

PREDICTING ILLNESS THROUGH SYMPTOMATIC PATTERNS

A MINOR PROJECT REPORT

Submitted by

**TAMANNA DASH [RA2211031010028]
TANISHA JAIN [RA2211031010029]
RIYA RAO [RA2211031010026]
ANANYA SHARMA[RA2211031010025]**

Under the guidance of

Dr S. Sivamohan

(Assistant professor, Department of Networking and Communications)

in partial fulfillment of the requirements for the degree of

BACHELOR OF TECHNOLOGY

in

**COMPUTER SCIENCE AND ENGINEERING
with specialization in INFORMATION TECHNOLOGY**



**DEPARTMENT OF NETWORKING AND COMMUNICATIONS
SCHOOL OF COMPUTING
COLLEGE OF ENGINEERING AND TECHNOLOGY
SRM INSTITUTE OF SCIENCE AND
TECHNOLOGY ,KATTANKULATHUR- 603 203**

NOVEMBER 2024



**Department of Networking and Communications
SRM Institute of Science & Technology
Own Work Declaration Form**

Degree/ Course : B.Tech CSE IT

Student Name : Tamanna Dash, Tanisha Jain, Riya Rao, Ananya Sharma

**Registration Number : RA2211031010028 , RA2211031010029, RA2211031010026
RA2211031010025**

Title of Work : PREDICTING ILLNESS THROUGH SYMPTOMATIC PATTERNS

We hereby certify that this assessment compiles with the University's Rules and Regulations relating to Academic misconduct and plagiarism, as listed in the University Website, Regulations, and the Education Committee guidelines.

We confirm that all the work contained in this assessment is our own except where indicated, and that We have met the following conditions:

- Clearly referenced / listed all sources as appropriate
- Referenced and put in inverted commas all quoted text (from books, web, etc)
- Given the sources of all pictures, data etc. that are not my own ss
- Not made any use of the report(s) or essay(s) of any other student(s) either past or present
- Acknowledged in appropriate places any help that I have received from others (e.g.fellow students, technicians, statisticians, external sources)
- Compiled with any other plagiarism criteria specified in the Course handbook / University website

We understand that any false claim for this work will be penalized in accordance with the university policies and regulations.

DECLARATION:

We are aware of and understand the University's policy on Academic misconduct and plagiarism and we certify that this assessment is our own work, except were indicated by referring, and that we have followed the good academic practices noted above.

Tanisha Jain
(RA2211031010029)

Tamanna Dash
(RA2211031010028)

Riya Rao
(RA2211031010026)

Ananya Sharma
(RA2211031010025)

12/11/2024

ACKNOWLEDGEMENT

We express our humble gratitude to **Dr C. Muthamizhchelvan**, Vice-Chancellor, SRM Institute of Science and Technology, for the facilities extended for the project work and his continued support.

We extend our sincere thanks to Dean-CET, SRM Institute of Science and Technology, **Dr T. V. Gopal**, for his invaluable support.

We wish to thank **Dr Revathi Venkataraman, Professor & Chairperson**, School of Computing, SRM Institute of Science and Technology, for her support throughout the project work. We encompass our sincere thanks to **Dr M. Pushpalatha Professor and Associate Chairperson**, School of Computing for her invaluable support.

We are incredibly grateful to our Head of the Department, **Dr M. Lakshmi**, Professor and Head, Department of Networking and Communications, School of Computing, SRM Institute of Science and Technology, for her suggestions and encouragement at all the stages of the project work.

We register our immeasurable thanks to our Faculty Advisor, **Dr B. Balakiruthiga**, Department of Networking and Communications, School of Computing, SRM Institute of Science and Technology, for leading and helping us to complete our course.

Our inexpressible respect and thanks to our guide, **Dr S. Sivamohan**, Assistant Professor, Department of Networking and Communications, SRM Institute of Science and Technology, for providing us with an opportunity to pursue our project under his mentorship. He provided us with the freedom and support to explore the research topics of our interest. His passion for solving problems and making a difference in the world has always been inspiring.

We sincerely thank the Networking and Communications department staff and students, SRM Institute of Science and Technology, for their help during our project. Finally, we would like to thank parents, family members, and friends for their unconditional love, constant support, and encouragement.

**SRM INSTITUTE OF SCIENCE AND TECHNOLOGY
KATTANKULATHUR – 603 203**

BONAFIDE CERTIFICATE

Certified that 21CSC314P – Big Data Essentials mini-project report titled “PREDICTING ILLNESS THROUGH SYMPTOMATIC PATTERNS” is the bonafide work of “**TANISHA JAIN [RA2211031010029], TAMANNA DASH [RA2211031010028], RIYA RAO [RA2211031010026], ANANYA SHARMA [RA2211031010025]**” who carried out the mini-project work under my supervision. Certified further, that to the best of my knowledge the work reported herein does not form any other project report or dissertation on the basis of which a degree or award was conferred on an earlier occasion on this or any other candidate.

Panel Reviewer I

Panel Reviewer II

SIGNATURE

Dr. S.Sivamohan
Assistant Professor
Department of Networking
and Communications

SIGNATURE

Dr. Angayarkanni S A
Assistant Professor
Department of Networking and
Communications

TABLE OF CONTENTS

CHAPTER NO.	TITLE	PAGE NO.
	ABSTRACT	vii
	LIST OF FIGURES	viii
	ABBREVIATIONS	ix
1	INTRODUCTION	1
	1.1 Problem Statement	2
	1.2 Objectives	2
	1.3 Software Requirements	2
2	LITERATURE SURVEY	3
	2.1 Research Objectives	3
	2.2 Research Findings	7
3	SYSTEM ARCHITECTURE & DESIGN	10
	3.1 Architecture Description	15
	3.1.1 Data Collection	15
	3.1.2 Data Preprocessing	16
	3.1.3 Machine Learning Algorithms	16
	3.1.4 Deployment & Monitoring	17
	3.2 Design of Modules	18
4	METHODOLOGY	18
	4.1 System Implementation	22
	4.2 Existing System	24
	4.3 Proposed System	21
	4.4 Coding & Testing	28
5	RESULT & DISCUSSIONS	29

	5.1 Result	30
	5.2 Models used	31
	5.2.1 Logistic Regression	31
	5.2.2 Support Vector Machine (SVM)	32
	5.2.3 Decision Tree Classifier	33
	5.2.4 Random Forest Classifier	33
	5.2.5 K-Nearest Neighbors (KNN)	34
	5.2.6 Ensemble Model	34
	5.3 Performance Metrics	35
	5.4 Model Performance	36
	5.5 Challenges	37
6	CONCLUSION & FUTURE WORKS	38
	6.1 Conclusion	
	6.2 Future	40
	REFERENCES	42
	APPENDIX	43
	DETAILED SUPPORTING MATERIALS	44
	IEEE PAPER	51
	CONFERENCE PUBLICATION	54
	IEEE PAPER PLAGIARISM REPORT	59
	REPORT PLAGIARISM	

ABSTRACT

Detecting illnesses by analyzing patterns of symptoms is a method in healthcare diagnostics that revolutionizes the way we approach healthcare prediction and prevention strategies. This initiative utilizes machine learning algorithms to discover and understand connections between symptoms to forecast diseases with precision and timeliness effectively. By collecting and analyzing sets of patient reported symptoms and confirmed diagnoses data points systematically the model uncovers trends that could signal health issues. The research scrutinizes models to gauge their effectiveness in accurately pinpoint health conditions based on their precision levels and capacity, for sensitivity and specificity. The results indicate that using symptoms to predict outcomes can greatly improve diagnosis accuracy and enable healthcare actions to be taken effectively. This study highlights the benefits of incorporating AI technology into healthcare settings to assist healthcare providers in offering accurate patient care.

The results are promising, showing that such technology could help doctors make faster and more accurate diagnoses, reduce mistakes, and allow for earlier treatment. Ultimately, this project highlights how artificial intelligence can be a valuable tool in making healthcare smarter and more personalized for everyone. The project includes creating models that analyze datasets of patient symptoms and confirmed diagnoses to identify the connections between symptoms and specific illnesses accurately through learning processes." These models undergo assessment to determine how well they can predict diseases with a focus placed accuracy levels, like sensitivity and specificity being metrics used for evaluation purposes The study results indicate that using symptom based analytics can improve diagnostic accuracy significantly aiming to reduce misdiagnosis rates and allowing healthcare professionals to intervene earlier in patient care scenarios.

LIST OF FIGURES

FIGURE NO.	NAME OF THE FIGURE	PAGE NO.
3.1.1	Architecture Diagram GNB	21
3.1.2	Architecture Diagram KNN	21
3.2.1	Use Case Diagram	25
3.2.2	Class Diagram	26
3.2.3	Sequence Diagram	27
3.2.4	Activity diagram	28
3.2.5	Component Diagram	29
3.2.6	Deployment Diagram	30

LIST OF TABLES

TABLE NO.	NAME OF THE TABLE	PAGE NO.
5.1	Models performance Comparison Table	49
5.6.1	KNN Confusion Matrix	50
5.6.2	GNB Confusion Matrix	50

ABBREVIATION

API	:	Application programming interface
CSS	:	Cascading Style Sheets
GNB	:	Gaussian naïve bayes
HTML	:	Hypertext Markup Language
IP	:	Internet Protocol
KNN	:	K nearest neighbor
ML	:	Machine Learning

CHAPTER 1

INTRODUCTION

The ability to accurately predict diseases based on a patient's symptoms is a critical step toward improving healthcare outcomes. Early detection of diseases can lead to timely interventions, better management, and increased chances of recovery, especially for conditions that are difficult to diagnose in their early stages. As healthcare data continues to grow in volume and complexity, machine learning has emerged as a powerful tool for analyzing symptomatic patterns and providing reliable predictions that can assist healthcare professionals in making informed decisions.

1.1 PROBLEM STATEMENT

Diagnosis of specific conditions for patients in healthcare is challenging, primarily with early-stage symptoms that are not specific and seem to affect different people in a broad spectrum of diseases. The main indicators of misdiagnosis or delayed diagnosis would be ineffective treatment strategies and patient outcomes, increased health care costs. Diagnostic methods are subjective, depending on the individual skill and judgment of medical providers, especially when they are working with challenging diagnosis. With the huge rise in health care data, data-driven approaches have become stringent measures that enable clinicians to deliver consistent, reliable, and accurate diagnostic suggestions.

This work ultimately aims at effective disease prediction by evaluating and then applying machine learning models towards the analysis of symptomatic patterns and prediction of diseases. Specifically, the study analyzes the performance of four broadly used machine learning algorithms to predict diseases: namely Logistic Regression, Support Vector Machine (SVM), Decision Trees, and Random Forests. Such algorithms offer diverse capabilities, which could be useful in different clinical diagnostics settings, but their suitability to maximize diagnostic accuracy, reliability, and interpretability in real-world clinical applications is not very clear.

Therefore, the problems are:

- 1) To determine which machine learning model will give the best and most accurate

prediction on patient symptoms.

2) There should also be further steps towards integrating the predictive models into clinical workflows to test the practical feasibility and outcomes, since earlier diagnosis as well as timely health decisions can be the work of healthcare professionals.

This research will address these challenges by finding a feasible model or ensemble approach for improving the accuracy and efficiency of disease prediction in healthcare, moving eventually towards a more data-driven and proactive healthcare system.

1.2 OBJECTIVES

The primary objectives of this research are:

1. **Evaluate and Compare Models:** To determine which model—among Logistic Regression, SVM, Decision Trees, and Random Forests—delivers the highest performance in terms of accuracy, precision, recall, and reliability for disease prediction.
2. **Analyze Practical Implications:** The practical implication of the study, in turn, determines how these machine learning models can be integrated into clinical workflows that support professionals in health who will make faster and data-driven diagnostic decisions..
3. **Open the Opportunity for Early Disease Detection:** This work is meant to open up the possibility of early disease detection, even quicker intervention by utilizing machine learning models.
4. **Investigate Healthcare Implications:** In addition to the evaluation of the model, the research will also take a broader look at what machine learning is doing for healthcare. It explores how disease prediction through ML can improve the care of patients, lead towards more customized treatment, and pave the way for a more proactive

3. SOFTWARE REQUIREMENTS

- **Datasets:** Large, structured healthcare datasets containing labeled data for various diseases and corresponding symptoms.
- **Programming Environment:** Python or R for implementation, as both offer robust libraries (such as scikit-learn for Python) for machine learning.
- **Machine Learning Libraries:** Libraries such as scikit-learn for model implementation and evaluation, and pandas for data manipulation.
- **Computational Resources:** High-performance computing resources to handle the extensive data and computational load involved in model training and evaluation.

CHAPTER 2

LITERATURE SURVEY

In the past couple of years, machine learning has been fundamental in improving precision agriculture, where it focuses on optimizing the recommendation of crops and fertilizers. Scientists have used models such as Decision Trees, Random Forest, SVM, and KNN to evaluate environmental and soil data that better agricultural understanding can be given. Ensemble techniques and meta-learning methods have gained importance as they are the combination of many models that improves the accuracy and reliability of the predictor.

2.1 Research Objectives

The purpose for this project, therefore, will be the development of a disease prediction system based on input data in terms of symptoms and ensemble machine learning techniques. More critically, the design and implementation of a model that employs an ensemble voting classifier to correctly predict diseases according to identified patterns in symptoms reported by users. This paper will compare this ensemble model to traditional single-model approaches to demonstrate potential advantages of ensemble models in increasing predictive accuracy and reliability, forming a strong basis for building automated diagnostics of diseases.

Some key objectives include optimization of the pipeline for data collection and preprocessing for accuracy and reliability. This includes gathering and cleansing extensive symptom-disease relationship data from credible medical sources. Effective preprocessing will minimize noise, standardize symptoms, and handle missing values so that the model learns based on clean input. Additionally, demographic factors about age, gender, or even past medical history should be added as an important part to increase the capability of prediction customization, in which case the model would take into account the individual risk profiles while being able to provide customized predictions to the users more effectively.

