# Capstone Project 3

## Cardiovascular Risk Prediction

### Satyajit Sahoo

# Understanding Problem Statement:

- We are given with a data set where one of the variables represent whether a person may suffer from heart disease in next 10 years. Of size(3390,16)

- Our aim is to predict the same.

- The target variable is TenYearsCHD and it is binary where:

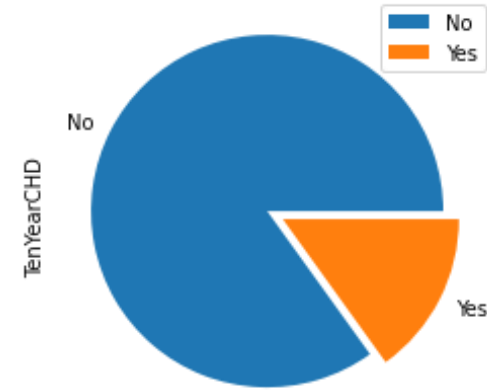    0  means **NO**
    1 means  **YES**

# Variables Involved :

- ID
- Age
- Education
- Sex
- Is_smoking
- Cigs per days

- BPMeds
- Prevalent stroke
- Prevalent Hypertensive
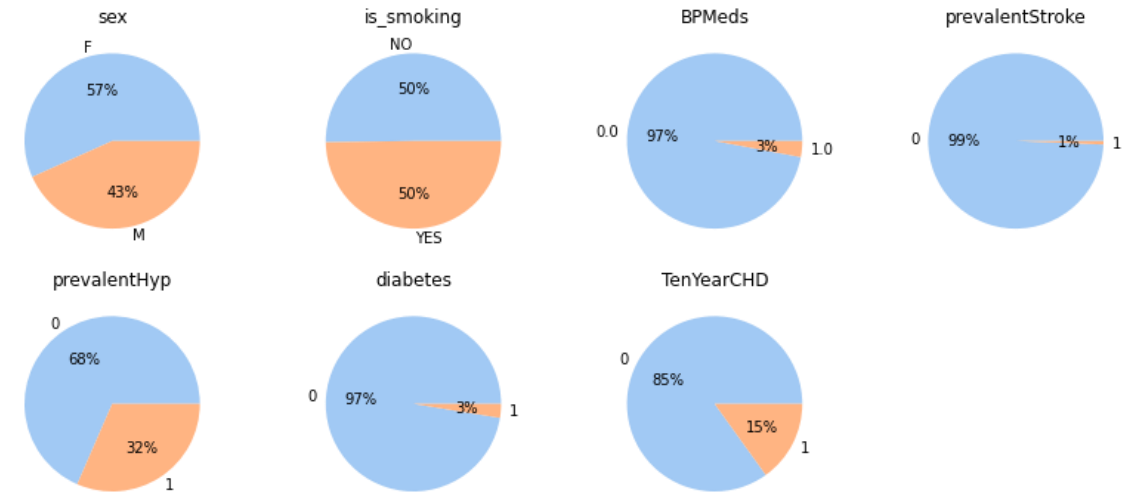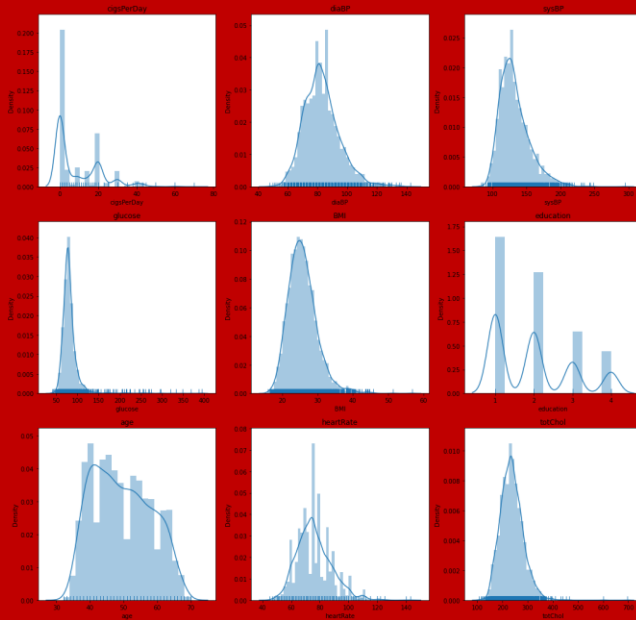- Diabetes
- Total Cholestrol
- sysBP

