

CARDIOVASCULAR RISK PREDICTION



PROJECT SUMMARY

Data Description

Demographic:

Sex: male or female("M" or "F")

Age: Age of the patient;(Continuous -

Although the recorded ages have been truncated to whole numbers, the concept of age is continuous)

Behavioral

is_smoking: whether or not the patient is a current smoker ("YES" or "NO")

Cigs Per Day: the number of cigarettes that the person smoked on average in one day.(can be

considered continuous as one can have any number of cigarettes, even half a cigarette.)

Medical (*history*)

BP Meds: whether or not the patient was on blood pressure medication (Nominal)

Prevalent Stroke: whether or not the patient had previously had a stroke (Nominal)

Prevalent Hyp: whether or not the patient was hypertensive (Nominal)

Diabetes: whether or not the patient had diabetes (Nominal)

Tot Chol: total cholesterol level (Continuous)

Sys BP: systolic blood pressure (Continuous)

Dia BP: diastolic blood pressure (Continuous)

BMI: Body Mass Index (Continuous)

Heart Rate: heart rate (Continuous -

In medical research, variables such as heart rate though in fact discrete, yet are considered continuous because of large number of possible values.)

PROBLEM STATEMENT

The Cardiovascular data consists of set of columns that speaks about the patient is suffering from heart disease or not

first we have to load the dataset and by understanding the dataset with the Meaningfull features and rows we perform the model

There are **3390** rows and **16** columns with in the dataset

FEATURES

age: Speaks about the age of the patient

education : shows the education the the patient done in four different phases 1. schooling,2.intermediate,3.UG,4.PG

****sex**:**Represents whether the patient is male or female

is_smoking:Represent whether the patient is smoking or not

cigsPerDay:If smoking then how many cigars the patient can take per a day



- ▶ **BPMEDS**: is there bp in the patient
- ▶ **PREVALENTSTROKE**: is there any previous stroke in the patient
- ▶ **diabetes**: Is the patient suffering with diabetes

TOTCHOL :Total cholesterol

- ▶ **SYSBP** :systolic blood pressure
- ▶ **DIABP** :Diastolic blood pressure
- ▶ **BMI** :body mass index
- ▶ **HEARTRATE** :represent the heart rate of the patient
- ▶ **glucose** :glucose levels of the patient
- ▶ **TENYEARCHD**:Predicting ten year coronary heart disease with in the patient



Finding the co related features

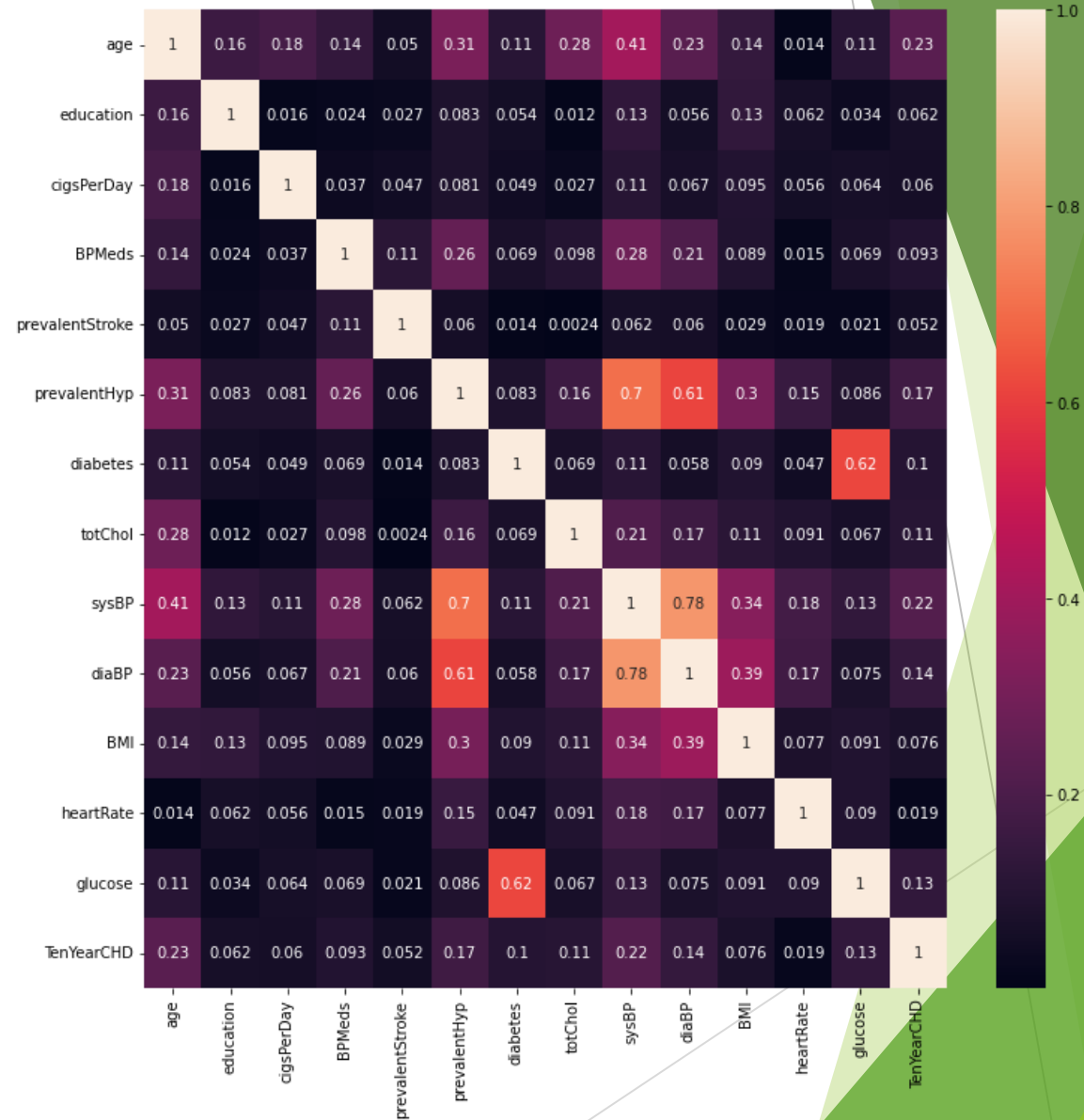
Here we can see that different features co related with each other

