# Disease Prediction using Machine Learning

[1]**Palle Pramod Reddy,** [2]**Dirisinala Madhu Babu,** [3]**Hardeep Kumar and**

[4]**Dr.Shivi Sharma**

[1,2,3]*Students, School of Computer Science and Engineering, Lovely Professional University, Jalandhar (Punjab), India*

[4] *Assistant Professor, School of Electronics and Electrical Engineering, Lovely Professional university, Jalandhar (Punjab) India*

*Abstract—*

The "Disease Prediction" method, which is concentrated on predictive modeling, it predicts the user's disease based on the symptoms that the user provides as input. The method examines the user's symptoms as input and returns the disease's likelihood as an output. Disease prediction is accomplished using the random forest classifier.

*Keywords—Random Forest, Chronic Disease,*

I.    INTRODUCTION

When anyone is currently afflicted with an illness, they must see a doctor, which is both time consuming and costly. It can also be difficult for the user if they are not near to doctors and hospitals because the illness cannot be identified. So, if the above procedure can be done using an automated software that saves time and money, it could be better for the patient, making the process go more smoothly. There are other Heart Disease Prediction Systems that use data mining methods to analyze the patient's risk level. Disease Predictor is a web-based system that predicts a user's disease based on the symptoms they have. Data sets from various health-related websites have been obtained for the Disease Prediction system. The consumer will be able to determine the likelihood of a disease based on the symptoms given using Disease Predictor. People are always curious to learn new things, particularly as the use of the internet grows every day. When an issue occurs, people often want to look it up on the internet. Hospitals and physicians have less access to the internet than the general public. When people are afflicted with an illness, they do not have many options. As a result, this system can be beneficial to people. Chronic illness is a disease that lasts a long time or takes a long time to heal, and many chronic diseases cannot be cured but can only be managed with daily treatments. India, like all other nations, is undergoing significant social and economic shifts, which is causing a fast rise in the frequency of cardiovascular disease.

Many advanced countries, including India, are dealing with a wide range of chronic diseases, mostly cardiovascular disease and diabetes, which could have deep consequences for global health, security, and economy. The rapid urbanization and economic growth of today's world has resulted in a wide range of lifestyles. Chronic diseases are now a problem in all nations, with chronic disease afflicting one-third of the population in each. Chronic disease care is more expensive, and it is difficult for those who are sick. In the medical field, a huge number of chronic disease datasets are gathered and processed, and data mining aids in disease early detection. Cardiovascular disease, diabetes, liver disease, Alzheimer's disease, and Parkinson's disease are the most high-priced diagnosis diseases.

It's a major challenge in the medical or healthcare industries to offer the highest quality services to all patients, and only those who can afford it can benefit from it. There is a vast amount of healthcare data available that is not being mined in a more efficient and reliable manner to uncover secret knowledge for successful decision-making. The proposed framework employs data mining techniques to detect Chronic diseases early.

Machine learning is the process of programming computers to improve their output based on examples or previous data. The study of computer systems that learn from data and experience is known as machine learning. Training and Testing are the two stages of the machine learning algorithm. Prediction of a disease based on the signs and medical history of the patient Machine learning has been a stumbling block for decades.

Machine Learning technology provides a strong forum in the medical sector for efficiently resolving healthcare issues.

## II. RESEARCH OBJECTIVE

There is a require to groundwork and evolve a system that will enable end users to predict chronic diseases without having to visit a physician or doctor for diagnosis. To identify various diseases by observing the symptoms of patients and applying various Machine Learning Models techniques. There is no proper procedure for handling text and structured data. Both structured and unstructured data would be considered by the proposed framework.

Machine Learning can improve the accuracy of predictions.

## III. LITERATURE REVIEW

The study for the best medical diagnosis mining technique was performed by K.M. Al-Aidaroos, A.A. Bakar, and Z. Othman. For this study, the authors compared Nave Baeyes to five other classifiers: LR, KStar (K*), Decision Tree (DT), Neural Network (NN), and a basic rule-based algorithm (ZeroR). The efficiency of all algorithms was evaluated using 15 real-world medical problems from the UCI machine learning repository (Asuncion and Newman, 2007). In the experiment, NB outperformed the other algorithms in 8 of the 15 data sets, leading to the conclusion that the predictive accuracy results in Nave Baeyes are superior to other techniques. Darcy A. Davis, Nitesh V. Chawla, Nicholas Blumm, Nicholas Christakis, and Albert-Laszlo Barabasi discovered that treating chronic illness at a global level is neither time nor cost effective. As a result, the authors performed this study in order to forecast potential disease risk. CARE (which uses only a patient's medical history and ICD-9-CM codes to predict possible disease risks) was used for this. Based on their own medical history and that of similar patients, CARE incorporates collective filtering approaches with clustering to predict each patient's greatest disease risks. ICARE, an iterative version that integrates ensemble principles for improved efficiency, has also been defined by the authors.

