* color spaces are employed. Feature extraction is achieved using GLCM.

In the work of Vijai Singh et al., [15], various plant leaf samples, such as rose/beans (bacterial disorder), lemon (sunburn disorder), banana (early scorch), and beans (fungal), captured with a digital camera. Green regions are considered as the background using a thresholding algorithm. The segmentation process involves the application of genetic algorithms, and features are extracted using color co-occurrence. The Minimum Distance Criterion, followed by an SVM classifier, is utilized for classification, achieving an average accuracy of 97.6%.

Sa'ed Abed et al., [13], enhanced the quality of input samples through a scaling and stretching (min-max linear) process. The creation of an HIS model is completed, followed by segmentation. A combination of Euclidean distance and K-means clustering is employed for sample segmentation. Feature extraction and classification are accomplished using GLCM and SVM, respectively.

Arya et al., [16], transformed input RGB images and converted them to the HIS format. Segmentation of components is carried out using Otsu's method. In Nema et al.'s [17] study, 81 images were included in the database, and analysis was conducted in the Lab color space. Leaf disease segmentation was performed using kmeans clustering, and SVM was employed for disease classification. Statistical measures such as mean, median, mode, and standard deviation were utilized to document their findings.

VidyashreeKanbur et al., [18], developed a model for leaf disease detection using multiple descriptors. The model exhibited superior performance when tested on a local leaf database but awaits evaluation on publicly available datasets.

In the investigation conducted by Pushpa et al. [19], the Indices Based Histogram technique was employed to effectively segment unhealthy regions of the leaf. Notably, the authors outperformed alternative segmentation techniques such as slice segmentation, polygon approximation, and mean-shift segmentation.

In the study led by Kaleem et al. [20], a meticulous pre-processing approach was applied, involving resizing the images to 300*300 dimensions, eliminating background noise, and enhancing brightness while adjusting contrast. The segmentation was executed using K-means clustering, and essential features were extracted utilizing Statistical GLCM. For the classification of leaf disorders, a SVM classifier was employed.

These studies showcase the versatility of segmentation techniques and the importance of adept pre-processing methods in achieving accurate results in plant leaf disease detection.

III . METHODOLOGY

In recent times, the detection of plant leaf diseases has garnered increasing attention, with a focus on leveraging vision and digital signal processing computer mechanisms. The application of these techniques involves a systematic approach comprising several stages [23]. Initially, the image undergoes pre-processing to enhance its quality and clarity. Subsequently, the process includes segmentation, where relevant portions of the image are delineated for further analysis. Following segmentation, features are extracted to capture essential characteristics related to the disease manifestation. Finally, a classification step is employed to categorize the disease class based on the extracted features. This comprehensive framework demonstrates the integration of advanced technologies in the domain of plant pathology, showcasing a methodical and automated approach for disease detection and classification on plant leaves.



Data Acquisition: In the study of plant leaf diseases, image acquisition serves as a crucial foundation for advancing research methodologies. Traditionally, data collection heavily leaned on online image repositories or controlled laboratory environments [6]. However, a contemporary and increasingly influential trend in this field involves the adoption of live image capturing directly from agricultural settings. This innovative approach addresses the limitations of relying solely on pre-existing datasets by providing realtime insights into the dynamic nature of plant diseases in their natural environment. By capturing live images from the field, researchers can obtain contextually rich data, allowing for a more nuanced understanding of disease progression, environmental influences, and the overall health of plants. This shift towards live image capturing represents a significant enhancement in research practices, enabling more accurate and up-to-date analyses in the study of plant leaf diseases.

Preprocessing: In the domain of computer vision applied to plant leaf diseases, image preprocessing assumes a pivotal role in refining input data for subsequent analysis. Essential techniques include smoothing and blurring, like Gaussian blurring, to reduce noise and create a uniform appearance, and contrast enhancement through methods such as histogram equalization to reveal finer details. Color space conversion, like transforming images to grayscale, simplifies analysis, while normalization ensures consistent scaling for fair comparisons. Resizing and cropping standardize input dimensions, aiding computational efficiency, and rotation/flipping during data augmentation improve model generalization [24]. Image registration aligns images for accurate comparison, while histogram stretching enhances pixel intensity distribution. Thresholding and edge detection are vital for segmentation, distinguishing diseased and healthy areas. Artifact removal eliminates unwanted distortions, contributing to cleaner images. Overall, these preprocessing techniques collectively optimize input data for more accurate and efficient computer vision analysis of plant leaf diseases.

Segmentation: This techniques in computer vision are critical for accurate plant leaf disease prediction. By precisely delineating regions of interest within leaf images, these techniques play a pivotal role in isolating diseaserelated patterns for focused analysis. Common segmentation methods include thresholding, which separates different regions based on pixel intensities, and edge detection, such as the Canny edge detector, highlighting boundaries for more accurate delineation. Region-based segmentation groups pixels with similar attributes, aiding in the identification of coherent structures relevant to disease manifestation. Watershed segmentation, utilizing gradient magnitudes, further refines the partitioning of distinct regions [19]. The effectiveness of these segmentation techniques directly influences the success of subsequent feature extraction and predictive modeling, providing a foundation for reliable and

automated plant leaf disease prediction in computer vision applications.

Feature Extraction: After Segmentation, Feature extraction stands as a pivotal stage within computer vision applications dedicated to predicting plant leaf diseases, encompassing the identification and retrieval of pertinent information from leaf images. The inclusion of various features, such as shape, texture, and color characteristics, proves essential in encapsulating distinct patterns linked to different diseases [15]. Feature extraction methods commonly involve histogram analysis for capturing color distribution, texture analysis utilizing techniques like cooccurrence matrices, and the use of contour-based shape descriptors. Through the conversion of raw image data into a discriminative feature set, this process lays the groundwork for the subsequent application of machine learning algorithms, facilitating accurate disease prediction. The meticulous selection of informative features holds paramount importance in constructing resilient models capable of effectively discerning between healthy and diseased leaves, thereby significantly contributing to the overarching success of computer vision-driven prediction systems in the field of plant pathology.

Classification: Classification plant leaf disease prediction is a crucial phase that involves assigning predefined labels or categories to processed image data based on learned patterns and features [5]. Leveraging machine learning algorithms, classification models are trained on extracted features from leaf images to discern between healthy and diseased plants. Commonly employed algorithms include support vector machines, decision trees, and neural networks. The success of classification hinges on the quality of features extracted during earlier stages, as well as the robustness and generalization capabilities of the chosen machine learning model. By effectively categorizing plant leaves into distinct classes, classification models contribute to accurate and automated prediction of leaf diseases, aiding in early detection and targeted mitigation strategies in the field of plant pathology [1].

Summary of computer vision techniques: The integration of computer vision techniques, especially in the domains of machine learning and deep learning, is pivotal for elevating the precision of plant leaf disease detection. Following Table 1 describes the preprocessing, segmentation and feature extraction techniques used in various plant diseases and Table 2 shows the performance measures based on the computer vision techniques applied.

Department of Computer Science, A.V.C. College (Autonomous), Mayiladuthurai – 609 305.

Authors	Plant Diseases Identified	Preprocessing & Segmentation Techniques	Feature Extraction Techniques
[1]	Tomato Leaf Diseases	K-Means Clustering & Contour Tracing	Discrete Wavelength Transform, GLCM, PCA
[2]	Black gram Leaf Diseases	Cropping , Resizing (512x512), rotation, mirror symmetry, illumination correction, shifting/translation, and noise injection	CNN
[3]	[3]Apple, Corn, Grapes, Potato, TomatoGray Scale conversion, Smoothening usi Gaussian Filter & Otsu Thresholdin		GLCM
[4]	Medicinal plant leafTulsi, Peppermint,Denoising, Sobel[4]Bael, Lemon, Balm, Catnip, Steviafilter for edge		Texture feature, Gray Level Runlength matrix feature, Multispectral feature
[5]	Tomato Leaf Diseases	Bilateral Filtering	EPO- MobileNet

Table 1. Analysis of computer vision techniques used on plant leaf diseases

Auth ors	Model Used	Precision (%)	Recall (%)	F1- Score (%)	Accuracy (%)
[1]	[1] SVM 9		97.8	91.5	89
[1]	KNN	97.8	97.5	95.2	97.3
[1]	CNN	99.5	99.5	98.8	99.09
[2]	ResNet1 8	98.46	98.37	98.40	99.39
[3]	Random Forest (Apple	92	-	91	91

	Leaf)					
[3]	[3] [3] (Corn Leaf)		-	98	94	
[3]	Random Forest (Gapes Leaf)	99	-	99	95	
[3]	Random Forest (Potato Leaf)	99	-	99	98	
[3]	Random Forest (Tomato Leaf)	91	-	93	87	
[4]	Multi- Layer Perceptr on	99.1	99	99	99.01	
[5]	OMNCN N	98.5	98.92	98.5	98.7	

Table 2. Analysis of performance measures achievedusing ML and DL models based on the computer visiontechniques

In the analysis of computer vision techniques applied to plant leaf diseases, the focus encompasses a range of methodologies, including image preprocessing, feature extraction, segmentation, and the integration of deep learning approaches like Convolutional Neural Networks (CNNs).

Simultaneously, a comprehensive examination is undertaken to evaluate the performance measures achieved through the implementation of machine learning (ML) and deep learning (DL) models. Metrics such as accuracy, precision, recall, and F1 score are meticulously assessed, providing a detailed perspective on the effectiveness of these models in accurately diagnosing and classifying plant leaf diseases. This dual-pronged analysis offers valuable insights into the efficacy of computer vision techniques in the field of plant pathology.

IV. CONCLUSION

In conclusion, computer vision plays a pivotal role in revolutionizing the field of plant disease prediction. The integration of advanced image processing techniques and machine learning and deep learning algorithms has enabled more accurate, efficient, and timely identification of diseases affecting plant health. By leveraging computer

Department of Computer Science, A.V.C. College (Autonomous), Mayiladuthurai – 609 305.

PROCEEDINGS 180

vision, researchers and practitioners can automate the analysis of large datasets of plant images, leading to early detection and targeted intervention strategies.

The preprocessing steps, including background noise removal, color space conversions, and segmentation techniques, contribute to enhancing the quality of input data for subsequent analysis. Feature extraction methodologies, encompassing shape, texture, and color characteristics, enable the development of robust models capable of discriminating between healthy and diseased plants. The utilization of techniques like 5-fold crossvalidation ensures the reliability and generalization of predictive models. Overall, the synergy between computer vision and plant pathology holds immense promise for sustainable agriculture, facilitating proactive measures to mitigate the impact of diseases on crop yield and food security. The continual advancements in this interdisciplinary field pave the way for innovative solutions, fostering resilience in agricultural practices and contributing to the broader goal of global food sustainability.

