

Capstone Project

Cardiovascular Risk Prediction

by

Roshan Jamthe

Points for Discussion

- Problem Statement
- Data Summary
- Data Visualization
- Feature engineering
- Feature selection
- Train Test datasets
- ML Models' performance
- Best model
- Feature Importance
- Conclusions

Problem Statement

The dataset is from an ongoing cardiovascular study on residents of the town of Framingham, Massachusetts. The classification goal is to predict whether the patient has a 10-year risk of future coronary heart disease (CHD).

The dataset provides the patients' information. It includes over 3,000 records and 15 attributes. Each attribute is a potential risk factor. There are demographic, behavioral, and medical risk factors.

Data Summary

Demographic:

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

Age: Age of the patient;(Continuous - 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)

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)

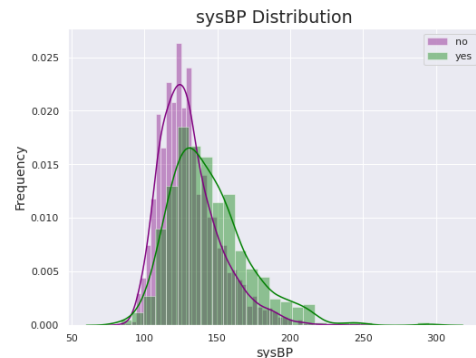
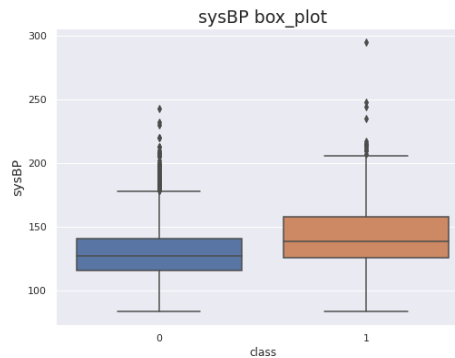
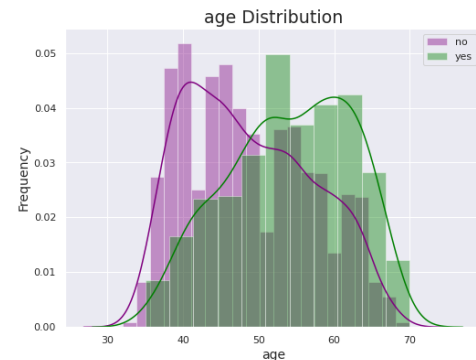
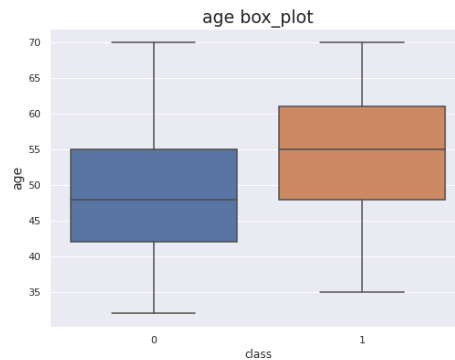
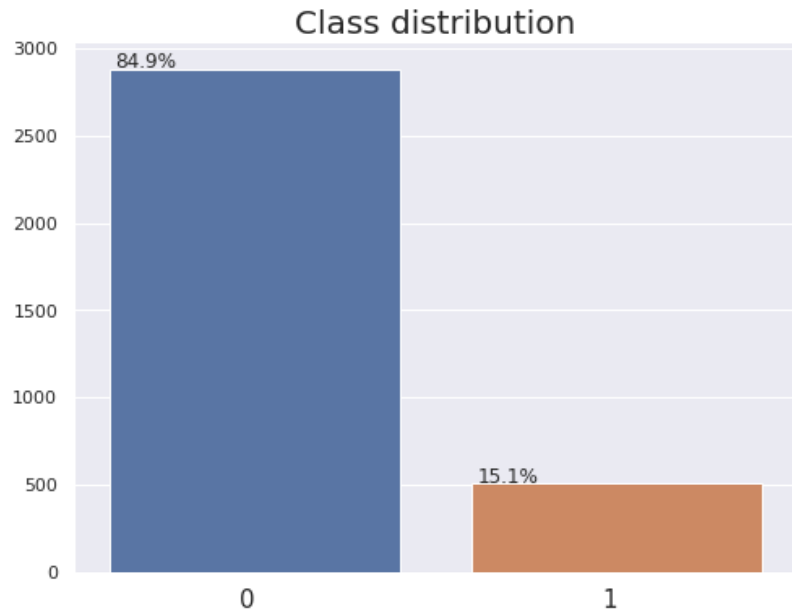
Heart Rate: heart rate (Continuous - In medical research, variables such as heart rate though in fact discrete, yet are considered continuous because of large number of possible values.)

