

www.ijprems.com editor@ijprems.com

# DISEASE PREDICTION USING MACHINE LEARNING

# Dhaman Kovachi<sup>1</sup>, Abhay Pendariya<sup>2</sup>, Himanshu Chandrakar<sup>3</sup>,

# Prof. Om Prakash Barapatre<sup>4</sup>

<sup>1,2,3</sup>B. Tech in Computer Science and Engineering (CSE), Department of Computer Science & Engineering , Bhilai Institute Of Technology Raipur, India.

<sup>4</sup>Guide, Department of Computer Science & Engineering, Bhilai Institute Of Technology Raipur, India.

### ABSTRACT

This project focuses on the prediction of diseases using machine learning, with the goal of improving early detection and enhancing diagnostic accuracy in healthcare systems. By leveraging advanced machine learning algorithms such as random forests, decision trees, and support vector machines (SVM), the system effectively analyzes patient data to identify patterns and assess the risk of various diseases. These algorithms are capable of recognizing subtle correlations within complex medical data, enabling early identification of potential health issues before they progress to more severe stages. In order to optimize the performance of these models, techniques like feature selection and data preprocessing are implemented to improve the quality of input data and ensure more accurate predictions. Additionally, addressing common challenges such as data imbalances, where certain conditions may be underrepresented, is crucial to avoid biased outcomes. The models are trained on extensive healthcare datasets, incorporating a diverse range of variables, such as demographic information, medical histories, lifestyle factors, and test results. This approach demonstrates significant potential for automating the diagnostic process, allowing for quicker and more reliable identification of diseases. Moreover, it can assist in personalized treatment planning, offering tailored solutions based on individual patient profiles. The integration of machine learning in healthcare systems can lead to improved decision-making, enhanced patient outcomes, and more efficient public health management, marking a crucial step toward the future of precision medicine. Keywords : Machine Learning , Random Forest, Support Vector Machine, Artificial Intelligence, Naive Bayes.

### 1. INTRODUCTION

### 1.1) Overview of the Project

The healthcare sector faces critical challenges in early disease detection, which is vital for timely intervention and improved health outcomes. Traditional methods of disease detection often rely on manual analysis by medical professionals, which can be time-consuming, resource-intensive, and prone to human error. With advancements in machine learning, there is an opportunity to enhance disease detection accuracy and efficiency by automating the process. This research focuses on developing a machine learning-based system that utilizes advanced algorithms such as Random Forests, Decision Trees, and Support Vector Machines (SVM) to predict and classify diseases based on patient data. By integrating techniques like data preprocessing and feature selection, the system ensures accurate predictions while addressing challenges such as data imbalance. The ultimate goal is to create a reliable framework for early diagnosis, enabling better health outcomes and reducing healthcare costs.

### 1.2) Objectives

The primary objectives of this research are:

**1.2.1**) **Develop a Predictive System:** Build a robust machine learning framework for disease prediction and classification.

**1.2.2**) Utilize Advanced Algorithms: Employ state-of-the-art algorithms to analyze complex medical datasets for accurate disease detection.

**1.2.3**) **Optimize Data Processing:** Implement data preprocessing techniques to enhance input quality and model performance.

**1.2.4) Improve Early Detection:** Provide tools that enable early diagnosis, aiding in timely treatment and improved health outcomes.

**1.2.5**) **Support Public Health**: Contribute to public health efforts by improving diagnostic efficiency and healthcare resource utilization.

# 2. LITERATURE SURVEY

### 2.1) Hybrid Disease Prediction Model Using RF, LSTM, and SVM

K. Gaurav, A. Kumar, P. Singh, A. Kumari, M. Kasar\*, T. Suryawanshi (2023) developed a hybrid disease prediction model combining Random Forest (RF), Long Short-Term Memory (LSTM), and Support Vector Machines (SVM).



This modular approach enhances prediction accuracy and reliability while offering scalability to incorporate additional diseases and data types. Their work highlights the potential of integrating diverse machine learning algorithms to create robust, customizable systems for disease detection.

### 2.2) Disease Prediction Using Random Forest Algorithm

Y. Deepika Reddy, V. Sandhya Rani, S.K.Sathyanarayana (2023) proposed a disease prediction system utilizing the Random Forest algorithm, treating disease prediction as a classification problem. By applying feature selection to patient symptom data, the model effectively predicts potential diseases, demonstrating the efficacy of Random Forest in handling classification tasks for medical diagnosis.