REFERENCES

- Harakannanavar, S. S., Rudagi, J. M., Puranikmath, V. I., Siddiqua, A., &Pramodhini, R. (2022). Plant leaf disease detection using computer vision and machine learning algorithms. Global Transitions Proceedings, 3(1), 305-310.
- [2] Talasila, S., Rawal, K., Sethi, G., & Sanjay, M. S. S. (2022). Black gram Plant Leaf Disease (BPLD) dataset for recognition and classification of diseases using computer-vision algorithms. Data in Brief, 45, 108725.
- [3] Joshi, K., Awale, R., Ahmad, S., Patil, S., &Pisal, V. (2022). Plant leaf disease detection using computer vision techniques and machine learning. In ITM Web of Conferences (Vol. 44, p. 03002). EDP Sciences.
- [4] Naeem, S., Ali, A., Chesneau, C., Tahir, M. H., Jamal, F., Sherwani, R. A. K., &UI Hassan, M. (2021). The classification of medicinal plant leaves based on multispectral and texture feature using machine learning approach. Agronomy, 11(2), 263.
- [5] Ashwinkumar, S., Rajagopal, S., Manimaran, V., & Jegajothi, B. (2022). Automated plant leaf disease detection and classification using optimal MobileNet based convolutional neural networks. Materials Today: Proceedings, 51, 480-487.
- [6] Krithika, P., &Veni, S. (2017, March). Leaf disease detection on cucumber leaves using multiclass support vector machine. In 2017 International Conference on Wireless Communications, Signal Processing and Networking (WiSPNET) (pp. 1276-1281). IEEE.
- [7] Meena, R., Joshi, S., &Raghuwanshi, S. (2022, December). An Automated System for Rice Plant Diagnosis Using Deep Learning. In International Conference on Communication and Intelligent Systems (pp. 373-390). Singapore: Springer Nature Singapore.
- [8] Pooja, V., Das, R., &Kanchana, V. (2017, April). Identification of plant leaf diseases using image processing techniques. In 2017 IEEE Technological Innovations in ICT for Agriculture and Rural Development (TIAR) (pp. 130-133). IEEE.
- [9] Mishra, B., Nema, S., Lambert, M., &Nema, S. (2017, March). Recent technologies of leaf disease detection using image processing approach—A review. In 2017 International Conference on Innovations in Information, Embedded and Communication Systems (ICIIECS) (pp. 1-5). IEEE.
- [10] Shaikh, R. P., & Dhole, S. A. (2017, April). Citrus leaf unhealthy region detection by using image processing technique. In 2017

International conference of Electronics, Communication and Aerospace Technology (ICECA) (Vol. 1, pp. 420-423). IEEE.

- [11] Harakannanavar, S. S., Shridhar, H., Premananda, R., Jambukesh, H. J., &Prashanth, C. R. (2022). Integrated Analysis of Tomato Plant leaf disease disorder using improved Machine Learning approach. Journal of Positive School Psychology, 1288-1297.
- [12] Dhaware, C. G., &Wanjale, K. H. (2017, January). A modern approach for plant leaf disease classification which depends on leaf image processing. In 2017 International Conference on Computer Communication and Informatics (ICCCI) (pp. 1-4). IEEE.
- [13] Esmaeel, A. A. (2018, April). A novel approach to classify and detect bean diseases based on image processing. In 2018 IEEE Symposium on Computer Applications & Industrial Electronics (ISCAIE) (pp. 297-302). IEEE.
- [14] Singh, V., &Misra, A. K. (2015, March). Detection of unhealthy region of plant leaves using image processing and genetic algorithm. In 2015 International Conference on Advances in Computer Engineering and Applications (pp. 1028-1032). IEEE.
- [15] Singh, V., &Misra, A. K. (2017). Detection of plant leaf diseases using image segmentation and soft computing techniques. Information processing in Agriculture, 4(1), 41-49.
- [16] Arya, M. S., Anjali, K., &Unni, D. (2018, January). Detection of unhealthy plant leaves using image processing and genetic algorithm with Arduino. In 2018 International Conference on Power, Signals, Control and Computation (EPSCICON) (pp. 1-5). IEEE.
- [17] Nema, S., & Dixit, A. (2018, December). Wheat leaf detection and prevention using support vector machine. In 2018 International Conference on Circuits and Systems in Digital Enterprise Technology (ICCSDET) (pp. 1-5). IEEE.
- [18] Kanabur, V., Harakannanavar, S. S., Purnikmath, V. I., Hullole, P., &Torse, D. (2020). Detection of leaf disease using hybrid feature extraction techniques and CNN classifier. In Computational Vision and Bio-Inspired Computing: ICCVBIC 2019 (pp. 1213-1220). Springer International Publishing.
- [19] Pushpa, B. R., AV, S. H., & Ashok, A. (2021, July). Diseased leaf segmentation from complex background using indices based histogram. In 2021 6th International Conference on Communication and Electronics Systems (ICCES) (pp. 1502-1507). IEEE.
- [20] Kaleem, M. K. (2021). A modern approach for detection of leaf diseases using image processing and ml based svm classifier. Turkish Journal of Computer and Mathematics Education (TURCOMAT), 12(13), 3340-3347.
- [21] Prakash, R. M., Saraswathy, G. P., Ramalakshmi, G., Mangaleswari, K. H., &Kaviya, T. (2017, March). Detection of leaf diseases and classification using digital image processing. In 2017 international conference on innovations in information, embedded and communication systems (ICIIECS) (pp. 1-4). IEEE.
- [22] Narmadha, R. P., & Arulvadivu, G. (2017, January). Detection and measurement of paddy leaf disease symptoms using image processing. In 2017 International Conference on Computer Communication and Informatics (ICCCI) (pp. 1-4). IEEE.
- [23] Sultana, N., Rahman, T., Parven, N., Rashiduzzaman, M., &Jabiullah, I. (2020). Computer vision-based plant leaf disease recognition using deep learning. International Journal of Innovative Technology and Exploring Engineering (IJITEE), 9(5), 622-626.
- [24] Nanehkaran, Y. A., Zhang, D., Chen, J., Tian, Y., & Al-Nabhan, N. (2020). Recognition of plant leaf diseases based on computer vision. Journal of Ambient Intelligence and Humanized Computing, 1-18.

A CLOUD BASED SECURE ELECTRONIC HEALTH HISTORY FRAMEWORK USING FERNET AND FULLY HOMOMORPHIC ENCRYPTION

N. SUBHALAKSHMI

Research Scholar, S.T.E.T Women's College (Autonomous), (Affiliated to Bharathidasan University, Tiruchirappalli) Sundarakkottai, Mannargudi - 614016, Thiruvarur Dt., Tamil Nadu, India subha.stet@gmail.com Mobile No. 9095598954

Dr. M.V. SRINATH

Research Supervisor, S.T.E.T Women's College (Autonomous), (Affiliated to Bharathidasan University, Tiruchirappalli) Sundarakkottai, Mannargudi - 614016, Thiruvarur Dt., Tamil Nadu, India sri_induja@rediffmail.com Mobile No. 9710944476

Abstract— Effectively managing health data and keeping Electronic Health History (EHHs) is a significant challenge. Physicians need patient's medical records, including multimedia big data analytics, to diagnose patients. Online EHH allows patients to efficiently manage their health records. This can be efficiently achieved by using cloud computing to store EHH. Cloud data storage offers a significant Quality of Service (QoS) aspect from a security perspective, but in this research EHHs are encrypted using Fully Homomorphic Encryption (FHE) and Fernet algorithm before being stored in the cloud environment. The proposed system in this work is divided into three modules: EHH storage after NLP; key management and authentication; and data storage on the Amazon Web Services (AWS) cloud. The authentication module uses a novel EHH algorithm, which performs both encryption and decryption. The text file saved in local drive, is encrypted and saved to Cloud bucket. The comparative analysis of performance metrics such as Bandwidth, memory usage and level of protection through different encryption techniques determines the OoS.

Keywords— Fully Homomorphic Encryption, Fernet, Electronic Health History, NLP, Amazon Web Service (AWS)

I. INTRODUCTION

India's healthcare system is the country's biggest challenge. Technology is used sparingly in healthcare maintenance [1]. Before making a meaningful patient diagnosis, doctors must have a complete understanding of the patient's medical history, but getting this information can be challenging. A patient consults numerous doctors over his or her lifetime. For their medical needs, the majority of people consult many doctors, including general practitioners, dermatologists, orthopaedic surgeons, and cardiologists [2].Every single doctor requires the patients' medical records. All offices also include information about health and treatments scattered around them. Health care is only provided by appointment and may be challenging; therefore, the doctor should be informed of any disease test results. There are times when we are unaware of what information is essential and lack access to all of the patient's medical records, which could threaten their lives.

Big data in the healthcare sector refers to an endless set of extremely large and intricate medical data. Therefore, those medical records cannot be managed using conventional gear or software. The transfer of patient information or other records to healthcare providers is necessary for the e-health service depending on the situation. Due to the increased benefits of cloud computing in terms of cost, storage, and scalability, data providers and organizations are increasingly outsourcing their information from local servers to remote cloud servers[3].

Healthcare cloud is a platform that allows server-based communication amongst service providers, clinics, dispensaries, and many other patients. As with cloud computing, there are many problems and disadvantages, the most important and common of which is security problems. Lawful and political issues (private patient document, responsibility, etc.), data guard, privacy safety (protection of users' personal data), deficiency of transparency, lack of safety standards, and software licensing are all included. Therefore, innovative approaches are required to ensure EHH safety and privacy [4].

Storing medical data on the cloud computing platform permits users to store their data conveniently. Personal data is first encoded by the user and directed to cloud storage. The encoded data is reorganized by the cloud platform after being stripped of important new data. Encoded cloud data is typically managed by a third party. Subsequently, the integrity, privacy, and safety of health record data are less secure than private storing systems. Consequently, it is significant to ensure that unauthorized third parties can access or modify encoded cloud data. Cloud platforms may use cryptographic methods with private keys generated by third parties wherever third parties are responsible for cloud storage. Another problem is that unauthorized users can search data, resulting in data leakage to strangers.

