# Capstone Project-3

# Cardiovascular Risk Prediction

**Pratik M Gumble**

# Point for Discussion

# Problem Statement

The classification goal is to predict whether the patient has a 10-year risk of future coronary heart disease (CHD).

# Introduction

Cardiovascular disease is a major health burden worldwide in the 21st century. The information which is gathered by data analysis of hospitals is utilizing by applying different blends of calculations and algorithms for the early-stage prediction of Cardiovascular ailments. Machine Learning is one of the slanting.

# Introduction

We will use a dataset to understand how to build different classification models in python from scratch. The models that will be introduced in this project are:

Logistic Regression
Random Forest
XG Boost
K-Nearest Neighbor
SVM

After we build the models using training data, we will test the accuracy of the model with test data and determine the appropriate model for this dataset.

# Dataset

The dataset is from an ongoing cardiovascular study on residents of the town of Framingham, Massachusetts. The dataset provides the patients' information. It includes over 4,000 records and 15 attributes.

# Dataset Variable Information

Each attribute is a potential risk factor. There are both demographic, behavioral, and medical risk factors.

**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)

# Dataset Variable Information

## 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.)

# Dataset Variable Information

## 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)

# Dataset Variable Information

## Medical(current)

• **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).

• **Glucose:** glucose level (Continuous)

# Dataset Variable Information

Predict variable (desired target)

• 10-year risk of coronary heart disease CHD

(binary: "1", means "Yes", "0" means "No") - DV

# Dataset Inspection & Processing

After data inspection we found that there are 3390 rows × 17 columns available in the dataset so we drop id and education columns and keep 15 columns for analysis.

## NaN Values Processing

As number of count of NaN values in cigsPerDay(22),BPMeds(44) totChol(38) and BMI(14),in are not more than 20% so we replace these values with zero instead of mean, median or mode.
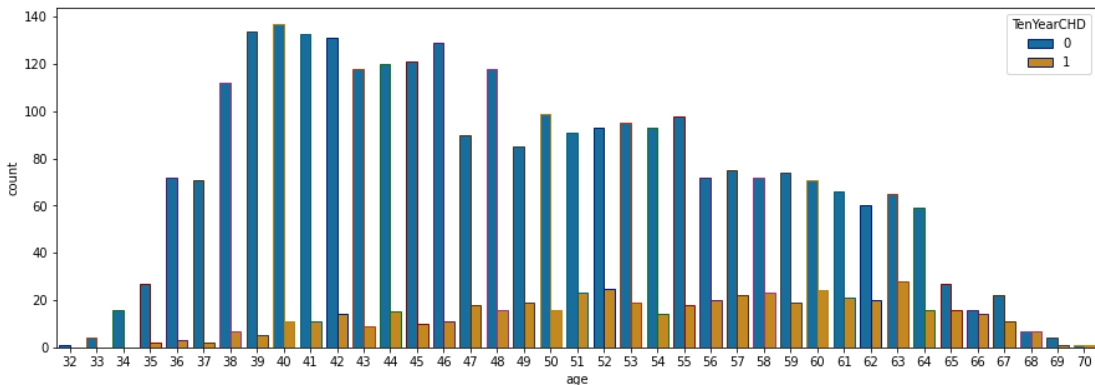Where Number of NaN values in glucose is 304 so we replace with it mean value.
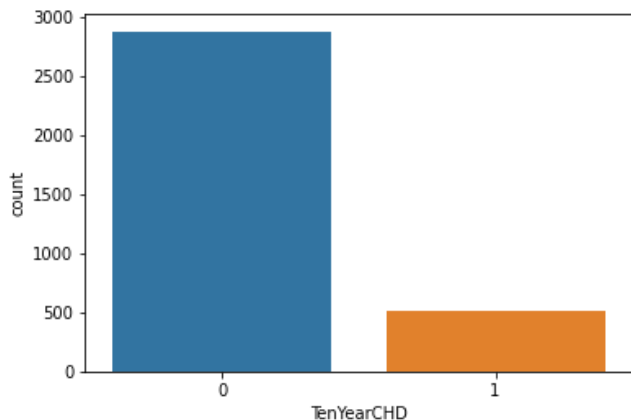
# Data Encoding

Before we start predicting, an important step to do is to convert our sex and is_smoking feature, which is a string, into integer.

M will be converted to 1 and F will be converted to 0.

