

Disease Prediction and Health Guidance Systems

Sadhana Mitra¹, Sonali Kumari², Sony Jena³, S. Nagendra Chary⁴, S.Ramadoss⁵

^{1,2,3,4}*Student, Dept of CSE, School of computing, Bharath Institute of Higher Educations and Research, Chennai, Tamil Nadu*

⁵*Asst. Professor/Dept of CSE, School of computing, Bharath Institute of Higher Educations and Research, Chennai, Tamil Nadu*

Abstract- The Disease Prediction and Health Guidance System is a machine learning-based application designed to predict potential diseases based on user-inputted symptoms. It utilizes a structured dataset of symptoms and corresponding diseases for model training. The previous system employed Decision Tree, Random Forest, and Naïve Bayes classifiers, while the current system uses K-Nearest Neighbors (KNN) and Support Vector Machine (SVM), improving classification accuracy. Additionally, the system provides users with personalized health recommendations, including medication, workouts, and dietary guidance. The frontend is developed using HTML, CSS, JavaScript, and Bootstrap, while Python with Flask powers the backend. By integrating machine learning-driven predictions, this system enhances accessibility to preliminary medical guidance, reducing unnecessary doctor visits and promoting early detection of health conditions. The previous system achieved an accuracy of 85.6%, whereas the current system has improved to 91.2%.

I.INTRODUCTION

Healthcare is a critical aspect of human well-being, and early disease detection plays a vital role in effective treatment and prevention. Many individuals face difficulties in recognizing symptoms and assessing their severity, leading to delayed medical intervention. The advancement of machine learning (ML) has paved the way for innovative solutions in disease prediction and healthcare assistance. Traditional diagnostic methods are often time consuming, costly, and dependent on the availability of medical professionals. In contrast, Machine learning based systems can analyze vast amounts of medical data and provide quick, preliminary diagnoses, improving accessibility and efficiency.

The Disease Prediction and Health Guidance System utilizes multiple machine learning algorithms, including K-Nearest Neighbors (KNN), Support Vector Machine (SVM), Random Forest, Gradient Boosting, and Naïve Bayes, to enhance prediction

accuracy. The system is designed to accept user-input symptoms, process the data using trained ML models, and generate disease predictions along with personalized healthcare recommendations. By leveraging Machine learning driven analytics, this system aims to empower individuals with valuable health insights, facilitating early detection and preventive care.

The rapid advancement of Machine learning (ML) has paved the way for innovative solutions in disease prediction and healthcare assistance. Machine learning driven systems are transforming the healthcare industry by analyzing vast amounts of medical data, identifying patterns, and providing quick, data-driven insights. Unlike traditional diagnostic methods that require extensive consultations, laboratory tests, and expert evaluations, Machine learning based systems can offer preliminary diagnoses in real time, reducing the burden on healthcare professionals while improving accessibility for patients.

By integrating Machine learning driven analytics, this system empowers individuals with valuable health insights, enabling them to take proactive measures toward disease prevention and management. It serves as a digital health assistant, providing accessible and efficient preliminary diagnoses, thus reducing the necessity for frequent doctor visits for minor ailments. Additionally, it aids healthcare professionals by acting as a supporting tool for symptom analysis, streamlining the diagnostic process, and optimizing patient care. The fusion of ML, and healthcare not only enhances disease prediction accuracy but also promotes early detection, leading to better health outcomes and a more efficient healthcare ecosystem.

II.LITERATURE SURVEY

In [1] Most people live long, healthy lives but are too busy to visit doctors for minor symptoms. Many