- diaBP
- BMI
- Heart Rate
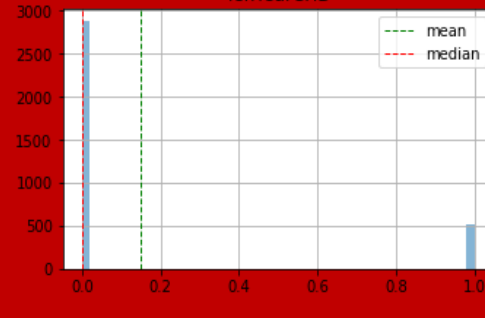- Glucose
- Ten Year CHD
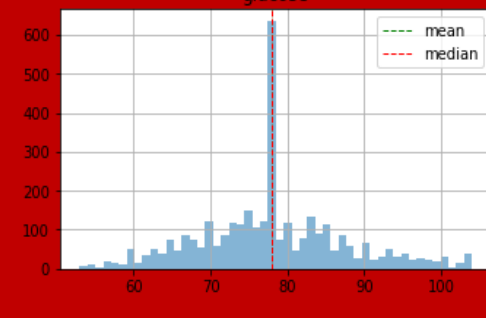
# Exploratory Data Analysis :

## Null Values

On exploring we see that these variables have null values in them.



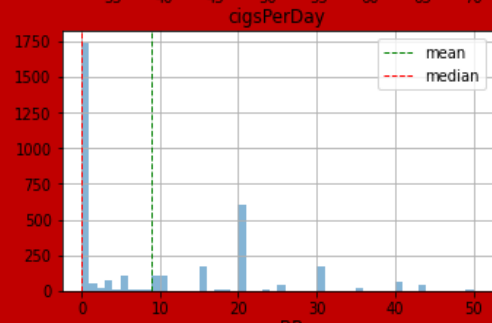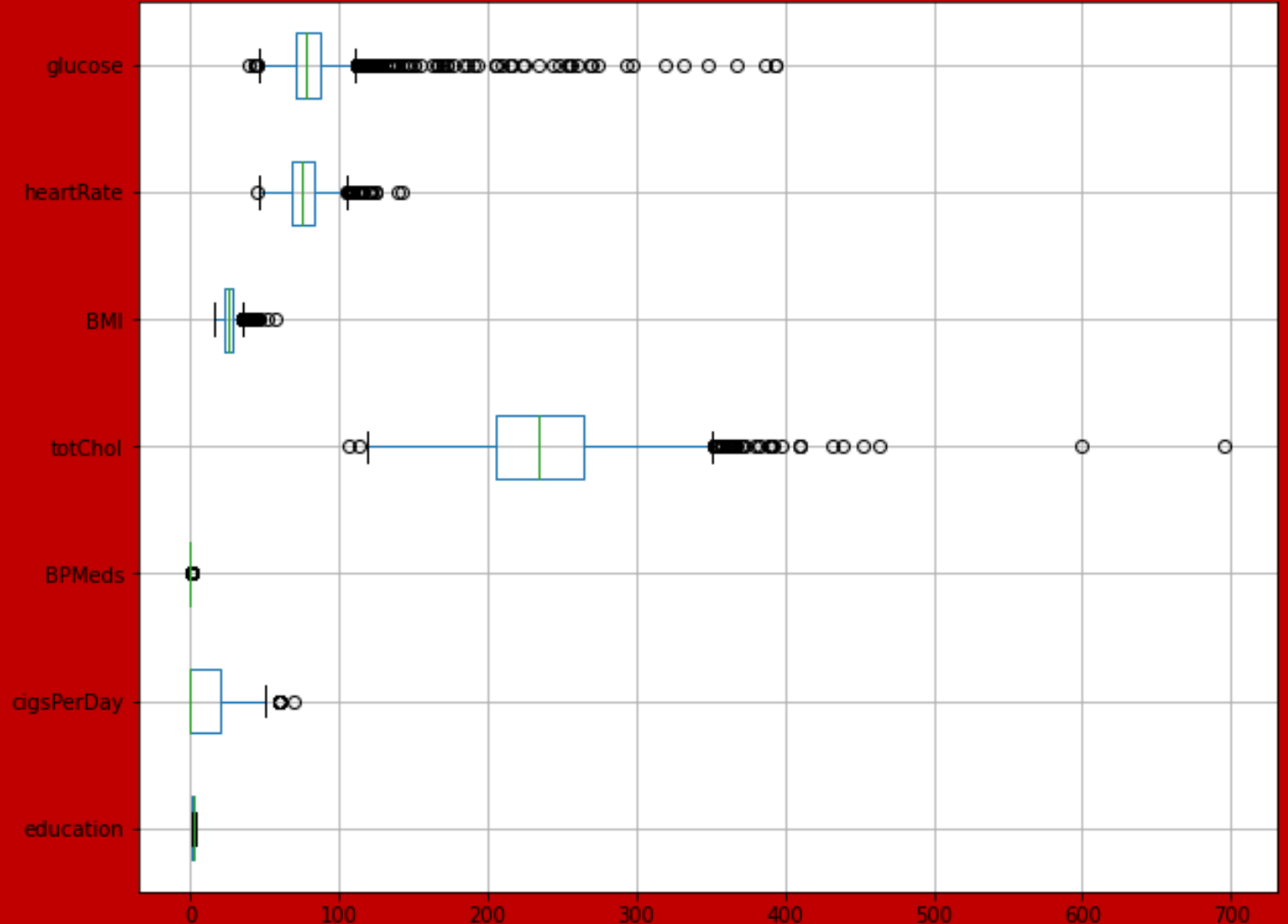| | value_counts | percentage |
|---|---|---|
| education | 87 | 2.57 |
| cigsPerDay | 22 | 0.65 |
| BPMeds | 44 | 1.30 |
| totChol | 38 | 1.12 |
| BMI | 14 | 0.41 |
| heartRate | 1 | 0.03 |
| glucose | 304 | 8.97 |