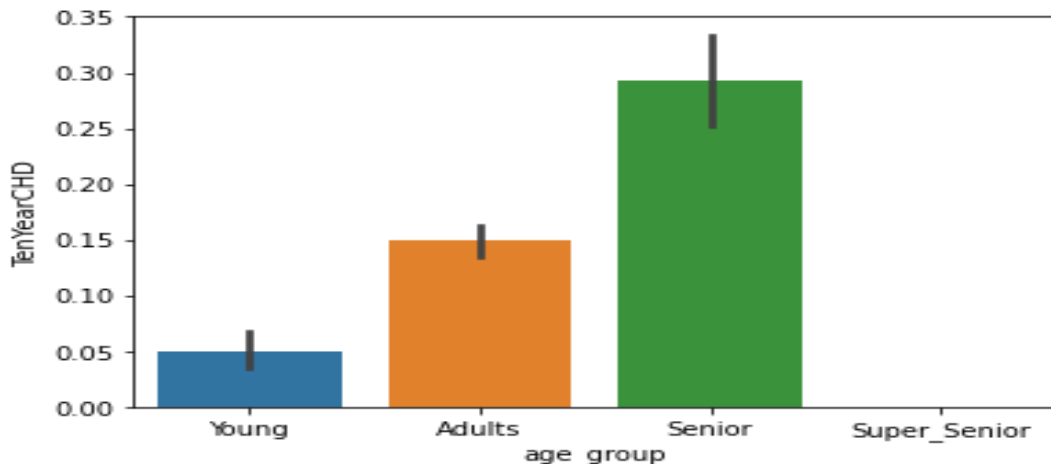Same for YES will be converted to 1 and NO will be converted to 0.

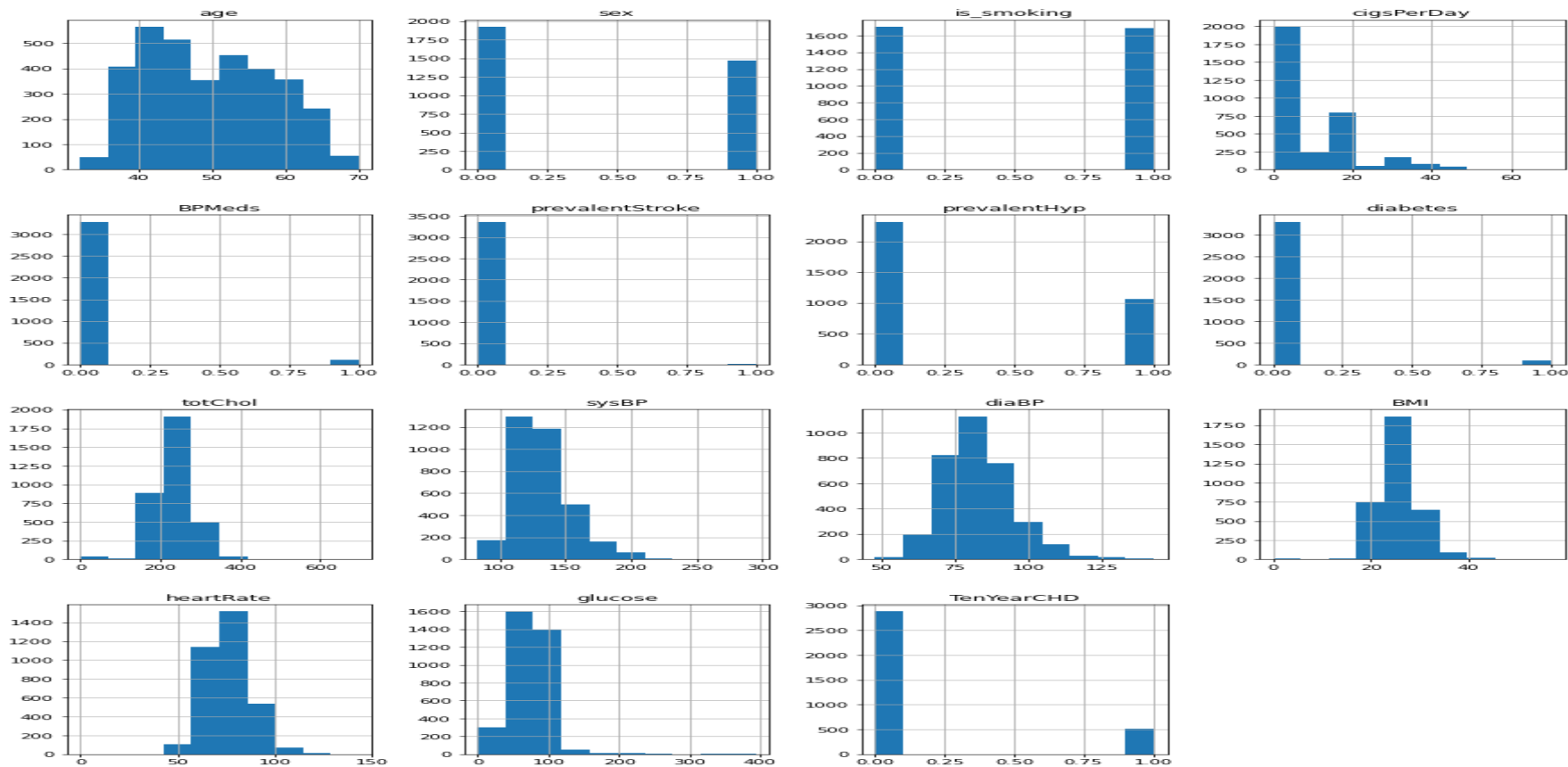We replaced this columns with same as a "sex" and "is_smoking".