This text delivers a protected data access scheme that delivers greater security and integrity when storing and sharing sensitive health care data in the cloud. This system consists of three different cloud platform techniques to enhance the security, privacy, and integrity of users' personal data. First, patient health histories are divided as Sensitive Data and Non-Sensitive Data. This is made likely by Natural Language Processing (NLP). The second technique involves storing data in Amazon Web Services (AWS) [5], which makes advantage of both cloud data and user permissions. Data controllers grant this authorization to access patients' medical history, which is typically hosted in the cloud. The third is homomorphic encryption and a novel model EHH is proposed, which is an option in case data owners modify shared data on cloud-based platforms. Enhanced security for cloud computing platforms where encrypted data is changed without the original data's awareness and authorized users recover and decode the data. This is a privacy mechanism that maintains the privacy of sensitive data over the user's verification process.

This work provides a brief overview to big data and its role in storing medical data. The usage of big data is attracting attention to the architecture and technology will continue to help us cope with the rapid growth of data in the healthcare applications. Here, we focus on the novel design and development of intelligent and safe healthcare information systems using NLP (spell check, capitalisation, tokenisation and POS tagging and notability detection) with EHH. Large amounts of data from the medical history can be handled using advanced security techniques. The novelty is in training the best data safety layers and storage practices to maintain privacy and security. The data are stored in cloud with EHH encryption algorithm to access the medical history anywhere. A proposed three-layer model seems to be a more efficient big data system for diagnosing diseases and determining patients' health histories as needed. Finally, comparative analysis of QoS parameters such as bandwidth, memory usage and level of protection for proposed EHH with the existing system determine the best fit model for securing the medical history of patients.

II. LITERATURE REVIEW

Mohammad et al., [6] proposes about the data security. According to their analysis, both the customer and the receiver side of a cloud are vulnerable to risk when transmitting and receiving data, so in order to provide security and other services in the future, the authors developed a flexible system of distribution that would allow a person to store and access data with significant ease. In contrast to shared access, remote access in the cloud computing environment could expose a person to minimal danger and data theft. Similarly, the F. Bracciet al., [7] has additionally analysed the statistics safety and garage in cloud in which they are not unusual place existences of statistics robbery, statistics breaches and cloud statistics in accessibility are analysed.

S. Neelaetal,[8] deals about allowing permissions to limit admission to data. Second, it provides a user-centric method to proactively avoid unpredicted user actions on the cloud side. This exhibits some degree of resistance at complex levels as it cannot deal with active threats, including new and future threats.

S. Narulaetal [9] train cloud NLP with techniques in a map-reduce cloud setting. The distributed cloud conforms to the NLP processes to appropriately categorize the dataset. Then combine the NLP of each node in the cloud. To discover the best dataset among the huge datasets in the algorithm's NLP data file, split several iterations, combine the consequences of the iterations, and do a few more iterations to get the best dataset which is difficult to implement.

I. Saeedetal.[10], focuses on computing cloud in data security and access control where three key facilities namely key data protection, key currency, and key validation are key to the dynamic data exchange of vigorously altering groups. To realize these three functions, the authors provide a dynamic group key record practice built on a central key generation center.

Chantamit-o-pas et al., [11] uses fully homomorphic encryption to improve the security of applications where complex data is found. Encrypted data is treated as encrypted input, and accurate fluctuations to encrypted data are made deprived of decrypting the data. The contents of the information are transparent to the user while it is being modified. This FHE scheme determines that isolated sharing of personal information does not compromise the confidentiality of personal information.

C. Weng et al., [12] uses four methods to create Homomorphic Encryption (HE) schemes. They are key generation algorithms, encryption algorithms, decryption algorithms, and additional evaluation algorithms. FHE involves two basic homomorphisms. They are, Multiple Homomorphic Encryption Algorithm and Additive Homomorphic Encryption Algorithm using homomorphic principles for multiplication and addition. Homomorphisms are only supported in addition and multiplication by HE algorithms. FHE consists of finding an encryption algorithm that allows the encrypted data to contain any amount of addition and multiplication algorithms. This work uses the full-encryption symmetric homomorphic algorithm along with the several modifications.

III. METHODOLOGY

A. FULLY HOMOMORPHIC ENCRYPTION

This paper simply uses the full-encryption symmetric homomorphic algorithm proposed by *F. Zhao*[13]. The FHE algorithm is the standard and most trusted algorithm in cryptography, especially for text encryption, where keys are generated and 128bits are reduced to 16/32/64 bits for fast processing and balanced reading of blocks[14]. The following stages were used to construct the FHE algorithm in the proposed work.

Steps for Encryption:

Step 1: Key generation: Firstly set-of-spherical keys are gained from the Cipher key;

Step 2: Loading of the array is implemented upon the plaintext / block array;

Step 3: Next, the number one spherical – secrets introduced to the initialized array;

Step 4: Next, country manipulation of 9 rounds is completed;Step 5: Finally, the 10tharray manipulation spherical is completed earlier acquisition of cipher text;

Step 6: Finally, the array received is copied through the received ciphered - text / encrypted information. Thus, the encryption is performed in which the 128bits collection is transformed as 16bytes for FHE because it plays nicely as 16bytes.

Decryption

For the decryption the inverse capabilities are performed

Step 1: Initial decryption spherical is completed;

Step 2: Nine-complete decryption rounds are followed

Step 3: Finally, the Xor Round key is completed to decrypt the cipher text.

Thus, the set of rules is advanced and implemented upon the datasets to teach and check the accuracy of the advanced encoding and interpreting rounds in FHE.

B. FERNET ALGORITHM

Fernet uses plaintext for encryption and decryption in a manner similar to the FHE algorithm. It offers rotation of keys produced by "MultiFernet." [15].

Encryption and Decryption:

The encryption key of fernet follows the below steps **Step1:** Generate a key.

Step2: Assign a key value to the selected variable.

Step3: Convert plaintext to ciphertext.

To **decrypt** the encrypted text, Fernet performs the ciphertext-to-plaintext conversion as an inverse function, and the output is displayed as a "string" value of bytes.

C. Reasons for Adopting FHE and Fernet

Although there are various symmetric and asymmetric encryption and decryption algorithms, FHE is the most extensively used and standardized algorithm [16].Provides users with keys that are more secure than others. Therefore, by combining AES and Fernet, this work aims at advanced encryption and decryption of text as input to NLP-based encoder/decoder models, something that has not been attempted before in existing research.

Dataset

The Patient information contains medical. demographic, and over-all data from the record. About 7,000 patients were treated for heart, cancer and tumour diseases. These records comprise particulars about patients namely name, age, gender, family history of illness, laboratory test results, physical examinations by physicians, and screening results. Severity disease changes into acute, chronic, and chronic conditions. Most patients wanted to identify their health status and get prescriptions from an isolated place [17]. Therefore, the data should be stored on the cloud platform. Security and privacy are challenges as patient data is sensitive and non-sensitive. It uses identity encryption together with FHE and Fernet to overcome security and privacy problems. Figure 1 shows the sample dataset consisting of medical history of 23-year-old patients stored as text file.

Input : Medical History (data)

<u>SUBJECTIVE</u>., This 23-year-old white female presents with complaint of allergies. She used to have allergies when she lived in Seattle but she thinks they are worse here. In the past, she has tried Claritin, and Zyrtec. Both worked for short time but then seemed to lose effectiveness. She has used Allegra also. She used that last summer and she began using it again two weeks ago. It does not appear to be working very well. She has used over-the-counter sprays but no prescription nasal sprays. She does have asthma but does not require daily medication for this and does not think it is flaring up_MEDICATIONS: , Her only medication currently is Ortho Tri-Cyclen and the Allegra ALLERGIES: , She has no known medicine allergies_OBJECTIVE. Vitals: Weight was 130 pounds and blood pressure 124/78. HEENT: Her throat was mildly erythematous without exudate. Nasal mucosa was erythematous and swollen. Only clear drainage was seen. TMs were clear. Neck: Supple without adenopathy Lungs: Clear_ASSESSMENT:, Allergic rhinitis,PLAN:,1. She will try Zyrtec instead of Allegra again. Another option will be to use loratadine. She does not think she has prescription coverage so that might be cheaper.,2. Samples of Nasonex two sprays in each nostril given for three weeks. A prescription was written as well.

Figure 1 Sample Dataset Consisting of Medical History of a 23-year-old Patient

System Architecture

The proposed architecture diagram of the Electronic Health History (EHH) is shown in figure 2. In this model, users transmit and store encrypted data in the cloud. Both provide safe data storage as well as data security during transmission. The information is handled by cloud computing service providers, but it is not readily accessible in plain text[15].



Figure 2 Proposed Framework EHH Encryption

Natural Language Processing`

print(d)

AWS Cloud Implementation

This comprises of the subsequent steps: spellchecking, capitalization, tokenization, positional tagging and noun entity recognition.

Algorithm1: EHH Encryption Algorithm #the Fernet module is #cryptopackage from cryptography. Fernet import Fernet #generate key key=Fernet.generate_key() #assign key value to variable f=Fernet(key) #Convert plaintext to ciphertext token=f.encrypt(b"Welcome to EHH Encryption of health data") #showciphertext print(token) #decrypttheciphertext d=f.decrypt(token) #showplaintext

The text file is collected after NLP then it is saved in local drive initially. The medical history data is stored as S3 bucket. The setup for the cloud storage of EHH model is demonstrated in figure 3. It comprises of upload and download options. We have interfaces for both upload and download. When we click "upload," the file is encrypted and the data is saved in a bucket as S3. This file will be in an encrypted format. Following this procedure, the content will be downloaded in encrypted format as seen in figure 4.

When using the download option, the same S3 file will be decrypted and stored in original format. Private key is used to encrypt all public information. Key is created and stored in local after loading every time.