lack medical knowledge, making consultations time-consuming. Machine learning can help by creating a medicine recommendation system that predicts diseases and suggests appropriate medicines based on symptoms entered by users. an ML-driven system to provide accurate medicine recommendations, improving accessibility and efficiency in healthcare. The study analyzed Naïve Bayes and KNN for medical applications. It found that Naïve Bayes, a probabilistic classifier, was highly effective for medical text classification due to its ability to handle categorical data efficiently. Meanwhile, KNN, a distance-based algorithm, performed well in clustering symptoms, making it useful for grouping similar medical cases. This highlights their complementary strengths in different healthcare-related tasks, improving disease classification and patient diagnosis by [2]. The study in the Journal of Machine Learning in Medicine examined Gradient Boosting for predictive healthcare. It emphasized how Gradient Boosting enhances weak models by sequentially correcting errors, improving accuracy on structured medical datasets. This method is particularly effective for handling complex healthcare data, making predictions more reliable. Its ability to reduce bias and variance makes it valuable for disease prediction, patient risk assessment, and medical decision-making by [3]. The study, presented at the International Conference on Computational Intelligence, compared KNN and SVM for disease classification. It found that SVM, a powerful classifier for high-dimensional data, performed better in complex medical datasets due to its margin optimization. However, KNN, a simpler algorithm, was more effective for smaller datasets, where its instance-based learning approach provided accurate disease classification with minimal computational complexity, making it suitable for limited medical records by [4]. The study explored ML-based diagnosis systems, emphasizing their efficiency in improving medical decision-making. It highlighted challenges in interpretability, as complex ML models often lack transparency. The research stressed the need for explainable ML to build trust in healthcare applications while ensuring accurate, data-driven diagnoses. Future trends focus on enhancing model transparency, robustness, and seamless integration with clinical workflows for better patient outcomes by [5]. The study highlighted that neural networks achieve high accuracy in medical diagnosis due to their ability to learn complex patterns from data. However, they

require extensive training data to generalize well and avoid overfitting. This limitation arises because deep learning models depend on large datasets for effective feature extraction. Despite their accuracy, challenges include data availability, computational cost, and the need for proper model tuning to ensure reliable healthcare predictions by [6]. The study found that CNN and RNN deep learning models outperformed traditional machine learning techniques in medical data analysis. CNN excelled in image-based diagnosis, while RNN was effective for sequential data like patient records. However, these models required significantly higher computational power due to complex architectures and large datasets. Despite their superior accuracy, challenges included processing costs, training time, and the need for specialized hardware like GPUs for optimal performance by [7]. The research highlighted that decision trees are simple, easy to interpret, and useful for medical decision-making. However, they often suffer from overfitting and lower accuracy compared to ensemble methods like Gradient Boosting. Gradient Boosting enhances predictions by combining multiple weak models, improving accuracy and robustness. While decision trees provide transparency, ensemble methods offer superior performance in handling complex medical datasets, making them more suitable for predictive healthcare applications by [8]. The study found that combining multiple AI models, such as ensemble learning or hybrid approaches, significantly improved prediction accuracy in healthcare applications. By leveraging the strengths of different algorithms, these models enhanced reliability and robustness. However, this approach required higher computational resources due to increased processing complexity, data requirements, and model training time. Despite these challenges, integrated ML models proved effective for more precise and comprehensive medical diagnosis and treatment predictions by [9]. The study examined reinforcement learning (RL) as a method to optimize treatment recommendations by analyzing patient history. RL algorithms learned from past medical decisions, adapting dynamically to improve future treatments. This approach personalized healthcare by selecting the most effective interventions based on patient responses. Despite its potential, challenges included the need for extensive data, computational power, and ensuring ethical, reliable decision-making in clinical applications for better patient outcomes by [10].

III. METHODOLOGY

KNN is a supervised learning algorithm that categorizes data points based on the majority vote of their 'K' nearest neighbors. The algorithm follows these steps:

1. *Data Collection and Preprocessing:* Medical datasets containing patient attributes such as age, blood pressure, glucose levels, cholesterol, and other diagnostic features are collected. Data preprocessing involves handling missing values, normalizing data, and selecting relevant features for classification. The dataset used for this study has been verified and validated by Bharath Medical Hospital to ensure accuracy and reliability in disease prediction.
2. *Choosing the Optimal Value of K:* The choice of K plays a crucial role in model performance. If K is too small, the model becomes sensitive to noise, leading to overfitting. If K is too large, the model may misclassify points due to excessive smoothing. The optimal K value is determined using cross-validation methods.
3. *Distance Calculation:* The similarity between the test sample and existing cases is determined using distance metrics. The most commonly used distance metric in KNN is the Euclidean distance, given by:

$$d(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

where x and y are two data points with features each. Other distance measures such as Manhattan distance:

$$d(x, y) = \sum_{i=1}^n |x_i - y_i|$$

and Minkowski distance:

$$d(x, y) = \left(\sum_{i=1}^n |x_i - y_i|^r \right)^{\frac{1}{r}}$$

can also be used depending on the dataset characteristics.