### 2.3) Classification Algorithms for Diabetes Diagnosis

Samin Poudel (2022) applied classification algorithms from the scikit-learn and AutoGluon libraries to diagnose diabetes. Implemented on AWS SageMaker, the study leveraged models such as Naïve Bayes, SVM, K-Nearest Neighbors (KNN), LightGBM, and XGBoost to improve prediction accuracy. This research highlights the effectiveness of combining multiple classification algorithms for medical diagnostics.

### 2.4) A Web-Based Disease Prediction System

Anjali Bhatt, Shruti Singasane, Neha Chaube (2022) developed a web application for disease prediction using machine learning. The system uses Gradient Boosting and Random Forest classifiers to identify heart, liver, and diabetes diseases based on user inputs. By offering features like disease trends by age group and tailored recommendations, this application demonstrates how machine learning can provide accessible and actionable insights for end-users.

### 2.5) Machine Learning Techniques for Disease Detection

Nareen O. M. Salim & Adnan Mohsin Abdulazeez (2021) reviewed various machine learning techniques used for human disease detection, employing classifiers like SVM, k-NN, and CNN. Their analysis highlights the strengths and limitations of supervised, unsupervised, and reinforcement learning approaches, showcasing the flexibility of machine learning in analyzing patient data for accurate disease predictions.

### 2.6) Predictive Modeling for Disease Detection

Palle Pramod Reddy, Dirisinala Madhu Babu, Hardeep Kumar and Dr.Shivi Sharma (2021) proposed a disease prediction system based on predictive modeling. Using Random Forest as the primary classifier, their approach processes symptom data to return accurate disease likelihoods, emphasizing the importance of robust algorithms for reliable disease predictions.

### 2.7) Symptom-Based Disease Probability Estimation

Dr C K Gomathy, Mr. A. Rohith Naidu (2021) developed a system using Naïve Bayes, linear regression, and decision tree models to estimate disease probabilities. Their work focuses on diseases such as Diabetes, Malaria, Jaundice, Dengue, and Tuberculosis, leveraging biomedical data to enable early detection and improved patient care. This study demonstrates how combining multiple models can enhance disease prediction capabilities.

### 3. METHODOLOGY

### 3.1 Technologies Used for Model Development

**3.1.1 Visual Studio Code:** A versatile and user-friendly code editor by Microsoft, supporting multiple languages and extensions for efficient code writing, debugging, and management.

**3.1.2 Python:** A high-level programming language known for its simplicity, readability, and extensive libraries, widely used for data analysis, machine learning, and application development.

**3.1.3 Scikit-learn (sklearn):** A Python machine learning library offering tools for data preprocessing, model building, evaluation, and algorithms like Random Forest, SVM, and clustering methods.

**3.1.4 NumPy:** A library for numerical computing, enabling efficient operations on multidimensional arrays and matrices, with functionalities like linear algebra and random number generation.

**3.1.5 Pandas:** A data manipulation library providing structures like DataFrames for handling, cleaning, and transforming structured data efficiently.

**3.1.6 Tkinter:** A Python library for developing GUI applications, offering customizable widgets and cross-platform compatibility for interactive desktop applications.

**3.1.7 Statistics:** A built-in Python module for basic statistical calculations, including measures of central tendency and data dispersion.

### 3.2 Machine Learning Models Used

3.2.1 Random Forest: An ensemble learning method combining decision trees to improve accuracy and generalization, suitable for both classification and regression tasks.

LIPREMS	INTERNATIONAL JOURNAL OF PROGRESSIVE RESEARCH IN ENGINEERING MANAGEMENT	e-ISSN : 2583-1062
	AND SCIENCE (IJPREMS)	Impact
www.ijprems.com	(Int Peer Reviewed Journal)	Factor :
editor@ijprems.com	Vol. 04, Issue 12, December 2024, pp : 692-697	7.001

3.2.2 Naïve Bayes: A probabilistic model based on Bayes' theorem, effective for text classification tasks like spam detection and sentiment analysis.

3.2.3 Support Vector Machine (SVM): A supervised learning algorithm that maximizes class separation margins, capable of handling linear and non-linear data with kernel functions.