Glucose: glucose level (Continuous)

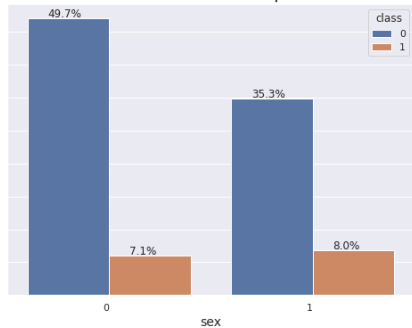
Predict variable (desired target):

TenYearCHD: 10-year risk of coronary heart disease CHD(binary: “1”, means “Yes”, “0” means “No”) - DV

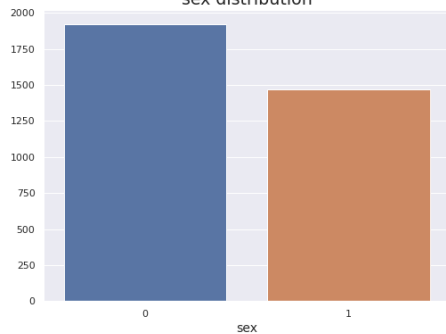
Data Visualization



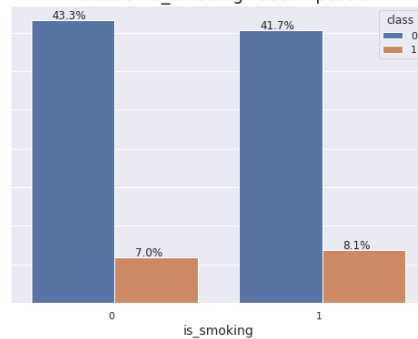
Feature "sex" decomposed



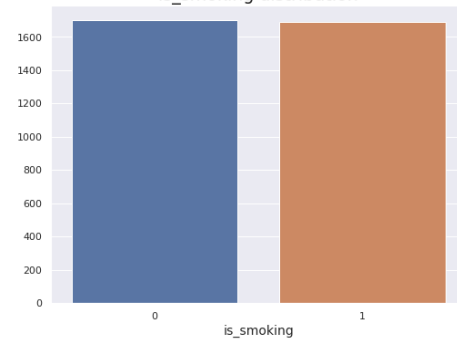
sex distribution



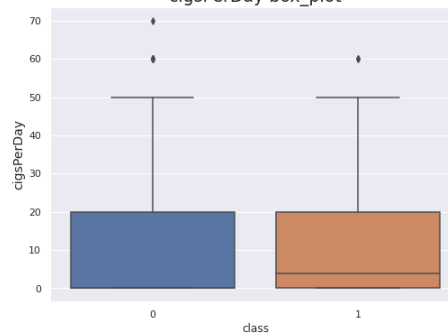
Feature "is_smoking" decomposed