These cutting-edge systems don't need any advanced knowledge and can predict a wide range of medical conditions in a single run. ICARE's remarkable potential risk coverage means more precise early alerts for thousands of illnesses, several years ahead of time. When used to its full extent, the CARE system can be used to investigate a wider range of disease backgrounds, raise previously unconsidered questions, and facilitate discussions regarding early detection and prevention.

This research paper was written by JyotiSoni, Ujma Ansari, Dipesh Sharma, and SunitaSoni to provide a survey of existing techniques of information discovery in databases using data mining techniques that are used in today's medical research, specifically in Heart Disease Prediction. A number of experiments have been carried out to compare the performance of predictive data mining

techniques on the same dataset, and the results show that Decision Tree outperforms, with Bayesian classification having comparable accuracy to Decision Tree in some cases, but other predictive approaches such as KNN, Neural Networks, and Classification based on Clustering underperform.

Shadab Adam Pattekari and Asma Parveen conducted a study to predict heart diseases using the Decision Tree Algorithm, in which the consumer provides data that is compared to a qualified set of values. As a result of this study, patients were able to provide basic information that was compared to data, and heart disease was expected. M.A.NisharaBanu and B. Gomathy analysed the various types of heart-related problems using medical data mining techniques such as association rule mining, grouping, and clustering I. The aim of a decision tree is to show any possible outcome of a decision. To achieve the best result, various rules are devised. The criteria used in this study were age, sex, smoking, being overweight, drinking alcohol, blood sugar, heart rate, and blood pressure. The risk level for various parameters is saved with their ids ranging from 1 to 100. (1-8). The standard level of prediction is represented by IDs less than 1, whereas higher IDs other than 1 represent higher risk levels. The pattern in the dataset is studied using the K- means clustering method. The algorithm divides the data into k groups. The closed cluster is allocated to each point in the dataset. Each cluster centre is recalculated as the average of the cluster's points.

## PROPOSED SYSTEM

We have mixed structured and unstructured data in the healthcare fields to determine disease risk in this project. The use of a latent factor model to recreate missing data in medical records obtained from online sources. We could also assess the major chronic diseases in a specific area and population using statistical information. We consult hospital experts to learn about useful features when dealing with structured data. In the case of unstructured text files, we use the randrom forest algorithm to automatically select features.
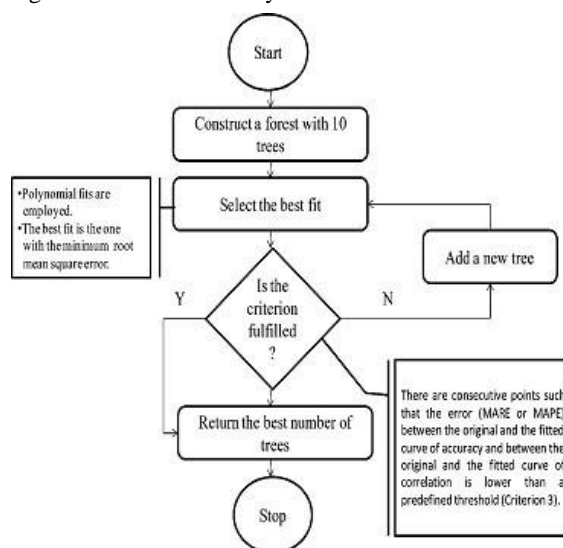


**Fig 1:- System Model**

## A. Data collection

Data collection has been done from the internet to identify the disease here the real symptoms of the disease are collected i.e. no dummy values are entered. The symptoms of the disease are collected from different health related websites.

### Data Preprocessing

Before feeding the data into the Prediction model, followingdata cleaning and preprocessing steps are performed

● Checking null values and filling using forward fill method

● Converting data into different cases

● Standardizing the data using mean  and standard deviation

● Splitting the dataset into training and testing sets

## B. Building Model

Many methods are used to perform data mining. Machine learning is one of the approaches. Random forest Machine learning strategies include grouping, clustering, summarization, and many others. Since classification techniques are used in this project, classification is one of the data mining processes in this

phase of categorical data classification. And this step is divided into two phases: training and testing. In the training phase, predetermined data and associated class labels are used for classification. The training stage is often referred to as supervised learning. The preparation and testing phases of the classification process are depicted in the diagram. In the training process, training tuples are used, and in the test data phase, test data tuples are used, and the classification rule's accuracy is calculated. Assume that the classification rule's accuracy on testing data is sufficient for the rule to be used for classification of unmined data.

