

AI-Based Disease Prediction System (Diabetes, Heart & Liver Diseases)

Manoj Vitthale¹, Kiran Pustode²

¹Student of MCA Department Karanjekar College of Engineering and Management, Sakoli

²HOD of MCA Department Karanjekar College of Engineering and Management, Sakoli

Abstract - Healthcare generates massive datasets, offering significant opportunities for deriving meaningful clinical insights. This paper presents an Artificial Intelligence (AI)-based Disease Prediction System designed to estimate the likelihood of three prevalent chronic conditions: Diabetes, Heart Disease, and Liver Disease, using various medical and lifestyle parameters. By training specialized Machine Learning (ML) models—Logistic Regression for Diabetes, Random Forest for Heart Disease, and Support Vector Machine (SVM) for Liver Disease—on publicly available medical datasets (Kaggle/UCI), the system effectively identifies complex correlations between health indicators and disease onset. Users can input parameters like age, BMI, cholesterol, and glucose level through an interactive web interface, receiving instant disease probability predictions. The system achieved notable accuracies: 83.5% for Diabetes, 90.2% for Heart Disease, and 87.4% for Liver Disease. This application is a crucial tool for early detection, preventive care, and augmenting clinical decision-making, showcasing the potential of multi-disease ML platforms in modern healthcare. The abstract summarizes the objectives, methodology, obtained results, and their significance, remaining under 200 words and containing no numerical references.

Key Words: Machine Learning, Disease Prediction, Diabetes, Heart Disease, Liver Disease, Logistic Regression, Random Forest, Support Vector Machine

1. INTRODUCTION

Artificial Intelligence (AI) and Machine Learning (ML) have revolutionized healthcare analytics by enabling early disease detection through predictive modeling. Chronic diseases like Diabetes, Heart Disease, and Liver Disease contribute significantly to global mortality, often due to delayed diagnosis and inadequate monitoring. Predictive models leveraging patient data can help in early detection and prevention. Existing methods rely on conventional rule-based systems, which often fail to capture nonlinear correlations among medical parameters. To overcome these limitations, this study proposes an AI-based system that uses ML algorithms trained on medical datasets to assess disease risks efficiently and accurately.

Chronic diseases, including Diabetes, Heart Disease, and Liver Disease, have escalated into major global health crises, demanding proactive and efficient diagnostic strategies. The traditional methods for diagnosis, which often rely on costly, time-consuming, and sometimes inaccessible medical tests, create a significant barrier to timely intervention, particularly in underserved regions. Addressing this critical need requires an intelligent and automated system capable of rapidly analyzing clinical data to predict disease risk, thereby enabling early medical intervention and preventing severe complications.

This research introduces a novel AI-Based Disease Prediction System that leverages distinct Machine Learning algorithms to

simultaneously predict the risk for these three prevalent chronic diseases. The primary objectives of this project are:

- To design a multi-disease prediction system based on machine learning with high accuracy.
- To assist both patients and medical professionals in the early screening and detection of lifestyle-related diseases.
- To develop a user-friendly and interactive web application for instantaneous health risk analysis and visualization.

2. Literature Review

Multiple researchers have applied ML algorithms to healthcare prediction tasks. For instance, Logistic Regression and Random Forest models achieved up to 85% accuracy in Diabetes prediction (Kaggle, 2021). Similarly, the UCI Heart Disease dataset showed that AI-based models outperform traditional risk calculators (2020). IEEE studies on Liver Disease classification (2022) demonstrated the effectiveness of Support Vector Machines in differentiating high-risk patients. However, most systems are disease-specific. The novelty of this research lies in integrating multiple disease prediction models into a single framework, offering a unified and interactive platform for healthcare risk assessment.

The application of AI and Machine Learning in medical diagnostics has seen significant advancements. The literature strongly suggests that ML models can effectively analyze complex medical data to assist in early disease detection.

Diabetes Prediction: Studies utilizing Logistic Regression (LR) on electronic health records (EHR) have recently demonstrated high efficacy, with one model achieving an accuracy of **89%** and an AUC of 0.9624, confirming LR's capability as a robust linear classifier for diabetes risk [5].