4. *Classification Rule:* After computing the distances, the K closest neighbors are identified. The class of the majority among these neighbors determines the classification of the new data

point. The decision function for classification is given by:

$$f(x) = \arg \max_c \sum_{i=1}^k I(y_i = C)$$

where is an indicator function that returns 1 if the neighbor belongs to class and 0 otherwise.

5. *Prediction and Decision Making:* The new data point is assigned the class label that appears most frequently among its K nearest neighbors. This classification is then used to predict diseases or recommend health interventions.

Applications in Disease Prediction

KNN is effectively used for predicting various diseases, including:

1. *Diabetes Prediction:* By analyzing features like blood sugar levels, BMI, and age, KNN can classify patients as diabetic or non-diabetic.
2. *Heart Disease Prediction:* Using parameters like blood pressure, cholesterol levels, and heart rate, KNN predicts the likelihood of cardiovascular diseases.
3. *Cancer Diagnosis:* KNN helps in identifying cancerous tissues by analyzing features from medical imaging and biopsy reports.

Applications in Health Guidance

Beyond disease prediction, KNN assists in providing personalized health recommendations:

1. *Dietary Suggestions:* Based on patient history and existing dietary habits, KNN suggests tailored nutrition plans.
2. *Exercise Recommendations:* By analyzing physical activity levels and medical conditions, KNN provides personalized workout guidelines.
3. *Preventive Health Measures:* KNN can identify at-risk individuals and suggest preventive health screenings based on historical data.

Advantages of KNN in Healthcare

1. *Simplicity:* KNN is easy to implement and does not require complex training processes.
2. *Adaptability:* The algorithm can be applied to various types of medical datasets.

3. *Non-Parametric Nature*: KNN does not assume any specific distribution of the data, making it suitable for diverse healthcare applications.

III. PROPOSED SYSTEM

The diagram represents a machine learning-based disease prediction and drug recommendation system. Below is a detailed explanation of each step in the process:

Dataset Import:

The process starts with importing a dataset that contains medical records, including symptoms, diseases, and prescribed medications. This dataset serves as the foundation for training the machine learning models.

Data Cleaning (Pandas):

The imported dataset is often raw and may contain missing values, duplicate records, or inconsistencies. Using the Pandas library in Python, data is cleaned to ensure that it is structured and ready for analysis. Cleaning involves handling missing values, standardizing formats, and filtering out unnecessary data.

Training of Disease, Symptoms, and Prescriptions:

Once the dataset is cleaned, the system trains on the relationships between symptoms, diseases, and prescribed medications. This step involves creating a structured model where symptoms are mapped to potential diseases, and corresponding prescriptions are identified.

Visualizing:

The processed data is visualized to gain insights into patterns, such as common symptoms associated with specific diseases. Visualization techniques like graphs, heatmaps, or bar charts help in understanding trends in the dataset.

Importing Machine Learning Classifier Modules:

Various machine learning algorithms are used for disease classification to enhance prediction accuracy. K-Nearest Neighbors (KNN) is a distance-based classifier that identifies diseases by comparing a patient's symptoms with similar cases in the dataset. Random Forest (RF) is an ensemble learning method that improves accuracy by combining multiple decision trees, reducing overfitting, and enhancing robustness. Support Vector Machine

(SVM) is a classification algorithm that separates diseases by mapping symptom features into a high-dimensional space and finding the optimal hyperplane for classification.

Checking Accuracy:

After training the classifiers, their accuracy is evaluated to determine which model performs best. Metrics such as precision, recall, F1-score, and accuracy are used to measure their effectiveness in predicting diseases.

SVM Selection

Based on the accuracy results, the Support Vector Machine (SVM) classifier is selected as the final model for disease prediction and drug recommendation. SVM is often preferred for its ability to handle high-dimensional medical data effectively. It works by finding the optimal hyperplane that best separates different disease classes, ensuring high accuracy. Additionally, SVM's robustness against overfitting makes it suitable for reliable medical diagnosis and treatment recommendations.

