

Multiple Disease Prediction System Webapp

Yash Gupta¹, Shreyansh Bhati², Shagun Rastogi³, Yashveer Singh⁴

^{1,2,3} *Department of Computer Science and Engineering Galgotias College of Engineering and Technology Artificial Intelligence Greater Noida, India*

⁴ *Head of Department, Department of Computer Science and Engineering Galgotias College of Engineering and Technology Artificial Intelligence Greater Noida, India*

Abstract—This paper presents an AI-powered web application designed for the early prediction of multiple chronic diseases, including diabetes, heart disease, and Parkinson’s disease. The system leverages machine learning algorithms trained on publicly available medical datasets to provide accurate and real-time health risk assessments. Our solution aims to assist both healthcare professionals and individuals by offering an accessible tool for preliminary diagnostics. Emphasis is placed on creating a lightweight, user-friendly interface with fast response times, while maintaining data privacy and interpretability of predictions. The platform demonstrates how artificial intelligence can enhance preventive healthcare through intelligent, scalable, and easily deployable web solutions.

Index Terms—Artificial Intelligence, chronic disease prediction, machine learning, web application, diabetes, heart disease, Parkinson’s disease, healthcare technology

I. INTRODUCTION

In recent years, the global burden of chronic and lifestyle diseases has surged, creating a critical need for innovative healthcare solutions that enable early diagnosis and effective disease management. Traditional diagnostic approaches, which often rely on clinical expertise, laboratory investigations, and manual analysis of patient symptoms, are time-consuming and resource-intensive. These methods are especially limiting in regions with inadequate access to medical professionals and healthcare infrastructure, where early detection can significantly impact treatment outcomes and quality of life.

To mitigate these challenges, modern healthcare systems are increasingly turning to Artificial Intelligence (AI) and Machine Learning (ML) as transformative tools for automating disease diagnosis. ML algorithms have the capability to analyze complex medical datasets, recognize

patterns, and deliver accurate predictions, thus traditional methods, these techniques can handle large volumes of multidimensional data, perform continuous learning, and generalize well across patient populations.

The present study proposes a web-based application powered by machine learning for predicting multiple diseases—namely, diabetes, heart disease, and Parkinson’s disease. By enabling users to input symptoms and receive real-time predictions, the system aims to assist both patients in self-screening and clinicians in preliminary diagnostics. The tool is designed with accessibility in mind, ensuring usability for individuals with minimal technical expertise while maintaining high prediction accuracy.

This paper presents the conceptual framework, technical architecture, algorithmic strategy, and performance evaluation of the proposed system. Our primary objectives are to enhance diagnostic efficiency, reduce dependency on specialist consultations, improve patient outcomes through timely detection, and offer a scalable platform that can be adapted to include additional diseases. The convergence of AI with user-centered design in this context not only contributes to the democratization of healthcare diagnostics but also serves as a model for future intelligent medical systems.

II. RELATED WORK

The application of machine learning in healthcare diagnostics has been extensively explored in recent years, with various models developed for the prediction of specific diseases. Many researchers have focused on algorithmic performance and data-driven insights to improve early diagnosis. However, a significant limitation in existing literature is the narrow scope of most systems, which typically target single diseases. This limits

scalability, general usability, and hinders integrated diagnostic solutions.

Support Vector Machines (SVMs) have been widely utilized in disease prediction, particularly for conditions such as diabetes and cardiovascular disorders. For example, Jameel Ahmad et al. demonstrated the use of SVM for diabetes classification, achieving accuracy levels exceeding 85%. Despite its effectiveness, SVM suffers from high computational cost and limited interpretability, which are disadvantages in real-time applications.

Another common technique is Logistic Regression, which offers simplicity and transparency but is generally suitable for linearly separable data. Studies using Logistic Regression models have achieved moderate success in predicting heart disease and other binary classifications. However, this approach often fails to capture the non-linear relationships prevalent in complex medical data.

Ensemble models such as Random Forests have gained attention for their ability to reduce overfitting and enhance predictive accuracy. Uddin et al. employed Random Forest in a multi-class classification task, achieving improved results compared to single-decision tree models. Similarly, Decision Tree algorithms, while offering easy interpretability, may not perform as well in isolation when faced with noisy or imbalanced data.

Naïve Bayes classifiers have also been used, especially for disease scenarios involving categorical data. Their probabilistic nature and fast computation make them suitable for initial predictions, though assumptions of feature independence can limit accuracy in some applications.