1. Sysbp and diabp features are highly co related with each other with .78 percentage among others

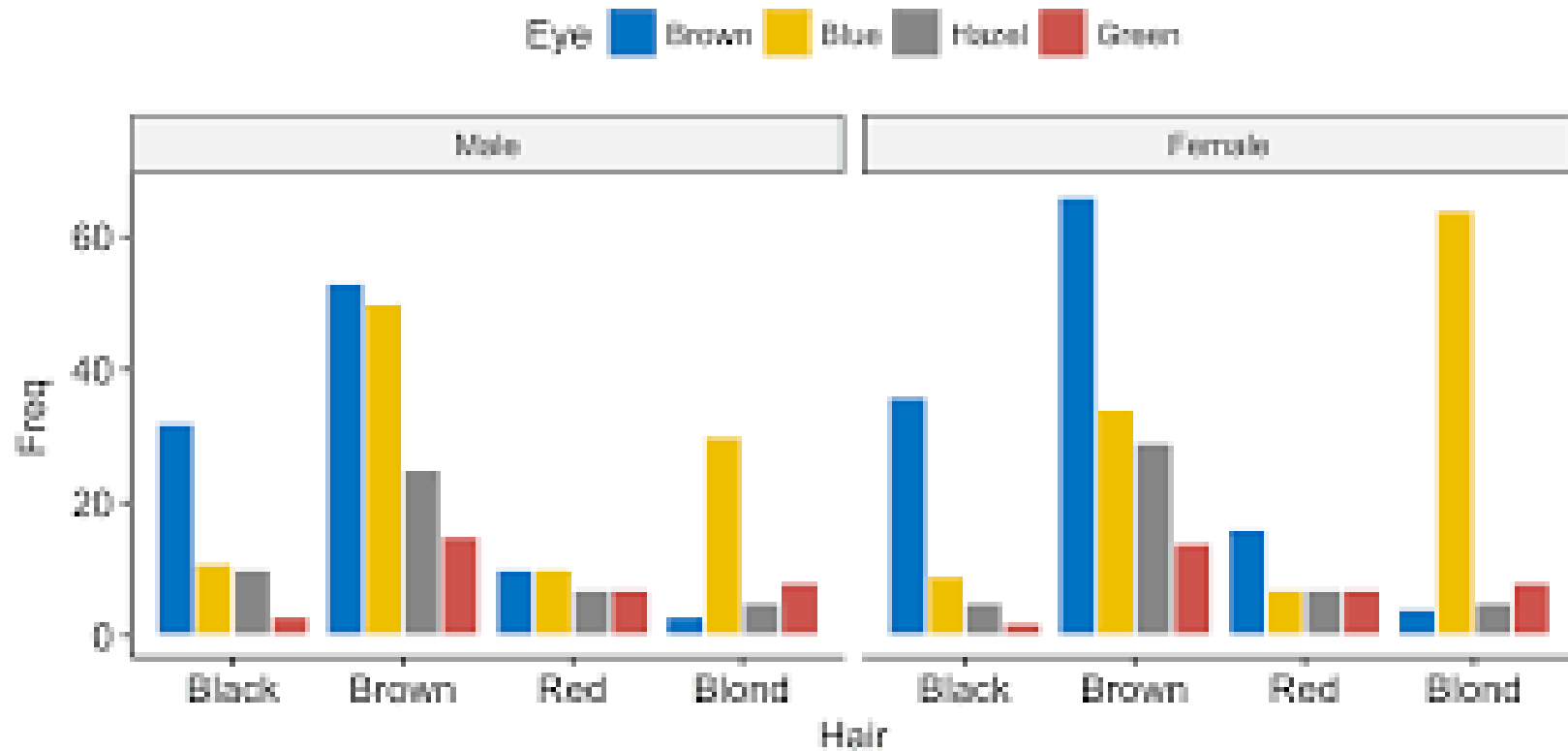
2. There is no significant correlation between independent and dependent variables

3. but there is high correlation between independent variables

4. Also some of the other features that are correlated with each other with some percentage



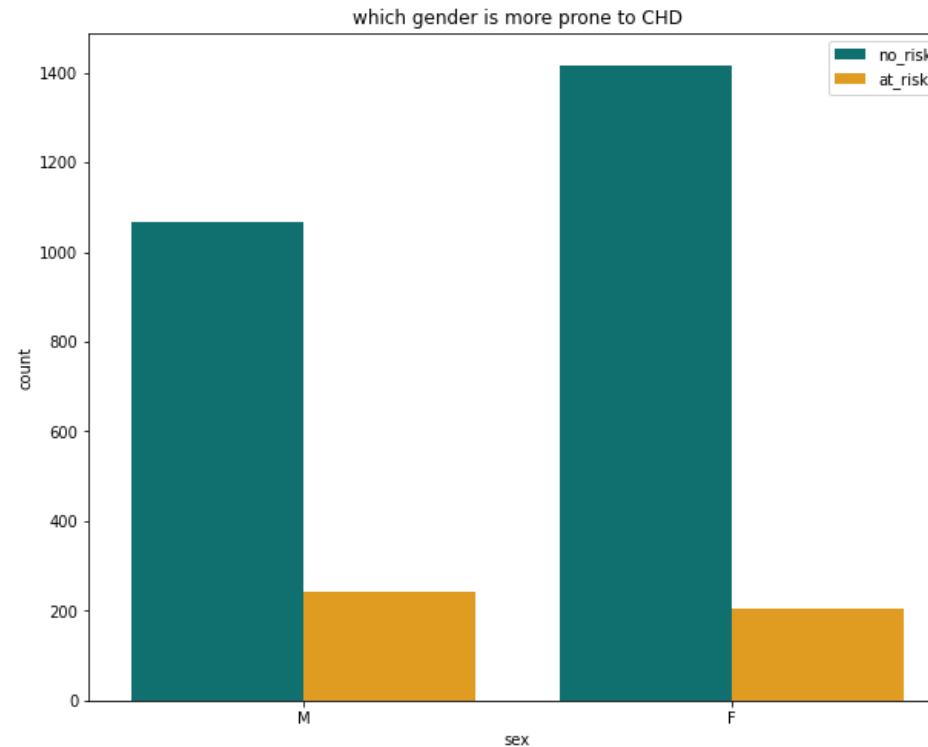
HERE WE HAVE DONE SOME CATEGORICAL VISUALISATION ON THE DATASET



Gender prone too heart disease

By the graph we intend that females are less effective of disease when compared to males