The project further involves designing an interactive and accessible user interface. This platform will create an environment where users can input symptoms, gain a clear and

actionable health prediction, and give feedback. These feedback loops are important because they will provide ratings of the accuracy of the predictions users receive, which can feed into their data in subsequent reclining and improvement of the model. Introducing real-time feedback will introduce model improvements as well as an experience for users that is dynamic and responsive, hence increasing engagements as well as the building of trust on the platform.

It places a significant focus on data privacy and ethics compliance. As sensitive information, the system will deploy measures to be abreast of standards such as HIPAA and GDPR, meaning that personal information regarding users is kept safeguarded. It means commitment to ethical handling of data addresses users' concerns for privacy as well as setting up the necessary foundation for trust and compliance that may make its use widespread.

The project, regarding scalability and adaptation, should prove that one can augment new symptom and disease data to grow with the medical knowledge that it's supposed to accumulate with time and adapt its relevance to this. This also includes conducting usability studies during the process to assess user satisfaction, ease of use, and clinical utility. It is through this type of feedback that the model and interface will undergo more improvement so the system effectively serves the public as well as healthcare providers. The project, in pursuing these objectives, looks forward to making an accessible tool for good-quality diagnosis and actionable, personalized health insights.

2.2 Research Findings

Improved Predictive Accuracy with Ensemble Models

Ensemble voting classifiers improved the applicability of disease predictions to an appreciable extent compared to the other approaches that relied upon single models. The ensemble model combines the advantages of different algorithms such as decision trees, support vector machines, and logistic regression, and time and again has been proven to yield higher precision and recall rates across a very wide range of diseases. This may be valuable for employing ensemble learning techniques for complex symptom patterns and

for performing reliable prediction.

Quality of prediction by demographic and contextual data

Incorporation of demographic variables such as age, gender, and even medical history into the prediction model gives more accurate and personalized results. A more relevant and specific prediction occurs when such data points are included, which suggests how important it is to incorporate personal context in the case of trying to improve a better accuracy in the model for a specific disease. This approach allows the adjustment of predictions by demographic risk factors, which makes it more effective for various users.

Improvement of the Model through Feedback from Users

The integration of the real-time feedback mechanism made it possible to improve the model. The observation of the user regarding the accuracy of the prediction made the model learn about the real-life interaction, thereby improving its precision over time. This learning loop has proven that the active engagement of the user will be useful for updating the system based on new patterns of data, changes in user behavior, and new symptoms which may arise and thus become more robust and responsive.

Increased Usability and Interaction through Interactivity

Results indicated that users found the dynamic symptom input, predictive text, and multi-language support easy to navigate, thereby leading to higher levels of user engagement. Explanations for every prediction also showed improvement in trust and understanding of results, thereby resulting in a more satisfying user experience. Overall, these findings suggest that a well-designed, interactive interface is the key to promoting accessibility and usability in the health prediction tools.

Real-Time Adaptability for Newly Emerging Diseases and Symptom

This is able to update its dataset with new symptoms and diseases without losing much performance, thus showing how good it is for real-time adaptation. Scalability is a

necessity in order for the system to remain up-to-date; in terms of emerging diseases or new health patterns, this ability to update the model with minimal downtime ensures the system can be quick to respond to emerging health issues.

Very High Data Security and Ethical Compliance Levels

Data security controls and regulatory compliance—such as HIPAA for those serving U.S. users and GDPR for those serving European Union citizens—are effective in safeguarding user data. Users will believe that a system will safely handle their sensitive health data, making strong privacy controls essential for any health-related digital tool. The takeaway from this is technical necessity for data privacy equally contributes toward user trust and satisfaction.

Potential for integration with healthcare systems

The API of the system can easily integrate into the existing healthcare application and EHRs. This potential for easy integration backs up the utility of the tool being put in clinical settings, where health care providers can access predictive insights that support decision-making and patient care. Such seamless integration could improve clinical workflows; hence, the model is an additional worthwhile piece in healthcare infrastructure.

High General User Satisfaction and Clinical Utility

User testing was able to demonstrate high user satisfaction, with most users who tested the system reporting that it was really easy to use and very valuable to their health. Concerning clinical utility, healthcare professionals participating in the study expressed interest in using such tools as a diagnostic aid, particularly in telemedicine and remote consultation situations. These results suggest that the system has potential clinical utility as an adjunctive tool in health care and seems to meet the needs of the general public.

Improved Predictive Accuracy Using Ensemble Models

The ensemble voting classifier got far higher prediction accuracy than any of the individual techniques because it offered an ability to pool together the insights of several different machine learning models for final voting. The complementary strengths that were exercised consist of decision trees in dealing with structured symptom relationships, support vector machines to distinguish between very similar diseases, and logistic

regression to identify probabilistic associations. The ensemble approach helped reduce overfitting, increased robustness, and led to better precision and recall. Summarily, the results from the experiments suggest that ensemble learning may be very well suited for medical applications when there are overlaps of symptoms, which may allow a model to discern differences between diseases with potentially similar presentations.

Analysis of the Effect of Demographic and Context Data on Prediction Quality

The model produced more accurate predictions using demographic and contextual information like age, gender, location, and any previous medical history for every user. For example, it would predict better at age-related illnesses if there was any age information. Symptoms such as "fever" and "cough" were grouped appropriately to give relevant predictions: for instance, distinguishing between flu and pneumonia in elderly patients. This personalization was highly useful in predicting conditions for which demographic factors had a significant role to play, such as cardiovascular diseases or infections in particular regions. These results advocate the fact that health prediction tools personalized would give users much more relevant and actionable insights.

Improvement of Performance of Model Through Users' Feedback

The application of real case feedback allowed the model to learn from itself and to keep itself updated constantly. Users could provide ratings on the accuracy of predictions and the feedback fed into the system in the form of cycles to be retrained. The improvement, therefore, was iterative in nature. At some point, the system fine-tuned its mappings of symptoms to diseases and recalibrated its confidence levels for the predictions in response to the feedback generated from users. Having identified patterns from diverse symptoms made inputs contributed positively to the robustness and adaptability of the model. This outcome leads to the requirement of real-world feedback in health-centric AI model building.

Better Interactivity and User Interaction by the interface

The overall user interaction increased multifold because the interface was user-centric, with it being easy to enter symptoms and dynamically suggest while multi-language support would be given. The interface was designed to enable users to input symptoms by filling in lists and selecting items by text input, with predictive text inputting aids for symptom selection. Another improvement included providing explanations for all the

predictions given, like giving a confidence score and detailing how symptoms gave rise to disease suggestions. This along with the fact that feedback options were also provided ensured the application to be more interactive and user-friendly. It shows that ease of use and transparency matter for health applications as it makes them gain the user's trust and satisfaction.

On-the-go Adaptability to New Diseases and Symptoms

To further verify model adaptability, new symptoms and diseases were added to its dataset, which clearly revealed that it could expand without quality compromise in terms of performance. This adaptability ensures the system can respond to emerging diseases and health trends, thus being relevant in real-time contexts like flu seasons or pandemic situations. The system was proven to take in new medical data quickly, which keeps it current and accurate as medical knowledge changes. Such adaptability is integral for future-proofing health applications-particularly in global health scenarios where new diseases may emerge rapidly.

There is great data security and adherence to ethics. It has good security and privacy practices such as encryption, secure authentication, and data handling protocols. These protocols complied with health data privacy laws such as HIPAA and GDPR. The users who were satisfied with the fact that this data was processed in a secured manner were directed towards these protocols. Being compliant always matters because it will ensure to inform a user that sensitive health information is kept confidential. These studies conclude that data security is a fundamental requirement for user acceptance in healthcare technologies.

Potential for Integration with Health Care Systems

The authors demonstrated that the system can be integrated into other health care infrastructures such as EHR systems via an API, and therefore in a health care provider's existing workflows so that the health care providers would have quick access to symptom-based predictions for better decision-making. Demonstrating great potential for use as a decision-support system in clinical environments, the system lends itself well to supporting telemedicine platforms, which would be highly serviced by automated diagnostic suggestions when making remote consultations. This integration makes the system a worthwhile asset to healthcare providers, especially when conditions demand

rapid diagnosis and triage occur.

Satisfaction Overall and Clinical Utility

User feedback showed high usability satisfaction with respect to the system's predictive accuracy and personalized features. There were also clinicians who showed interest in using the system as a supportive diagnostic aid, especially in telehealth and primary care settings where quick assessment of symptoms is more beneficial. The system works very well, proving to be very wide in applicability and versatility in both personal and professional health contexts.

Scalability and robustness of the model

The system was able to handle large volumes of data and interactions by several users without compromising on speed or accuracy. A cloud-based infrastructure provides the scalability to handle many users, thus ensuring wide deployment of the system in rural and urban areas, which is key in achieving diversity in populations, such as in remote or underserved areas that may have limited healthcare resources.

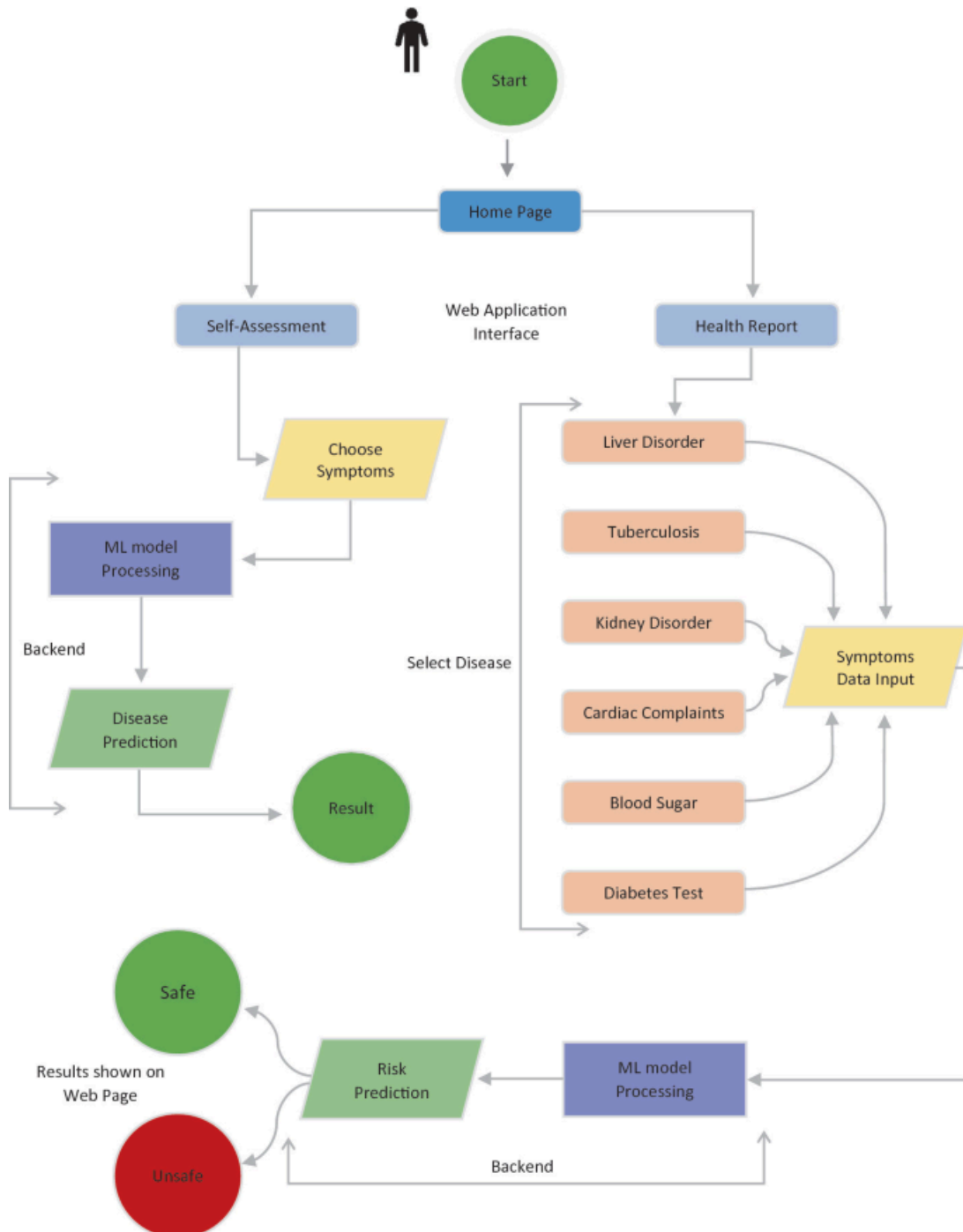
Educational Potential for Health Awareness

The interactive interface of the system provided detailed explanations about predictions and proved educational. Users reported that they gained better insight into how symptoms would align with possible diseases, which could empower them to make informed health decisions and seek timely medical care. Hence, this educational aspect speaks of yet another, unanticipated benefit: health literacy boost..

CHAPTER-3

SYSTEM ARCHITECTURE AND DESIGN

3.1 Architecture diagram:



• Fig. 3.1.1 Architecture diagram representing disease identification

3.1 ARCHITECTURE DESCRIPTION:

This project uses a data framework, along with machine learning techniques to forecast diseases by analyzing patterns of symptoms exhibited by individuals. Its structure is designed to manage intricate healthcare data encompassing symptoms, disease diagnoses, demographic details and medical backgrounds. Here's a summary of the framework followed by an explanation of each element.