# Null -  Value treatment

Deciding which value to impute in place of null values.

Mean

Mode
Median

# Outlier treatment

# Logistic Regression Data Preparation

## Verifying linear dependency



## Removing Multicollinearity

| | feature | VIF |
|---|---|---|
| 0 | education | 4.599856 |
| 1 | sex | 1.969307 |
| 2 | cigsPerDay | 1.760078 |
| 3 | prevalentHyp | 1.621659 |
| 4 | glucose_age_chol | 6.468307 |

# Preparing data for SVCs

# Preparing data for Ensemble tree models.

# Handling Data Imbalance.

```
Before X_train : (2712, 5)
Before Y_train :
0     2303
1      409
Name: TenYearCHD, dtype: int64
After X_smote : (3810, 5)
After Y_smote :
0     1905
1     1905
Name: TenYearCHD, dtype: int64
```
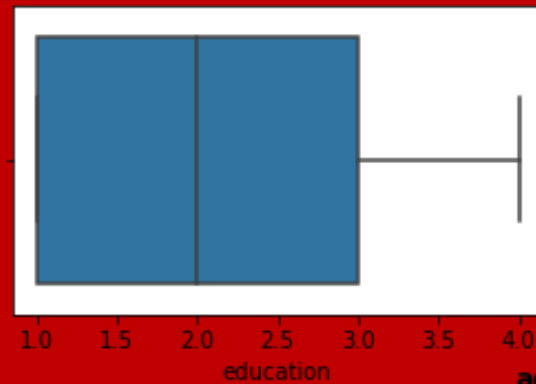
# Logistic Regression Implementation