- As the green represents the no risk of disease and yellow represents at risk

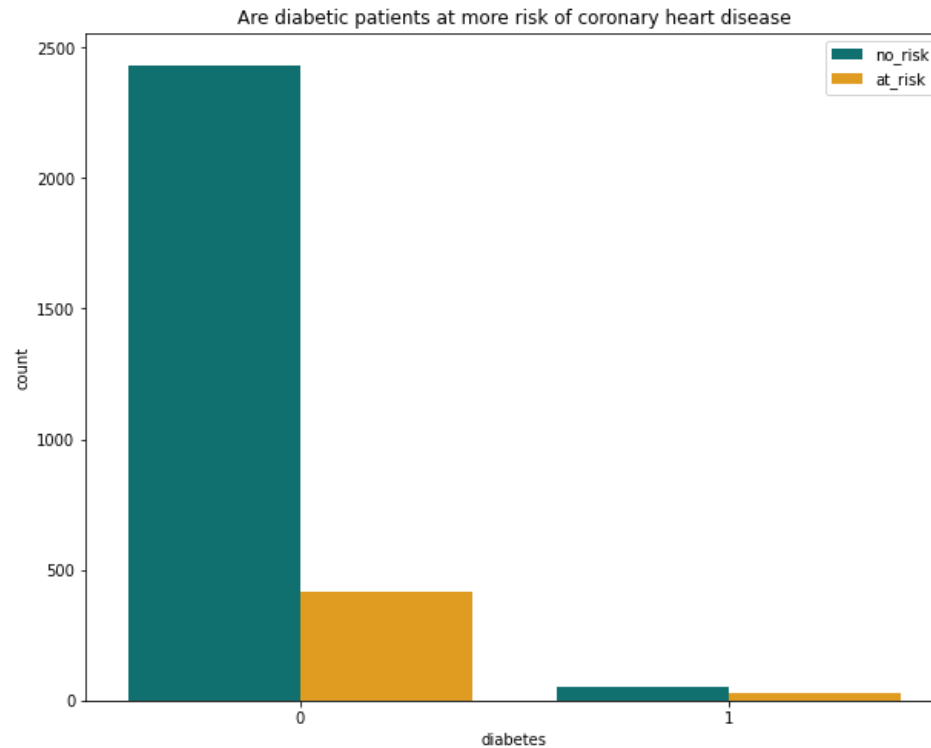


Are diabetic patients at more risk of coronary heart disease

we represent the diabetic patients at more risk of coronary heart disease

0 indicates tenyearCHD occurs and 1 indicates not

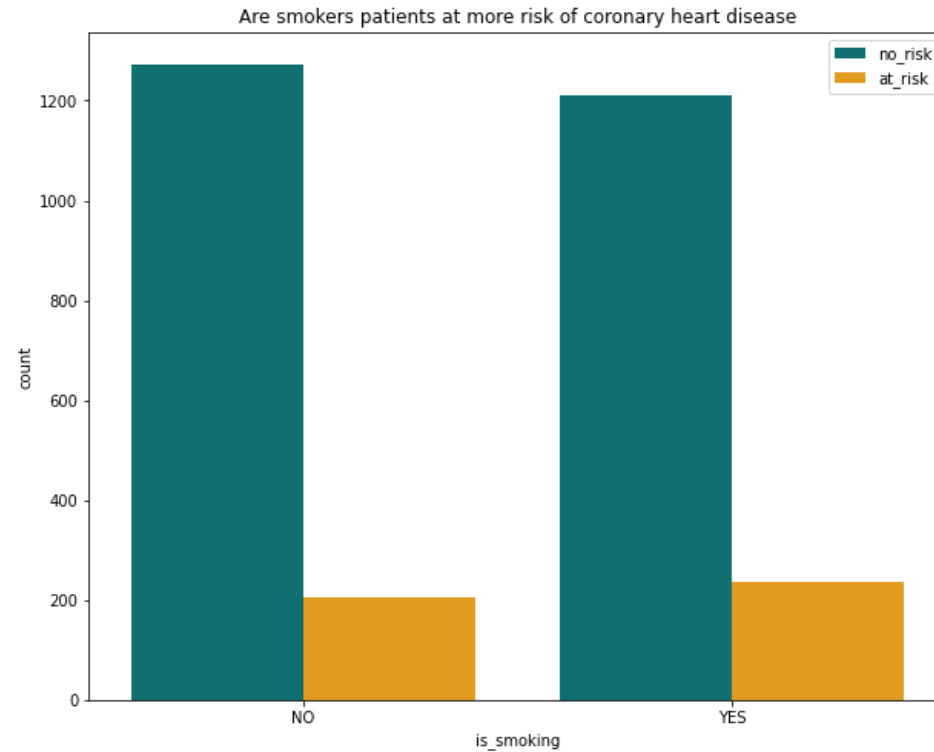
here we can see that diabetes patients are more seriously occurs with CHD