-> C (0 127.0.0.1:8000/socs#/								C A	• •		1
NLP 🛅 Data Engineering 🛅 CV 🛅 DL 🛅 M	6. 🛅 Statistics & Proba	B MLOPS	E LeetCode	🗄 Spark Hadoop	E Meas	E GNN	E Others	🛅 Image Cluster	ing 🖻	CBIR	
FastAPI 🏧 🤐											
penepi juon											
default											1
POST /upload Upiced											~
Post Alexandroid Developed											
Poor Poor Poor Poor											
Schemas											
Rody unload unload post											
anna-ahnan-ahnan-hnar a											
HTTPValidationError >											
ValidationError >											

Figure 3 Local Cloud Setup API

UP 🛅 Data Engineering 🛅 CV 🛅	DL 🛅 ML 🛅 Statistics & Proba. 🛅 M	ILOPS 🛅 LeetCode 🛅 Spark Hadoos	p 🛅 kõres 🛅 GN	N 🗄 Others 🗎 Image Clustering 🛗 CBR
Services Q. Scorch		[Option+5]		All. Data
mazon S3 ×	Amazon 53 🗦 Buckets			
ackets ccess Points bject Lambda Access Points	Account snapshot Storage laws provides visibility into sto	nige usage and activity trends. Learn more [e	View Storage Lens dashboard
ulti-Region Access Points Itch Operations coss analyzer for 53	Buckets (4) and Buckets ere containers for data stored in S	a Lease more	Copy ARN	Impty Delete Create bucket
	Q, Find buckets by nome			< 1 > @
ock Public Access settings for is account	Name A	AWS Region 9	Access V	Creation date 9
srage Lens	O eridata 🎙	US East (N. Virginia) us-east-1	Bucket and objects not public	August 20, 2021, 12:46:41 (UTC+05:30)
indoards	O arldatabackup	US East (N. Virginia) us-east-1	Bucket and objects not public	August 21, 2021, 13:03:05 (UTC+05:30)
ature spotlight 🚺	O artikeying	US East (N. Virginia) us-east-1	Objects can be public	September 28, 2021, 15:18:36 (UTC+05:30)
	O secure-private-big-data	Asia Pacific (Mambai) ap-south-1	Objects can be public	November 21, 2022, 18:01:09 (UTC+05:30)

Figure 4 S3 files stored as encrypted

IV. RESULTS AND DISCUSSION

The final proposed system comprises of the aws cloud implementation, performance metrics such as level of protection, memory usage and bandwidth. Figure 5 displays the file's contents.

185

<pre>#ipip install soluppop #ipip install https://sl-us-west-2.amssonaws.com/ai2-s2- solppep/relamses/v0.4.0/nm_nes_bobds_md=0.4.0.tar.gs [pip install https://dl-us-west-2.amsdonaws.com/ai2-s2-</pre>
scimpacy/releases/v0.5.0/en_core_sci_am=0.5.0.tar.gs
inport pandas as pd
<pre>med_transcript = pd.read_cov("/content/sample_data/stamples.cov", index_col= 0) med_transcript.info()</pre>
med_transcript.head()
delana faandan oora faana Dataffaanata
Int64Index: 4999 entries. 0 to 4998
f Column Non-Null Count Dtype
0 description 4999 non-null object
1 medical_appectatty 4999 non-null object
3 transcription 4966 non-null object
4 keywords 3931 non-null object
memory usage: 234.3+ KB
med_transcript.dropna(subset=['transcription'], inplace=True)
<pre>med_transcript_mmail = med_transcript.sample(n=100, replace=raise, random_sta te=32)</pre>
med_transcript_small.info() med_transcript_small.bed()
and require the man of the second s
Let's take one transcription to see how we can work with NER;
a provide entity recognition (NER) is a subtank of natural language processing
* classify named entities mentioned in unstructured text into pre- tional sector of the sector of
mample_transcription = med_transcript_mmali['transcription'].iloc[0]
print(mample_transcription[:1000]) # prints just the first 1000 characters
HIBTORY OF PRESENT ILLNESSI, The patient is well known to me for a history
or iron-derigiency anemia due to chronic blood load from colitin. We corrected her hematocrit last year with intravenous (VV) iron. Ultimately,
she had a total proctocolectomy done on 03/14/2007 to treat her colitis. Her
course has been very complicated since then with needing multiple surgeries for removal of hematoma. This is partly because she was on anticoagulation
for a right arm deep venous thrombosis (DVT) she had early this year.
complicated by emptic phintitis. Chart was reviewed, and I will not reitorate

Figure 5 Content of the File

Figure 6 shows the throughput [18] comparison between DES,AES, RSA, 3DES, FHE, Fernet and EHH encryption algorithm. The time taken is measured in terms of number of bytes /sec.



Figure 6 Throughput Comparison of Encryption algorithms

The file is running at different sizes using AES, DES, FHE, RSA, 3DES and proposed EHH algorithms. An overall examination of these processes is evaluated based on parameters such as execution time, memory relative to implementation, and throughput. Python is cast-off to create a conventional comparison model EHH. The proposed encryption strategy aims to determine which of the aforementioned encryption techniques is the fastest and most secure. This allows users to modify the encryption algorithm that best suits the type of information being encrypted. The results show that the proposed EHH system of algorithms raises the productivity of encryption algorithms by firmly encrypting data in short period of time. Time-Based Analysis and Evaluation of Encryption and Decryption analysis shows that EHH is faster than all other algorithms, and that the

algorithm EHH system outperforms all other algorithms. This works numerous times faster than Key generation must be finished before encrypting data. This method is implemented with the cooperation of the patients and the cloud service provider. Compared to prior encryption processes, the EHH process takes 10-15% less time to encrypt files. The inverse of encryption time is known as decryption time. This is the premeditated time it takes for the decryption algorithm to create the original text from the encrypted text. By analyzing the decoding time, the hybrid system of algorithms performs better than all other algorithms. Figure 7 show the comparison of encryption times for AES, DES, FHE, and Proposed EHH algorithms.



Figure 7 Comparison of Encryption time of AES, DES, FHE, and proposed EHH Algorithms

V. CONCLUSION

As the demands of big data grow, sustaining the integrity of the big data, cloud and its users turn out to be a top significance. There are numerous safety approaches available in the cloud. We use encryption methods to ensure cloud storage security in medical data history. With the increasing use of big data to store numerous data, safety and confidentiality have become a major daily challenge. This work proposes building an encryption model that leverages the influence of NLP and text mining with data encryption algorithms EHH to protect the portion of the data in documents that actually needs protection. This method encrypts attributed at the content of the text itself, which speed sup processing. The proposed technique outperforms the existing methods in terms of memory usage, throughput, and encryption time, according to a comparison of the performance of the AES, DES, Fernet, FHE, and proposed EHH encryption algorithms. Compared to prior approaches, this methodology achieves EHH encryption and decryption in less time. Therefore, the proposed EHH classification of

algorithms is helpful for modern needs. This structure generates random AES and DES keys and encrypts the accepted data files. The main advantage of this proposed system is saving the confidential medical history in cloud which will help in emergency period.

In the future, the industry based on cloud computing will be greatly hindered, and the challenge remains that users will lose data and information stored in the cloud. According to this research's findings, cloud design and taxonomies make it simpler for outsiders to access the cloud and steal sensitive data; as a result, the author believes that modifying the approach to architecture is the best way to address this problem. Likewise, users should choose remote access with two levels of encryption to safeguard their data.

REFERENCES

- K. Sudheep and S. Joseph, "Review on Securing Medical Big Data in Healthcare Cloud," 2019 5th International Conference on Advanced Computing & Communication Systems (ICACCS), 2019, pp. 212-215.
- [2] T. Javid, M. Faris, H. Beenish and M. Fahad, "Cybersecurity and Data Privacy in the Cloudlet for Preliminary Healthcare Big Data Analytics," 2020 International Conference on Computing and Information Technology (ICCIT-1441), 2020, pp. 1-4
- [3] Z. Tang, "A Preliminary Study on Data Security Technology in Big Data Cloud Computing Environment," 2020 International Conference on Big Data & Artificial Intelligence & Software Engineering (ICBASE), 2020, pp. 27-30.
- [4] F. Wang, H. Wang and L. Xue, "Research on Data Security in Big Data Cloud Computing Environment," 2021 IEEE 5th Advanced Information Technology, Electronic and Automation Control Conference (IAEAC), 2021, pp. 1446-1450.
- [5] Zhiying Wang, Nianxin Wang, Xiang Su and Shilun Ge, "An empirical study on business analytics affordances enhancing the management of cloud computing data security", *International Journal of Information Management*, vol. 50, pp. 387-394, 2020.
- [6] Mohammad Anwar Hossain, Ahsan Ullah, Newaz Ibrahim Khan and Md Feroz Alam, "Design and Development of a Novel Symmetric Algorithm for Enhancing Data Security in Cloud Computing", *Journal of Information Security*, vol. 10, no. 04, pp. 199-236, 2019.
- [7] F. Bracci, A. Corradi and L. Foschini, "Database security management for healthcare SaaS in the Amazon AWS Cloud," 2012 IEEE Symposium on Computers and Communications (ISCC), 2012, pp. 000812-000819.
- [8] S. Neela, Y. Neyyala, V. Pendem, K. Peryala and V. V. Kumar, "Cloud Computing Based Learning Web Application Through Amazon Web Services," 2021 7th International Conference on Advanced Computing and Communication Systems (ICACCS), 2021, pp. 472-475.
- [9] S. Narula, A. Jain and Prachi, "Cloud Computing Security: Amazon Web Service," 2015 Fifth International Conference on Advanced Computing & Communication Technologies, 2015, pp. 501-505.
- [10] I. Saeed, S. Baras and H. Hajjdiab, "Security and Privacy of AWS S3 and Azure Blob Storage Services," 2019 IEEE 4th International Conference on Computer and Communication Systems (ICCCS), 2019, pp. 388-394.
- [11] J Chantamit-o-pas, Pattanapong and Goyal, Madhu, "Prediction of Stroke Using Deep Learning Model", 2017, pp. 774-781.
- [12] C. Weng, C. Friedman, C. Rommel and J. Hurdle, "A Two-Site Survey of Medical Center Personnel's Willingness to Share Clinical Data for Research: Implications for Reproducible Health NLP Research," 2018 IEEE International Conference on Healthcare Informatics Workshop (ICHI-W), 2018, pp. 78-79.