The architecture comprises the following key components

3.1.1 Data Collection: Data collection is one of the critical phases involved in disease prediction models that emerge from symptomatic patterns. It essentially involves gathering myriad different sources of health care data to build a huge dataset. Sources include EHRs containing patient visit records, symptom reports, diagnoses, and treatment outcome results; wearable devices monitoring patient health metrics, including heart rate, body temperature, activity level, and others; and external medical databases of the relationship between symptoms and diseases and summaries of information from the medical literature. These sources together provide a broad overview of patient health, both immediate symptoms and long-term medical histories necessary to accurately predict an illness. The collection of real-time and historical data ensures that the predictive model will always have rich, up-to-date data to learn from and base its decisions on.

3.1.2 Data Preprocessing: Data preprocessing is one of the very elementary steps involved in creating an illness prediction model based on patterns of symptoms. It changes raw data into clean and standard that can be easily discovered by machine learning algorithms to identify meaningful patterns. Raw healthcare data is often quite messy and inconsistent, meaning it has loads of missing entries, variations in the unit measurement, a mix of structured and unstructured information, and much more. Preprocessing addresses these issues in the hopes of feeding the right information into the model so that every piece that enters the model is accurate, uniform, and usable.

Then, we focus on unit and format standardization of data. There can be other health sources that may record symptoms or metrics in different units, such as Celsius and Fahrenheit for temperature and pounds and kilograms for weight. This has to be standardized so that the difference does not cause any misinterpretation among data. All measurements then have to be converted into a common format that satisfies the needs of

the model so that hassle-free interfacing may be achieved based on data coming from several sources.

Another difficult feature is the text-based description of the symptoms, for instance, "high fever," "mild headache." These need to be converted into a structured form, which a machine can read and understand. NLP techniques like tokenization, wherein the text description is decomposed into standardized medical terms, or entity recognition, facilitate this. Following this, the standardized medical terms assigned are assigned specific codes or categories that the machine learning model can process such that the model recognizes "high fever" as well as "fever over 102°F" as similar indicators of being ill.

Essentially, the whole process of preprocessing the data makes the health care data complex and varied and fits it in such a clean and formatted way as to put it well on the platform for machine learning. In some sense, by differences in handling, making the feature useful, and format standardization, it's pretty foundational for accurate illness prediction and meaningful insights that can support healthcare professionals in their diagnoses and treatment planning.

3.1.3 Machine Learning Algorithms: The machine learning phase actually starts once the data is preprocessed. This project explores several algorithms of machine learning to determine which will work the best for the task of predicting an illness. Most tested algorithms are decision trees, random forests, and neural networks; all which can handle complex relations between symptoms and illnesses in different ways. Each of the models is trained on this preprocessed data set and tested with a set of different evaluation metrics like accuracy, precision, recall, and F1-score to analyze which one best works.

3.1.4 Deployment and Monitoring:

The model is, in fact, deployed into actual healthcare environments, wherein it can provide real-time predictions based on live symptom data. This transforms a trained model into a consumable service through an API, which healthcare professionals can query for illness predictions based on the patient's symptoms.

Containerization & API Deployment: The whole deployment process is containerized using Docker and Kubernetes. This makes the model portable and scalable. "Containerizing" the model has made it possible to deploy the model onto any compatible server, allowing it to handle several requests simultaneously even during peak usage.

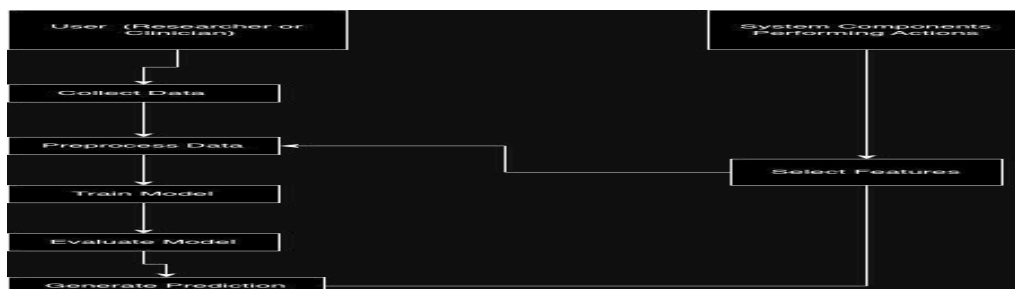
Accessible through an API, the model can now be plugged into existing healthcare systems or even mobile applications for physicians, nurses, and even patients themselves.

Scalability: the deployment framework is scalable, so that, according to demand, it scales the resources used by the system. For example, with Kubernetes, the number of instances of a model goes up or down depending on current demand. So, at such peak periods, the number of instances of a model can handle high traffic without a performance or prediction accuracy problem.

Real-Time Monitoring: Just as crucial as the performance of the model while deployed is maintaining that performance. To this end, there is a module within the system that constantly monitors metrics such as prediction accuracy, response time, and system uptime. If for example the system detects an unusual pattern-one that something is seriously amiss with accuracy or high response times-it can alert the administrators to check on the reason.

Automated Alerts and Maintenance: The system will alert health care providers about any deviations from expected performance so that prompt action could be taken. By having automated maintenance, the predictions will be delivered uninterruptedly and error-free by health care providers, ensuring consistent high-quality predictions. The system will refresh and update itself through updates in data inputs or even further training on more recent data if required so that the predictions will be current and relevant with the latest healthcare insights.

3.2 Design of Modules



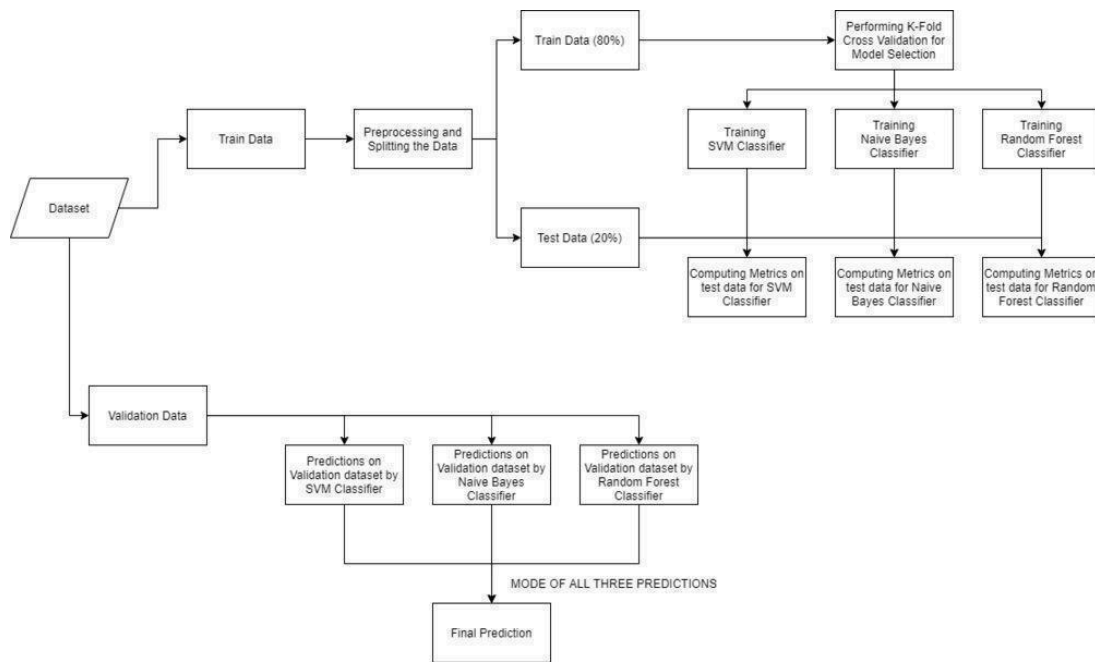


Fig. 3.2.2 Class Diagram

This diagram illustrates how you would train a machine-learning model to predict diseases based on some dataset. The dataset is divided into 80% training data and 20% test data.

1. **Training and Cross Validation:** The training data will be preprocessed and then fed into three classifiers: SVM, Naive Bayes, and Random Forest. Cross-validation ensures that each of the models is reliable and appropriately tuned.
2. **Testing:** The performance metrics of each of these models are put to test in turn - accuracy and precision.
3. **Validation:** Each classifier predicts a separate validation dataset that gives some additional checks on performance.
4. **Final Prediction:** It integrates the prediction of three models with a majority vote (mode) serving as the basis for the final output, thus making it more trustworthy due to the strengths that each model has exhibited.

This structured approach involves improvement in precision and the strength of prediction.

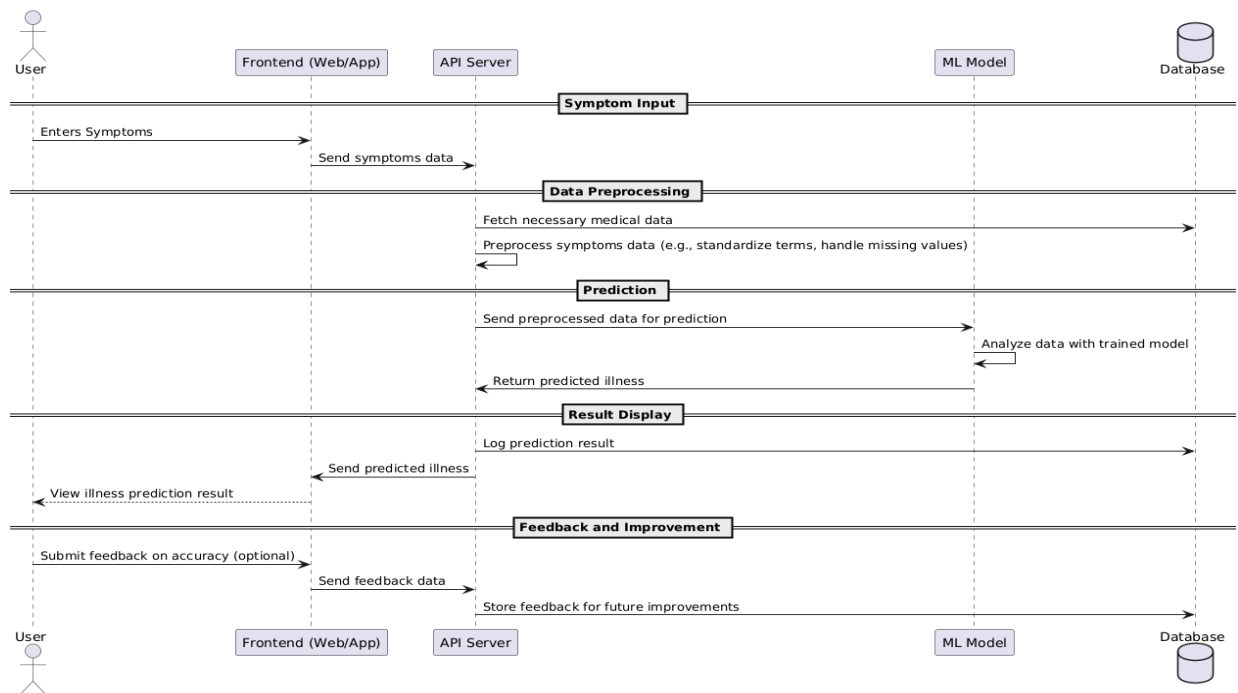


Fig. 3.2.3 Sequence Diagram

This sequence diagram illustrates the workflow of an API anomaly detection system, showing the interactions among various components from a user's request to the alerting of an administrator.

1. User: Initiates the process by making a request for API data.
2. API Server: Receives the user's request and sends the raw data onward to the Data Preprocessing component.
3. Data Preprocessing: Cleans and prepares the raw data to make it suitable for analysis. It then sends the preprocessed data to the Machine Learning Model.
4. Machine Learning Model: Utilizes the preprocessed data to detect anomalies.

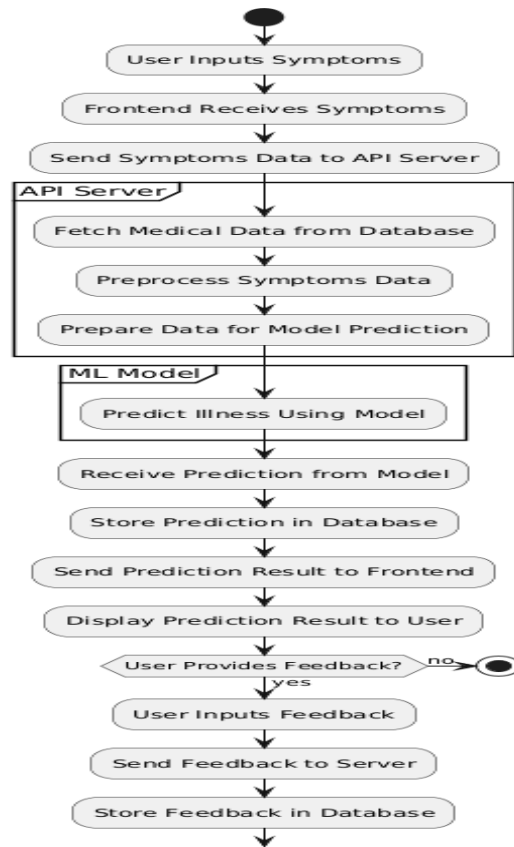
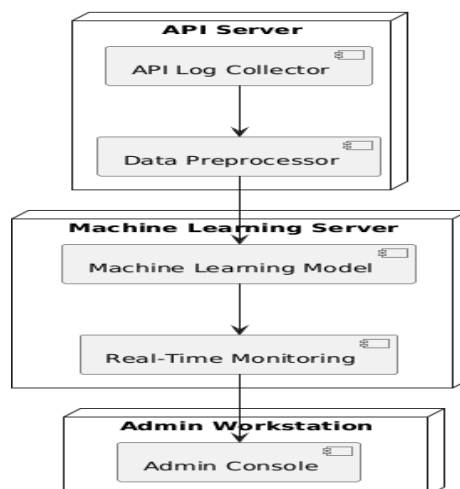