# Logistic Regression Implementation (Contd..)

| Classification Report of train set | | | | | |
|---|---|---|---|---|---|
| | 0 | 1 | accuracy | macro avg | weighted avg |
| precision | 0.639166 | 0.635697 | 0.638643 | 0.637432 | 0.637771 |
| recall | 0.908081 | 0.238313 | 0.638643 | 0.573197 | 0.638643 |
| f1-score | 0.750255 | 0.346667 | 0.638643 | 0.548461 | 0.587897 |
| support | 1621.000000 | 1091.000000 | 0.638643 | 2712.000000 | 2712.000000 |

| Classification Report of test output | | | | | |
|---|---|---|---|---|---|
| | 0 | 1 | accuracy | macro avg | weighted avg |
| precision | 0.642361 | 0.607843 | 0.637168 | 0.625102 | 0.628717 |
| recall | 0.902439 | 0.231343 | 0.637168 | 0.566891 | 0.637168 |
| f1-score | 0.750507 | 0.335135 | 0.637168 | 0.542821 | 0.586319 |
| support | 410.000000 | 268.000000 | 0.637168 | 678.000000 | 678.000000 |

# Gaussian Naïve Bayes Implementation

# Support Vector Classifier Implementation (Contd..)

| Classification Report of train set | 0 | 1 | accuracy | macro avg | weighted avg |
|---|---|---|---|---|---|
| precision | 0.616153 | 0.696822 | 0.628319 | 0.656487 | 0.650925 |
| recall | 0.919637 | 0.243798 | 0.628319 | 0.581718 | 0.628319 |
| f1-score | 0.737910 | 0.361217 | 0.628319 | 0.549563 | 0.575537 |
| support | 1543.000000 | 1169.000000 | 0.628319 | 2712.000000 | 2712.000000 |

| Classification Report of test output | 0 | 1 | accuracy | macro avg | weighted avg |
|---|---|---|---|---|---|
| precision | 0.631944 | 0.676471 | 0.638643 | 0.654208 | 0.650398 |
| recall | 0.916877 | 0.245552 | 0.638643 | 0.581214 | 0.638643 |
| f1-score | 0.748201 | 0.360313 | 0.638643 | 0.554257 | 0.587440 |
| support | 397.000000 | 281.000000 | 0.638643 | 678.000000 | 678.000000 |

# Support Vector Classifier with balanced loss Implementation (Contd..)

## Classification Report of train set

|  | 0 | 1 | accuracy | macro avg | weighted avg |
|---|---|---|---|---|---|
| precision | 0.616153 | 0.696822 | 0.628319 | 0.656487 | 0.650925 |
| recall | 0.919637 | 0.243798 | 0.628319 | 0.581718 | 0.628319 |
| f1-score | 0.737910 | 0.361217 | 0.628319 | 0.549563 | 0.575537 |
| support | 1543.000000 | 1169.000000 | 0.628319 | 2712.000000 | 2712.000000 |

## Classification Report of test output

|  | 0 | 1 | accuracy | macro avg | weighted avg |
|---|---|---|---|---|---|
| precision | 0.631944 | 0.676471 | 0.638643 | 0.654208 | 0.650398 |
| recall | 0.916877 | 0.245552 | 0.638643 | 0.581214 | 0.638643 |
| f1-score | 0.748201 | 0.360313 | 0.638643 | 0.554257 | 0.587440 |
| support | 397.000000 | 281.000000 | 0.638643 | 678.000000 | 678.000000 |

# LGBM Implementation

# LGBM Implementation

## Classification Report of train set

|  | 0 | 1 | accuracy | macro avg | weighted avg |
|---|---|---|---|---|---|
| precision | 0.612679 | 0.704156 | 0.626475 | 0.658418 | 0.652481 |
| recall | 0.921018 | 0.244068 | 0.626475 | 0.582543 | 0.626475 |
| f1-score | 0.735854 | 0.362492 | 0.626475 | 0.549173 | 0.573403 |
| support | 1532.000000 | 1180.000000 | 0.626475 | 2712.000000 | 2712.000000 |

## Classification Report of test output

|  | 0 | 1 | accuracy | macro avg | weighted avg |
|---|---|---|---|---|---|
| precision | 0.604167 | 0.627451 | 0.607670 | 0.615809 | 0.614195 |
| recall | 0.901554 | 0.219178 | 0.607670 | 0.560366 | 0.607670 |
| f1-score | 0.723493 | 0.324873 | 0.607670 | 0.524183 | 0.551816 |
| support | 386.000000 | 292.000000 | 0.607670 | 678.000000 | 678.000000 |