Are smokers patients at more risk of coronary heart diseases

WE compared that smoking effects the CHD then we noptices that

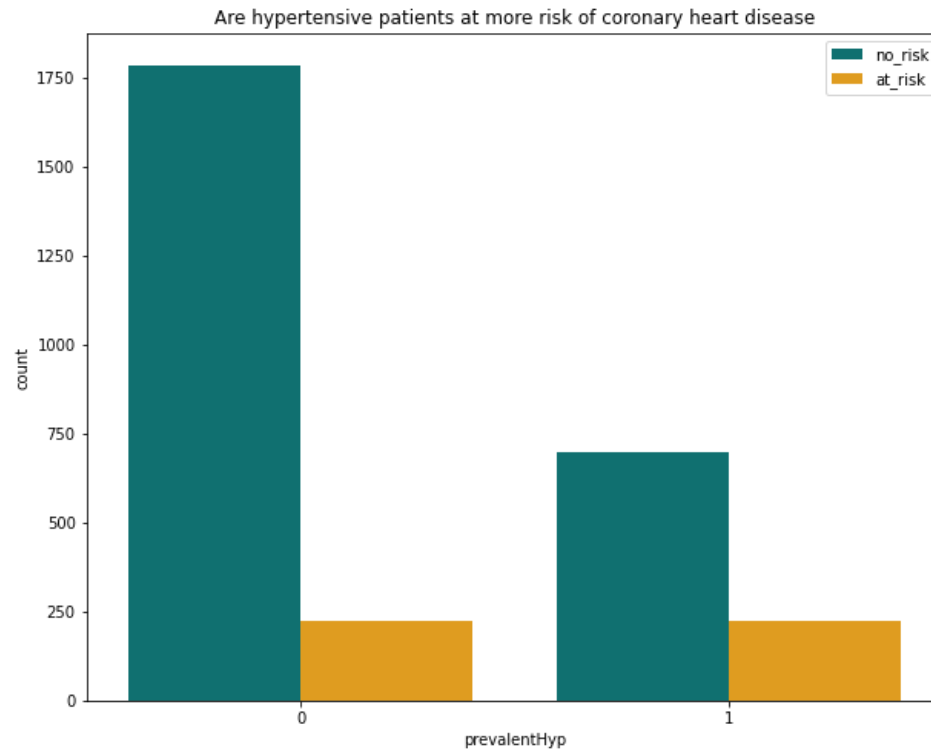
there is no relation of smoiking causes tha CHD by this graph



Are hypertensive patients at more risk of coronary heart disease

by this hypertensive patients at not at all more risk of coronary heart disease by prevlenthyp

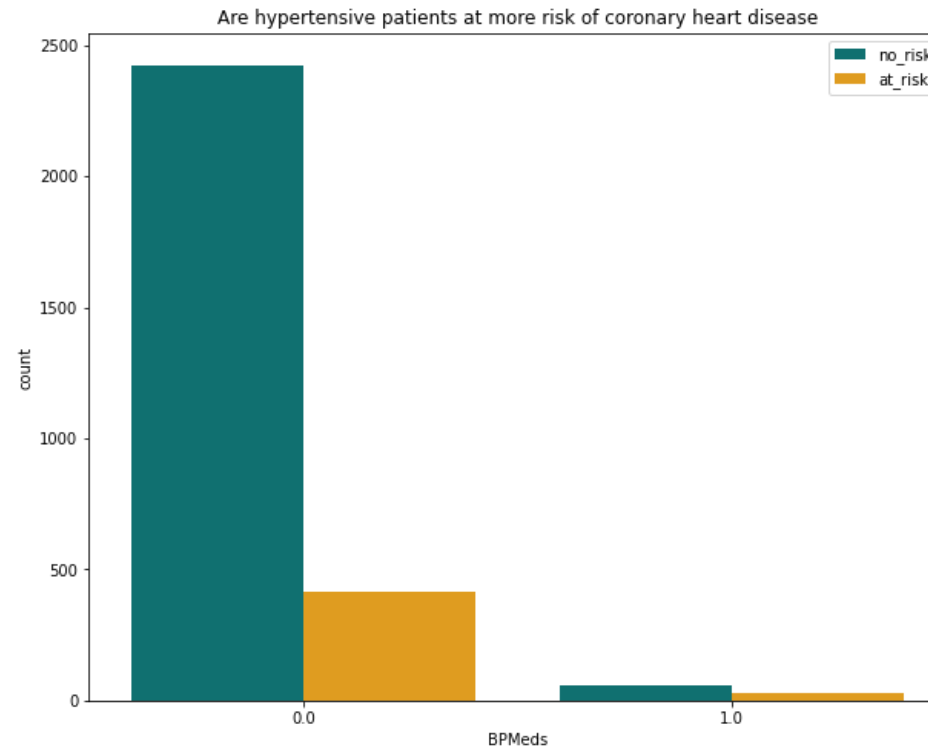
it doesn't occurs difference among them



Are patients with blood pressure on medication at more risk of coronary heart disease