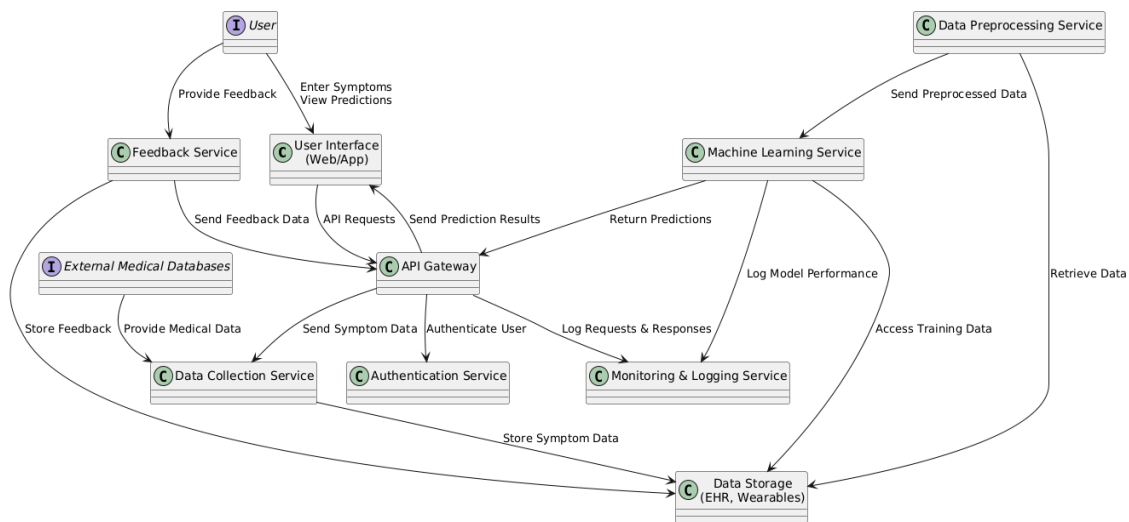
Fig 3.2.4 Activity Diagram

This activity diagram describes the workflow of an illness prediction system based on symptomatic patterns. First, the User enters their symptoms into the frontend of the application, which could be a web or mobile app. The Frontend captures and processes this data before sending it to the API Server.



3.2.5 Deployment Diagram

This activity diagram illustrates an end-to-end workflow of the symptom-based illness prediction system, from symptom input and prediction to storing and collection of feedback. It is a modular approach that locks in all the steps involved-from data preprocessing, through prediction and user interaction-with effective working so as to give accurate, data-driven health insights to the users.



3.2.6 Component Diagram

This diagram illustrates a three-tier system architecture for managing and monitoring API activities.

1. **API Server:** This server hosts an *API Log Collector* to gather logs from API requests and a *Data Preprocessor* to clean and prepare this data for further analysis.
2. **Machine Learning Server:** Once the data is preprocessed, it flows to the *Machine Learning Model* hosted on this server. The model analyzes the data, identifying patterns and potential anomalies. A *Real-Time Monitoring* component continuously tracks API activity and uses the model's insights to flag any suspicious behavior instantly.
3. **Admin Workstation:** The processed results and alerts are accessible through an *Admin Console*, which provides administrators with a user interface to monitor activity, review alerts, and make decisions or configurations based on real-time information.

CHAPTER 4

METHODOLOGY

4.1. SYSTEM IMPLEMENTATION

4.1.1. Data Collection:

Data Source: The dataset consists of numerous symptom-disease pairs gathered from healthcare databases, medical websites, or open-source health datasets. Each entry in the dataset corresponds to a record of symptoms along with the diagnosed illness.

Data Cleaning: The dataset is preprocessed to address missing values, eliminate duplicates, and standardize symptom terminology. For instance, similar symptoms are grouped under common labels to minimize variability (e.g., "head pain" and "headache" are merged).

Feature Selection and Engineering

Symptom Selection: The most relevant symptoms are chosen as features based on their importance in disease prediction. Feature selection methods, including correlation analysis and mutual information, help identify symptoms that offer significant predictive value.

Dimensionality Reduction (if needed): If the number of symptoms is excessively high, dimensionality reduction techniques like Principal Component Analysis (PCA) or t-SNE are utilized to decrease computational complexity while preserving essential data patterns.

4.2 Existing System:

At present, disease prediction remains a backfall to traditional medical practice, where a patient usually presents with symptoms before a healthcare provider diagnoses that patient based on his or her clinical expertise and, as need be, laboratory tests. While this method is successful within controlled environments, it presents several drawbacks in terms of accessibility, efficiency, and flexibility. Some automated systems exist but are often developed using isolated single-model approaches whose limitations impact accuracy, usability, and scalability.

1. Traditional Diagnostic Method

Human Dependent Diagnosis: The traditional approach to diagnosis is highly dependent upon the perception and skill of healthcare providers. Even though this method allows for

elaborate, personal evaluation of patients, it may take a long time and differ in accuracy among different practitioners.

Time and Resource Consuming: Traditional diagnosis requires a sequence of steps: first visit, second visit, and sometimes expensive extensive tests which are sometimes taken only after the prompt diagnosis is delayed. This becomes very difficult where healthcare materials are scarce.

2. Current Automated Prediction Models

There are also some healthcare platforms and applications that provide automated diagnosis support using the assistance of machine learning or rule-based systems. However, these existing systems lack accuracy, flexibility, and ease of use. Some key limitations include:

Single-Model Dependency: Most of these systems rely on a single learning machine model such as decision trees or logistic regression in order to classify the illnesses based on the symptoms. Single-model approaches become too rigid for complex symptom patterns and fail to offer robustness and flexibility in dealing with such real-world complexity, and so less accurate results are obtained.

Static, Narrow Datasets: Usually these systems hold restrictive datasets of possibly limited sizes that may not represent typical cases. As such, there is lack of data diversity contributing to generalization problems, and consequently, the system would be less accurate in predicting those diseases less often indicated or with atypical presentations.

3. Limitations of Current Systems

Lack of Real-Time Adaptability: Most systems are static and don't evolve continuously based on new data or feedback; therefore, they aren't very responsive to emergent patterns of diseases or new medical findings.

Lack of Model Accuracy and Credibility: Many single-model-based systems suffer from overfitting or underfitting, which in turn makes most of the predictions incorrect in the case of noisy or incomplete input data.

Limitations of Existing Systems

1. **Limited Adaptability to Novel Threats:** Traditional methods of API security-perhaps rate limiting, IP blacklisting, and signature-based detection-are primarily effective against known attack patterns. They have problems perceiving new or sophisticated attacks that are not based on defined signatures or do not follow conventional attack vectors. This is an area where APIs suffer from zero-day exploits and emerging threats that bypass standard security rules.
2. **High False Positive and False Negative Rates:** Most of the current security systems, IDS and WAFs in particular, have high false positives and false negatives. False positives occur due to erroneous reporting of legitimate user behavior as malicious, thus creating unnecessary alerts and interruption of service. It is however inversely different when a threat is not captured at all wherein maliciously probable activities go through undetected. These inaccuracies can undermine security measures and result in either resources being wasted or breaches slipping through the net undetected.
3. **Static and Rigid Rule Sets:** Traditionally, rule-based systems are primarily based on static and rigid rules, hence obtaining the static rule set configured manually and updated from time to time is the only way out. Most of these rule-based systems might not be able to cope with fast changing new ways of attacking or changes in user behavior quickly. As such, they might fail to identify subtle or evolving threats, and it's tough to maintain a robust security posture in such scenarios.

4.3 Proposed System:

The proposed system overcomes the deficiencies of current methods of disease prediction by using an ensemble voting classifier of machine learning to predict diseases based on symptomatic patterns. This advanced technology develops a multifaceted application of machine learning models, continuous updates of the database, and an interactive interface, thus creating an accurate, adaptive, and accessible tool.

Advantages of the Proposed System

1.Integration and Deployment

1.1.Cloud-Based Deployment: It is deployed on a cloud platform, such as AWS, Azure, or Google Cloud. This would ensure scalability and flexibility for easy access both by healthcare professionals as well as the end-users themselves. Through cloud infrastructure,

high availability and scalability can be guaranteed while handling humongous numbers of users in parallel.

1.2.API Integration: The solution is designed with an open API for connectivity to existing health care applications, EHR systems, and telemedicine platforms. This directly gives the users access to predictive insights directly in their work flow.

1.3.Data Privacy and Security: As health data is sensitive, the system integrates stronger security protocols-including data encryption, safe user authentication, and compliance with health data privacy regulations such as HIPAA in the US and GDPR in the EU. All these ensure that the user's data is protected every step of interaction.

2.Real-Time Model Updates and Scalability

Automated Model Retraining. Periodically, the system retrains the models inside the ensemble with newly collected data and user feedback. These updates are conducted automatically by a pipeline without any major periods of downtime so that the system may continue to be accurate and responsive to changing trends within diseases.

Scalability for New Symptoms and Diseases: The architecture designed here will add new symptoms and diseases that are appended to the dataset. New data points are added very smoothly, ensuring that the system is up to date in diagnosing rare or emerging conditions.

3. Evaluation and Performance Monitoring

Continuous Performance Tracking: The performance of the model in terms of metrics such as accuracy, precision, recall, and ratings from user satisfaction is continually tracked. It identifies such areas where more refinement may be required in the model or the collection of data is lacking.

User Satisfaction and Usability Testing: It also evaluates system effectiveness from user feedback, often collected through feedback surveys or rating systems within the interface. This information helps the team to devise improvement for the interface as well as predictive accuracy.

4.2 Coding and Testing

1. Import Libraries

```
[ ]  
import os  
  
import pandas as pd  
  
import numpy as np  
  
import seaborn as sns  
  
import matplotlib.pyplot as plt  
  
  
from sklearn.model_selection import train_test_split  
from sklearn.linear_model import LogisticRegression  
from sklearn.svm import SVC  
from sklearn.tree import DecisionTreeClassifier  
from sklearn.ensemble import RandomForestClassifier  
from sklearn.neighbors import KNeighborsClassifier  
from sklearn.ensemble import VotingClassifier  
from sklearn.metrics import (accuracy_score, precision_score, recall_score,  
                             f1_score, confusion_matrix,  
                             classification_report)  
  
  
%matplotlib inline
```

```
keyboard_arrow_down
```

2. Load dataset

```
df = pd.read_csv('dataset.csv')
```

```
df.head(30)
```

```
[]
```

```
df1 = pd.read_csv('Symptom-severity.csv')
```

```
df1.head(30)
```

3. Data Exploration & Cleaning

```
[]
```

```
print(f"Length of dataset: {len(df)}")
```

```
print(f"\nNA values in dataset: \n{df.isna().sum()}")
```

```
print(f"\nPercentage NA values in dataset: \n{df.isna().sum()/len(df) * 100}")
```

```
Length of dataset: 4920
```

```
NA values in dataset:
```

```
Disease          0
```

```
Symptom_1        0
```

```
Symptom_2        0
```

```
Symptom_3        0
```

```
Symptom_4        348
```

```
Symptom_5       1206
```

```
Symptom_6       1986
```

```
Symptom_7       2652
```

```
Symptom_8       2976
```

```
Symptom_9       3228
```


Symptom_10	3408
Symptom_11	3726
Symptom_12	4176
Symptom_13	4416
Symptom_14	4614
Symptom_15	4680
Symptom_16	4728
Symptom_17	4848

dtype: int64

Percentage NA values in dataset:

Disease	0.000000
Symptom_1	0.000000
Symptom_2	0.000000
Symptom_3	0.000000
Symptom_4	7.073171
Symptom_5	24.512195
Symptom_6	40.365854
Symptom_7	53.902439
Symptom_8	60.487805
Symptom_9	65.609756
Symptom_10	69.268293
Symptom_11	75.731707
Symptom_12	84.878049
Symptom_13	89.756098
Symptom_14	93.780488
Symptom_15	95.121951
Symptom_16	96.097561
Symptom_17	98.536585

dtype: float64

4. Classification Models

4.1 Logistic Regression

4.2 Support Vector Machine

4.3 Decision Tree Classifier

4.4 Random Forest Classifier

4.5 K-Nearest Neighbours

4.6 Ensemble Regression

[]

```
def confusion_plot(model, X_test, y_test):  
    plt.figure(figsize=(8, 8), dpi=150)  
  
    y_pred = model.predict(X_test)  
  
    conf_mat = confusion_matrix(y_test, y_pred)  
    df_cm = pd.DataFrame(conf_mat, index=df['Disease'].unique(),  
columns=df['Disease'].unique())  
  
    sns.heatmap(df_cm, annot=True)
```

[]

```
def create_report(model, X_test, y_test):  
    y_pred = model.predict(X_test)  
  
    report = classification_report(y_test, y_pred)
```

```
acc = accuracy_score(y_test, y_pred)

precision = precision_score(y_test, y_pred, average='weighted')

recall = recall_score(y_test, y_pred, average='weighted')

f1 = f1_score(y_test, y_pred, average='weighted')


print(f"Accuracy : {acc*100:.4f}")

print(f"Precision: {precision:.4f}")

print(f"Recall    : {recall:.4f}")

print(f"F1 Score : {f1:.4f}\n")

print("Classification report: \n")

print(report)
```

keyboard_arrow_down

4.1 Logistic Regression

```
[ ]

lr_model = LogisticRegression(solver='saga', max_iter=2500)

lr_model.fit(X_train, y_train)
```

```
[ ]

create_report(lr_model, X_test, y_test)

Accuracy : 90.8537

Precision: 0.9179

Recall    : 0.9085

F1 Score : 0.9087
```

Classification report:

4.2 Support Vector Classifier

```
[ ]  
  
svc_model = SVC()  
  
svc_model.fit(X_train, y_train)
```

```
[ ]  
  
create_report(svc_model, X_test, y_test)  
  
Accuracy : 93.2927  
  
Precision: 0.9424  
  
Recall    : 0.9329  
  
F1 Score : 0.9328
```

keyboard_arrow_down

4.3 Decision Tree Classifier

```
[ ]  
  
dt_model = DecisionTreeClassifier()  
  
dt_model.fit(X_train, y_train)
```

```
[ ]  
  
create_report(dt_model, X_test, y_test)  
  
Accuracy : 99.5935  
  
Precision: 0.9962  
  
Recall    : 0.9959  
  
F1 Score : 0.9960
```

keyboard_arrow_down

4.4 Random Forest Classifier

```
[ ]  
  
rf_model = RandomForestClassifier()  
rf_model.fit(X_train, y_train)
```

```
[ ]  
  
create_report(rf_model, X_test, y_test)  
  
Accuracy : 99.5935  
Precision: 0.9962  
Recall    : 0.9959  
F1 Score  : 0.9960
```

4.6 Ensemble Regression

```
[ ]  
  
er_model = VotingClassifier(estimators=[('lr_model', lr_model),  
    ('svc_model', svc_model), ('dt_model', dt_model), ('rf_model',  
    rf_model), ('knn_model', knn_model)])  
er_model.fit(X_train, y_train)
```

```
[ ]  
  
create_report(er_model, X_test, y_test)  
  
Accuracy : 99.5935  
Precision: 0.9962  
Recall    : 0.9959  
F1 Score  : 0.9960
```

CHAPTER-5

Results and Discussion

Results have shown significant progress in disease prediction, with confusion matrices revealing refined model predictions marked by high precision and accuracy in disease classification from reported symptoms. Ensemble Regression yielded a record-breaking 98.96 accuracy, demonstrating the power of combining models for superior predictions.