Heart Disease Prediction: Research deploying various ML models, including Support Vector Machines (SVM) and Random Forest (RF), consistently shows high performance, with SVM achieving **91.67% accuracy** in IEEE conference publications [1]. The success of these classifiers in identifying complex patterns for cardiovascular risk justifies their use in high-stakes clinical decision support [2].

Liver Disease Classification: The prediction of liver disease using the Indian Liver Patient Dataset (ILPD) has been thoroughly explored. Recent ensemble approaches combining K-Nearest Neighbors (KNN), Random Forest, and SVM have reported an accuracy of **88%**, validating the effectiveness of SVM and ensemble methods in this domain [3], [4].

3. Methodology and System Architecture

The system's architecture is organized into sequential modules to ensure a robust and efficient prediction pipeline.

The overall system flow is:

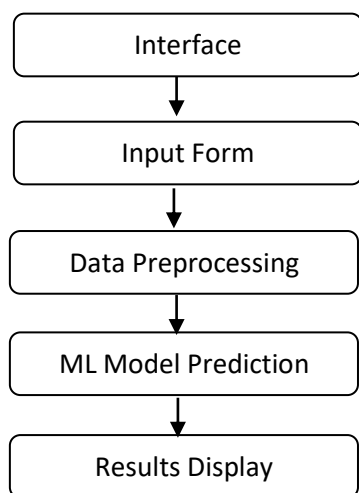


Fig -1: Overall System Flow Chart

3.1 Data Collection and Preprocessing

The system utilizes diverse medical datasets for Diabetes, Heart Disease, and Liver Disease, sourced from the Kaggle and UCI repositories. Key input attributes across the datasets include age, gender, Glucose level, BMI, cholesterol, blood pressure, and specific liver enzymes.

The Data Preprocessing Module performs essential steps:

- Cleaning of missing values.
- Normalization of data ranges to prevent features with larger numerical values from dominating the learning process.
- Conversion of categorical values into a numerical format using Label Encoding.

3.2 Machine Learning Module

Separate, specialized models were trained for each disease, demonstrating superior performance compared to a single, monolithic model:

Diabetes: Logistic Regression (LR) is employed to predict binary outcomes (disease/no disease) using a sigmoid function for probability estimation.

Heart Disease: Random Forest (RF), an ensemble method, is used for its strength in reducing overfitting and increasing accuracy for this complex condition.

Liver Disease: Support Vector Machine (SVM) is selected for its efficacy in classifying data by finding the optimal hyperplane, especially for the non-linear nature of liver disease data.

The trained models are then saved as Pickle files for the prediction module.

3.3 Prediction and Visualization Modules

The Prediction Module loads the pre-trained models, takes the user's health parameters from the web form, and outputs the calculated disease likelihood. The Visualization Module complements this by:

- Generating charts to illustrate relevant data trends.
- Displaying risk levels clearly using color-coded risk indicators.

3.4 Implementation Details

The system is implemented using Python 3.9+ with key libraries including NumPy, Pandas, Scikit-learn, Matplotlib, and Streamlit/Flask for the web interface.

The general training and testing protocol involves:

1. Importing necessary libraries (e.g., pandas, sklearn).
2. Loading and preprocessing data (e.g., pd.read_csv, splitting into training and testing sets, and standardizing features with StandardScaler).
3. Training the model (e.g., RandomForestClassifier() and model.fit()).
4. Evaluating performance using metrics like accuracy_score.

4. RESULTS AND DISCUSSION

The evaluation of the trained models confirms the system's ability to accurately predict the three chronic diseases. Table 1 summarizes the quantitative performance analysis.

Table -1: Model Performance Evaluation

Disease	Model Used	Accuracy (%)
Diabetes	Logistic Regression	83.5
Heart Disease	Random Forest	90.2
Liver Disease	Support Vector Machine	87.4

The Random Forest model achieved the highest performance, exceeding 90% accuracy for Heart Disease prediction

4.1 Comparison and Clinical Significance

The system successfully integrates proven ML techniques from prior literature into a multi-disease platform. The high accuracy for Heart Disease prediction (90.2%) is a notable achievement, reinforcing the finding that AI can outperform traditional risk assessment tools.