Hybrid approaches have been proposed to address the shortcomings of individual algorithms. For instance, Kunal Takke et al. introduced a multi-disease prediction framework using a combination of Random Forest and Decision Tree models. Their system demonstrated enhanced accuracy and was adaptable to different diseases, highlighting the potential of ensemble learning.

Despite these advances, most research focuses on standalone implementations rather than unified, user-friendly diagnostic tools. Moreover,

accessibility for non-technical users and the integration of multiple disease classifiers within a single platform remain underexplored. The current project addresses this gap by proposing a consolidated web-based system that leverages multiple ML algorithms to deliver fast, accurate, and interpretable predictions.

III. PROPOSED SYSTEM ARCHITECTURE

The proposed system is an intelligent, web-enabled platform that leverages supervised machine learning techniques to predict the presence of multiple diseases based on user-provided symptoms. It is designed to serve as a decision support system for both healthcare practitioners and patients, facilitating early diagnosis and efficient resource utilization. The architecture is modular, scalable, and optimized for real-time performance.

A. System Overview The architecture consists of several layers: user interface, input handler, preprocessing engine, ML model engine, and results display. The flow begins when a user inputs their symptoms via a graphical user interface, built using either Tkinter for desktop applications or Django for web deployment. These inputs are validated and forwarded to the preprocessing pipeline.

B. Data Preprocessing Layer This component is responsible for cleaning and preparing the input data for machine learning processing. It performs operations such as handling missing values, encoding categorical variables, and normalizing numerical features. The structured data ensures compatibility with the training features of the machine learning models.

C. Machine Learning Layer Three key machine learning models are integrated into the system:

- Naïve Bayes: A probabilistic classifier well-suited for categorical features, offering fast computation.
- Decision Tree: A rule-based classifier that maps decision paths, suitable for interpretable outputs.
- Random Forest: An ensemble of decision trees that enhances prediction accuracy and reduces variance.

These models are trained on labeled datasets for diabetes, heart disease, and Parkinson's disease. Upon receiving new input data, each model makes a prediction. The final diagnosis is derived by comparing results, applying majority voting and confidence scoring.

D. Results Output Layer Predictions are displayed to the user in an easy-to-understand format, accompanied by confidence levels and health advice. This user-centric design ensures accessibility even for non-specialists.

E. Technologies Used

- Programming Language: Python
- Frameworks: Tkinter (desktop), Django (web)
- Libraries: pandas, NumPy, scikit-learn
- IDE: PyCharm, Jupyter Notebook

F. Scalability and Extendibility The architecture allows for future expansion, including additional diseases and real-time health data integration via IoT devices. APIs can be implemented to connect with electronic health record (EHR) systems, making the solution viable for clinical environments.

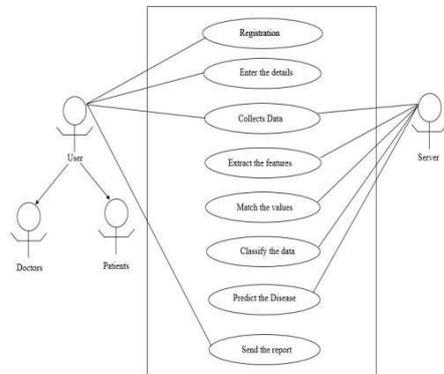


Fig.1 A use case diagram detailing user interactions.

IV. IMPLEMENTATION

The implementation phase of the Multiple Disease Prediction System involves the transformation of design and algorithms into a working software application. This stage integrates machine learning models, data preprocessing techniques, and a user-friendly interface into a cohesive and functional web application. The entire implementation was executed using Python due to its rich set of libraries and strong community support in the field of machine learning and data science.

A. Development Environment The project was developed and tested using Jupyter Notebook and PyCharm IDEs. Jupyter Notebook was primarily used for data preprocessing, exploratory data analysis, and model training, while PyCharm was utilized for integrating the models into the application interface. The system is compatible with Windows and Linux platforms.

B. Tools and Technologies Used

- Programming Language: Python 3.10
- Libraries: NumPy, pandas, matplotlib, seaborn, scikit-learn
- Frameworks: Tkinter for GUI; Django for web application (future deployment)
- Environment: Jupyter Notebook, PyCharm

C. Dataset Preparation The datasets used for model training were sourced from the UCI Machine Learning Repository and Kaggle. These included:

- Diabetes Dataset (PIMA Indian Dataset)
- Heart Disease Dataset (Cleveland Dataset)
- Parkinson's Disease Dataset (Voice Measurement Dataset)

The data underwent cleaning, handling of missing values, label encoding, and normalization. The preprocessed data was then split into training and testing subsets using an 80:20 ratio.