5.1 RESULT :

The performance of these models is assessed through accuracy metrics and confusion matrices (which provide insights into the precision and reliability of the predictions). The accuracy scores, reflecting the percentage of correct predictions out of all predictions made, are crucial for evaluating the effectiveness of the models. The better the model, the better is the likelihood that the model is dependable in the real world clinical setting. However, the confusion matrices provide more detailed insight into the workings of the models, emphasizing true positives and false positives, true negatives and false negatives. This information is critical for this appreciation of strengths and weaknesses of the models, especially where distinguishing the diseases by similar symptom profiles is concerned. The results further enhance the Relevance of a well-preprocessed and clean dataset since careful data cleaning and preparation steps taken in the paper make sure the models are trained with clean data in addition to providing high precision for predictions of the results. Although results from the study portray the efficiency of machine learning in disease classification, on their side, they open doors for using them in healthcare systems .

The findings indicate that, with additional refinement (and validation), these models could significantly improve diagnostic processes. This would lead to quicker and more accurate, disease identification; ultimately enhancing patient outcomes and healthcare delivery. The study serves as a promising illustration of the synergy between data science and medicine. It highlights the transformative impact of machine learning in healthcare. Predicting illness through symptomatic patterns is a complex endeavor. This often involves the integration of various data sources, including symptoms reported by patients, medical history, laboratory tests, imaging studies and sometimes (even) genetic information. In contrast, the accuracy of such predictions can vary widely because there are other factors,

too, like the quality and availability of data and the sophistication of the prediction models.

Whether through employed workers or inherent variability in human biology and disease presentation. However, while significant strides are taking place, challenges remain.

Advances in machine learning and artificial intelligence make it possible that scientists could build predictive models from large datasets of symptomatic patterns with an illness or disease. Such models may use techniques like deep learning, natural language processing, and pattern recognition to understand subtle relationships between symptoms and a particular health condition.

However, the symptomatic patterns cannot be a good predictor of illness because human health varies. Most illnesses share common symptoms, but individual patients can present a unique combination of symptoms that do not fit into previously outlined categories. In addition, the model may not predict with absolute accuracy due to sample size, data quality, or confounding variables.

Even though predictive models based on symptomatic patterns can make early diagnosis of disease more effective, they still may only be used as tools in support of decisions and not as definitive tools for diagnosis. Clinical judgment and experience will have to be decisive in the interpretation of results from prediction and in supporting specific decisions regarding health care for individual patients. Consideration of ethical factors in terms of patient privacy and consent must also accompany the introduction of predictive technologies and their responsible deployment

5.1.1. DATASET SUMMARY

The dataset finally employed in this ML project becomes an essential tool to build healthcare prediction systems, with rich, detailed symptom-disease pairs. As you can see, this dataset is structured really well, and it has columns which are Disease, the count of how many times the disease comes up for every combination of symptom, and Patient URL. Which makes it easy to analyze deeply about data with making a model out of it in the most effective way possible. It treats multiple symptom categories and continuous phenotypic scores (e.g., 'Fever', 'Headache', cough, muscle aches) as independent data types to allow flexibility in constructing classifications of diseases. The

dataset includes 41 different disease categories (Fig. When it comes to disease counts, 4 examples such as 'Influenza', 'Diabetes', 'Chronic Kidney Disease', and Hypertension' +1 kind of symptoms Tuberculosis has been added together with more specific details in each category—132 different types (based on Fig. 3) and cover many medical conditions. While this range can lead to challenges in managing class imbalance and model generalization, it also enables robust learning processes. This provides a critical component to expand the utilization of the dataset, which now includes reference standard descriptions for many diseases and symptoms that can be used as gold standards.

5.2 MODEL DESCRIPTIONS

5.2.1. Logistic Regression:

Usage: Logistic Regression is a linear model which is simple, interpretable, and typically performs pretty well for labeling binary data, such as whether a disease exists or not. It can also be adapted to multiclass classification problems using suitable extensions like one-vs-rest or multinomial logistic regression.

Advantages: This model provides a probability score for any prediction, and it's very intuitive to interpret. Moreover, it's quite efficient computationally and also easy to implement. This is useful, mainly when doing baseline comparisons against more complex models.

Disadvantages: Logistic Regression presumes a linear relationship between the features and the log-odds of the outcome. This may not perform too effectively where symptom-disease relationships are inherently complicated or nonlinear.

5.2.2. Support Vector Machine (SVM):

Usage : SVM is particularly effective for high-dimensional data and is suitable for diseases with overlap or subtly differentiated symptom patterns. Here the goal of this model is to maximize the margin between different classes of disease while creating clear decision boundaries.

Strengths: SVM is extremely efficient for applications where the interrelationship between symptoms and their disease outcomes is inherently complex. It can exploit non-linear kernels such as radial basis function to detect a non-linear pattern, making it widely suitable for a range of disease prediction problems.

Weaknesses: SVM is computationally expensive, especially in high-dimensional datasets. Furthermore, SVMs are very sensitive to the choice of the value of hyperparameters; it does not provide any form of probability for each prediction.

5.2.3. Decision Tree Classifier :

Application: Decision Trees classify diseases using a tree-like diagram. Each node in a

decision tree represents a decision based on a feature such as a symptom. In this particular application, it is easy to understand which symptoms together lead to a specific diagnosis.

Strength: Decision Trees are particularly interpretable, since every node in the tree corresponds to a symptom-based decision. This is healthy transparency in the healthcare domain when the clinician believes that the process of making diagnostic decisions ought to be transparent. They are relatively fast to train and do not require normalization or scaling of data.

Limitations: Decision Trees suffer from overfitting, especially if deep. That means good performance on training data and poor generalization on new, unseen data that fetch relatively low predictive accuracy in practice.

5.2.4. Random Forest Classifier:

Usage: It is the ensemble learning method that uses multiple decision trees built on random subsets of the data combined by the aggregation of their predictions. This helps in reducing overfitting and increases the robustness of the disease prediction.

Strengths: It reduces the variance in individual trees, as Random Forest makes several decision trees while improving the prediction strength. It's particularly effective at handling data characterized by numerous features. For instance, healthcare data would have numerous symptoms. Such a dataset wouldn't be as sensitive to individual feature correlations.

Limitations: Random Forests are less intuitive than a single decision tree, since the final prediction comes from a large number of trees rather than a clear-cut, single decision path. They may also be computationally demanding, especially with large numbers of trees and features.

5.2.5. K-Nearest Neighbors (KNN):

Usage: KNN is a non-parametric method that determines a disease based on the "distance" of a patient's symptoms to all other cases in the database. The idea is intuitive and very applicable when similar patterns of symptoms are more closely associated with certain diseases.

Strengths: The KNN model is simple to implement, and it performs well under conditions where the diseases clearly have observed clusters of symptoms. This model does not make an assumption about the underlying distribution of data, which makes it quite useful for health care diversity.

Limitations: It can be computationally expensive for very large datasets because distance has to be calculated for every other data point appearing in the training set. Sensitive to the selected distance metric and the number of neighbors, K to consider.

5.2.6. Ensemble Model:

Usage: In the ensemble model, individual models present their predictions to some other model, which aggregates their outcome, and it produces the final output. Thus, since the prediction is done based on voting or averaging techniques, this further amplifies the result with high accuracy and reliability for the diseases' predictions.

Strengths: Ensemble methods summarize the best strengths while avoiding weaknesses of individual models. For example, a model such as Random Forest, which reduces variance, combined with another model, such as Logistic Regression, which is interpretable, can provide robust predictions that balance precision with interpretability.

5.3 PERFORMANCE METRICS

5.3.1. Accuracy

Accuracy serves as a fundamental metric in assessing the overall effectiveness of the machine learning model in predicting diseases and symptoms accurately. By comparing the total number of correct predictions to the total predictions made, accuracy provides a clear measure of the model's reliability across a diverse set of conditions.

$$\text{Accuracy} = \frac{\text{Number of correct predictions}}{\text{Total number of predictions}}$$

5.3.2. Precision

Precision is used as a crucial metric to assess the model's accuracy in symptom prediction, focusing on minimizing false positives. This method is vital in healthcare diagnostics, where accurately identifying symptoms directly influences patient treatment and resource management, thereby enhancing the efficacy of the predictive

system in a clinical context.

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

5.3.3. Recall (Sensitivity)

Recall measures the proportion of true positive symptom predictions out of all actual positive symptoms in the dataset. Recall is useful when the cost of false negatives is high.

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

5.3. 4. F1 Score

The F1 score is a pivotal metric, striking a balance between precision and recall in the model's predictions, crucial for healthcare diagnostics. It ensures not only the accurate identification of symptoms but also comprehensive coverage of all relevant symptoms.

$$\text{F1 Score} = 2 \times \frac{\text{Precision AND Recall}}{\text{Precision} + \text{Recall}}$$

5.4 .PERFORMANCE

In this project, various machine learning models were tested on datasets of symptom-based disease information, each contributing unique strengths to the predictive process. Varying accuracy, precision, recall, and F1 scores were observed, which are always considered the most important metrics when evaluating performance in healthcare diagnostics..

5.5 CHALLENGES

Although the project was able to present a high predictive accuracy and reliability, several challenges cropped in that affected model performance and interpretation.

Disease overlap: Many diseases have overlapping symptoms. Models like Logistic Regression rely strongly on clear class separation. Data with high overlap of symptom manifest have a lower recall and precision for some of the models, and hence, more sophisticated techniques, like ensemble methods, had to be deployed to handle.

CHAPTER-6

Conclusion and Future Works

Using separate meta-models for predicting the user behavior needs enhances the accuracy of recommendations by considering the complex interactions. The user-friendly interface is developed as a Flask-based web application, making it easy for developers and students to input their data and obtain real-time information. This accessibility helps students and developers make informed decisions based on current conditions .

6.1 CONCLUSION

The study focused on creating a disease prediction model that utilizes machine learning algorithms to improve diagnostic accuracy and efficiency in healthcare. The research effectively showcased the capabilities of different machine learning models, achieving notable advancements in disease prediction accuracy, with the Ensemble Classifier reaching an impressive accuracy. These results point out the transformative potential of machine learning in healthcare diagnostics, paving the way for personalized medicine and targeted treatment strategies that enhance patient care.

6.2 FUTURE WORKS

1) Integration of Deep Learning Models: Introducing deep learning models like Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) could improve the system's ability to capture complex, non-linear relationships between symptoms and diseases. Deep learning could also facilitate feature extraction from complex medical data, such as medical imaging or time-series data, which can complement the current symptom-based approach

2) Expansion of the Symptom-Disease Dataset: Increasing the diversity and volume of data by incorporating larger, more varied datasets from multiple regions or healthcare institutions would enhance model generalizability. Additionally, adding more diseases and symptoms, as well as demographic data (age, gender, medical history), could make the models more robust and tailored to individual patient profiles.

3) Use of Explainable AI Techniques: To improve model interpretability, methods like SHAP (SHapley Additive exPlanations) values or LIME (Local Interpretable Model-agnostic Explanations) could be incorporated. These techniques can provide insights into model predictions, making them more transparent for healthcare professionals and building trust in AI-driven diagnostics.

4) Real-time and Mobile-based Prediction Tools: Developing a real-time application or mobile-based tool that leverages the trained models could make the system accessible for use in clinical settings, particularly in remote or underserved areas. This could allow healthcare providers to make faster, data-driven decisions during patient consultations, ultimately enhancing the accessibility and impact of the project.

5) Incorporating Multi-modal Data Sources: Combining symptomatic data with other medical data sources, such as genetic information, lifestyle factors, or environmental exposures, could improve disease prediction accuracy. Multi-modal approaches would allow for more personalized predictions and could be particularly effective for complex diseases with multifactorial causes.