# EDA



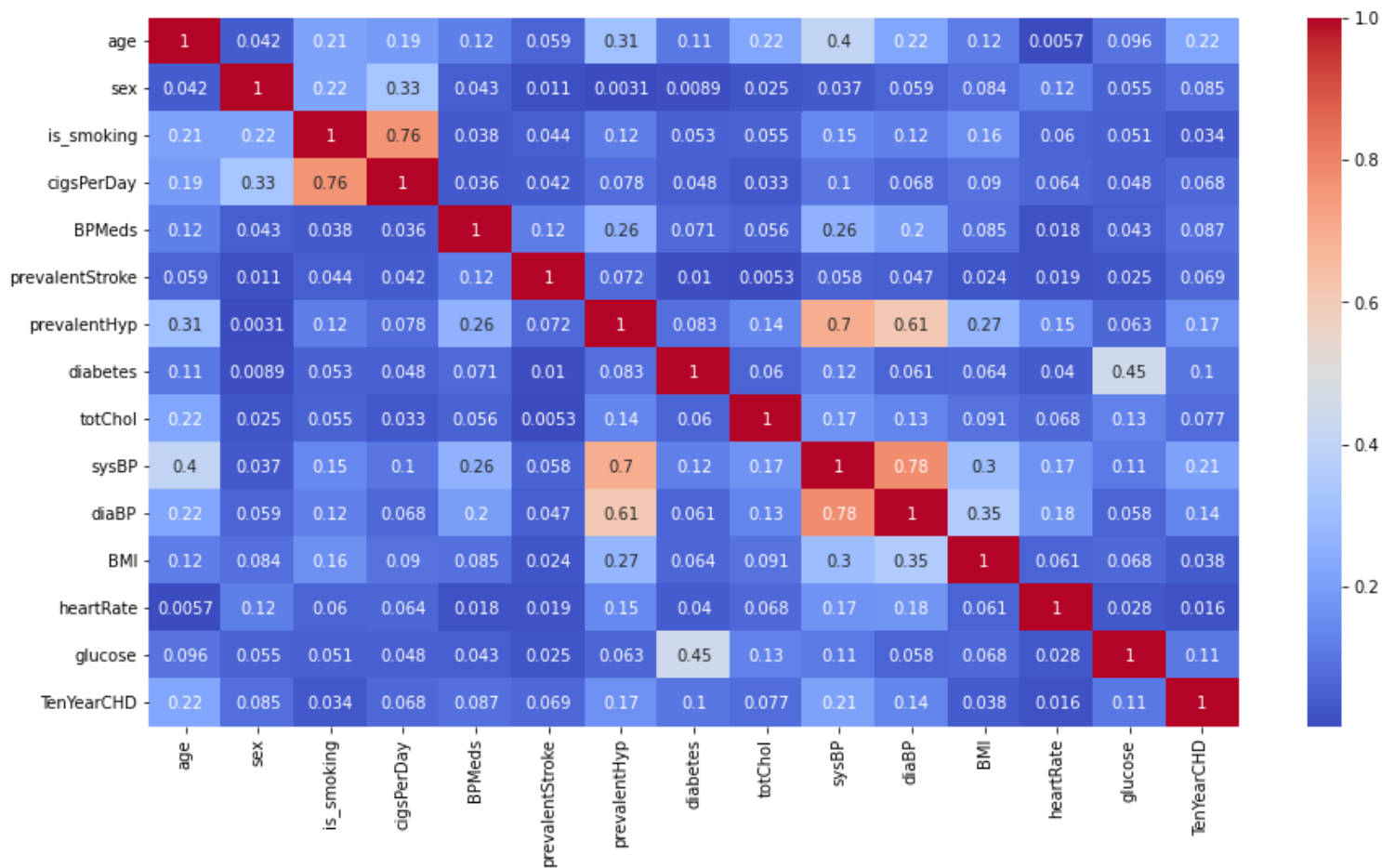| 0 | 2879 |
|---|------|
| 1 | 511 |

From the graph and statistical analysis it is clear that most of the people in the dataset have not 10-year risk of coronary heart disease CHD(2879 peoples) and 511 peoples have 10-year risk of coronary heart disease CHD
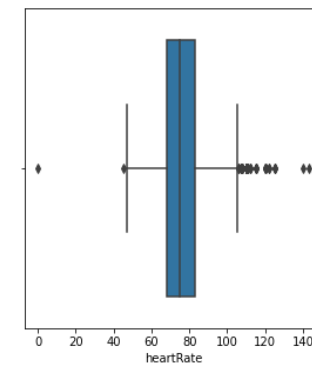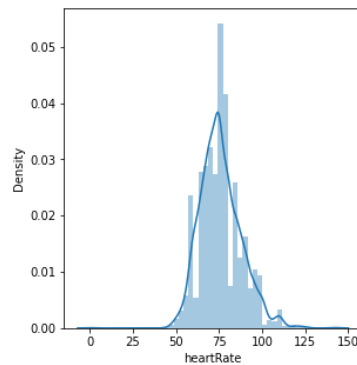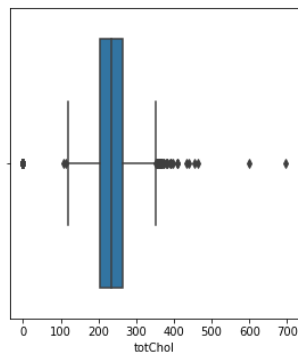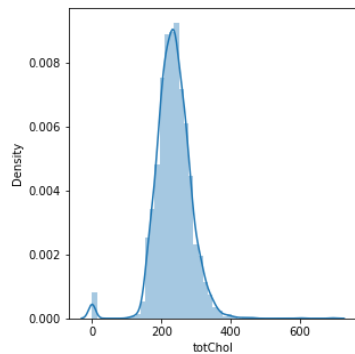
# EDA

# EDA

# EDA

# EDA

# Outliers Removal

## Interquartile Range (IQR)

After removing outliers using IQR, the data contains (2439, 15)

## Z-Score

After using Z-score to detect and remove outliers, the number of records in the dataset is (3046, 15).
As the number of records available is higher after Z-score, we will proceed with z score.

# Imbalanced Dataset

Sampling technique comes to save us and deal with imbalanced data.
There are two sampling techniques available to handle the imbalanced data:

Under Sampling
Over Sampling

Here we choose Oversampling technique as Unlike under-sampling, this method leads to no information loss.

Now to overcome data imbalance we apply Oversampling.
In oversampling we use SMOTE.
It stands for Synthetic Minority Oversampling Technique.

Original dataset shape Counter({0: 2646, 1: 400})
Resample dataset shape Counter({0: 2646, 1: 2646})

# Implementation

Now Lets implement 5 models on our dataset:

1. **Logistic Regression**

2. **Random Forrest**

3. **XGBoost**

4. **K-NN**

5. **SVM**

# 1. Logistic Regression

| Train accuracy | Test accuracy | Train precision | Test precision | Train recall | Test recall | Train f1 score | Test f1 score | Train ROC-AUC | Test ROC-AUC |
|---|---|---|---|---|---|---|---|---|---|
| 67% | 68% | 67% | 65% | 70% | 71% | 68% | 68% | 67% | 68% |

# 2. Random Forest

| Train accuracy | Test accuracy | Train precision | Test precision | Train recall | Test recall | Train f1 score | Test f1 score | Train ROC-AUC | Test ROC-AUC |
|---|---|---|---|---|---|---|---|---|---|
| 100% | 90% | 100% | 87% | 100% | 93% | 100% | 90% | 100% | 90% |

