



**INTERNATIONAL JOURNAL OF MULTIDISCIPLINARY
ADVANCED SCIENTIFIC RESEARCH AND INNOVATION
(IJMASRI)**

ISSN: 2582-9130

IBI IMPACT FACTOR 1.5

DOI: 10.53633/IJMASRI

RESEARCH ARTICLE

DIABETES AND HEART DISEASE PREDICTION

Harshit Pandey¹ and Sachin Garg²

^{1,2} *Department of Information Technology Maharaja Agrasen Institute of Technology, Rohini, Delhi-110086,*

Abstract

Machine Learning (ML), one of the most well-known applications of Artificial Intelligence, is revolutionising the field of research. In this work, machine learning is utilised to determine whether or not a person has cardiac disease. Many people suffer from cardiovascular diseases (CVDs), which claim the lives of people all over the world. On a monthly basis, a large amount of patient-related data is maintained. The information gathered can be used to predict the occurrence of future diseases. Machine learning may be used to determine whether a person has a cardiovascular illness based on particular characteristics such as chest discomfort, cholesterol levels, gender, and other factors. Logistic Regression (86.3% accuracy), Naïve Bayes (86.3% accuracy), Random Forest (93.6% accuracy), and K-Nearest Neighbour (87.8% accuracy) are some of the machine learning algorithms that we are using to predict cardiac disease.

Keywords: Diabetes Disease, Heart Disease, Machine Learning

Introduction

Human body is made up of various organs, all of which have their own functions. Heart is one such organ which pumps blood throughout the body and if it does not do so, the human body can have fatal circumstances. One of the main reasons for mortality today is having a heart disease (Mohan 2019). So, it becomes necessary to make sure that

our cardiovascular system or any other system in the human body for that matter must remain healthy. Unfortunately, people all around the world have been facing cardiovascular diseases. Any technology that can help diagnose these diseases before much damage is done will prove as helpful in saving people's money and more importantly their lives. Data mining techniques can be useful in predicting heart diseases. Predictive models can be made by finding previously unknown patterns and trends in databases and using the obtained information (Bhatla

646

and Jyoti 2012) Data mining means to extract knowledge from large amounts of data (Patel 2015). Machine learning is a technology which can help to achieve diagnosis of heart disease before much damage happens to a person. As an emerging field in science and technology, machine learning can classify whether a person might be suffering from a heart disease or not. The techniques and algorithms can be directly used on a dataset for creating some models or to draw vital conclusions, and inferences from the dataset. Common attributes used for heart disease are Age, Sex, Fasting Blood Pressure, Chest Pain type, Resting ECG (test that measures the electrical activity of the heart), Number of major vessels colored by fluoroscopy, Threst Blood Pressure (high blood pressure), Serum Cholestrol (determine the risk for developing heart disease), Thalach (maximum heart rate achieved), ST depression (finding on an electrocardiogram, trace in the ST segment is abnormally low below the baseline), painloc (chest pain location (substernal=1, otherwise=0)), Fasting blood sugar, Exang (exercise included angina), smoke, Hypertension, Food habits, weight, height and obesity (Mythili et al., 2013).

Litreature Review

Using the UCI Machine Learning dataset, a lot of research has been done. Using different machine algorithms we got different levels of accuracy which are explained as follows.

Avinash Golande and others studied various different ML algorithms that can be used for classification of heart disease. Research was carried out to study Decision Tree, KNN and K-Means algorithms that can be used for classification and their accuracy were compared (Avinash Golande and Pavan Kumar 2015). This research concludes that accuracy obtained by Decision Tree was highest, further it was inferred that it can be made efficient by combination of different techniques and parameter tuning.

Fahd Saleh Alotaibi has designed a ML model comparing five different algorithms (Fahd

Saleh Alotaibi 2019). Rapid Miner tool was used which resulted in higher accuracy compared to Matlab and Weka tool. In this research the accuracy of Decision Tree, Logistic Regression, Random forest, Naive Bayes and SVM classification algorithms were compared. Decision tree algorithm had the highest accuracy.

(Anjan Nikhil Repaka, et al., 2019) proposed a system in that uses NB (Naive Bayesian) techniques for classification of dataset and AES (Advanced Encryption Standard) algorithm for secure data transfer for prediction of disease.

(Theresa Princy *et al.*, 2016), executed a survey including different classification algorithms used for predicting heart disease. The classification techniques used were Naive Bayes, KNN (KNearest Neighbour), Decision tree, Neural network and accuracy of the classifiers was analyzed for different numbers of attributes

Nagaraj M Lutimath, et al., has performed the heart disease prediction using Naive bayes classification and SVM

(Support Vector Machine). The performance measures used in analysis are Mean Absolute Error, Sum of Squared Error and Root Mean Squared Error, it is established that SVM emerged as superior algorithm in terms of accuracy over Naive Bayes (Nagaraj 2019).

Proposed Models

Logistic Regression

Logistic Regression is a classification algorithm mostly used for binary classification problems. In logistic regression instead of fitting a straight line or hyper plane, the logistic regression algorithm uses the logistic function to squeeze the output of a linear equation between 0 and 1. There are 13 independent variables which makes logistic regression good for classification.