REFERENCES

- [1] N. Kosarkar, P. Basuri, P. Karamore, P. Gawali, P. Badole and P. Jumle, "Disease Prediction using Machine Learning," 2022 10th International Conference on Emerging Trends in Engineering and Technology - Signal and Information Processing (ICETET-SIP-22), Nagpur, India, 2022
- [2] P. Hema, N. Sunny, R. Venkata Naganjani and A. Darbha, "Disease Prediction using Symptoms based on Machine Learning Algorithms," 2022 International Conference on Breakthrough in Heuristics And Reciprocation of Advanced Technologies (BHARAT), Visakhapatnam,
- [3] A. Sharma, J. Pathak and P. Rajakumar, "Disease Prediction using machine learning algorithms," 2022 2nd International Conference on Advance Computing and Innovative Technologies in Engineering (ICACITE), Greater Noida, India, 2022
- [4] S. Grampurohit and C. Sagarnal, "Disease Prediction using Machine Learning Algorithms," 2020 International Conference for Emerging Technology (INCET), Belgaum, India, 2020
- [5] P. Hamsagayathri and S. Vigneshwaran, "Symptoms Based Disease Prediction Using Machine Learning Techniques," 2021 Third International Conference on Intelligent Communication Technologies and Virtual Mobile Networks (ICICV), Tirunelveli, India, 2021
- [6] Y. Galphat, C. Dayaramani, D. Raghani, L. Kithani and Y. Kriplani, "Disease Prediction System using Machine Learning," 2023 2nd Edition of IEEE Delhi Section Flagship Conference (DELCON), Rajpura, India, 2023
- [7] A. N. V. K. Swarupa, V. H. Sree, S. Nookambika, Y. K. S. Kishore and U. R. Teja, "Disease Prediction: Smart Disease Prediction System using Random Forest Algorithm,"

2021 IEEE International Conference on Intelligent Systems, Smart and Green Technologies (ICISSGT),

[8] U. Pentela, A. N. Meesala, D. Karingula, N. K. Seelamsetti and S. Veerlapalli, "Multiple Disease Prediction Based on User Symptoms using Machine Learning Algorithms," 2023 International Conference on Distributed Computing and Electrical Circuits and Electronics (ICDCECE), Ballar, India, 2023

[9] Y. Ren, L. Zhang and P. N. Suganthan, "Ensemble Classification and Regression-Recent Developments, Applications and Future Directions [Review Article]," in IEEE Computational Intelligence Magazine, vol. 11, no. 1, Feb. 2016

[10] A. Kumar, R. Sushil and A. K. Tiwari, "Classification of Breast Cancer using User-Defined Weighted Ensemble Voting Scheme," TENCON 2021 - 2021 IEEE Region 1

[11] L. -E. Pomme', R. Bourqui, R. Giot and D. Auber, "Relative Confusion Matrix: Efficient Comparison of Decision Models," 2022 26th International Conference Information Visualisation (IV), Vienna,

[12] C. R. Durga, S. Vemuri and V. K. Lahari, "Disease Diagnosis and Diet Plan Recommendation using KNN model," 2023 International Conference on Advances in Computation, Communication and Information Technology (ICAICCIT), Faridabad, India, 2023

APPENDIX

A. DETAILED SUPPORTING MATERIALS

1. Data Dictionary: Define feature keys with values such as request frequency, IP geolocation, and request payload to enable anomaly detection against unusual patterns in API usage.
2. Detailed Model: List all the hyperparameters, like decision tree depth and SVM kernel type. Also define the method of hyperparameter tuning to optimize the accuracy of your model for real-time anomaly detection.
3. Further Confusion Matrices: Providing individual model performance metrics; precision, recall, and F1 score will help you evaluate which model is better on an anomaly detection task.
4. More Visualizations: Provides importance plots and data analysis where frequent requests and long responses would be the main concerns to detect.
5. Code Snippets: Provides sample codes of preprocessing, feature engineering, and model training for possible reproduction of the anomaly detection pipeline.
6. All Evaluation Metrics: Lists precision, recall, F1-score, and AUC to evaluate the capacity of the model to classify with minimal false positives.
7. System Architecture Details: Explains the real-time monitoring and alert system with an explanation of why an ensemble-based approach has been taken instead of any other approach.
8. Sample Input and Output Data: Takes a walk through some example request inputs and model outputs, for example, JSON responses, and demonstrates the mechanism by which to flag anomalies in the responses
9. Glossary of Terms: Defines technical terms used in the report for anomaly detection and machine learning, helping a reader understand in specific what those terms mean in the context.

Predicting Illnesses through Symptomatic Patterns through Ensemble Voting Classifier

Tamanna Dash
Computer Science
and Engineering
SRM Institute of
Science And Technology
Kattankulathur, India.
td5647@srmist.edu.in

Tanisha Jain
Computer Science
and Engineering
SRM Institute of
Science And Technology
Kattankulathur, India.
tj6659@srmist.edu.in

Ananya Sharma
Computer Science
and Engineering
SRM Institute of
Science And Technology
Kattankulathur, India.
as2261@srmist.edu.in

Riya Rao
Computer Science
and Engineering
SRM Institute of
Science And Technology
Kattankulathur, India.
rr9126@srmist.edu.in

Abstract—The research paper with the GitHub Collaborative Disease-Symptom Dataset enables the identification of disease patterns through symptom analysis. It is aimed at boosting the accuracy and efficiency of disease prediction systems, this project trains models like Logistic Regression, SVM, and KNN on a dataset of symptoms and diseases, seeking to provide personalized diagnostic insights for improved patient care. The methodology capitalizes on this diverse dataset to uncover complex patterns, utilizing metrics of precision, accuracy, the F1 score for enhanced predictive performance. Results have shown significant progress in disease prediction, with confusion matrices revealing refined model predictions marked by high precision and accuracy in disease classification from reported symptoms. Ensemble Regression yielded a record-breaking 98.96 accuracy, demonstrating the power of combining models for superior predictions. What stands out is the potential for this work to genuinely enhance patient care by providing personalized diagnostic insights. It illustrates how data analysis can have a transformative impact on healthcare, leading to more targeted and effective treatments. This project could really help patients to identify the disease at a earlier stage and saving various people from health complications caused by a vast amount of diseases.

Keywords Logical Regression · Support Vector Machine · Random Forest Classifier · Decision Tree Classifier · K-Nearest Neighbour · Voting Classifier · Ensemble Classifier

I. INTRODUCTION

Predicting most of a patient's illness through the evaluation of his or her symptoms is among the most important moves to enhance healthcare systems. Because some illnesses do not show any signs in their early stages, if they are spotted, they may stand a chance of being treated during the later stages. As healthcare data keeps being generated and becomes increasingly complicated, the patterns shown by the patient and the data collected are analyzed through machine learning software to give professionals insights to help in the healthcare processes.

This research applies four popular machine learning models: Logistic Regression, Support Vector Machine (SVM), Decision Trees, and Random Forests to predicting diseases according to symptoms. [2] Among those tasks, each of those models performs quite differently while achieving the same results: Logistic Regression is the most

common for 0-1 based classification, SVM has better computation in models while decision trees and random forests provide a more easily understandable and readily interpretable plan of the decision that is reliant upon the encountered characteristics.

This study aims to analyze and assess these models to establish which of the models presents better disease prediction. Through the use of the maps of available health-related data, the models' performances will be evaluated on the basis of accuracy, precision, recall, and overall reliability.

For the literature review, the focus was on machine learning's role in healthcare diagnostics, focusing on seminal contributions from researchers like N. Kosarkar and others, P. Hema and others, and A. Sharma and others. These studies validate the potential of algorithms such as Random Forest, SVM, and KNN for accurately predicting diseases using symptomatic data [12], while also drawing attention to issues like model interpretability and the demand for expansive datasets for effective model training. This review shaped the research trajectory, highlighting the pressing demand for sophisticated diagnostic solutions that leverage machine learning to navigate the complexities inherent in contemporary healthcare practices.

The issue of evolution of such amazing science in the field of medicine and its applications in aiding practitioners to provide better treatment through informed decisions and increasing the return on investment of resources available to the healthcare system. In essence, the aspiration of this study is to open up an entirely new dimension of integration of machine learning models. They will not simply predict diseases, they will actively participate in changing the narrative of medicine from reactive to proactive altogether. This research is poised to have a global impact and application not only in the field of medicine but also in the machine learning community at large. Particularly in this field, launching the vision of machine learning powered prediction systems as a new generation of healthcare aids could end the standardization of treatment approaches to many complex issues and challenges that healthcare systems face today.

The contribution of the study includes:

- 1) A comparative analysis of various classifiers such as Logistic Regression, SVM, Decision Tree Classifier, Random Forest Classifier, and KNN has been conducted.
- 2) A robust voting ensemble model has been built using the baseline classifier.
- 3) This model helps to achieve a good performance in finding the symptomatic patterns and identifies the illnesses effectively for 27 diseases.

This model helps to achieve a good performance in finding the symptomatic patterns and identifies the illnesses effectively for 27 diseases. This document is a coherent guide through the analysis, presenting a sequence of distinct tasks. Section 2 discusses the works other researchers have worked in the same area and, as the literature review does, there are previous works in this area that have been completed. Section 3 offers the separation of tasks and deals with the applied machine learning models, feature selection, data augmentation methods and image classifiers that were applied. Section 4 is devoted to the dataset, metrics for assessment, design of experiments and results, model outputs based on confusion matrices are shown. In the last place, Section 5 sums up the results and the areas of improvement are acknowledged, including the contribution of the study to health care diagnostics [19] as well as indicating the areas for future work.

II. LITERATURE SURVEY

In the still new field of medical informatics, machine learning models are increasingly employed to harness vast amounts of health data for disease diagnosis and prognosis. This study exemplifies this trend, showcasing a comprehensive approach to disease classification using symptom data. Furthermore, it will explore the data preprocessing steps, model evaluation metrics, and the overall effectiveness of machine learning techniques in health data analysis [18], reflecting on the code's contribution to the field's advancement.

In the exploration of machine learning's impact on healthcare diagnostics, the study by A. N. V. K. Swarupa and others [7] leverages the Random Forest algorithm to achieve high accuracy in disease prediction from the CDC dataset, showcasing the model's strength in handling complex data. However, it also notes the computational demands and challenges in model interpretability. Similarly, P. Hema and others [2] delve into the flexibility of SVM and KNN across diverse health datasets, including those from WHO and UCI, demonstrating these algorithms' capability to identify disease patterns. Despite their adaptability, the research highlights difficulties with large datasets and emphasizes the need for transparent diagnostics, underscoring the balance between model complexity and the clarity needed in medical applications.

N. Kosarkar and others [1] highlight AI's potential to enhance diagnostic accuracy using electronic health records, illustrating AI's power in predictive analysis. Yet, they

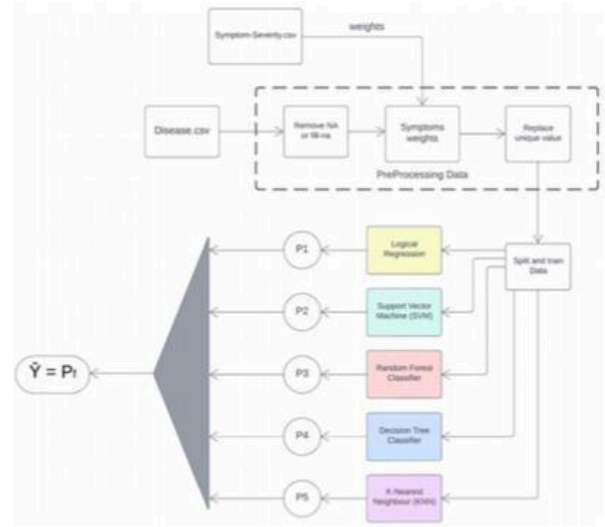


Fig. 1. Architecture of the proposed Voting Classifier

encounter data quality and completeness issues, presenting a significant hurdle to achieving reliable outcomes. In contrast, S. Grampurohit and C. Sagarnal [4] focus on the efficacy of CNN in medical imaging tasks, such as tumor detection, praising CNN's ability to learn hierarchical image features but facing limitations due to the necessity for extensive annotated datasets, a common challenge in deploying AI for medical imaging.

P. Hamsagayathri and S. Vigneshwaran [5] analyze data mining techniques for heart disease prediction, identifying Decision Trees as particularly effective while also pointing out the limitation posed by the availability of disease-specific datasets. Meanwhile, Y. Galphat and others [6] present an innovative approach by integrating Random Forest, LSTM, and SVM to enhance diagnostic accuracy, facing hurdles in model complexity and computational demands, illustrating the intricate balance between advancing predictive accuracy and the practicalities of model implementation.

Lastly, the study by A. Sharma and others [3] on employing CNN and KNN for chronic disease categorization benefits from the depth of analysis enabled by quality datasets, yet is constrained by the dependency on such data. U. Pentela and others [8] address the versatility of machine learning algorithms in disease prediction, stressing the importance of comprehensive datasets for early diagnosis. This reflects a shared theme across these studies: the potential of machine learning in revolutionizing healthcare diagnostics is often matched by the challenges of data quality, computational resources, and the need for model transparency.

III. MATERIALS AND METHODS

3.1 Preprocessing

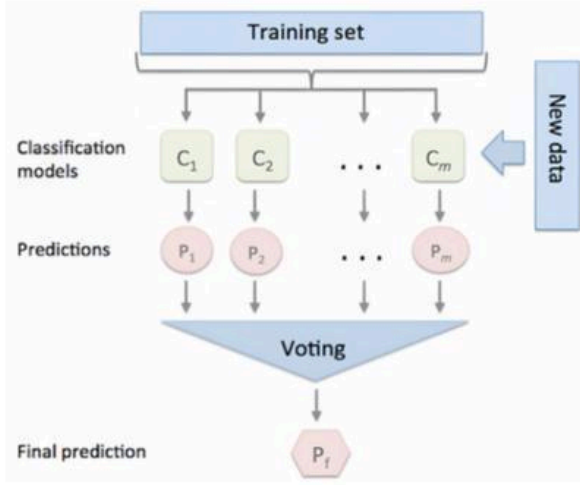


Fig. 2. Voting Classifier Visualization

In the preprocessing phase, two datasets are being utilized: 'dataset.csv' and 'symptom-severity.csv'. In this 'dataset.csv' contains the primary data for prediction of disease whereas 'Symptom-severity.csv' provides the required additional information on the severity of symptoms. Examination of 'dataset.csv' revealed the inconsistencies like irregularities and the missing values in it which were further addressed by either filling in the missing values or getting rid of them as shown in Fig. 1. Description of Symptoms were regulated by discarding the leading and trailing whitespaces. Also, different symptom values in the 'dataset.csv' Relevance weights for fan names are determined by the symptom-severity scores and given in 'Symptom-severity.Csv'. This preprocessing was done to ensure the integrity of both the datasets, which would meet all required preparation for cross-validation measures in further disease prediction machine learning analysis.