#### 3.3) Code Snippet

**3.3.1) Importing Libraries** 

import numpy as np import pandas as pd from sklearn.preprocessing import LabelEncoder from sklearn.model\_selection import train\_test\_split from sklearn.svm import SVC from sklearn.naive\_bayes import GaussianNB from sklearn.ensemble import RandomForestClassifier import tkinter as tk import statistics

### **3.3.1.1)** Preparing the Dataset for Training:

# Reading the train.csv by removing the last column since it's an empty column
DATA\_PATH = "Training.csv"
data = pd.read\_csv(DATA\_PATH).dropna(axis=1)
# Encoding the target value into numerical values using LabelEncoder
encoder = LabelEncoder()
data["prognosis"] = encoder.fit\_transform(data["prognosis"])
X = data.iloc[:, :-1] # Features
y = data.iloc[:, -1] # Target

**3.3.1.2)** Splitting the Dataset and Train the Model:

# Train-test split X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, y, test\_size=0.2, random\_state=24) # Initialize and train models final\_svm\_model = SVC() final\_nb\_model = GaussianNB() final\_rf\_model = RandomForestClassifier(random\_state=18) # Train models on the entire dataset final\_svm\_model.fit(X, y) final\_nb\_model.fit(X, y) final\_rf\_model.fit(X, y)

#### 3.1.2.3) Prepare for Prediction

```
# Prepare symptom index dictionary
symptoms = X.columns.values
symptom_index = {symptom: index for index, symptom in enumerate(symptoms)}
# Prepare prediction classes
data_dict = {
    "symptom_index": symptom_index,
    "predictions_classes": encoder.classes_
}
```



www.ijprems.com

# INTERNATIONAL JOURNAL OF PROGRESSIVE **RESEARCH IN ENGINEERING MANAGEMENT AND SCIENCE (IJPREMS)** (Int Peer Reviewed Journal)

e-ISSN: 2583-1062 Impact **Factor:** 7.001

Vol. 04, Issue 12, December 2024, pp : 692-697



#### 3.1.2.4) Submit Button Functionalities

```
# GUI Part
def on_submit():
    symptoms_input = input_box.get() # Get the input text (symptoms)
    if symptoms_input.strip() == "":
         output_box.delete(1.0, tk.END)
         output_box.insert(tk.END, "Please enter symptoms.")
         return
    predictions = predictDisease(symptoms input)
    print("Predictions:", predictions) # Debug: Print predictions to the console
    output_box.config(state=tk.NORMAL) # Enable text box for editing
    output_box.delete(1.0, tk.END) # Clear the output box
output_box.insert(tk.END, f"Random Forest Prediction: {predictions['rf_model_prediction']}\n")
    output_box.delete(1.0, tk.END) # Clear the output box
    output_box.insert(tk.END, f"Naive Bayes Prediction: {predictions['naive_bayes_prediction']}\n'
    # output_box.insert(tk.END, f"SVM Prediction: {predictions['svm_model_prediction']}\n")
output_box.insert(tk.END, f"Final Prediction: {predictions['final_prediction']}\n")
    output_box.config(state=tk.DISABLED) # Revert text box to readonly
    root.update() # Force an immediate update to the window
```

### 4. RESULTS

The machine learning model developed for symptom-based disease prediction demonstrated reliable accuracy and practical utility in supporting healthcare diagnostics. The model utilized algorithms such as Support Vector Machine (SVM), Random Forest, and Neural Networks, evaluated through metrics like precision, recall, and computational efficiency.

LIPREMS	INTERNATIONAL JOURNAL OF PROGRESSIVE	e-ISSN :
	<b>RESEARCH IN ENGINEERING MANAGEMENT</b>	2583-1062
	AND SCIENCE (IJPREMS)	Impact
www.ijprems.com	(Int Peer Reviewed Journal)	Factor :
editor@ijprems.com	Vol. 04, Issue 12, December 2024, pp : 692-697	7.001

**4.1)** Accuracy: The model achieved high accuracy in predicting diseases based on symptom data, showcasing its ability to distinguish between overlapping or complex symptom patterns.

**4.2) Precision and Recall:** The model demonstrated high precision and recall, effectively minimizing false positives and false negatives, ensuring dependable diagnostic suggestions.

**4.3**) User Interface: The system provided a user-friendly interface for inputting symptoms, offering ranked disease predictions with associated confidence levels to aid medical professionals and patients.

**4.4) Real-World Testing:** The model was tested using diverse symptom datasets, including data from varied demographic and clinical contexts. It maintained consistent performance, highlighting its robustness and scalability.

The successful development and testing of this system underline its potential to complement medical diagnostics, especially in telemedicine and resource-constrained settings. While not a replacement for professional judgment, it serves as a valuable tool to enhance healthcare accessibility and decision-making.