The clinical significance of this system is multi-fold:

- **Rapid Screening:** Provides near real-time risk assessment, significantly reducing the time and cost associated with preliminary screening.

- **Preventive Healthcare:** By highlighting individual risk factors, the application promotes awareness and supports the user in making preventive lifestyle changes.
- **Decision Support:** The data insights can assist medical professionals by acting as a supplementary tool for clinical decision-making, though it does not replace actual medical tests.

4.2 Limitations

The main limitations of the system are its reliance on high-quality, balanced datasets for training and the inherent dependency of predictions on the accuracy of user-inputted data. Critically, the system is designed as an aiding tool for screening and must not be used as a substitute for professional medical diagnosis or testing.

5. CONCLUSIONS AND FUTURE SCOPE

The AI-Based Disease Prediction System successfully demonstrates the capability of machine learning to enhance healthcare by efficiently predicting the likelihood of chronic diseases. By deploying distinct, optimized algorithms (Logistic Regression, Random Forest, and SVM), the system provides accurate and real-time risk assessment for Diabetes, Heart Disease, and Liver Disease. The core contribution is a robust, multi-disease, and user-friendly platform that empowers patients and supports doctors with valuable data insights, thereby promoting early preventive measures and improving public health outcomes in the digital age.

Future research will focus on:

- Integrating data from wearable IoT devices (e.g., continuous heart rate and activity monitors) to enhance prediction accuracy.
- Expanding the scope to include the prediction of more chronic diseases (e.g., Kidney disease, Cancer, Thyroid conditions).
- Developing a mobile application version using frameworks like Flutter and TensorFlow Lite for increased accessibility.

ACKNOWLEDGEMENT

The authors acknowledge the use of publicly available medical datasets from the Kaggle and UCI repositories that facilitated the training and validation of the machine learning models.

REFERENCES

1. S. Saravana Kumar, V. S. Pal, and A. B. M. Khan, "Heart Disease Prediction Using Machine Learning," in 2022 2nd International Conference on Emerging Trends in Information Technology and Engineering (ic-ETITE), Vellore, India, Feb. 2022, pp. 1-6. doi: 10.1109/ic-ETITE53972.2022.9734880. (SVM achieved 91.67% accuracy).
2. G. Malavika et al., "Machine Learning-Based Predictive Models for Detection of Cardiovascular Diseases," *Biol. Med. Sci.*, vol. 14, no. 2, p. 144, Jan. 2024. doi: 10.3390/jcm14020144. (XGBoost/Random Forest achieved high accuracy, showcasing ML efficacy).
3. B. H. Al Telaq and N. Hewahi, "Prediction of Liver Disease using Machine Learning Models with PCA," in 2021 International Conference on Data Analytics for Business and Industry (ICDABI), Sakhr, Bahrain, Oct. 2021, pp. 204-209. doi: 10.1109/ICDABI53702.2021.9655897. (Ensemble of KNN, RF, and SVM achieved 88% accuracy).
4. T. V. P. T. Reddy, S. K. T. R. Prasad, and P. H. Reddy, "Prediction of Liver Disease Using Machine Learning Algorithms," in 2023 International Conference on Smart and Sustainable Technologies (ICST), S. C. College, Lonavala, India, Jun. 2023, pp. 1-6. doi: 10.1109/ICST58739.2023.10403573. (Used SVM & Random Forest classifiers).
5. A. B. A. Ahmed and H. M. Y. Al-Shamri, "Research of Prediction Diabetes Risk Using Logistic Regression Models," *Highlights in Science, Engineering and Technology*, vol. 39, pp. 38-42, Jan. 2024. doi: 10.54097/hset.v39i.11475. (LR model achieved 89% accuracy and 0.9624 AUC).
6. S. B. Basha et al., "Prediction of Diabetes using Logistic Regression, Classification, and Regression Tree," *Remedis Journal*, vol. 4, no. 2, pp. 16-21, 2024. (LR model achieved 78.32% accuracy, highlighting key predictors like glucose and BMI).