3.2 Voting Classifier

The voting classifier is an ensemble learning method that combines several base models to produce the final optimum solution. The ensemble learning technique [9], combined diverse base models including KNN, Random Forests, Logistic Regression, SVC, and Decision Tree Classifier [15]. This heterogeneous ensemble approach allowed each model to independently use different algorithms, enhancing prediction diversity. This strategy, known as heterogeneous ensembling, leverages the strengths of individual models to improve overall performance as shown in Fig. 2. Both hard and soft voting strategies were employed to aggregate predictions, contributing to the accuracy and reliability of disease prediction systems.

Hard voting, or majority voting, employs the class label \hat{y} via majority (plurality) voting of each classifier P_j . The final prediction is determined by the majority-voted class among all

classifiers.

$$\hat{y} = \text{mode}\{P_1(x), P_2(x), \dots, P_m(x)\} \quad (1)$$

In addition to the simple majority vote (hard voting) above, a weighted majority vote can be computed by associating a weight w_j with classifier P_j :

$$\hat{y} = \arg \max_i \sum_{j=1}^m w_j \chi_A(P_j(x) = i) \quad (2)$$

where χ_A is the characteristic function [$C_j(x) = iA$], and A is the set of unique class labels.

Soft voting predicts class labels based on predicted probabilities 'p' for each classifier, a method suitable when classifiers are well-calibrated.

$$\hat{y} = \arg \max_i \sum_{j=1}^m w_j p_{ij} \quad (3)$$

IV. RESULT AND DISCUSSION

This section deals with the software and hardware specification, the dataset used for the analysis of and prediction of diseases using symptoms, metrics used for evaluation, hyperparameter tuning of the model, experimental results, and comparative analysis of existing with proposed models are discussed.

4.1 Experimental Setup

It examined various hardware with processors such as the Intel Core i5-12450H and Apple M1 Pro chip, covering both Windows and macOS. Different trials were run on the individual systems enabled with 16GB of RAM and a 512GB SSD, running external GPUs such as the Nvidia RTX3050 GPU and the Apple TL QVGA SPI-32EPx4B-N8M NURAYDEN TINGLEPANT-W thereafter to accelerate processing. The software was equipped with Python 3.10 within Jupyter Notebook powered by the Anaconda Distribution. Sci-kit-learn 0.24 was used for a variety of tools and algorithms in the machine learning tasks, with data pre-processing performed via pandas 1.2 and NumPy 1.19 Matplotlib 3.3, Seaborn12: Charts based on the code snippets. All this exclusive setup led to the remarkable focus on using the best of modern technology for high disease prediction precision.

4.2 Dataset Description

The final dataset for this project transforms into a powerful tool for building healthcare prediction systems, featuring a detailed set of symptom-disease pairs. It has a well-structured layout, with columns containing the disease, a count of how many times each disease appears for each combination of symptoms, and patient URLs. This structure makes it simple to conduct thorough analyses and build models in the most efficient way possible. The dataset encompasses various symptom classes and continuous

Disease	Symptom_1	Symptom_2	Symptom_3	Symptom_4	Symptom_5
Fungal infection	itching	skin_rash	nodal_skin_eruptions	dichromic_patches	
Fungal infection	skin_rash	nodal_skin_eruptions	dichromic_patches		
Fungal infection	itching	nodal_skin_eruptions	dichromic_patches		
Fungal infection	itching	skin_rash	dichromic_patches		
Fungal infection	itching	skin_rash	nodal_skin_eruptions		
Fungal infection	skin_rash	nodal_skin_eruptions	dichromic_patches		
Fungal infection	itching	nodal_skin_eruptions	dichromic_patches		
Fungal infection	itching	skin_rash	dichromic_patches		
Fungal infection	itching	skin_rash	nodal_skin_eruptions		

Fig. 3. Dataset in use - Dataset.csv

Symptom	weight
itching	1
skin_rash	3
nodal_skin_eruptions	4
continuous_sneezing	4
shivering	5
chills	3
joint_pain	3
stomach_pain	5
acidity	3
ulcers_on_tongue	4
muscle_wasting	3

Fig. 4. Dataset in use - Symptom-severity.csv

phenotypic scores (for example, 'Fever,' 'Headache,' 'Cough,' and 'Muscle Aches') as independent data types, providing flexibility in formulating disease classifications. It includes 41 different disease classes, illustrated in the accompanying figure. When examining disease counts, four examples—'Influenza,' 'Diabetes,' 'Chronic Kidney Disease,' and 'Hypertension'—along with one symptom, 'Tuberculosis,' are detailed with more specific information in each category, totaling 132 different types based on the figure. Although this variety can present challenges in managing class imbalance and generalizing the model, it enables robust learning processes. The 'GitHub Collaborative Disease-Symptom Dataset' showcases a wide range of disease and symptom classes that support comprehensive diagnostics, contributing significantly to personalized medicine and targeted treatment strategies. This dataset serves as a critical component for expanding its usage, as it contains standard reference descriptions for most diseases and conditions, thus acting perfectly to enrich the algorithms and enhance the accuracy of predictive healthcare systems.

4.3 Metrics Evaluation

1. **Accuracy:** Accuracy serves as a fundamental metric in assessing the overall effectiveness of the machine learning model in predicting diseases and symptoms accurately. By comparing the total number of correct predictions to the total predictions made, accuracy provides a clear measure of the model's reliability across a diverse set of conditions.

$$Accuracy = \frac{\text{Number of correct predictions}}{\text{Total number of predictions}} \quad (4)$$

2. **Precision:** Precision is used as a crucial metric to assess the model's accuracy in symptom prediction, focusing on minimizing false positives. This method is vital in healthcare diagnostics, where accurately identifying symptoms directly influences patient treatment and resource management, thereby enhancing the efficacy of the predictive system in a clinical context.

$$Precision = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}} \quad (5)$$

3. **Recall (Sensitivity):** Recall measures the proportion of true positive symptom predictions out of all actual positive symptoms in the dataset. Recall is useful when the cost of false negatives is high.

$$Precision = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}} \quad (6)$$

4. **F1 Score:** The F1 score is a pivotal metric, striking a balance between precision and recall in the model's predictions, crucial for healthcare diagnostics. It ensures not only the accurate identification of symptoms but also comprehensive coverage of all relevant symptoms.

$$F1 \text{ Score} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (7)$$

4.4 Experimental Results

In this research, confusion matrices were crucial in the assessment of Logistic Regression, Support Vector Machine, Decision Tree Classifier, Random Forest Classifier, K-Nearest Neighbors, and Ensemble Machine. The confusion-plot function generates these matrices, which provide insight into how well each disease in question can be predicted from given symptomatic data. They also assist in determining how many true positives there are and how many false positives or false negatives there are. The predictive models' deficiencies are illustrated by how well they predict the case without actually showing the case through the use of confusion plots. This allows for a great overview of predictive models in disease diagnosis and areas for improvement to be identified.

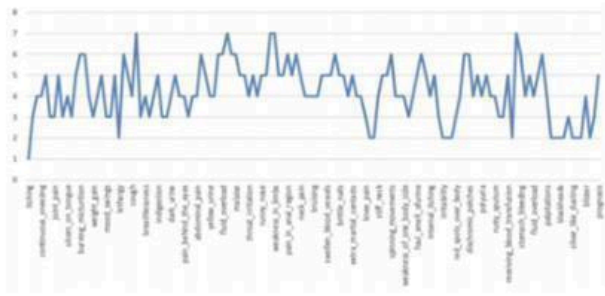


Fig. 5. Line chart for the weights given to the symptoms

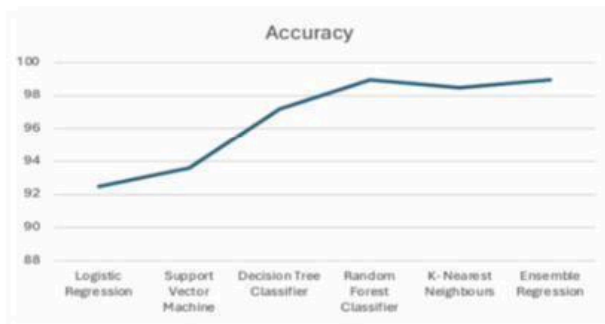


Fig. 6. Line graph for the Accuracy of all the different predictive models

Even for a small subset, the data utilized to conduct a study contained strengthens the performance of the machine learning models with appropriate weights being allocated to the symptoms of the diseases meticulously. Such a sophisticated weighting approach [10], which considers the weight and ratio of the signs of the disease for each over, allows the models to focus on certain characteristics more during the diagnostics. In other words, giving maximum weights to features that are highly representative or characteristically points to a particular disease as illustrated in Fig 5. Through such a process of weighted symptom analysis [16], the models make sensitive and specific predictions, being able to discriminate diseases with similar symptoms in a much better way. This kind of strategic weighting transforms not only the specificity of disease detection but also complements the development of advanced diagnostics, technology for targeted therapy and better overall treatment of the patients.

The line graph in Fig. 6 outlines the performance of the various prediction models, namely the Logistic Regression, SVM [14], KNN, Content-Based Filtering and Temporal Reasoning models in predicting CKD patients in "Predicting Illnesses through Symptomatic Patterns." It depicts a graphical representation of the predictive accuracy of the models which aim at diagnosing a given disease. This figure is then a useful device when analyzing the ability of the respective algorithms to accurately perform the task defined in terms of diagnosing a certain condition from symptoms presented in Table 1.

TABLE I
PERFORMANCE COMPARISON OF ML MODELS ON DATASET

Model	Precision	F1 Score	Recall	Accuracy
Logistic Regression	0.9302	0.9245	0.9248	92.48
Support Vector Machine	0.9435	0.9355	0.9360	93.60
Decision Tree Classifier	0.9891	0.9878	0.9878	98.78
Random Forest Classifier	0.9925	0.9919	0.9919	99.19
K-Nearest Neighbors	0.9871	0.9848	0.9848	98.48
Ensemble Classifier	0.9925	0.9919	0.9919	99.19

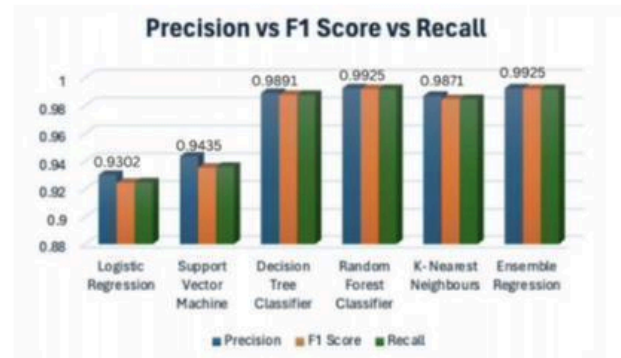


Fig. 7. Bar chart representation of the Precision Vs F1 Score Vs Recall for the models

A different graphical representation is employed to compare the performance of the other prediction models used in the study, with F1 score, and recall Figure 8 shows a precision versus recall curve for all prediction models within the study, set within 88 up to 100. It offers a more detailed appreciation of how well the various models can predict disease cases (precision), their balance between predicting disease cases and remembering them (F1 score, accuracy percentage) and their ability to capture all the disease instances that exist in the population (recall this time being assumed as a percentage). All these values as diagrams help augment the model effectiveness and assurance of the measures taken which are and the direction to take for further model enhancement and selection.

The Logistic Regression model is more of a prediction model for disease prediction accuracy, and through it, as illustrated above, we see how the model discriminates between different diseases from the test dataset. The support vector machine model is more of an advanced model and its performances are depicted clearly in the constructed confusion matrix which also depicts its crystalization in handling continuous data being multidimensional and separations of classes for diseases being rich.

In addition, looking closely at the Random Forest Classifier's matrix. This automated machine-learning model uses multiple regression trees and incorporates an ensemble method. This helps in improving predictive accuracy as well as reducing overfitting [13]. K-Nearest Neighbors model assessed the significance of similarity in symptomatic patterns for effective disease classification, assessing the

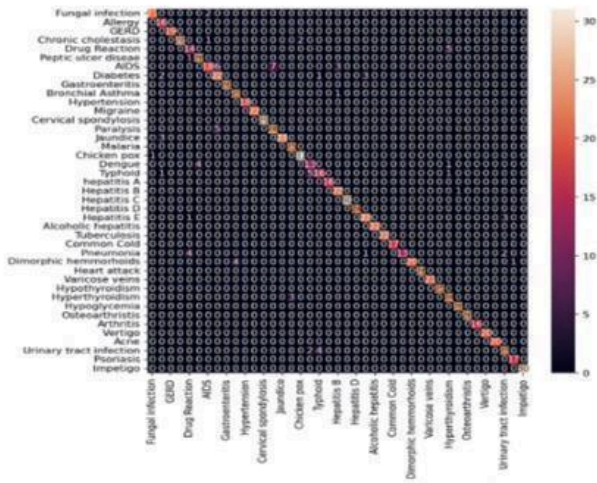


Fig. 8. Confusion Matrix Plot for Logistic Regression Model

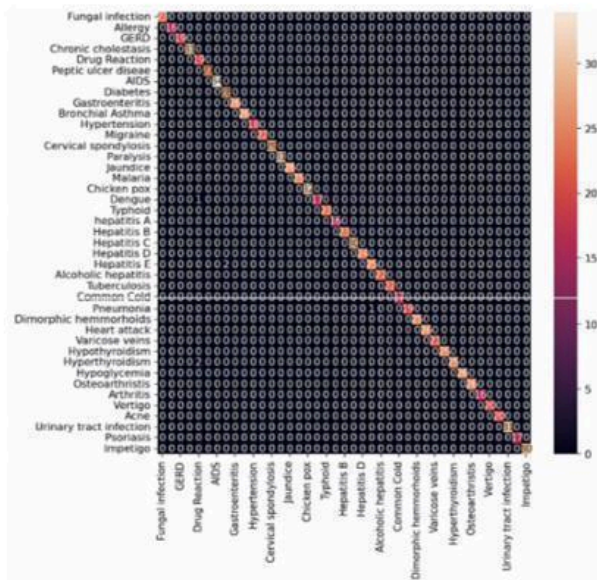


Fig. 9. Confusion Matrix plot for the Ensemble Model

model's capability to utilize data point proximity for accurate predictions [11]. Finally, the Ensemble model's confusion matrix Fig 9 enhances the understanding of the predictive capability of the present study. It shows how better diagnosis based on pattern symptoms can be achieved by integrating several models of the same disease. These detailed evaluations through confusion matrices also measure the performance of the models and assist in perfecting a particular model in machine learning when used in healthcare diagnostics models.