- [13] F. Zhao, C. Li and C. F. Liu, "A cloud computing security solution based on fully homomorphic encryption," 16th International Conference on Advanced Communication Technology, 2014, pp. 485-488.
- [14] J. Chen, "Cloud Storage Third-Party Data Security Scheme Based on Fully Homomorphic Encryption," 2016 International Conference on Network and Information Systems for Computers (ICNISC), 2016, pp. 155-159.
- [15] A. Olumide, A. Alsadoon, P. W. C. Prasad and L. Pham, "A hybrid encryption model for secure cloud computing," 2015 13th International Conference on ICT and Knowledge Engineering (ICT & Knowledge Engineering 2015), 2015, pp. 24-32. -16
- [16] Zhao, Feng; Li, Chao; Liu, Chun Feng, 'A cloud computing security solution based on fully homomorphic encryption', Global IT Research Institute (GIRI) 16th International Conference on Advanced Communication Technology (ICACT), 2015, pp. 485 488 -19
- [17] L. d. S. Dourado, R. S. Miranda, A. P. F. de Araujo and e. E. Ishikawa, "Performance Evaluation of Big Data Applications in Cloud Providers," 15th Iberian Conference on Information Systems and Technologies (CISTI), 2020, pp. 1-6.
- [18] S. Ayyub, and P. Kaushik, "Secure Searchable Image Encryption in Cloud Using Hyper Chaos", The International Arab Journal of Information Technology, 2019, 16(2), pp. 251-259. -27

187

Preserving health care data privacy using federated learning

N. Srinivasan¹ and Dr. S. Selvamuthukumaran² ¹Research Scholar and ² Professor & Head

¹Department of ComputerApplications, A.V.C. College of Engineering, Affiliated to Anna University, Chennai, Tamil Nadu, India. itsrini6@gmail.com

²Department of ComputerApplications, A.V.C. College of Engineering, Affiliated to Anna University, Chennai, Tamil Nadu, India. smksmk@gmail.com

Abstract - Healthcare Frameworksarestructured models that provide a foundation for designing, implementing, and managing various aspects of healthcare systems and services. These frameworks help healthcare organizations, policymakers, and professionals to address various challenges in healthcare delivery, quality improvement, patient safety, and healthcare management. In conventional centralised learning, all participating institutions transfer their dataset to a centralized location, where the machine learning or deep learning model is developed. In order to achieve good robustness and generalizability of model performance, ML and DL training strategies require traditional method. Different training datasets from various sites to be transferred and pooled into a "centralized location". Such data transferring process could raise practical concerns related to data security and patient privacy. Federated learning (FL)is a distributed collaborative learning model which enables the coordination of multiple collaborators without the need for sharing confidential data. This approach has potential to ensure data privacy, increasing the model performance, reducing the data transfer costs, improving scalability, among different institutions.

Keywords: -Healthcare, Federated learning, Decentralized Learning, Federated Machine Learning, Federated Deep Learning, Data privacy.

I. INTRODUCTION

In healthcare industries, Electronics Healthcare data lead to a new wave of a scientific and technological revolution. During the last few years, the number of intelligent devices has increased exponentially with the advent of Internet-of-Things (IoT). Many of these devices are integrated with multiple sensors and increasingly powerful hardware that enable them not only to collect but more importantly, to process data on an enormous scale. Artificial Intelligence (AI) has breached every aspect of our lives, which stimulates data analytics to become a powerful field for helping healthcare industries identify opportunities and avoid risks.

Some of the health care organization, providing efficient data solutions for other institutions. Health care data seems like a "double-sidedknife-edge", which brings a variety of patient personal information leakages if the patient data are not properly used and maintained [1] [2]. Therefore, many Healthcare organizations focused on data security and deployed new strategies to handle the variety of data.



Figure 1: A model of data and privacy bottleneck in AI.

For example, the United States' California Consumer Privacy Act (CCPA) and the European Union's General Data Protection Requirements (GDPR) consecutivelyframe the rules to strengthen the protection of patient data and privacy by standardizing the behaviour of enterprises [3]. Especially, CCPA claims that, the patients have the rights to protect their personal information to access by the third parties, which is a severer restriction on the sharing of individual information for commercial purposes.

II. DATA AND PRIVACY BOTTLENECK IN AI

The "data islands" bottleneck and the emphasis on data privacy and security, as shown in Figure 1, there are the two new challenges in AI [3]. That's why the federated learning is introduced in healthcare industry, which is the decentralized training approach with features and applications shown in Figures 2 and 3, is suited for their resolution are discussed below:

zDepartment of Computer Science, A.V.C. College (Autonomous), Mayiladuthurai – 609 305. PROCEEDINGS

ISBN: 978-81-967420-1-0



Figure 2: Practical usage of federated learning in personal healthcare domains.



Figure 3: Features and applications of federated learning.

- A. Data Islands bottleneck AI has experienced certain low points in its development process which are resultant because of the lack of excellent algorithms and computing power. While things go transversely, "data isolated islands" means data is stored, maintained, and isolated from each other in different departments. In most cases, "data isolated islands" is a big challenge to integrate the data scattered in various healthcare organizations, and probably at a huge cost.
- **B.** Privacy-preserving Problem With the development of big data, it has become a global harmony to focus on data privacy and security. Once the data is leaked, it may not only compromise individuals' privacy but also cause social panic. However, driven by economic advantages, Healthcare industries usually capture patient data from many sources, such as asking patient directly, tracking patient, and appending other sources of patient data to their own. Then the data are analysed and turned into knowledge. In the era of big data, the behaviour of individuals on the Internet is triggered into data, and the collection of these data may eventually lead to the disclosure of personal privacy. In terms of the frequent incidents of personal data leakage, personal data rights and

C. institutional data rights are not equal, where consume are passive while enterprises are active. Such issues can be resolved through strict data privacy regulations. As traditional machine learning exposes more and more of its drawbacks, finding new and secure effective ways to collect data becomes crucial. Various privacy-preserving enhancement techniques and privacy-preserving machine learning solutions should be proposed in succession.



Figure 4: Difference between Machine Learning and Federated Learning.

III. PRELIMINARIES OF FEDERATED LEARNING

Compared with traditional machine learning using centralized approaches, federated learning is a decentralized training approach (e.g., split learning and large-batch synchronous Stochastic Gradient Descent (SGD), etc.) which enables smartphones located at different geographical locations to collaboratively learn a machine learning model while keeping all the personal data that may contain private information on the device. The existing federated learning can be classified into three types, namely, horizontal (or sample-based) federated learning, vertical (or feature-based) federated learning and federated transfer learning [3]. Vertical federated learning and federated transfer learning have similar types of protocols- both involve at least two participants and can be used for privacy preserving machine learning algorithms.

In a nutshell, federated learning inherits most of the features of the general machine learning with a difference of decentralized training. Another difference is that federated learning maintains the users' privacy by not uploading sensitive data to a centralized server, which is only used for sharing global updates. This feature also increases efficiency by decentralizing the training process to many devices. The requirements and architecture of federated learning are briefly introduced in the next subsection. **A.Requirements of Federated Learning**: Federated learning allows designing machine learning systems without direct access to the training data. Similar to the evolution of computing, from mainframes to client-server setups, federated learning decentralizes the machine learning with privacy by default. The key features of federated learning are, 1) Performance improves with more data. 2) Models can be meaningfully combined. 3) Edge devices can train models locally.



Figure 5: Federated Learning System (Horizontal) [3].

- **B.** System Architecture of Federal Learning: In federal learning, each edge device trains the model with its data locally and sends the small update to the central server. A horizontal federated learning technique [3] is taken as an example, shown in Figure 5, with details as follows: 1) Train global model in the server. 2) Deploy global model to edge devices. 3) Optimize model from each edge device. 4) Upload locally trained model update. 5) Average the updatevalues and apply the average to the global model. 6) Repeat step 2 to step 5. The updates in the model contain the parameters and corresponding weights, and all these updates from various users are then averaged to improve the shared global model.
- C. Two Approaches of Sending Updates Sending the update to the server is a steppingstone of federated learning to success. Currently, there are mainly two ways of attaining this: Federated Stochastic Gradient Descent (FedSGD) and Federated Averaging (FedAvg).

FedSGD.FedSGD is inspired by SGD, which is a well-established approach in the field of statistical optimization. FedSGD is an extended SGD that

assumes there are k participants Pj $(j \in [1, k])$ of the training data, and n elements in the input data while forming the global objective function. When FedSGD is to be used, each edge device needs to send gradients or parameters to the server which averages gradient or parameters and applies to new parameters. Note that the FedSGD is simple than FedAvg but needs frequent

communication between devices and servers. FedAvg. In FedSGD, each client performs gradient descent on the deployed model by using the local data, then the server calculates the average of the resulting models. The FedAvg is designed by adding more computation to each client. Specifically, FedAvg iterates the local update multiple times before the averaging step. Different from FedSGD, FedAvg enables each edge device to train and update parameters by using gradient descent iteratively. Therefore, even though FedAvg has a higher requirement for the edge devices, it results in better performance than FedSGD.

IV. Challengesand Solutions: Federated learning

plugs the most obvious and gaping security issues in distributed machine learning by leaving the training data at its source. It protects the privacy of user-data in different ways for various situations, such as by using differential privacy and homomorphic encryption. For better understanding of the challenges in FL, we primarily focus on efficiency and accuracy

A. *Challenge*: How to hide updates?

In federated learning, only global updates are sent to the central server. However, the cloud is not trusted and still allows to steal sensitive information from the data owners. For example, Phong et. al. [4] demonstrated that even a small portion of the gradients obtained by the maliciouscloud, useful information leaked by these portions are still enough to be exploited by malicious-cloud. The attack usually increases neurons and the noise in the model.

Solutions: Fully Homomorphic Encryption (FHE) is an elegant solution for this challenge, and it aims to preserve the structure of ciphers such as that addition and multiplicative operations can be performed after the encryption. All operations in a neural network except for activation functions are sum and product operations which can be encoded using FHE. Activation functions are approximated with either higher degree polynomials, Taylor series, standard or modified Chebyshev polynomials that are then implemented as part of homomorphic encryption schemes. In practice, FHE seems theoretical, and additively homomorphic encryption [5] are widely used to evaluate non-linear functions in machine learning algorithms that require balancing the tradeoffs between data privacy and prediction accuracy. Recently, Phong et al. [4] built an enhanced system to guarantee that no information is leaked to the server. Inspired by [4], all asynchronous stochastic gradients can be encrypted using the somewhat homomorphic encryption and stored on the cloud server. Then, the encrypted gradients can be applied to neural networks, properties (addition where homomorphic and multiplication) enable the computation across the gradients.

B. *Challenge*: How to optimize communication and computation complexity?

In federated learning, to predict the next word for a smartphone user when a he/she is composing a message is one of the classical scenarios.

The main reason is that mobile devices have only sporadic access to power and network connectivity. Additionally, it is difficult to establish direct and stable communication channels among mobile devices, and authenticate locally other devices that arein-charge by the service provider. Thus, how to reduce communication and computational overheads decide whether federated learning can be employed in practice while settling the trade-offs between power consumption and local training.



Figure 6: Aggregation of Vertical Federated Learning via Homomorphic Encryption [4].