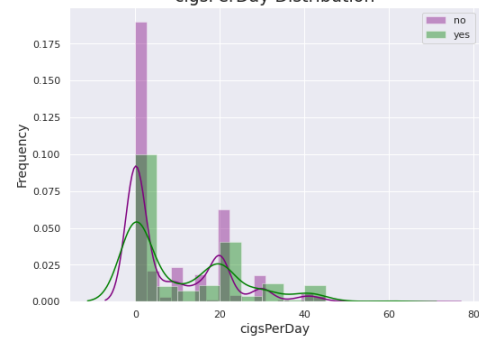
is_smoking distribution



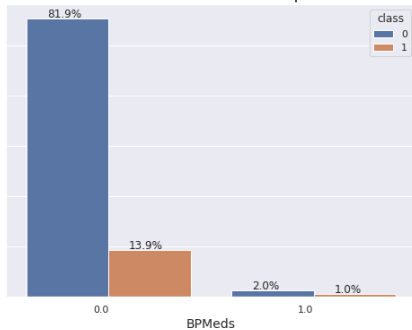
cigsPerDay box_plot



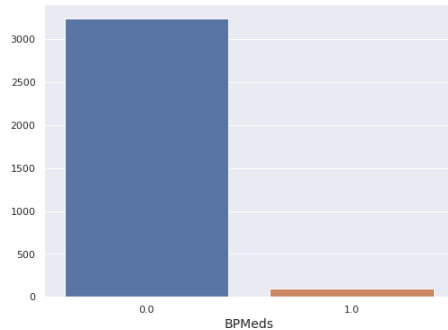
cigsPerDay Distribution



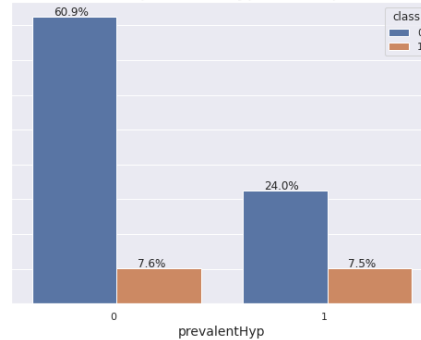
Feature "BPMeds" decomposed



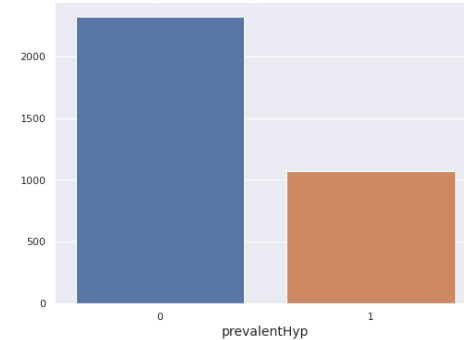
BPMeds distribution



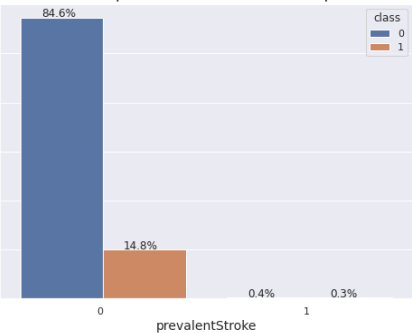
Feature "prevalentHyp" decomposed



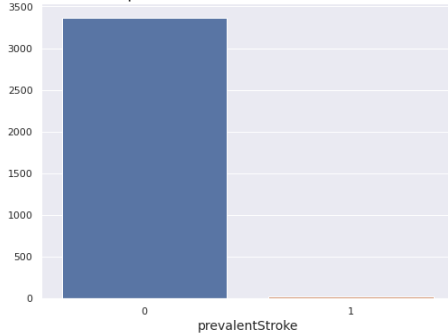
prevalentHyp distribution