Recommendation of Drugs:

After a disease is predicted, the system suggests suitable medications based on the trained model. The recommended drugs are derived from past medical records, ensuring that the prescriptions align with expert medical practices.

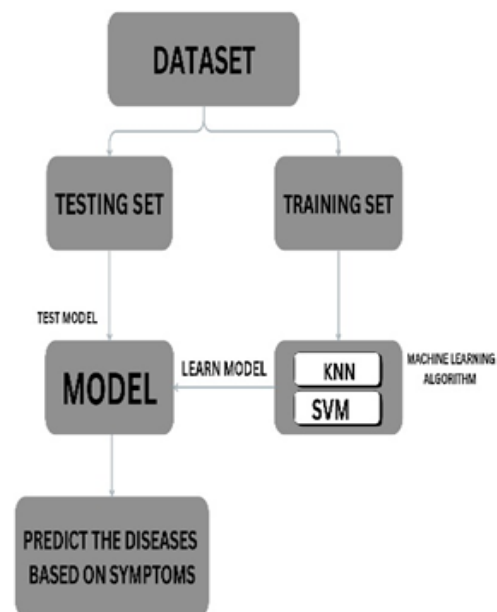


Fig 1. Block diagram

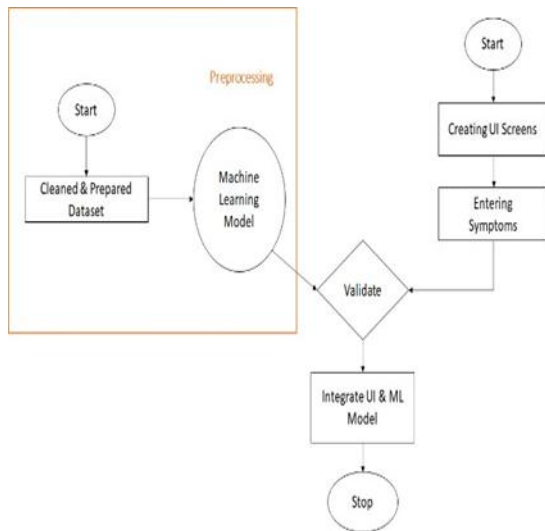


Fig 2. Designed architecture diagram

Fig.1. represents a machine learning-based system workflow, particularly focusing on preprocessing and integration with a user interface (UI). It starts with data preprocessing, where a cleaned and prepared dataset is fed into a machine learning model. After training, the model is validated and integrated with a UI system. The UI is designed for users to enter symptoms, and after validation, the machine learning model is incorporated into the system for prediction, leading to the final stop.

Fig 2. elaborates on the dataset processing for disease prediction. It divides the dataset into training and testing sets. The training set is used to train the model using machine learning algorithms like K-Nearest Neighbors (KNN) and Support Vector Machine (SVM). The trained model is then tested on the testing set, and ultimately, it predicts diseases based on the symptoms provided. This diagram focuses more on the machine learning model training and testing process.

IV. RESULT AND DISCUSSION

The images depict a web-based AI-powered health guidance system that predicts diseases based on user-inputted symptoms and provides recommendations regarding precautions, medications, workouts, and diet. The user interface consists of a search bar where users can enter symptoms such as itching, sleeping issues, or aches. Upon clicking the "Predict" button, the system processes the input and provides relevant results in categorized sections.

Fig 1, displays the initial user input phase of the system. The title "Health Care Center" is

prominently visible at the top. The interface presents a text box where the user has typed the symptoms "stomach pain, acidity" as input. Below the text box is the Predict button, which, when clicked, allows the AI system to analyze the symptoms and suggest a possible disease. This phase of interaction is crucial as it serves as the starting point where the user engages with the system by entering their health concerns.

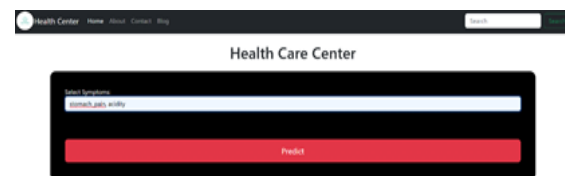


Fig 1. Received symptoms

Fig 2, shows the AI-powered health prediction system after it has analyzed the user's symptoms and predicted a disease. The system has displayed a pop-up box titled "Predicted Disease", where it identifies GERD (Gastroesophageal Reflux Disease) based on the symptoms entered by the user. Below the prediction interface, the system presents several buttons labelled Disease, Description, Precaution, Medications, Workouts, and Diets, suggesting that the user can explore more details about their predicted condition. The interface is designed with a clear layout, with a dark-themed input section and a contrasting red Predict button.