# 3. XG Boost

| Train accuracy | Test accuracy | Train precision | Test precision | Train recall | Test recall | Train f1 score | Test f1 score | Train ROC-AUC | Test ROC-AUC |
|---|---|---|---|---|---|---|---|---|---|
| 85% | 80% | 84% | 77% | 86% | 83% | 85% | 80% | 85% | 80% |

# 4. KNN Classifier

| Train accuracy | Test accuracy | Train precision | Test precision | Train recall | Test recall | Train f1 score | Test f1 score | Train ROC-AUC | Test ROC-AUC |
|---|---|---|---|---|---|---|---|---|---|
| 83% | 77% | 77% | 69% | 95% | 93% | 85% | 79% | 82% | 78% |

# 4. SVM Classifier

| Train accuracy | Test accuracy | Train precision | Test precision | Train recall | Test recall | Train f1 score | Test f1 score | Train ROC-AUC | Test ROC-AUC |
|---|---|---|---|---|---|---|---|---|---|
| 78% | 74% | 76% | 70% | 84% | 93% | 80% | 75% | 78% | 75% |

# Comparison

| Classifier | Train accuracy | Test accuracy | Train precision | Test precision | Train recall | Test recall | Train f1 score | Test f1 score | Train ROC-AUC | Test ROC-AUC |
|---|---|---|---|---|---|---|---|---|---|---|
| Logistic Regression | 67% | 68% | 67% | 65% | 70% | 71% | 68% | 68% | 67% | 68% |
| Random Forest | 100% | 90% | 100% | 87% | 100% | 93% | 100% | 90% | 100% | 90% |
| XGBoost | 85% | 80% | 84% | 77% | 86% | 83% | 85% | 80% | 85% | 80% |
| KNN | 83% | 77% | 77% | 69% | 95% | 93% | 85% | 79% | 82% | 78% |
| SVM | 78% | 74% | 76% | 70% | 84% | 93% | 80% | 75% | 78% | 75% |

It is quite evident from the results Random forest is the best model that can be used for cardiovascular risk prediction dataset since all the performance metrics (accuracy, precision, recall and roc-auc score) show a higher value for the random forest model !

# Predictive Model

Predictive model help us to understand possible future occurrences by analyzing our existance model.

Predictive models make assumptions based on what has happened in the past and what is happening now.

If incoming, new data shows changes in what is happening now, the impact on the likely future outcome must be recalculated, too.

In our predictive model we can see that Random forest classifier predict well as patient having **10-year risk of coronary heart disease or not.**

| Classifier | Actual | Prediction | Remarks |
|---|---|---|---|
| Logistic Regression | 1 | 1 | The Person have Cardiovascular Risk for next 10 years |
| Random Forest | 1 | 1 | The Person have Cardiovascular Risk for next 10 years |
| XG Boost | 1 | 1 | The Person have Cardiovascular Risk for next 10 years |
| KNN Classifer | 1 | 1 | The Person have Cardiovascular Risk for next 10 years |

