

# Prediction of Heart Disease Using Machine Learning

## Abstract

Heart disease is the one of the most common disease. This disease is quite common now days we used different attributes which can relate to this heart diseases well to find the better method to predict and we also used algorithms for prediction. Naive Bayes, algorithm is analyzed on dataset based on risk factors. We also used decision trees and combination of algorithms for the prediction of heart disease based on the above attributes. The results shown that when the dataset is small naive Bayes algorithm gives the accurate results and when the dataset is large decision trees gives the accurate results.

## 1. Introduction

The main topic is prediction using machine learning technics. Machine learning is widely used now days in many business applications like e commerce and many more. Prediction is one of area where this machine learning used, our topic is about prediction of heart disease by processing patient's dataset and a data of patients to whom we need to predict the chance of occurrence of a heart disease.

## 2. Literature Survey

[1]. M.A.Nishara Banu and B.Gomathy has given a paper named Disease Predicting system using data mining techniques. In this paper they talk about MAFIA (Maximal Frequent Item set algo- rithm) and K-Means clustering. As classification is important for prediction of a disease. The classification based on MAFIA and K-Means results in accuracy.

[2]. Wiharto and Hari Kusnanto have given a paper named Intelli- gence System for Diagnosis Level of Coronary Heart Disease with K-Star Algorithm. In this paper they exhibit an expectation framework for heart infection utilizing Learning vector Quantiza- tion neural system calculation The neural system in this frame- work acknowledges 13 clinical includes as information and pre- dicts that there is a nearness or nonattendance of coronary illness in the patient, alongside various execution measures.

## 3. Methodology

### 3.1 Data Pre-Processing

Cleaning: Data that we want to process will not be clean that is it may contain noise or it may contain values missing of we process we can't get good results so to obtain good and perfect results we need to eliminate all this, the process to eliminate all this is data cleaning. We will fill missing values and can remove noise by using some techniques like filling with most common value in missing place.

Transformation: This involves changing data format to one form to other that is making them most understandable by doing normalization, smoothing, and generalization, aggregation techniques on data.

Integration: Data that we need not process may not be from a single source sometimes it can be from different sources we do not integrate them it may be a problem while processing so integration is one of important phase in data pre-processing and different issues are considered here to integrate.

Reduction: When we work on data it may be complex and it may be difficult to understand sometimes so to make them understand- able to system we will reduce them to required format so that we can achieve good results.

### 3.2 ID3 Algorithm

To do this we have many machine learning algorithms out of which the more widely used methods are Naïve Bayes classification technique and decision tree construction, in this decision tree construction we have many algorithms one of which we took for this ID3 algorithm. The ID3 algorithm is one of the old algorithms which is used for building decision trees. In the process of building a decision tree it handles missing values and removes outliers. So we can build this decision tree even if the data is not cleaned well. Decision trees construct classification or regression models as a structure which is similar to a tree. It separates a dataset into fewer and fewer sub-sets while in the meantime a related decision tree is incrementally created. The last outcome is a tree with a choice point and a leaf point. A choice node has a minimum of 2 branches. Leaf nodes speak to a grouping or choice. The highest choice node in a tree which compares to the best indicator is called the root point. Choice trees can deal with both all-out and numerical information.

**Step 1:** If all occurrences in  $X$  are certain, then make YES node and end. On the off chance that all cases in  $X$  are negative, make a NO node and end. Generally select an element,  $B$  with qualities  $U_1, U_n$  and make a choice node.

**Step 2:** Partition the preparation occurrences in  $X$  into subsets  $X_1, X_2, \dots, X_n$  as indicated by the estimations of  $U$ .

**Step 3:** apply the calculation recursively to each of the sets  $A_i$ .

### 3.3 Naive-Bayes Classification

The Naive-Bayesian classifier relies upon Bayes' speculation with autonomy suppositions among attributes. A Naive-Bayesian output is definitely not hard to run, with no entrapped repetitive parameter estimation which makes it particularly supportive for broad datasets in spite of its effortlessness, the Naive Bayesian classifier generally completes its job shockingly good and is broadly used in light of the fact that it frequently outflanks high order techniques which are complex. The Naive Bayes treats every variable as independent which helps it to predict even if variables don't have proper relation.

### 3.4 K-means

k-means clustering is one of the clustering techniques used to cluster datasets based on nearest-neighbor. Here the data is clustered in  $k$  clusters based on a similarity between them. We also fill missing values of data using this k-means. Once we clustered the data every dataset will come into any one of the clusters by using these clusters. If we have missing values in the dataset we can fill those values as these are categorized into groups. Now as these missing values are all cleared we can apply different prediction techniques on this. For an example we can apply now as we know that for a dataset to be used for prediction in Naive Bayes need to be pre-processed we can use this data for prediction in Naive Bayes. By different combination of using these algorithms we can achieve good accuracy.