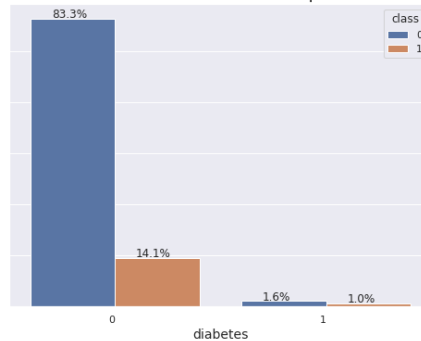
Feature "prevalentStroke" decomposed



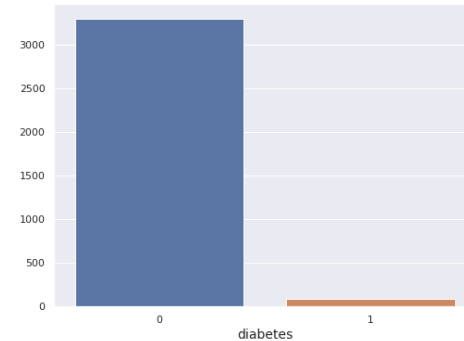
prevalentStroke distribution

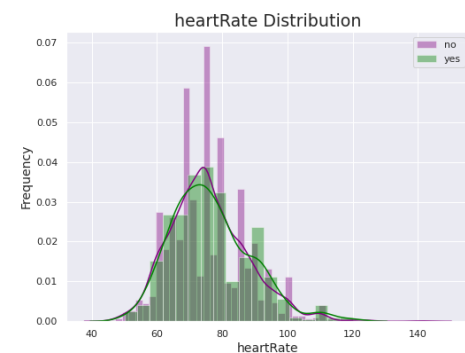
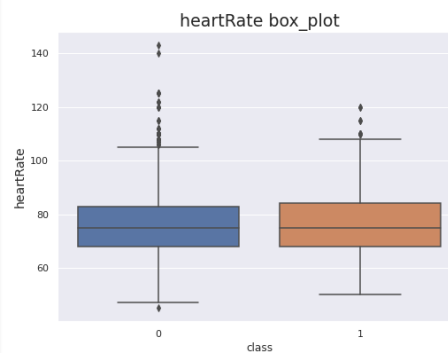
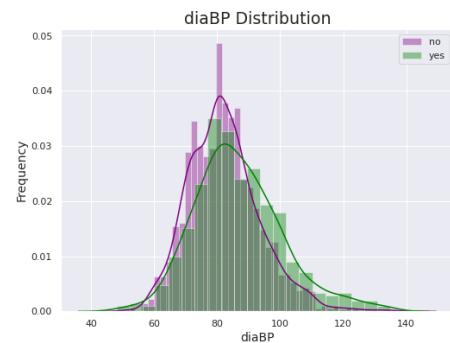
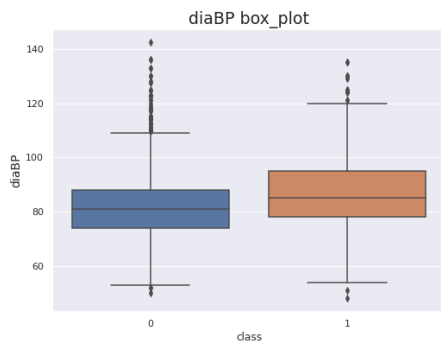
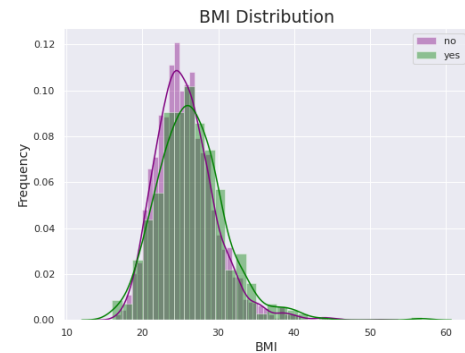
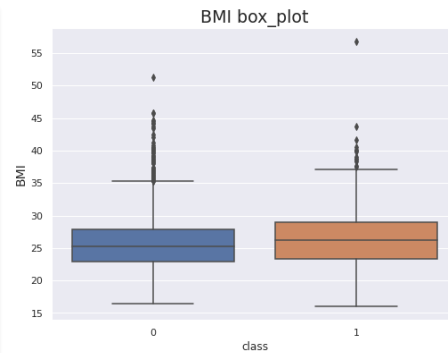
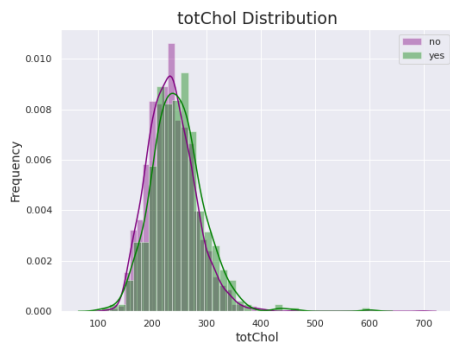
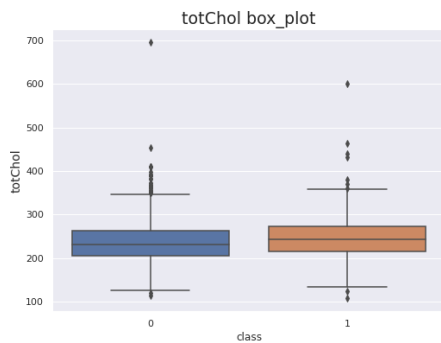