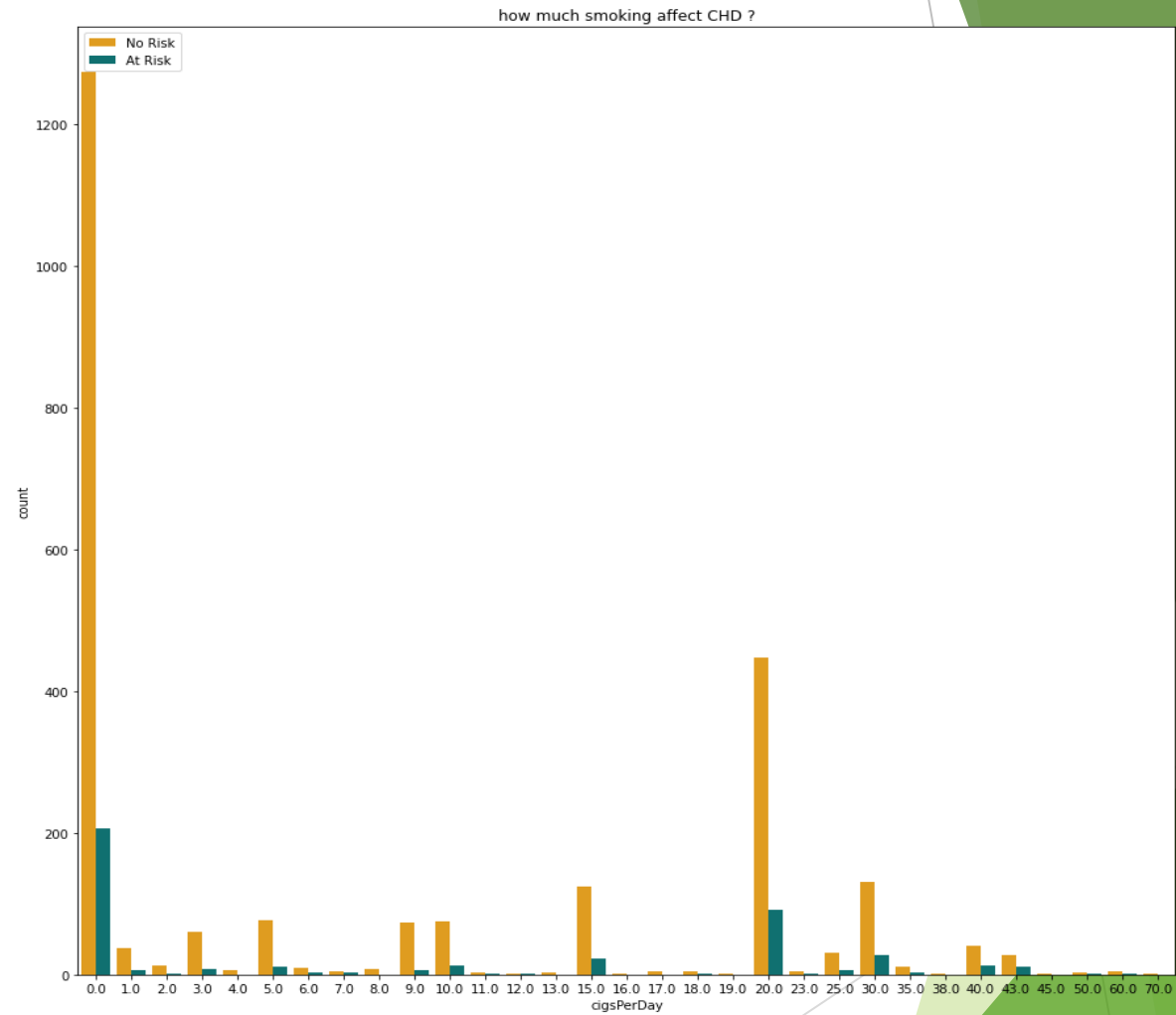
Patients with blood pressure can indicate the symptoms of occurring the CHD

There's almost equal chance of getting the CHD with BP



How much smoking affect CHD

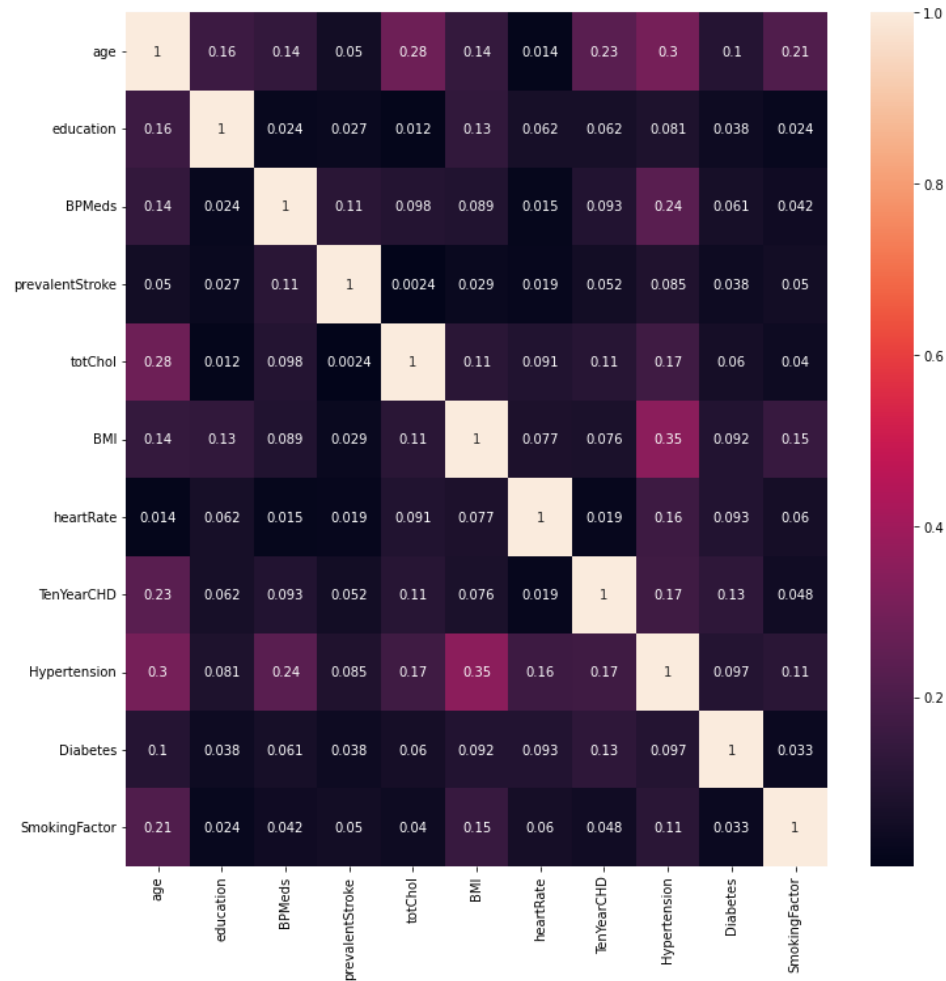
By understanding the graph we may know that smoking doesn't put any efforts to get effected by the heart disease