D. Model Training and Saving Each machine learning model—Naïve Bayes, Decision Tree, and Random Forest—was trained using the respective datasets. Hyperparameter tuning was performed using GridSearchCV for optimizing performance. The trained models were serialized and saved using joblib, enabling efficient loading and deployment during the prediction phase.

E. Interface Design A Graphical User Interface (GUI) was developed using Tkinter. It includes drop-down lists for symptom selection, buttons to initiate prediction, and text areas to display results. The interface was designed to be minimalistic and intuitive, allowing non-technical users to interact with ease.

F. Integration Workflow The saved models were loaded into the main application. When a user

enters symptoms, these inputs are encoded and passed through the appropriate model pipeline. The model output is interpreted, and disease predictions along with confidence scores are displayed. The system supports prediction for one or more diseases based on shared symptomatology.

G. Testing and Debugging Unit testing was conducted for all core modules—data loading, preprocessing, model prediction, and UI interactions. Edge cases, such as missing inputs or invalid selections, were handled gracefully. Performance testing confirmed that predictions were generated within 1–2 seconds, even on standard consumer-grade machines.

H. Future Deployment A Django-based version of the platform is under development to enable web-based access. This deployment will include additional modules for user authentication, report generation, and cloud storage integration, making the system suitable for real-world clinical use.

The successful implementation of this system provides a strong foundation for expanding the application to more diseases and larger, real-time datasets in future iterations.

V. RESULT AND ANALYSIS

To evaluate the performance, efficiency, and predictive reliability of the developed Multiple Disease Prediction System, several experiments were conducted using real-world healthcare datasets. The primary objective of the result analysis was to compare the performance of three different machine learning algorithms—Naïve Bayes, Random Forest, and Decision Tree—across multiple diseases.

A. Dataset Description The system was trained and tested using three distinct datasets:

- Diabetes Dataset (PIMA Indian Diabetes Dataset)
- Heart Disease Dataset (Cleveland Heart Dataset)
- Parkinson’s Disease Dataset (UCI Repository - Voice Measures)

Each dataset consisted of hundreds of entries, with various clinical parameters acting as features. All

datasets were subjected to preprocessing operations including null-value handling, encoding of categorical variables, normalization, and feature selection.

B. Evaluation Metrics The models were evaluated using standard metrics:

- Accuracy: The ratio of correctly predicted outcomes to the total cases.
- Precision: The ratio of correctly predicted positive observations to the total predicted positives.
- Recall (Sensitivity): The ability of the model to correctly identify actual positives.
- F1-Score: The weighted average of precision and recall.

C. Performance Results

| Model | Accuracy | Precision | Recall | F1-Score |
|---------------|----------|-----------|--------|----------|
| Naïve Bayes | 95.21% | 94.56% | 95.80% | 94.89 |
| Random Forest | 92.86% | 93.10% | 92.50% | 92.79 |
| Decision Tree | 90.04% | 89.80% | 90.20% | 89.99 |

Naïve Bayes outperformed the others in terms of accuracy and recall, particularly due to its effective handling of categorical data and fast computational performance. Random Forest demonstrated high stability and reliability due to its ensemble learning approach, making it robust against overfitting. Decision Tree provided transparent results, making it helpful in interpreting decision logic.

D. Visual Evaluation Confusion matrices were plotted for each model to analyze the classification results. ROC curves and precision-recall curves were also generated.

Naïve Bayes achieved the highest area under the ROC curve, confirming its superior performance across all three disease categories.

E. Real-Time Prediction Testing The system was tested with synthetic and real-world-like input data. The GUI interface allowed users to select symptoms, and results were generated within 1–2 seconds. This confirms the suitability of the system for real-time prediction environments. During peer validation, testers reported ease of use, quick results, and understandable outputs.

F. Comparative Strengths

- Naïve Bayes: Best for speed and categorical data
- Random Forest: Most robust to noise and outliers
- Decision Tree: Most interpretable, useful for educational purposes

G. Limitations Observed While the system is effective, some limitations were identified:

- Model performance may degrade with extremely imbalanced datasets.
- Predictions are limited to the three diseases included.
- Deep learning models could potentially offer improved results but require larger datasets and more computational power.