Solutions: Bonawitz et al. [6] discussed the problem of computing a multiparty sum in federated learning by leveraging the spirit of secure aggregation protocol. Inspired by the work of [6], it is concluded that multi-party computing (MPC) and FHE are two important approaches in federated learning, and the above-mentioned challenge in federated learning can be solved via FHE-based MPC. Specifically, compared with Garbled circuitbased MPC, FHE-based MPC can be executed in limited rounds. Therefore, to reduce the communication and computation overhead, a constant (at most 3) rounds threshold FHE-based MPC protocol can be designed under the common reference string (CRS) model against a semi-honest adversary by combining light-weight cryptographic primitives, e.g., secret sharing, authenticated encryption, and somewhat FHE. Additionally, FHE can guarantee the privacy and confidentiality of the updates, and threshold-FHE guarantees that the approach can tolerate users dropping out of the protocol in the recovery phase (see Figure 6).

Irrespective of the promising collaboration via federated learning, some attacks [7] have demonstrated that machine learning models

Remembered too much that the privacy of the user cannot be protected. An inference attack is one of those attacks. It implies that an attacker can infer sensitive information to which it has no granted access, by using prevailing common knowledge and authorized query results. The overview of inference attacks against collaborative learning is shown in Figure 7. To the most recent, new inference attacks [8], [9] emerge endlessly and show that information about individual training data can also be inferred from the model itself, and the most indirect way requires only the ability to query the model several times.



Figure 7: An illustration of the Inference Attacks against Collaboration Learning [9].

Notably, Orekondy et al. [8] proposed two likability attacks against decentralized learning to learning generalizable user-specific patterns in the model updates. This is an identification attack to associate a user profile with a model update and a matching attack to associate two model updates with each other. In addition, Melis et al. [9] designed and evaluated several inference attacks against collaborative learning. The authors showed that an adversary can infer the presence of exact data points leading to the exposure of sensitive information, however, for a certain subset of training data.

Solutions: To address the above-mentioned inference challenge, the most often used method is differential privacy [10] that provides efficient and statistical guarantees against learning for an adversary.

The common practice to utilize differential privacy is adding noise to the data to obscure sensitive items such that the other party cannot distinguish the individual's information. Therefore, it is impossible to restore the original data, which means inference attacks become ineffective. Notably, application specific trade-off between the privacy of the training data and accuracy of the resulting model is an open question, thus, how to choose the

191

C. Challenge: How to defend inference attacks?

parameters (e.g., ε) to control this trade-off is a central issue, but the discussions on this is out of the scope of this paper. As discussed in [9], record-level ε differential privacy is an elegant approach to constitute an obstacle to the success of membership inference whereas it cannot prevent property inference. To mitigate the risks of linkability attacks, according to various strategies of Orekondy et al. [8], it is required to reduce the distinctiveness in model updates by using calibrated domain-specific data augmentation. Such a technique can provide promising results in achieving privacy with minimal impact to the utility.

D. *Challenge*: How to prevent model poisoning attacks?

According to this study, a formidable challenge is the possibility of the existence of misbehaving clients introducing backdoor functionality [11], mounting Sybil attacks [12], or label flipping attacks [13] to poison the global model, often named as the poisoning attacks. It is difficult to assert which kind of poisoning attack is the most threatening attack because they happen in different scenarios. Contrary to inference attacks, poisoning attacks happen when the adversary can inject bad data into the model's training pool, and has a chance to learn something it shouldn't. The most common result of a poisoning attack is that the model's boundary shifts in some way (see Figure 8). In fact, Bagdasaryan et al. [11] showed that stealthy backdoor functionality can be introduced into the global model in the federated learning, and designed a new approach based on the model replacement. The idea of this attack is depicted in Figure 9. Specifically, the attacker compromises one or several participants; trains a model on the backdoor data using their new constrain-and-scale technique; submits the resulting model. After federated averaging, global model is replaced by the attacker's backdoored model



Figure 8: An exemplary illustration of Poisoning Attack [13].

Solutions: There are distinct solutions to prevent model poisoning attacks. Especially, to prevent the backdoor attack, Bagdasaryan et al. [11] is a competitive one, who analysed and evaluated several defences to suggest their approach for federated

learning by specifically combining anomaly detection, Byzantine-tolerant gradient descent, and participant-level differential privacy. Alongside, to resist Sybil attacks, Fung et al. [14] proposed a new defense approach to federated learning and named it FoolsGold. Additionally, to defend model poisoning, Bonawitz et al. [6] suggested using secure aggregation because the updates from each participant are invisible to the aggregator. However, to mitigate known risks, the mentioned solutions just target one particular type of attack that happened at a different place. Thus, it is hard to convince which kind of solution is best. Furthermore, integrating these solutions into an automatic predictable model to prevent poisoning attacks depending on the actual conditions is an open question. Detecting the types of attacks, and determining an accurate solution accordingly, can be a good strategy. To resist Sybil-based poisoning attacks, one of the known defences is suggested to assume that the training data can be explicitly observed or clients can be controlled. But how to apply to federated learning for these assumptions is another problem, because the server only touches the updates from each participants' interaction. To prevent backdoor attacks, it seems to be a candidate solution that can keep their backdoor attack into limits but at the expense of sacrificing the model's performance. To prevent data poisoning attacks, the approach of participant-level privacy is highly differential recommended. Specifically, participant-level differential privacy for federated learning relies heavily on two prior works: the FedAvg algorithm which trains deep networks on user-partitioned data, and the moments' accountant of Abadi et al. [15] which provides tight composition guarantees for the repeated application of the Gaussian mechanism combined with amplification-viasampling. Another feature of participant-level differential privacy is providing a required level of privacy to each participant.

ISBN: 978-81-967420-1-0



Figure 9: An exemplary illustration of the attacker's backdoored model [11].

V. PROMISING RESEARCH DIRECTIONS

Federated Learning can be a great fit for the resourceconstrained mobile devices, Internet-ofthings (IoT),

PROCEEDINGS

zDepartment of Computer Science, A.V.C. College (Autonomous), Mayiladuthurai – 609 305.

ISBN: 978-81-967420-1-0

industrial sensor applications, and other privacysensitive use cases. Some of the promising open issues for data integrity and privacy through federated

learning along with basic research directions are shown in Figure 10. Applications for protected data including on device item ranking, next-word prediction, and content suggestion based on federated learning are the major research aspects. Recently, Google released its first production-level federated learning platform to operate sensitive data in the privacy-preserving ways that covers many federated learning-based applications. However, many tradeoffs between performance and security are waiting for us to explore. How to train the data without counting on the computational resources while users need not trade their privacy for better services is a prompt problem. Once solved, an immediate and meaningful application is the computationally inexpensive privacy-preserving for Genome-Wide Association Study (GWAS) and smart healthcare armed with FHE under decentralized settings. In particular, taking into account the confidentiality of government and business data, the use of federal learning in the joint modeling between government and enterprises can establish a complete credit system. Furthermore, federated learning for finance applications via FHEbased MPC techniques is another research direction [3]. In particular, the financial industry can form a financial data alliance that needs the collaborative effort of all financial institutions. However, one of the important obstacles is that no one wants to share his/her data in an unrequited way while he/she also would like to collaborate with other financial institutes. Hence, how to collaborate while keeping personal sensitive information by using the FHEbased MPC protocol can be an important direction to follow. Besides, with modern networking of 5G and beyond, edge-cloud integration can certainly help the easier deployment of federated learning mechanisms. However, with the availability of different functions from the 5G or beyond, it is required to decide the location of the servers as well as plan for a function that will accommodate the global updates. Moreover, aspects of authentication also need to consider when the initial model is exchanged for localized operations.

VI. CONCLUSIONS

Federated learning is revolutionizing the way machine learning models are trained. In this paper, the existing challenges in federated learning are investigated and details of the corresponding solutions are additionally provided for each problem. Several solutions for the associated challenges in federated learning are discussed, such as how to hide updates, how to

optimize communication and computation complexity, how to defend inference attacks, and how to prevent model poisoning attacks. The discussions in terms of generalized methods can be followed to build fullyfledged solutions for resolving the privacy-protection of data via federated learning.

REFERENCES

- [1] D. A. Hahn, A. Munir, and S. P. Mohanty, "Securityand privacy issues in contemporary consumer electronics [energy and security]," IEEE Consumer ElectronicsMagazine, vol. 8, no. 1, pp. 95-99, December 2018.
- [2] W. Z. Khan, M. Y. Aalsalem, M. K.

Khan, and O. Arshad,"Data and privacy: Getting consumers to trust productsenabled by the internet of things," IEEE Consumer Electronics Magazine, vol. 8, no. 2, pp. 35–38, March 2019.

- [3] Q. Yang, Y. Liu, T. Chen, and Y. Tong, "Federated machine learning: Concept and applications," ACM Transactions on Intelligent Systems and Technology, vol. 10,no. 2, pp. 12:1–12:19, January 2019.
- [4] L. T. Phong, Y. Aono, T. Hayashi, L.

Wang, and S. Moriai, "Privacy-preserving deep learning via additively homomorphic encryption," IEEE Trans. Information Forensics and Security, vol. 13, no. 5, pp. 1333-1345, May2018.

[5] A. Acar, H. Aksu, A. S. Uluagac, and

M. Conti, "Asurvey on homomorphic encryption schemes: Theory and implementation," ACM Comput. Surv., vol. 51, no. 4, pp.79:1-79:35, January 2018.

- [6] K. Bonawitz, V. Ivanov, B. Kreuter, A.
 - Marcedone, H. B.McMahan, S. Patel, D. Ramage, A. Segal, and K. Seth,"Practical secure aggregation for privacy-preserving machine learning," in Proceedings of the 2017 ACM SIGSACConference on Computer and Communications Security, CCS 2017, Dallas, TX, USA, October 30 - November 03,2017, pp. 1175-1191.
- [7] A. Pyrgelis, C. Troncoso, and E. D.

Cristofaro, "Knock, who's there? membership inference on aggregatelocation data," in Proceedings of the 25th Annual Networkand Distributed System Security Symposium, NDSS 2018.

- [8] T. Orekondy, S. J. Oh, B. Schiele, and
 - M. Fritz, "Understanding and controlling user linkability in decentralizedlearning," arXiv Computing Research vol.abs/1805.05838,2018, Repository, https://arxiv.org/abs/1805.05838.
- [9] L. Melis, C. Song, E. D. Cristofaro, and

V. Shmatikov,"Inference attacks against collaborative learning," arXivComputing Research Repository, vol. abs/1805.04049,2018, https://arxiv.org/abs/1805.04049.