### 5. CONCLUSION AND FUTURE WORK

#### 5.1) Conclusion:

This study underscores the transformative potential of machine learning (ML) and artificial intelligence (AI) in disease prediction. The integration of advanced algorithms with real-time data enables early, accurate diagnoses, improving treatment outcomes, reducing costs, and enhancing patient quality of life. By uncovering complex patterns in diverse datasets, these models excel at identifying diseases at their nascent stages. This capability not only empowers healthcare providers with a proactive approach but also facilitates the adoption of personalized medicine, which tailors interventions to individual genetic, lifestyle, and environmental factors.

While ML and AI have shown immense promise, these technologies are not replacements for professional medical judgment but valuable tools that complement clinical expertise. The continued evolution of these systems will revolutionize disease prediction and management, contributing significantly to global public health.

#### 5.2) Future Work:

To further enhance the capabilities of AI-driven disease prediction systems, several areas of improvement and expansion are proposed:

**5.2.1) Integration of Genomics and Personalized Medicine:** Future models will incorporate genomic data to enable precise predictions of disease susceptibility. By combining genetic profiles with environmental and lifestyle factors, healthcare providers can deliver highly targeted and effective treatment plans.

**5.2.2) Real-Time Health Monitoring with Wearable Devices:** Wearables like smartwatches and fitness trackers can provide real-time health data, such as heart rate and oxygen levels. These devices, integrated with AI, will allow for continuous monitoring, early warning alerts, and proactive disease prevention, particularly for chronic conditions and acute medical events.

**5.2.3**) **AI-Powered Imaging and Diagnostics:** Advancements in AI algorithms will enable the detection of subtle abnormalities in medical imaging, such as X-rays and MRIs, that might elude human observation. This will facilitate early diagnosis and improve treatment outcomes for conditions like cancer and cardiovascular diseases.

**5.2.4) Predicting Disease Outbreaks:** AI systems can analyze diverse datasets, including travel patterns and environmental factors, to predict and mitigate outbreaks of infectious diseases. This capability will help public health agencies allocate resources effectively and implement preventive measures.

**5.2.5)** Holistic Health Assessments: Future systems will provide comprehensive health evaluations by considering genetic, lifestyle, environmental, and social determinants. These insights will enable healthcare providers to address risk factors preemptively and guide patients toward healthier choices.

**5.2.6 Automation and Efficiency in Healthcare:** AI-powered automation will streamline diagnostic processes, reduce manual workloads, and improve the efficiency of healthcare delivery. By reducing administrative burdens, medical professionals can focus more on patient care, leading to better outcomes.

**5.2.7 Addressing Ethical and Data Challenges:** Ensuring data security, fairness, and representation in AI systems is critical. Efforts to eliminate biases in training datasets and safeguard patient information will enhance trust and ensure equitable healthcare for all.

By exploring these future directions, disease prediction systems can evolve into indispensable tools for global healthcare, fostering a proactive, personalized, and efficient approach to medical diagnostics and disease management.



editor@ijprems.com

INTERNATIONAL JOURNAL OF PROGRESSIVE	e-ISSN :
RESEARCH IN ENGINEERING MANAGEMENT	2583-1062
AND SCIENCE (IJPREMS)	Impact
(Int Peer Reviewed Journal)	Factor :
Vol. 04, Issue 12, December 2024, pp : 692-697	7.001

#### 6. **REFERENCES**

- [1] K. Gaurav, A. Kumar, P. Singh, A. Kumari, M. Kasar\*, T. Suryawanshi (2023). Human disease prediction using machine learning techniques and real-life parameters. International Journal of Engineering (IJE).
- [2] Y. Deepika Reddy, V. Sandhya Rani, S.K.Sathyanarayana (2023). Human disease detection using ML. International Journal of Creative Research Thoughts (IJCRT).
- [3] Samin Poudel (2022). A study of disease diagnosis using machine learning. MDPI
- [4] Anjali Bhatt, Shruti Singasane, Neha Chaube (2022). Disease prediction using machine learning. International Research Journal of Engineering and Technology (IRJETS).
- [5] Nareen O. M. Salim & Adnan Mohsin Abdulazeez (2021). Human diseases detection based on machine learning algorithms: A review. International Journal of Scientific and Academic Research (IJSAB).
- [6] Palle Pramod Reddy, Dirisinala Madhu Babu, Hardeep Kumar and Dr.Shivi Sharma (2021). Disease prediction using machine learning. International Journal of Creative Research Thoughts (IJCRT).
- [7] Dr C K Gomathy, Mr. A. Rohith Naidu (2021). The prediction of disease using machine learning. International Journal of Science and Research Methodology (IJSREM)