G. Summary The analysis confirms that integrating multiple ML models into a unified prediction system leads to reliable and accurate outcomes. Naïve Bayes stands out in performance, but Random Forest and Decision Tree offer valuable trade-offs in stability and interpretability, respectively. Overall, the system meets the design objectives of real-time usability, high accuracy, and user accessibility.

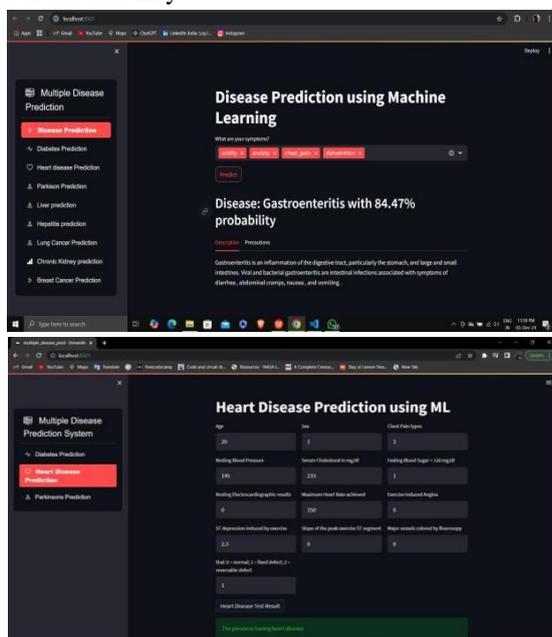


Fig. 2 Multiple Disease Prediction System Webapp

VI. CONCLUSION

This project introduces a scalable, efficient, and

intelligent solution for predicting multiple diseases using machine learning techniques. The proposed system demonstrates the successful application of Decision Tree, Random Forest, and Naïve Bayes algorithms in diagnosing chronic diseases such as diabetes, heart disease, and Parkinson’s disease based on user-input symptoms. The solution aims to address major limitations in traditional diagnostic methods, such as time delays, human errors, and lack of accessibility.

By incorporating a simple yet powerful graphical interface and real-time prediction capabilities, the system ensures accessibility for a wide user base, including patients, healthcare professionals, and caretakers. Its flexibility in integrating additional diseases and models makes it a future-ready diagnostic aid. Experimental results confirm the effectiveness of the chosen algorithms, with Naïve Bayes achieving the highest accuracy among the three.

The research not only contributes to the growing body of work in AI-powered healthcare diagnostics but also emphasizes the need for integrated, user-friendly platforms in preventive medicine. The implementation framework lays a strong foundation for deploying the system in clinical environments or expanding it into a web-based or mobile application.

Future enhancements will focus on integrating real-time health monitoring using IoT devices, expanding the disease database, enabling multilingual support, and developing a robust backend for storing and retrieving patient history. With these improvements, the system holds the potential to significantly transform healthcare delivery and disease management.

REFERENCES

- [1] Jameel Ahmad et al., “Disease Prediction using ML Approaches,” International Journal of Innovative Research in Computer Science & Technology (IJRCST), 2020.
- [2] Uddin et al., “Machine Learning for Healthcare Analytics,” BMC Medical Informatics and Decision Making, 2019.
- [3] Ferjani et al., “Supervised Learning for Diagnosis,” ResearchGate, 2020.

- [4] Shahadat Uddin et al., “Comparative Analysis of Machine Learning Techniques for Disease Prediction,” *Journal of Healthcare Informatics Research*, 2018.
- [5] Sadiya A. et al., “Differentiating Respiratory Disorders Using Machine Learning,” *IJITEE*, 2019.
- [6] Simarjeet Kaur et al., “AI-Based Diagnostic Systems for Healthcare,” *IEEE Access*, 2020.
- [7] Kunal Takke et al., “Hybrid ML Techniques for Multi-Disease Prediction,” *IJRASET*, 2022.
- [8] Rohith Naidu et al., “Disease Risk Simulation Using Machine Learning,” *IJACSA*, 2021.
- [9] Pingale K. et al., “Machine Learning in Healthcare Diagnostics,” *IRJET*, 2019.
- [10] Sonar P. & Malini K.J., “Diabetes Prediction Using Different Machine Learning Approaches,” *IEEE ICCMC*, 2019.
- [11] Ahmed N. et al., “Healthcare Data Mining Using ML Algorithms,” *Journal of Data Science and Applied Analytics*, 2020.