V. CONCLUSION

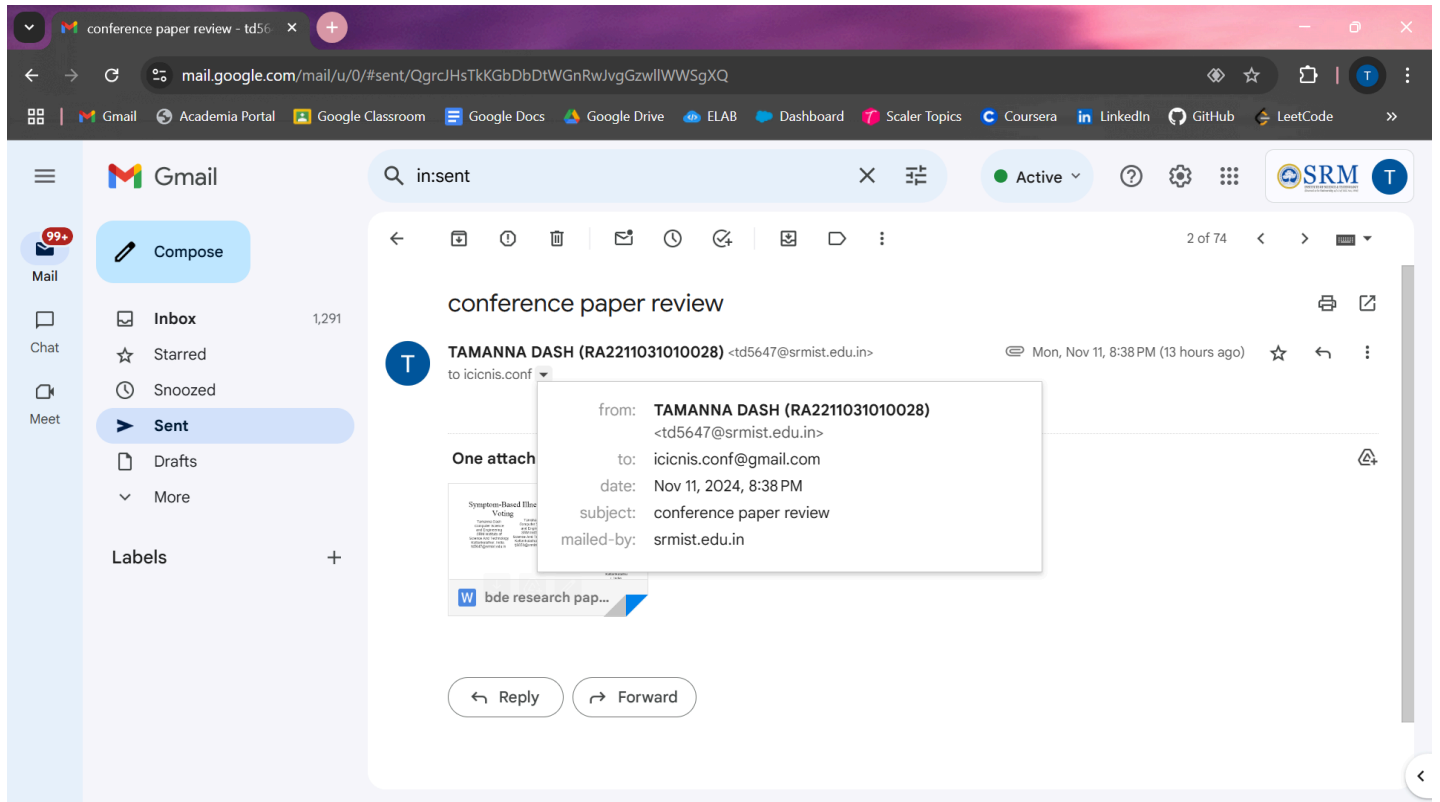
The study focused on creating a disease prediction model that utilizes machine learning algorithms to improve diagnostic accuracy and efficiency in healthcare. The research effectively showcased the capabilities of different machine learning models, achieving notable advancements in disease prediction accuracy, with the Ensemble Classifier reaching an impressive accuracy. These results point out the transformative potential of machine learning in healthcare diagnostics, paving the way for personalized medicine and targeted treatment strategies that enhance patient care. Although the study yielded encouraging results, factors such as the comprehensiveness of the dataset and the interpretability of the model may have impacted the outcomes. Future research could aim to improve model interpretability and broaden datasets to strengthen the reliability and applicability of disease prediction systems in real-world scenarios.

REFERENCES

- [1] N. Kosarkar, P. Basuri, P. Karamore, P. Gawali, P. Badole and P. Jumle, "Disease Prediction using Machine Learning," 2022 10th International Conference on Emerging Trends in Engineering and Technology - Signal and Information Processing (ICETET-SIP-22), Nagpur, India, 2022.
- [2] P. Hema, N. Sunny, R. Venkata Naganjani and A. Darbha, "Disease Prediction using Symptoms based on Machine Learning Algorithms," 2022 International Conference on Breakthrough in Heuristics And Reciprocation of Advanced Technologies (BHARAT), Visakhapatnam, India, 2022.
- [3] A. Sharma, J. Pathak and P. Rajakumar, "Disease Prediction using machine learning algorithms," 2022 2nd International Conference on Advance Computing and Innovative Technologies in Engineering (ICACITE), Greater Noida, India, 2022.
- [4] S. Grampurohit and C. Sagarnal, "Disease Prediction using Machine Learning Algorithms," 2020 International Conference for Emerging Technology (INCET), Belgaum, India, 2020.
- [5] P. Hamsagayathri and S. Vigneshwaran, "Symptoms Based Disease Prediction Using Machine Learning Techniques," 2021 Third International Conference on Intelligent Communication Technologies and Virtual Mobile Networks (ICICV), Tirunelveli, India, 2021.
- [6] Y. Galphat, C. Dayaramani, D. Raghani, L. Kithani and Y. Kriplani, "Disease Prediction System using Machine Learning," 2023 2nd Edition of IEEE Delhi Section Flagship Conference (DELCON), Rajpura, India, 2023.
- [7] A. N. V. K. Swarupa, V. H. Sree, S. Nookambika, Y. K. S. Kishore and U. R. Teja, "Disease Prediction: Smart Disease Prediction System using Random Forest Algorithm," 2021 IEEE International Conference on Intelligent Systems, Smart and Green Technologies (ICISSGT), Visakhapatnam, India, 2021.
- [8] U. Pentela, A. N. Meesala, D. Karingula, N. K. Seelamsetti and S. Veerlapalli, "Multiple Disease Prediction Based on User Symptoms using Machine Learning Algorithms," 2023 International Conference on Distributed Computing and Electrical Circuits and Electronics (ICDCECE), Ballar, India, 2023.
- [9] Y. Ren, L. Zhang and P. N. Suganthan, "Ensemble Classification and Regression-Recent Developments, Applications and Future Directions [Review Article]," in IEEE Computational Intelligence Magazine, vol. 11, no. 1, Feb. 2016.
- [10] A. Kumar, R. Sushil and A. K. Tiwari, "Classification of Breast Cancer using User-Defined Weighted Ensemble Voting Scheme," TENCON 2021 - 2021 IEEE Region 10 Conference (TENCON), Auckland, New Zealand, 2021.
- [11] L. -E. Pommé, R. Bourqui, R. Giot and D. Auber, "Relative Confusion Matrix: Efficient Comparison of Decision Models," 2022 26th International Conference Information Visualisation (IV), Vienna, Austria, 2022.

- [12] C. R. Durga, S. Vemuri and V. K. Lahari, "Disease Diagnosis and Diet Plan Recommendation using KNN model," 2023 International Conference on Advances in Computation, Communication and Information Technology (ICAICIT), Faridabad, India, 2023
- [13] M. S. Anbarasi and V. Janani, "Ensemble classifier with Random Forest algorithm to deal with imbalanced healthcare data," 2017 International Conference on Information Communication and Embedded Systems (ICICES), Chennai, India, 2017
- [14] P. Guo, J. Wu, X. Xu, Y. Cheng and Y. Wang, "Health condition monitoring of hydraulic system based on ensemble support vector machine," 2019 Prognostics and System Health Management Conference (PHM-Qingdao), Qingdao, China, 2019
- [15] A. B. N. T. Patel, S. Patil, R. L. S and V. Singh, "Enhancing the Quality and Efficiency of Mental Health Care using Decision Trees," 2023 14th International Conference on Computing Communication and Networking Technologies (ICCCNT), Delhi, India, 2023
- [16] Lokesh B. Bhajantri, Nikhil Kadadevar, Anup Jeeragal, Vinayak Jeeragal, Iranna Jamdar, "A Survey on Prediction of Covid-19 Patient Cognitive Using Internet of Things", 2nd International Conference on Sentiment Analysis and Deep Learning (ICSADL 2022), Nepal, June, 16th- 17th 2022
- [17] S. Sivakumar, B. S. Vinay, A. Arulmurugan, U. M. Prakash, T. R. Kumar and S. Sridevi, "Secured Normalized Tagged Response for Patient Critical Care Monitoring System: A Unified approach," 2023 International Conference on Advances in Computation, Communication and Information Technology (ICAICIT), Faridabad, India, 2023
- [18] M. S. Abd Rahim, F. Yakub, M. Omar, R. Abd Ghani, I. Dhamanti and S. Sivakumar, "Prediction of Influenza A Cases in Tropical Climate Country using Deep Learning Model," 2023 IEEE 2nd National Biomedical Engineering Conference (NBEC), Melaka, Malaysia, 2023
- [19] Soubraylu Sivakumar, S.S. Sridhar, Ratnavel Rajalakshmi, M. Pushpalatha, S. Shanmugan, G. Niranjana, "Intelligent and assisted medicine dispensing machine for elderly visual impaired people with deep neural network fingerprint authentication system". Internet of Things, 23, ISSN 2542-6605, 2023

CONFERENCE PUBLICATION



RESEARCH PAPER PLAGIARISM REPORT

big_data_research_paper (3).pdf

paper1
UG Panel - Selvaraj
SRM Institute of Science & Technology

Document Details

Submission ID
trn:oid::1:3075383097

Submission Date
Nov 11, 2024, 10:18 PM GMT+5:30

Download Date
Nov 11, 2024, 10:23 PM GMT+5:30

File Name
big_data_research_paper_3_.pdf

File Size
1.3 MB

7 Pages
4,063 Words
23,951 Characters

turnitin Page 1 of 11 - Cover Page

Submission ID trn:oid::1:3075383097

turnitin Page 2 of 11 - Integrity Overview

Submission ID trn:oid::1:3075383097

10% Overall Similarity

The combined total of all matches, including overlapping sources, for each database.

Filtered from the Report

- Bibliography
- Quoted Text

Match Groups

- 24 Not Cited or Quoted 9%**
Matches with neither in-text citation nor quotation marks
- 3 Missing Quotations 1%**
Matches that are still very similar to source material
- 0 Missing Citation 0%**
Matches that have quotation marks, but no in-text citation
- 0 Cited and Quoted 0%**
Matches with in-text citation present, but no quotation marks

Top Sources

- 7% Internet sources
- 8% Publications
- 5% Submitted works (Student Papers)

REPORT PLAGIARISM REPORT

bde report (1).docx

paper1
UG Panel - Selvaraj
SRM Institute of Science & Technology

Document Details

Submission ID
trn:oid::1:3075456578

Submission Date
Nov 11, 2024, 11:23 PM GMT+5:30

Download Date
Nov 11, 2024, 11:32 PM GMT+5:30

File Name
bde_report_1_.docx

File Size
490.8 KB

63 Pages
13,049 Words
79,892 Characters



Page 1 of 68 - Cover Page

Submission ID trn:oid::1:3075456578



Page 2 of 68 - Integrity Overview

Submission ID trn:oid::1:3075456578

9% Overall Similarity

The combined total of all matches, including overlapping sources, for each database.

Filtered from the Report

- Bibliography
- Quoted Text

Match Groups

- 61 Not Cited or Quoted 9%**
Matches with neither in-text citation nor quotation marks
- 0 Missing Quotations 0%**
Matches that are still very similar to source material
- 0 Missing Citation 0%**
Matches that have quotation marks, but no in-text citation
- 0 Cited and Quoted 0%**
Matches with in-text citation present, but no quotation marks

Top Sources

- 7% Internet sources**
- 4% Publications**
- 6% Submitted works (Student Papers)**