Naive Bayes

Naïve Bayes algorithm is based on the Bayes rule. The independence between the attributes of the dataset is the main assumption and the most important in making a classification. It is easy and fast to predict and holds best when the assumption of independence holds. Bayes' theorem calculates the posterior probability of an event (A) given some prior probability of event B represented by $P(A|B)$ [10] as shown in the following equation

$$P(A|B) = (P(B|A)P(A)) / P(B)$$

Random Forest

Random Forest algorithms are used for classification as well as regression. It creates a tree for the data and makes prediction based on that. Random Forest algorithm can be used on large datasets and can produce the same result even when large sets record values are missing. The generated samples from the decision tree can be saved so that it can be used on other data. In random forest there are two stages, firstly create a random forest then make a prediction using a random forest classifier created in the first stage.

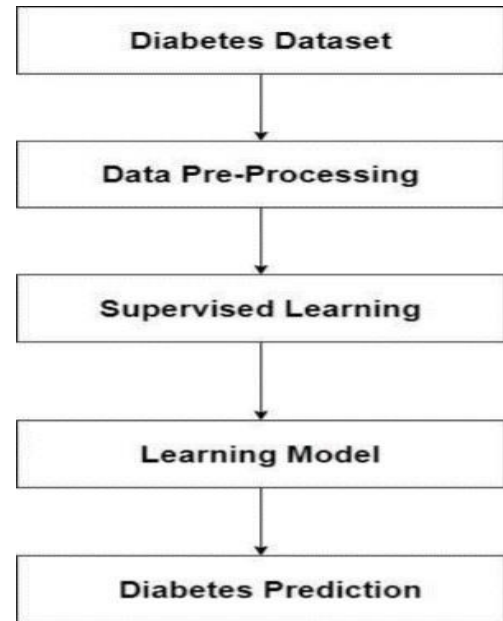
K-Nearest Neighbour

In K-NN algorithm a data point is taken whose classification is not available, then the number of neighbors, k is defined. After that k neighbors are selected according to the lowest Euclidian distance between the selected data points and their neighbors. The selected data point is then classified into a category, which is the same as the category which has the majority of neighbors among the K neighbors.

Dataset Collection And Processing

The dataset used in this experiment is the heart disease dataset which is the combination of 4 databases. Cleveland, Hungary, Switzerland, and Long Beach V. It contains 76 attributes, including the predicted attribute, but all published experiments refer to using a subset of 14 of them. The "target" field refers to the presence of heart disease in the patient. It is integer valued 0 = no disease and 1 = disease. The dataset is already

preprocessed and is available on the kaggle website.



Results:

The result have been collected of the used algorithms. The accuracy is different for the algorithm used.

	Model	Accuracy
0	Logistic Regression	86.341463
1	Naive Bayes	85.365854
2	Random Forest	93.658537
3	K-Nearest Neighbour	87.804878

Conclusion

Globally, heart disease is the leading cause of death. With the rising number of deaths due to heart disease, it is becoming increasingly important to build a system that can effectively and accurately forecast heart disease. The goal of the research was to discover the most effective machine learning system for detecting cardiac problems. Using the dataset present in kaggle website, this research examines the accuracy of K nearest neighbour, Logistic Regression, Random Forest, and Naive Baye's algorithms for predicting heart disease. Random forest model gives the highest accuracy among all other models which are used in this project. In the future, the study might be improved by creating a web application based on the Random Forest method and employing a larger dataset than the one used in this analysis, which would help to deliver better results and aid health professionals in successfully and efficiently forecasting cardiac disease.

References

1. .Mohan, S., Thirumalai, C and Srivastava, G.(2019). Effective heart disease prediction using hybrid machine learning techniques. *IEEE Access*, 7, 81542-81554.
2. Bhatla, N and Jyoti, K. (2012). An analysis of heart disease prediction using different data mining techniques. *International Journal of Engineering*, 1(8), 1-4.
3. Patel, J., Tejal Upadhyay, D and Patel, S. (2015). Heart disease prediction using machine learning and data mining technique. *Heart Disease*, 7(1), 129-137.
4. Mythili, T., Dev Mukherji, Nikita Padaila and Abhiram Naidu,(2013). "A Heart Disease Prediction Model using SVM- Decision Trees- Logistic Regression (SDL)", *International Journal of Computer Applications*, vol. 68, 16 April 2013.
5. Avinash Golande and Pavan Kumar, T(2019). "Heart Disease Prediction Using Effective Machine Learning Techniques", *International Journal of Recent Technology and Engineering*, Vol 8, pp.944-950,2019.
6. Fahd Saleh Alotaibi,(2019)"Implementation of Machine Learning Model to Predict Heart Failure Disease", (*IJACSA*) *International Journal of Advanced Computer Science and Applications*, Vol. 10, No. 6, 2019
7. Anjan Nikhil Repaka, Sai Deepak Ravikanti, Ramya G Franklin (2019)"Design And Implementation Heart Disease Prediction Using Naives Bayesian", *International Conference on Trends in Electronics and Information(ICOEI 2019)*.
8. Theresa Princy, R,J (2016). Thomas,'Human heart Disease Prediction System using Data Mining Techniques', *International Conference on Circuit Power and Computing Technologies*, Bangalore, 2016.
9. Nagaraj, M (2019). Lutimath, Chethan C, Basavaraj S Pol., 'Prediction Of Heart Disease using Machine Learning', *International journal Of Recent Technology and Engineering*, 8,(2S10), pp 474-477, 2