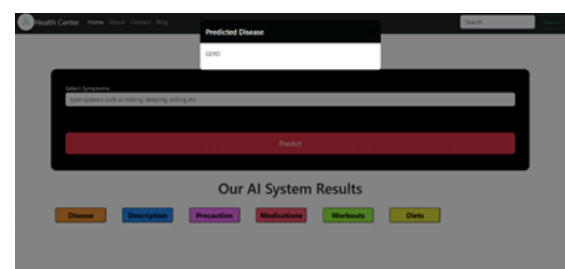


Fig 2: Predicted diseases

Fig 3, shows the disease description feature of the AI health system. After predicting GERD, the user has clicked on the "Description" button, triggering a pop-up box labelled "Description". The content in the box explains that GERD (Gastroesophageal Reflux Disease) is a digestive disorder that affects the lower esophageal sphincter. This feature is beneficial as it provides users with a brief understanding of the predicted disease, allowing them to make informed decisions about seeking medical attention.

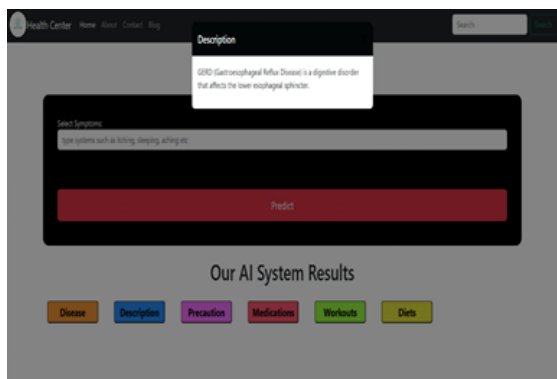


Fig 3: Description of disease

Fig 4, is a medications section, showing possible pharmaceutical treatments based on the predicted disease. The listed medications include Proton Pump Inhibitors (PPIs), H2 Blockers, Antacids, Prokinetics, and Antibiotics, which are commonly prescribed for conditions affecting the digestive system, such as acid reflux, ulcers, or bacterial infections.

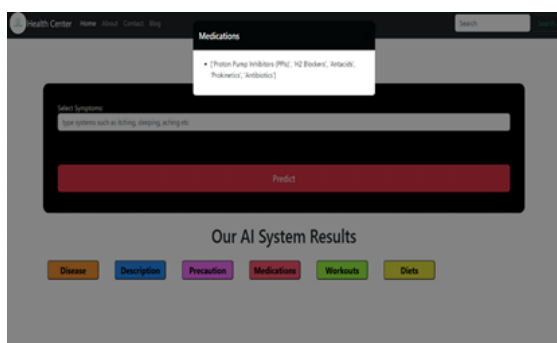


Fig 4: Medication

Fig 5, after medication the main prediction interface of the system categorizes its results into sections labeled Disease, Description, Precaution, Medications, Workouts, and Diets, each with color-coded buttons for easy navigation. This ML-driven tool aims to assist users in understanding their symptoms, suggesting possible conditions, and providing health recommendations, although a professional medical consultation is still advised for accurate diagnosis and treatment.

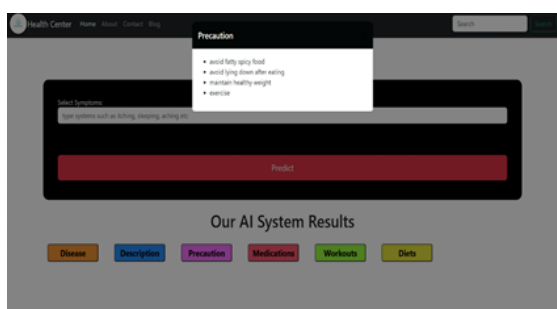


Fig 5: Precaution

Fig 6, is highlighted for displaying dietary recommendations for a specific health condition. The suggested foods include a low-acid diet, fiber-rich foods, ginger, licorice, and aloe vera juice, which are known for their digestive benefits and soothing effects on conditions like acid reflux or gastrointestinal issues.