## C.Prediction:

Prediction using Random Forest : -

Prediction done by Random Forest Model using Flask frame work model trained by training chronic disease dataset

## IV.    RESULTS AND CONCLUSION

| Model | Accuracy |
|---|---|
| Diabetes Model | 98.25 |
| Breast Cancer Model | 98.25 |
| Heart Disease Model | 85.25 |
| Kidney Disease Model | 99 |
| Liver Disease Model | 78 |

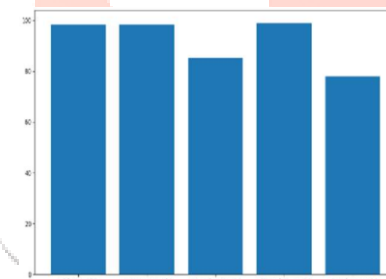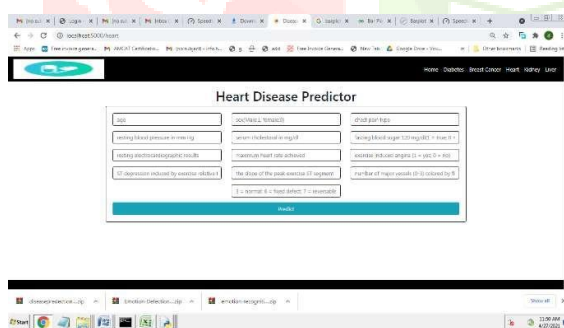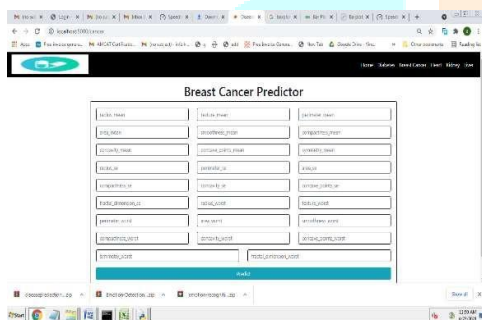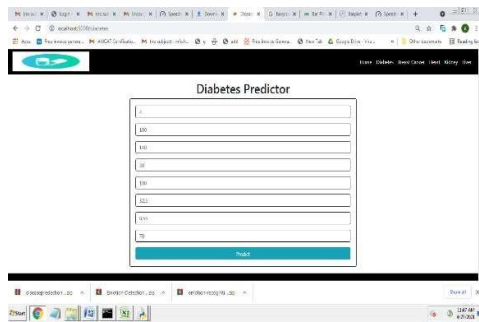**Table 1:- shows the accuracy achieved using random forest for each disease**



**Fig. 2 shows the accuracy or each model using Random forest classifier**

**Fig. 7:- Kidney Disease Prediction entry form**





Fig. 3 Home scr

**Fig. 4 :- Diabetes Prediction entry form**



**Fig. 5 :- Breast cancer Prediction entry form**



**Fig. 6 :- Heart Disease Prediction entry form**

**Fig. 7:- Liver Disease Prediction entry form**

## Conclusion

The aim of this project is to predict disease based on symptoms. The project is set up in such a way that the device takes the user's symptoms as input and generates an output, which is disease prediction. A prediction accuracy probability of 95% is obtained on average. The grails system was used to successfully incorporate Disease Predictor.

## REFERENCE

1) A.Davis, D., V.Chawla, N., Blumm, N., Christakis, N., & Barbasi, A. L. (2008). Predicting Individual Disease Risk Based On Medical History.

2) Adam, S., & Parveen, A. (2012). Prediction System For Heart Disease Using Naive Bayes.

3) Al-Aidaroos, K., Bakar, A., & Othman, Z. (2012). Medical Data Classification With Naive Bayes Approach. Information Technology Journal.

4) Darcy A. Davis, N. V.-L. (2008). Predicting Individual Disease Risk Based On Medical History.

5) JyotiSoni, Ansari, U., Sharma, D., & Soni, S. (2011). Predictive Data Mining for Medical Diagnosis: An Overview Of Heart Disease Prediction.

6) K.M. Al-Aidaroos, A. B. (n.d.).

K.M. Al-Aidaroos, A. B. (n.d.). 2012. *Medical Data sssClassification With Naive Bayes Approach* .

7) Nisha Banu, MA; Gomathy, B;. (2013). Disease Predicting System Using Data Mining Techniques.