Plotting the correlation matrix using heatmap

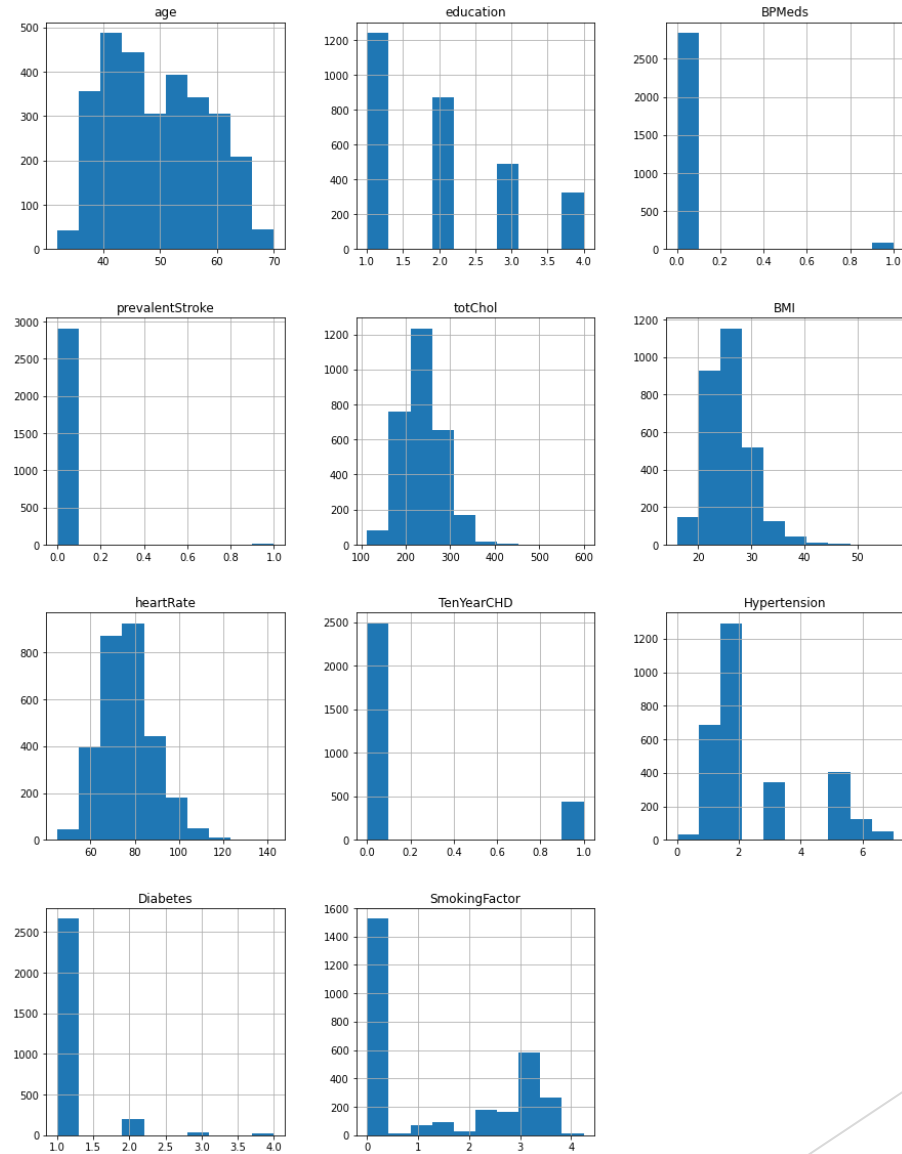
Here we can see that the multicollinearity doesn't exist between the columns

Every feature consists of its own correlation



plot histogram to see the distribution of the data

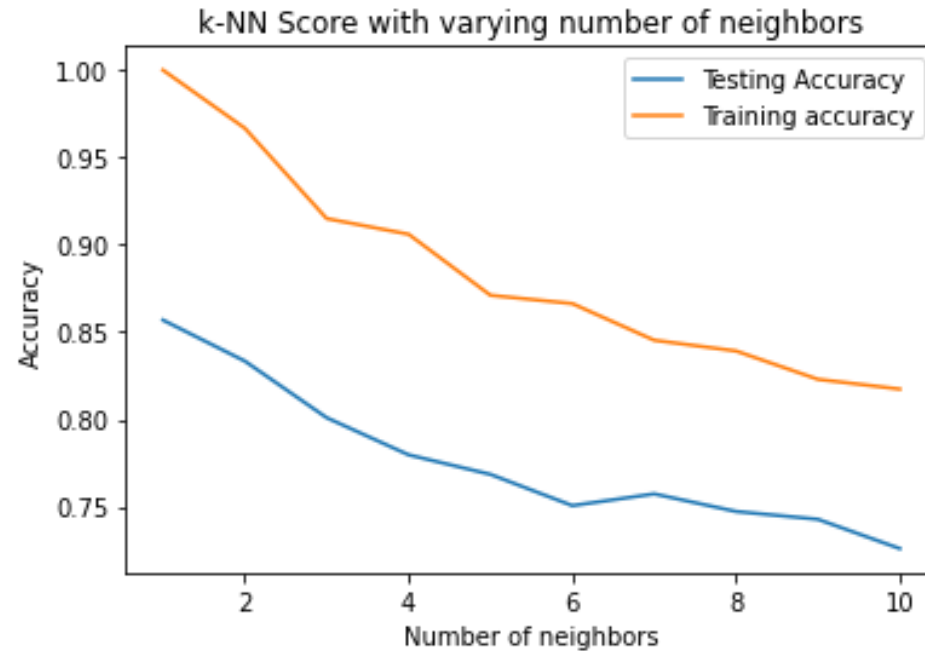
BPMeds and prevalentStroke does not help in explaining variance so we can remove those columns