- [10] I. Wagner and D. Eckhoff, "Technical privacy metrics:a systematic survey," ACM Computing Surveys (CSUR), vol. 51, no. 3, p. 57, July 2018.
- [11] E. Bagdasaryan, A. Veit, Y. Hua, D. Estrin, and V. Shmatikov, "How to backdoor federated learning," arXiv Computing ResearchRepository, vol. PROCEEDINGS 193

zDepartment of Computer Science, A.V.C. College (Autonomous), Mayiladuthurai – 609 305.

Puthal, and L. T. Yang, "Analytical model for sybil attack phases in internet ofthings," IEEE Internet of Things Journal, vol. 6, no. 1, pp. 379–387, February 2018.

[13] R. Tourani, S. Misra, T. Mick, and G.

 Panwar, "Security, privacy, and access control in information-centricnetworking: A survey," IEEE Communications Surveysand Tutorials, vol. 20, no. 1, pp. 566–600, Firstquarter2018.

[14] C. Fung, C. J. M. Yoon, and I. Beschastnikh, "Mitigating sybils in federated learning poisoning," arXivComputing Research Repository, vol. abs/1808.04866,2018, https://arxiv.org/abs/1808.04866.

B. McMahan,I. Mironov, K. Talwar, and L. Zhang, "Deep learning withdifferential privacy," in Proceedings of the ACM SIGSACConference on Computer and Communications Security, Vienna, Austria, October 24-28, 2016, pp. 308–318.

^[12] A. K. Mishra, A. K. Tripathy, D.

^[15] M. Abadi, A. Chu, I. J. Goodfellow, H.

CUSTOMER SEGMENTATION FOR ANALYSIS OF PREDICTION USING DATA MININGTECHNIQUES

A.Srilekha

Assistant Professor Department of Computer Science,

S.T.E.T Womens College (Autonomous), Sundarakottai, Mannargudi. Affiliated to Bharathidasan University,

Tiruchirappalli, Tamil Nadu.

stetcsdepartment23@gmail.com

P.Punitha

II M.Sc., Computer Science Department of Computer Science,

S.T.E.T Womens College (Autonomous), Sundarakottai, Mannargudi. Affiliated to Bharathidasan University,

Tiruchirappalli, Tamil Nadu.

Abstract

Customer relationship Management is a mediator between customer management activities in all stages of a relationship and business performance. In CRM, the big challenge is the customer retention. Customer is the soul of the any organization. To maintain the customers and should remain the loyal customers is the major role. So prediction of the customer behavior is essential, for that process data mining tools helps to predicting the customers. Customer segmentation is one of the important issues in customer relationship management. The objective of the thesis is segmenting the customers using data mining techniques such as classification and clustering and used for the prediction. Prediction is the practice of forecasting future customer behaviors. In case of segmentation, customer information is isolate customer characteristics and correlated with targeted customers. Customer analysis is carried out in all business. It refers to the process of understanding who our customers are, where they come rom and what motivates them to buy our product and services. The aim of the work is to perform automatic segmentation of customer based on buying behaviors using data mining techniques. In this information age, the companies faced major problem was retain the customers. It is methodology to build the long term profitable customers through targeted customers using customer segmentation. By analyzing the customer needs and behaviors, customers are segmented into groups. For this work, data mining techniques such as classification and clustering helps to segmenting the customers. Target customer analysis is used to analyses the customers to the suitable cluster who have similar purchase behavior. Customer retention and customer development deals with retaining the existing customers and maximizing the customer purchase value.

Keywords---CRM, Data mining techniques.

I. INTRODUCTION

Data mining has been described as the non trivial extraction of implicit, previously unknown information from data and it is the science of extracting useful information from large databases. In general, data mining is the process of analyzing data from different perspectives and summarizing it into useful information. It helps organizations to focus on the most vital information in their data repositories. It also envisages future trends and behaviours which may help organizations to make practical knowledge-driven decisions. Analysis offered by data mining system is most useful than that provided by the analysis of past events by decision support systems. The historical decision support systems respond to business queries very slowly whereas data mining systems answer very quickly. Data mining has attracted a great deal of attention due to the vast availability of huge amounts of data and the imminent need for turning such data into useful information.

Classification:

Classification is the process of creating a set of models that describe and distinguish data classes. These models can be used to predict the class of objects whose class label is unknown. The derived model is based on the analysis of a set of training data. Common applications of classification include credit card fraud detection, insurance risk analysis, bank loan approvals, etc

Clustering:

Clustering is the process of grouping the objects in a database based on the principle of maximizing intra-class similarity and minimizing interclass similarity. It differs from classification process in the sense that it uses predefined labels whereas clustering automatically predicts the class labels. Applications of clustering include market segmentation by identifying common behaviors of groups of people, discovering new types of stars in datasets of planetary objects, and so on.

Application of Predictive Analytics

- Analytical CRM
- Clinical Decision support systems
- Insurance
- Cross-sell
- Fraud detection
- Direct marketing

- Product or economy level prediction
- Text mining

Predictive models which identify patterns in historical and transactional data to determine various risks and opportunities. Forecasting models capture relationships between many factors to allow assessment of the risks potential associated with a particular of conditions, guiding decision making for candidate transactions.

II. LITERATURE REVIEW

EVOLUTION OF CRM:

The evolution of CRM can be traced to the pre-industrial era when agriculturists and artisans used to produce customized products as per their customers' requirements leading to direct interaction and subsequent bonding or a friendly relationship between them. After the advent of the industrial revolution this interaction went on the decline with middle men reducing the direct interaction between the sellers and customers. The relationship management development was dependent on contact management interaction between the seller and the customer as a result of which organizations started to evolve sales force automation (SFA) leading to the improvement in personal contact with customers and gaining information about them.

CRM IN ACQUIRING AND RETENTION OF CUSTOMERS

"Customer Relationship Management is a comprehensive strategy and process of acquiring, retaining, and partnering with selective customers to create superior value for the company and the customer. The above definition throws light on the importance of CRM as a means of acquiring new customers, for customer retention and for creating a loyal customer base while also emphasizing about role coordination amongst different departments of the organization for improving service value for the customers.

CRM AS A BUSINESS STRATEGY

"CRM is a business strategy to understand, anticipate and manage the needs of an organization's current and potential customers". The above definition throws light on the importance of CRM as a business strategy helping in achieving customer satisfaction leading to customer loyaltyand increased customer base which finally translates into increased profitability for organizations as they need to invest lesser resources for attracting new customers.

CRM TO MANAGE INFORMATION

"CRM is the process of carefully managing detailed information about individual customers and all customer touch points such as anything and any occasion that customer approaches the brand or product to maximize loyalty".

CRM FOR LONG TERM RELATIONSHIPS

"CRM can be viewed as an application of one-to-one marketing and relationship marketing, responding to an individual customer on the basis of what the customer says and what else is known about that customer".

CRM AS A TECHNOLOGY

"CRM is a term for methodologies, technologies, and ecommerce capabilities used by companies to manage customer relationships". The above definition throws light on the importance of use of technology in CRM for collecting and processing the information to understand customers' needs and wants.

IMPORTANCE AND CHARACTERISTICS OF CRM

The success of every organization depends not only on maximum utilization of resources but also making right investment decisions. The recent years had organizations getting much higher returns from investments in CRM leading to organizations realizing the importance of CRM and are also discovering the different characteristics of CRM.

IMPORTANCE OF CRM

There has been lot of research studies conducted by researchers which have highlighted the importance of adapting CRM practices in various organizations which are as follows -

MAINTAINING AND PROCESSING CUSTOMER DATABASE

The availability of customer information in the form of database which could be processed and analyzed with help of database software helped organizations to understand the different customer wants and needs.

SELECTING OF CUSTOMERS AND DEVELOPING RELATIONSHIPS

CRM is useful to identify and select customers who have potential of developing a long term healthy relationship with the organization.

DEVELOPING OF STRATEGIES FOR ACQUIRING CUSTOMERS

Once customers are selected, organizations can develop customized strategies for attracting new customers and helping in maintaining a competitive edge in the market by increasing the loyal customer base.

INCREASED SALES

Increased loyal base has potential for increased sales through repeat purchases, cross selling of products leading to increase in the profit of the organization, increased customer satisfaction. When the selected loyal customer base is targeted with the right strategies, customer satisfaction can be achieved leading to word of mouth publicity and attracting more customers further enhancing the sales of the organization and increase in profits.

FEEDBACK FOR ORGANIZATIONS

Organizations are actively focused on collecting information from the customer, which is useful for future CRM strategies in acquiring and effectively using customer and organization relationship information.

ISBN: 978-81-967420-1-0

IMPROVED MARKET KNOWLEDGE

The knowledge obtained from results of analyzing customer database is useful for the organization in determine the marketing strategies and tactics and collecting information of competitor's plan and strategy.

III.METHODOLOGY

The research focuses on the areas of research in CRM Segmentation namely the Association rules mining. The rules mined from traditional frequent association rule mining approaches reflect the regularities and general trends, and provide, interesting insight to the user.

This process can be experimented by the following aspects:

- Define the Problem
- Prepare the Data
- Explore the Data
- Build mining Models
- Explore and Validate the Models
- Deploy and Update the Models

Apriori Algorithm

Apriori is a rapid algorithm for generating the association rules. This algorithm uses a "bottom-up" methodology and contains two steps. In the first step, the task is to determine all large item sets having support value greater than minimum support. The candidates are generated from the priorvia a fast and clever cross-product function. Then candidates support value is determined. All candidates with support greater than min sup (Minimum Support Threshold) are placed in thenext L(Frequent item set). This process stops when L_n is empty.

Association is very helpful for decision-making and effective in marketing sector. The advantage of this algorithm is that it is very easyto implement and uses large item set property. The main drawbacks of the association rule algorithm are the following:

- Repeated scanning of the database
- Huge Candidate sets.
- Generate the frequent data item sets relying on minimum support
- Huge number of discovered rules
- Large storage space

IV.RESULTS AND DISCUSSION

Our results suggest the impact of IT infrastructure on superior CRM capability is indirect and fully mediated by human analytics and business architecture. We also find that CRM initiatives jointly emphasizing customer intimacy and cost reduction outperform those taking a less balanced approach. Overall, this paper helps explain why some CRM programs are more successful than others and what capabilities are required to support success. Customer relationship management suffers when it is poorly understood, improperly applied, and incorrectly measured and managed. This study reveals the combination of investment commitments in human, technological and business capabilities required to create a superior CRM capability. This exploratory study was conducted to analyze customer satisfaction and business intelligence activities of industry sector. The study proves that customer satisfaction is the major goal of the insurance sector. This is because customer satisfaction and business intelligence are major determinant of sustainable competitive advantage in the insurance market.