We reviewed different papers on heart disease prediction out of all prediction techniques and methods what everyone uses when it comes to prediction is Naive Bayes and decision trees. We have different methods one of which that we used here is ID3 algorithm. We took a medical data of heart disease patients from UCI machine learning repository one of the popular repositories to get data for machine learning experiments. It contains a record of nearly 300 patients. We performed both this Naive Bayes and ID3 techniques on this training data using R tool. In R tool we used some 3<sup>rd</sup> party libraries like e1071 for implementing Naive Bayes and part to construct decision tree. In the data set that we took for implementing this contains variables

**Table 1: data attributes for prediction**

Attribute	Values and meaning
Age1	Age in year
Gender1	Value 1 and 0 for male and female
Cp1	Pain in chest Yes/No
Trestbps1	blood pressure during resting
Chol1	Cholesterol of serum in mg/dl
Fbs1	blood sugar during fasting
Restecg1	Resting electrocardiographic results
Oldpeak1	ST depression induced by exercise relative to rest

<b>Slope1</b>	The slope of peak exercise of ST segment. Value 1: up sloping Value 2: flat Value 3: down sloping
<b>Ca1</b>	Number of major vessels (0-3) colored by fluoroscopy
<b>Thal1</b>	3 = normal; 6 = fixed defect
<b>Num1</b>	Diagonal of heart disease Value 0: No Risk; Value 1: Low Risk; Value 2: Risk; Value 3: High Risk; Value 4: Higher Risk

Number which indicates rate of getting heart attack on a scale of 0 to 4. What we observed in the results are out of different executions on different data sets both algorithms predict with a good accuracy but when comparing both in most of cases decision tree is giving a result which is has less probability. We also observed that as Naïve Bayes treats all members of class as independent it can get all probabilities i.e., probabilities occurrence of different values of class members.

The result of decision tree:

```
> fltTree<-rpart( Num1 ~ Age1 + Sex1 + Cp1 + Trestbps1 + Chol1 +Fbs1 +Restecg1 + Thalach1 + Exang1 + Oldpeak1 + Slope1 + Ca1 + Thal1,df)
> predict(fltTree,df_new,type="vector") 1
0.1478261
```

The result of Naïve Bayes is:

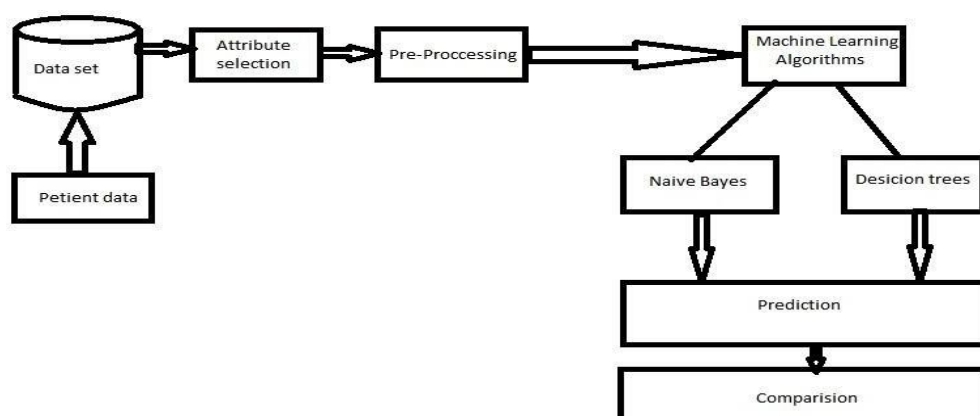
```
> model<-naiveBayes(Num1~.,data=df)
> predict(model,df_new,type="raw")
```

This above results of decision tree and Naïve Bayes both are applied on same dataset taking same training set and the actual answer is 0. We can observe that in Naïve Bayes is showing the actual results with more probability that means we can predict 0 as answer but decision tree shows it as one which is approximately correct but not exactly correct.

Decision tree:

The above decision tree is constructed in R on the dataset that we used in this. To achieve that in R we used C50 library.

## 4. Proposed System



By the above experiment what we say is as Naïve Bayes results and decision tree results may change so for every prediction we need not have a comparison of both the algorithms so get accurate results and in the same way if we use only a single algorithm which cannot pre-process data we even can't get good accuracy so it's better to have combination of algorithms like k-means, ID3 and k-means and Naive Bayes.

### 5. Conclusion

In this what we found is during small datasets in some other cases most of time decision trees direct us to a solution which is not accurate, but when we look at Naïve Bayes results we are getting more accurate results with probabilities of all other possibilities but due to guidance to only one solution decision trees may miss lead. Finally we can say by this experiment that Naïve Bayes is more accurate if the input data is cleaned and well maintained even though ID3 can clean it self it cannot give accurate results every time, and in this same way Naïve Bayes also will not give accurate results every time we need to consider results of different algorithms and by all its results if a prediction is made it will be accurate. But we can use Naïve Bayes consider variables as individual we can use combination of algorithms like Naïve Bayes and K-means to get accuracy.