# XGRFB Implementation

# XGRFB Implementation

## Classification Report of train set

|  | 0 | 1 | accuracy | macro avg | weighted avg |
|---|---|---|---|---|---|
| precision | 0.575771 | 0.748166 | 0.601770 | 0.661968 | 0.657328 |
| recall | 0.927922 | 0.238504 | 0.601770 | 0.583213 | 0.601770 |
| f1-score | 0.710611 | 0.361702 | 0.601770 | 0.536157 | 0.545548 |
| support | 1429.000000 | 1283.000000 | 0.601770 | 2712.000000 | 2712.000000 |

## Classification Report of test output
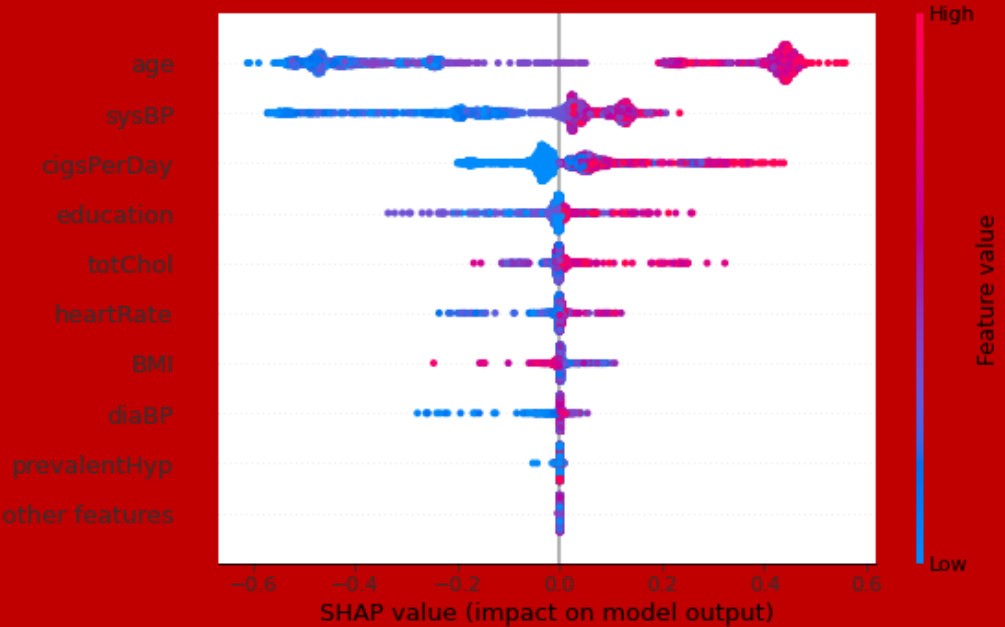
|  | 0 | 1 | accuracy | macro avg | weighted avg |
|---|---|---|---|---|---|
| precision | 0.574653 | 0.686275 | 0.591445 | 0.630464 | 0.626512 |
| recall | 0.911846 | 0.222222 | 0.591445 | 0.567034 | 0.591445 |
| f1-score | 0.705005 | 0.335731 | 0.591445 | 0.520368 | 0.533440 |
| support | 363.000000 | 315.000000 | 0.591445 | 678.000000 | 678.000000 |

# Neural Network Implementation

# Conclusion..

1. Logistic Regression can predict 61% of the negative values long with 35% of False     Negative predictions.

2. Gaussian Naive Bayes can predict 55% of the negative values along with 34% of False Negative predictions.

3. Support Vector Classifier *(without balanced loss function)* can predict 68% of the negative values with 37% of False Negative predictions.

4. Support Vector Classifier *(with balanced loss function)* can predict 70% of the negative values with 37% of False Negative predictions.

5.  LGBM can predict 63% of the negative values with 40% of False Negative predictions.

6. XGRFB can predict 69% of the negative values with 42% of False Negative predictions.

7. Random Forest Model can predict 66% of the negative values with 39 % False Negative predictions.