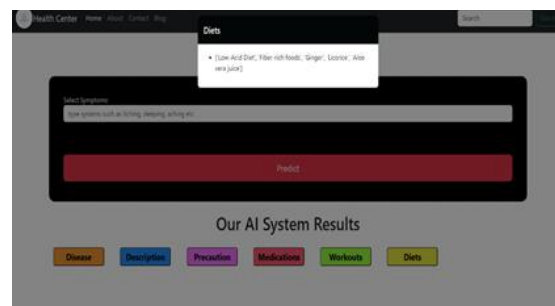


Fig 6: Diets

V. CONCLUSION

This study highlights the significance of machine learning in disease prediction and health guidance. The integration of machine learning algorithms such as K-Nearest Neighbors (KNN), Random Forest (RF), and Support Vector Machine (SVM) enables accurate classification of diseases based on patient symptoms. KNN is useful for identifying patterns by comparing symptoms with similar cases, while SVM efficiently classifies diseases by separating feature spaces using hyperplanes. Among these, Random Forest has shown superior accuracy and reliability due to its ensemble learning approach, reducing overfitting and improving classification robustness.

K-Nearest Neighbors is a promising machine learning approach for disease prediction and health guidance. Its ability to classify patient data based on similarity makes it a valuable tool in personalized healthcare. However, challenges such as computational efficiency and data preprocessing must be addressed to enhance its effectiveness. Future research can focus on optimizing KNN for large-scale medical datasets and integrating it with deep learning models for improved accuracy.

By incorporating these models, the system not only predicts potential diseases but also provides personalized medication recommendations, dietary suggestions, precautions, and workout plans. This AI-driven approach can enhance accessibility to

healthcare, offering preliminary diagnoses and guidance to individuals who may not have immediate access to medical professionals. However, challenges remain, such as improving dataset quality, handling complex symptom variations, and integrating real-time clinical data for better accuracy.

Future enhancements may include expanding the dataset, optimizing model performance, and integrating deep learning techniques for better generalization. While machine learning provides a powerful tool for disease prediction, it should complement, rather than replace, professional medical consultations to ensure accurate diagnosis and treatment.

REFERENCE

- [1] Satvik Garg Department of Computer science Jaypee University of Information Technology Solan, India 2021.
- [2] T. Venkat Narayana Rao, Anjum Unisa, Kotha Sreni TERNATIONAL JOURNAL OF SCIENTIFIC & TECHNOLOGY RESEARCH VOLUME 9, ISSUE 02, FEBRUARY 2020.
- [3] Binu Thomas and Amruth K John 2021 IOP Conf. Ser.: Mater. Sci. Eng. 1085 012011, 2019.
- [4] Rekha Nagar^{1*}, Yudhvir Singh^{2*} U.I.E.T (M.D.U), India International journal of Emerging Technologies and Innovative Research (www.jetir.org), ISSN:2349-5162. Vol.6 Issue 4, page no 471-474, April 2019.
- [5] Benjamin Stark¹, Constanze Knahl², Mert Aydin³, Karim Elish⁴ Department of Computer Science, Florida Polytechnic University, Lakeland, US 2019.
- [6] Recommendation system using Machine Learning Suhasini Parvatikar Computer Engineering, Assistant Professor SCOE, Kharghar, India.
- [7] Smitha Rani, Prakash Kumar, "Deep Learning Approaches for Medical Diagnosis," *International Conference on Artificial Intelligence in Healthcare*, 2020.
- [8] A. Gupta, R. Sharma, "Comparative Analysis of Machine Learning Algorithms for Disease Prediction," *Journal of Healthcare Informatics*, Vol. 12, Issue 3, 2021.
- [9] D. Patel, M. Verma, "Optimizing Healthcare Predictions Using Ensemble Learning," *International Journal of Data Science and Analytics*, 2022.
- [10] S. Choudhary, P. Mehta, "Reinforcement Learning for Personalized Treatment Plans," *IEEE Transactions on Computational Healthcare*, 2023.