ML model implementation KNN

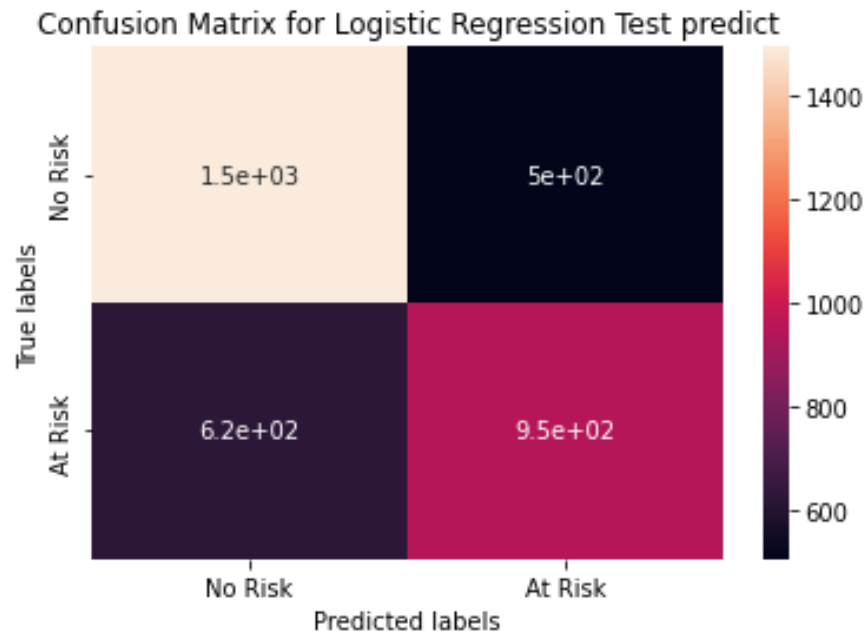
We can see that test score is increasing as number of neighbours increases.

Let's try to find best parameter for knn.