Comparison Results

The overall comparison of the cluster assignment instances are listed in the below table.

Name of	No of Cluster	Instances in	Instances in		
the Cluster	Cluster	Cluster1	Cluster 2		
K-Means	2	752	238		
Hierarchical	2	598	402		
EM	2	643	357		





Cluster Assignments

V. CONCLUSION

The aim of the work is use of data mining techniques in CRM. Needs/attitudinal segmentation are commonly employed through marketer search data in the industry to gain insight into the customer attitudes, wants, views, preferences, and opinions about the enterprise and the competition. In addition to external/marketresearch data, transactional data can also be used for the development of effective segmentation solutions. A value-based segmentation scheme allocates customers to groups according to their spending amount. It can be used to identify high-value customers and to prioritize their handling according to their measured importance. The clustering is used for grouping the customers.

The information obtained from the customers is helpful to the organization to full fill the needs and wants of their customers. The organization can also earn profit and maintain a long term relation with the customers, but they should ensure that the privacy of the customers, should not get affected during usage of the data. The data mining with BI is the most common trend in field of Research.

FUTURE WORK

The proposed model of data mining for CRM have some limitations, is an integration of the existing models of data mining and CRM, with a reference to the various techniques that could be adopted as well as the various applications of the results of the analysis. Now-a-days, Business intelligence for CRM applications provides a firm with actionable information from the analysis and interpretation of vast quantities of customer/market related data. Databases for business intelligence include customer demographics, buying histories, cross-sales, service calls, website navigation experiences and online transactions. Here, in future the concepts applicable to all types of industrial applications through the appropriate use of analytical methods and software, which inturn helps the firm gain competitive advantage by creating greater valuefor the customer.

REFERENCES

- Sheth, J.N., Parvatiyar, A., &Shainesh, G. (2001). Customer Relationship management: Emerging Concepts, Tools and Applications. New Delhi: *Tata McGraw-Hill Publishing Company Limited*.pp.6-7.
- [2] Callaghan, M., Shaw, R. (2001). Relationship orientation: Towards an antecedent m&odel of trust in marketing relationship. Massey University Auckland. NZ: Proceedings of the Australian and NewZealand Marketing Academy Conference. In S.Chetty& B. Collins(Eds.).pp.1-9.
- [3] CRM Community (2002). CRM Community Library Fundamentals.
- [4] Wahab, S., & Ali, J. (2010, Dec). The Evolution of Relationship Marketing (RM) Towards Customer Relationship Management (CRM): A Step towards Company Sustainability. *Information Management andBusiness Review*, 1(2), 88-96.
- [5] Gronroos, Christian (1995). Relationship marketing : The Strategy Continuum. *Journal of the Academy of Marketing Science*, Fall, 252-254.
- [6] Crosby, L. A., Kenneth, R. E., & Deborah, C. (1990). Relationship Quality in Services Selling – An Interpersonal Influence Perspective. *Journal of Marketing*, 52 (April), 21-34.
- [7] Janakiraman, B., &Gopal, R.K. (2006). Total quality management: text and cases. New Delhi: *Prentice Hall of India Private Limited*. (pp.205-215).
- [8] Ngai, E. W. T. (2003) Internet Marketing Research (1987-2000): A Literature Review and Classification. *European Journal of Marketing*, 37(1/2), 24-49.

- [9] Frazier, G. L., Spekman, R.E., & O'Neal, C. (1988). Justin- Time Exchange Systems and Industrial Marketing. *Journal of Marketing*, 52. (October), 52-67.
- [10] Wahab, S., & Ali, J. (2010, Dec). The Evolution of Relationship Marketing (RM)Towards Customer Relationship Management (CRM): A Step towards Company Sustainability. *Information Management and Business Review*, 1(2),88-96.
- [11] Gummesson, E. (2004). Return on relationships (ROR): Thevalue of relationship marketing and CRM in business-tobusiness contexts. *The Journal of Business and Industrial Marketing*, 19(2), 136-148.
- [12] Parvatiyar, A., & Sheth, J.N. (2001). Customer Relationship Management: Emerging Practice, Process, and Discipline. *Journal of Economic and Social Research.* 3(2), *Preliminary Issue* (2002), 1-34.
- [13] Payne, A., &Frow, P., (2005). A Strategic Framework for Customer Relationship Management. *Journal of Marketing, American Marketing Association*, 69 (October), 167–176
- [14] Kotler, P., & Keller, K. (2011). Marketing Management. (14th ed.). New Delhi: *Pearson Education*.
- [15] Peppers, D., Rogers, M., &Dorf, B. (1999). Is Your Company Ready for One-to-One Marketing?.*Harvard Business Review*, 77 (January–February), 151–160.
- [16] Payne, A., &Frow, P., (2005). A Strategic Framework for Customer Relationship Management. *Journal of Marketing, American Marketing Association*, 69 (October), 167–176
- [17] Nyarku, K.M. (2013). Assessing Customer Relationship Management (CRM) Practices at National Investment Bank (NIB) Ghana Limited: (A Study of the Cape Coast Branch). *International Journal of Advances in Management and Economics*, 2 (2, Mar.-April), 151-162.
- [18] Antar, J., &Gholamifar, D. (2006, June). Customer Relationship Management in Fashion Companies for men's wear (Master's dissertation, Jonkoping International Business School, 2006).
- [19] Greenberg, P. (2004). CRM at the speed of light. New Delhi: *Tata Mcgraw Hill.*
- [20] Reinartz, W., &Kumar, V. (2000). On the Profitability of Long-Life Customers in a Noncontractual Setting: An Empirical Investigation and Implications for Marketing. *Journal of Marketing*, 64 (October), 17–35.

DEEP LEARNING AND APPLICATIONS

Nagaraj. M

I year MCA, Department of Computer Applications, A.V.C.College of engineering, Mannampandhal, Mayiladuthurai-609305. mnagamurugan07@gmail.com

Abstract

Deep learning is a powerful multi-layer architecture that has important applications in image processing and text classification. This paper first introduces the development of deep learning and two important algorithms of deep learning: convolutional neural networks and recurrent neural networks. The paper then introduces three applications of deep learning for image recognition, image detection, and image forensics, as well as three text classification methods based on convolutional neural networks, recurrent neural networks, and other text classification methods. Finally, the development trend of deep learning in the field of text and image processing and the difficulties to be further researched are summarized and prospected.

---00000000----

BLUE EYES TECHNOLOGY

Nithiya Bharathi B^1 and , Sathiyadevi S^2

I year MCA Student Department of Computer Applications A.V.C. College of engineering, Mannampandhal, Mayiladuthurai - 609305 shiva4015p@gmail.com, sathiydevi0503@gmail.com

Abstract

The world of science cannot be measured in terms of development and progress. It showshow far human mind can work and think. It has now reached to the technology known as "Blue eyes technology" that can sense and control human emotions and feelings through gadgets. The eyes, fingers, speech are the elements which help to sense the emotion level of human body. This paper implements a new technique known as Emotion Sensory World of Blue eyes technology which identifies human emotions (sad. happy. excited or surprised) using image processing techniques by extracting eye portion from the captured image which is then compared with stored images of data base. After identifying mood the songs will be played to make human emotion level normal.

---00000000----

Coded Cryptosystem

A .Kavya and M. Madhavan

UG CS Student, Nehru Memorial College, Puthanampatti, Trichy- 621 007. kavyaanandhakumar004@gmail.com madimah2004@gmail.com

Abstract

With the advent of increasing the internet, the message transmitted through it must be protected so that an adversary could not recover the original message called plaintext. For that cryptography is being used. When the message is sent through the communication channel, it is assumed that the message transmitted after encrypted called cipher text may reach the receiver always error free. But, in practical life situation it is not so. This is because the communication channel is always noisy and the errors can happen at any time. Thus, the cipher text must be coded using some coding techniques for detecting and correcting the errors before it is transmitted. Once the coded cipher text reaches the receiver, the errors must be corrected and then decrypted so that original plain text is obtained. If the error is single bit error, it is corrected using Hamming code and if it has errors they are corrected using Bose-Chaudhuri-Hocquenghem code. Thus, this paper focus on cryptography along with coding theory and it is termed as coded cryptosystem.

---00000000----
International Conference

On

COMPUTIONAL INTELLIGENCE AND ITS APPLICATIONS (ICCIA - 2024)

About the College

The Anbanathapuram Vahaira Charities, an organization with a track record of more then two centuries was established in 1806 by the philanthropic members of the five families. In 1955, the A.V. Charities established the A.V.C. College to serve the cause of higher education. A.V.C. College (Autonomous), Mayiladuthurai is part of the history of Anbanathapuram Vahaira Charities (AVC). The College is affiliated to Annamalai University. It was founded in 1955. The College started its Post Graduate Programme in 1970 and celebrated its Silver Jubilee in 1980. The College started the Evening Section for women students in 1984. With its commendable performance, the College was conferred the 'Autonomous' status in 1987 and the same is extended till date. The College celebrated its Golden Jubilee in the year 2005. The College has celebrated its Diamond Jubilee in 2016.



About the Department

The College is one among the leading Institutions imparting quality Computer education from 1984. The College offers several Computer Science courses, namely B.Sc., B.C.A., M.Sc., and M.Phil. The B.Sc Computer Science (Aided) started in the year 1984, M.Sc., Computer S cience in the year 1987, B.Sc., Computer Science (Evening) in the year 1988, B.C.A., in the year 2001 and M.Phil., Computer Science in the year 2015. The main objective of the Department is to educate the student to meet out the growing demand for the Computer professionals in the dynamic IT industry and Research field. The Department has the Vision of Empowerment of rural poor through quality Computer education and the Mission of Providing venue to technological innovations through Computer and Internet education.

About the Conference

The International Conference on Computational Intelligence and its Applications aims to bring together researchers to exchange and share their experiences and research results on all aspects of Computational Intelligence. Computational Intelligence is a subfield of Artificial Intelligence (AI) that deals with the design and development of intelligent computer systems. Computational Intelligence is a branch of Artificial Intelligence that deals with creating algorithms and systems that can learn from data and make decisions based on what they have learnt. The primary applications of Computational Intelligence include Image Processing, Natural Language Processing, Data Mining, Big data Analytics, Artificial Intelligence and Robotics etc. By utilizing these applications, it has become possible for machines to replicate human behaviour. This is our first International Conference and the research articiles are invited from various Institutions. Many research scholars, research supervisors and PG students actively sent the research articles and the peer reviewed research papers were published in this edited volume.