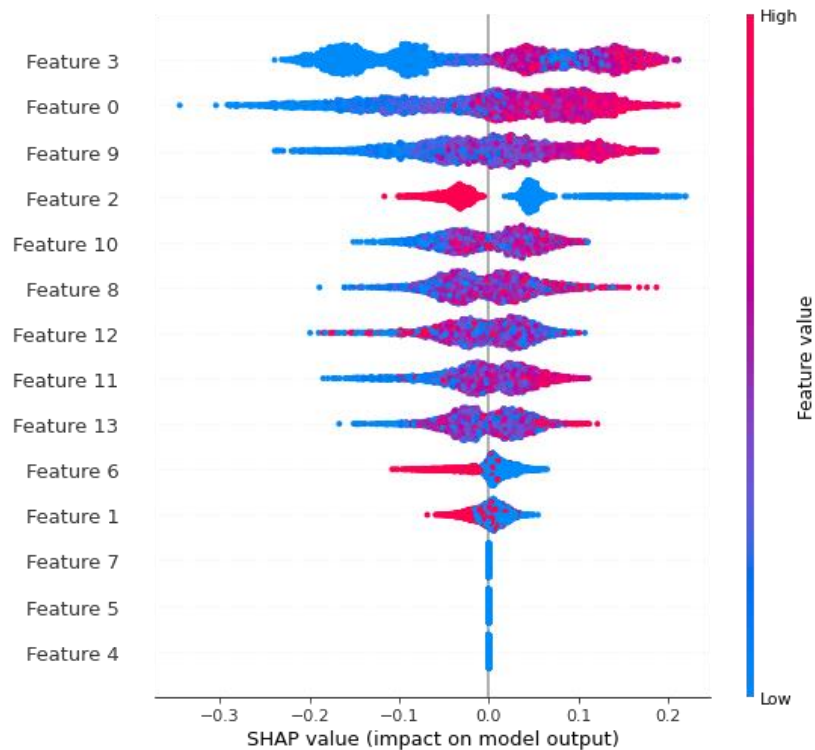| Classifier | Actual | Prediction | Remarks |
|---|---|---|---|
| Logistic Regression | 0 | 1 | The Person have Cardiovascular Risk for next 10 years |
| Random Forest | 0 | 0 | The Person does not have Cardiovascular Risk for next 10 years |
| XG Boost | 0 | 1 | The Person have Cardiovascular Risk for next 10 years |
| KNN Classifer | 0 | 1 | The Person have Cardiovascular Risk for next 10 years |

# Feature Importance

Feature (variable) importance indicates how much each feature contributes to the model prediction. Basically, it determines the degree of usefulness of a specific variable for a current model and prediction.

By analyzing variable importance scores, we would be able to find out irrelevant features and exclude them. Reducing the number of not meaningful variables in the model may speed up the model or even improve its performance.
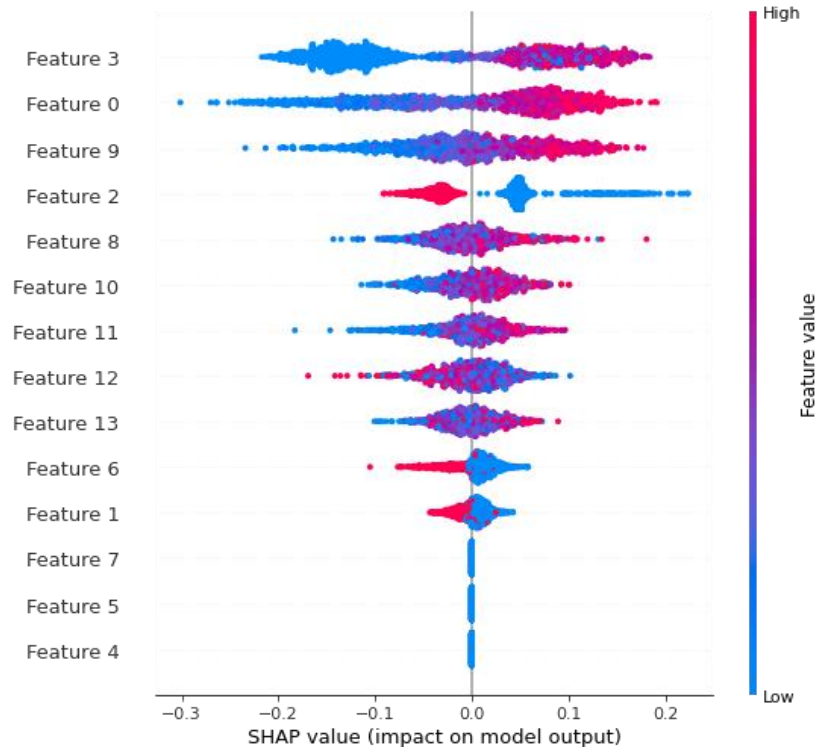
In almost all classifier  CigsPerDay, age, sysBP, diaBP, tatchol, heartrate, BMI, glucose are the important feature to determine person having 10 year risk for heart diastase or not. Out of that CigsPerDay, age and BP are the most important feature which decided the person having  chance of 10 year heart diseases or not.

# Feature Importance

# Conclusion

The main aim of this project is to compare the accuracy and other parameters like precision, recall ,f1 score, auc roc of all the classification algorithms to evaluate the risk of 10-year CHD using 14 features. After implementing four classification models and comparing their accuracy and other scores, we can conclude that for this dataset Random forest Classifier is the appropriate model to be used. Also out of all features CigsPerDay, age and BP are the most important feature which decided the person having chance of 10 year heart diseases or not.

Thank You