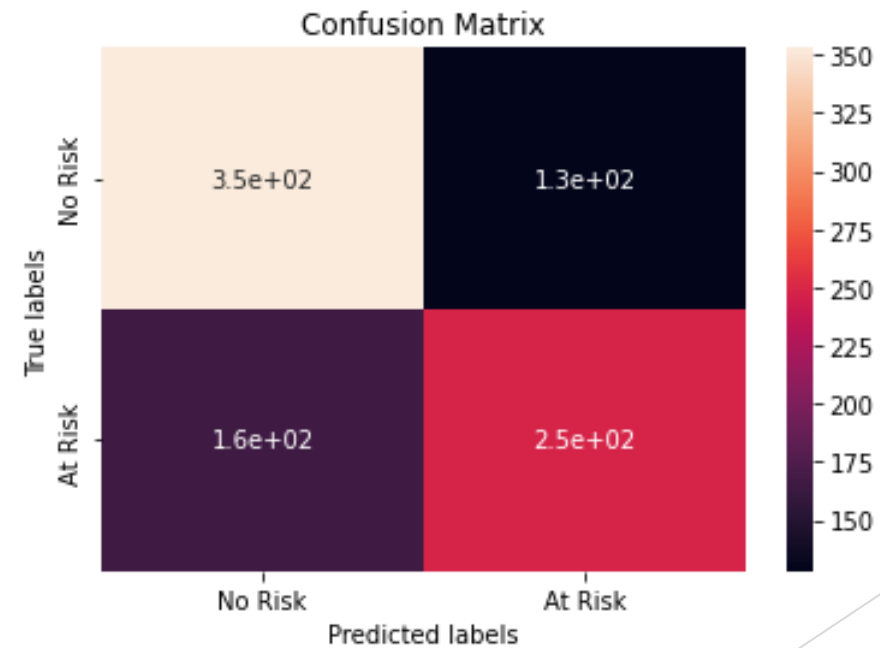
REGRESSION

confusion matrix for train data

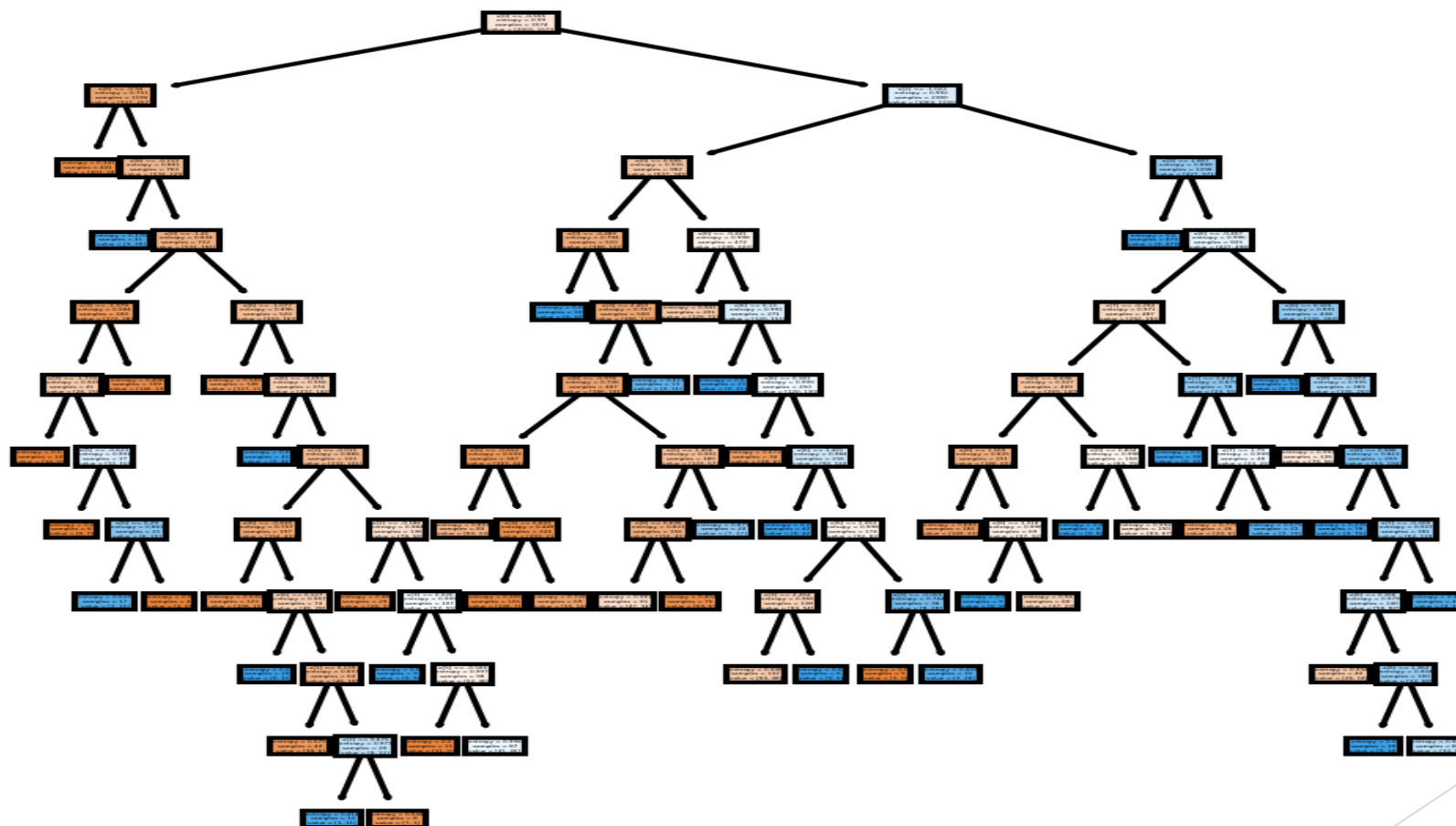


LOGISTIC

confusion matrix for test data

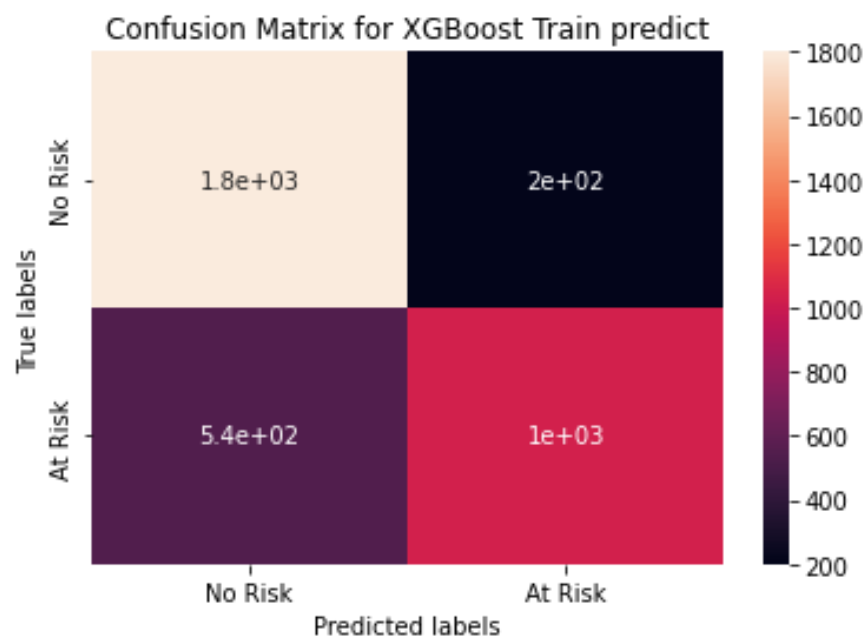


DECISION TREE CLASSIFIER

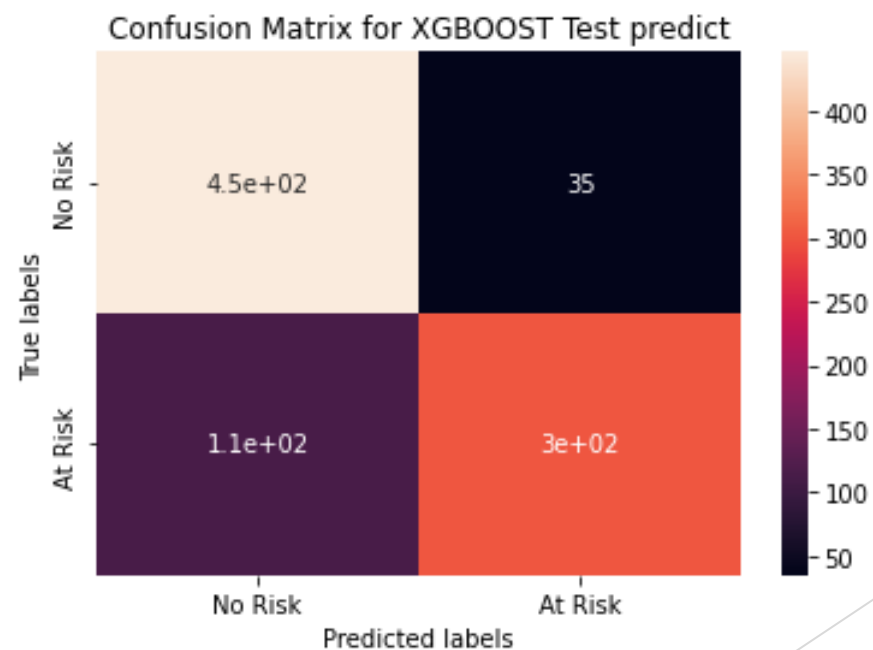


XGBOOST

CONFUSION matrix for Xgboost(train)



Confusion matrix for xgboost (test)



By selecting all the accuracies of the model gradientboosting classifier appears to be best accuracy

The accuracy of KNN : 0.86

The accuracy of Logistic regression : 0.68

the accuracy of decision tree : 0.79

The accuracy of randomforestclassifier : 0.88

the accuracy of XGBoost : 0.83

The accuracy of gradientboostingclassifier : 0.89

THANK YOU