Feature "diabetes" decomposed



diabetes distribution





Feature Engineering

Blood Pressures:

Sys_BP: systolic blood pressure (Continuous)

Dia_BP: diastolic blood pressure (Continuous)

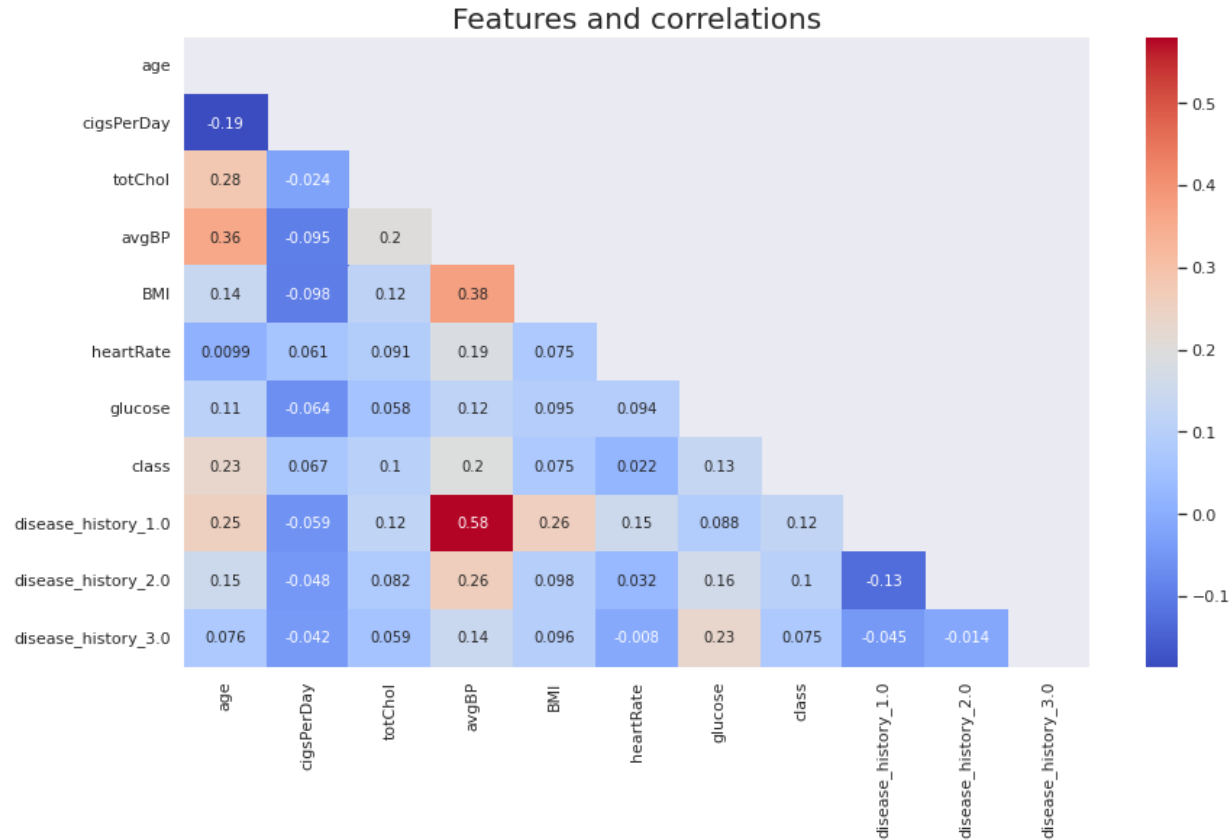
New feature,

avg_BP = average of sys_BP and dia_BP

Historical features:

disease_history = BPMeds + prevalentStroke + prevalentHyp + diabetes

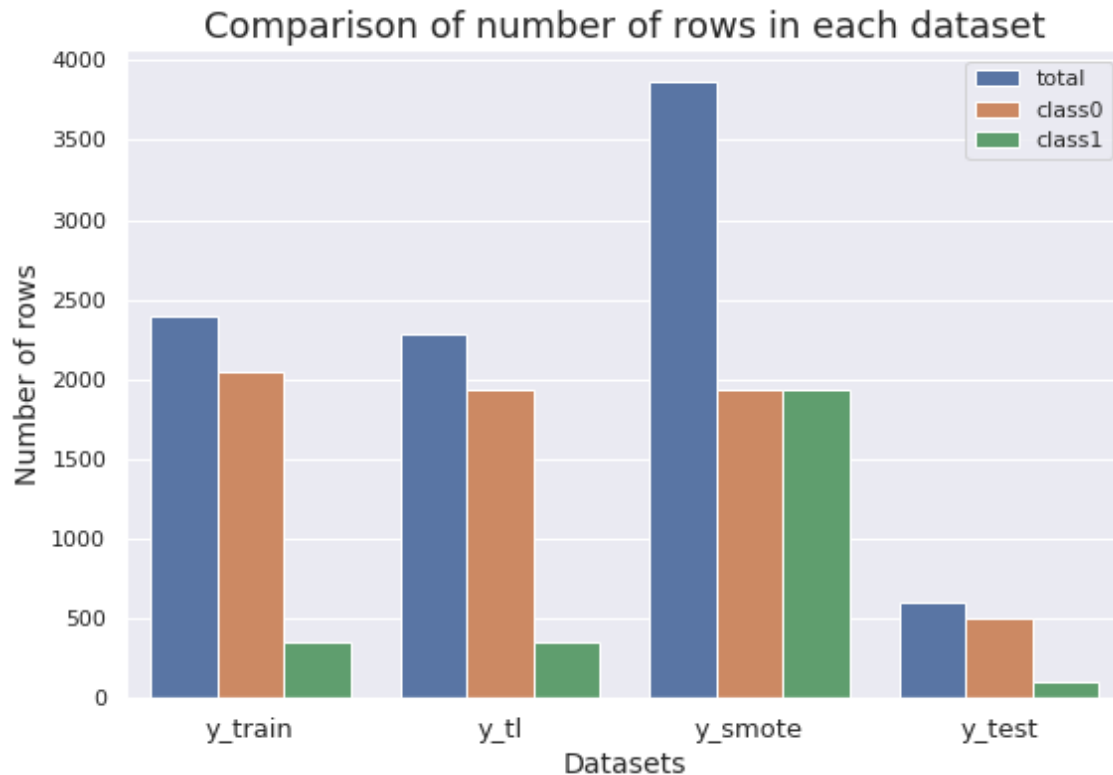
Feature Selection



Method

For feature selection we used Boruta algorithm which is build on top of Random Forest

Train Test datasets



Datasets Explanation

y_train :
Untreated imbalanced
train dataset

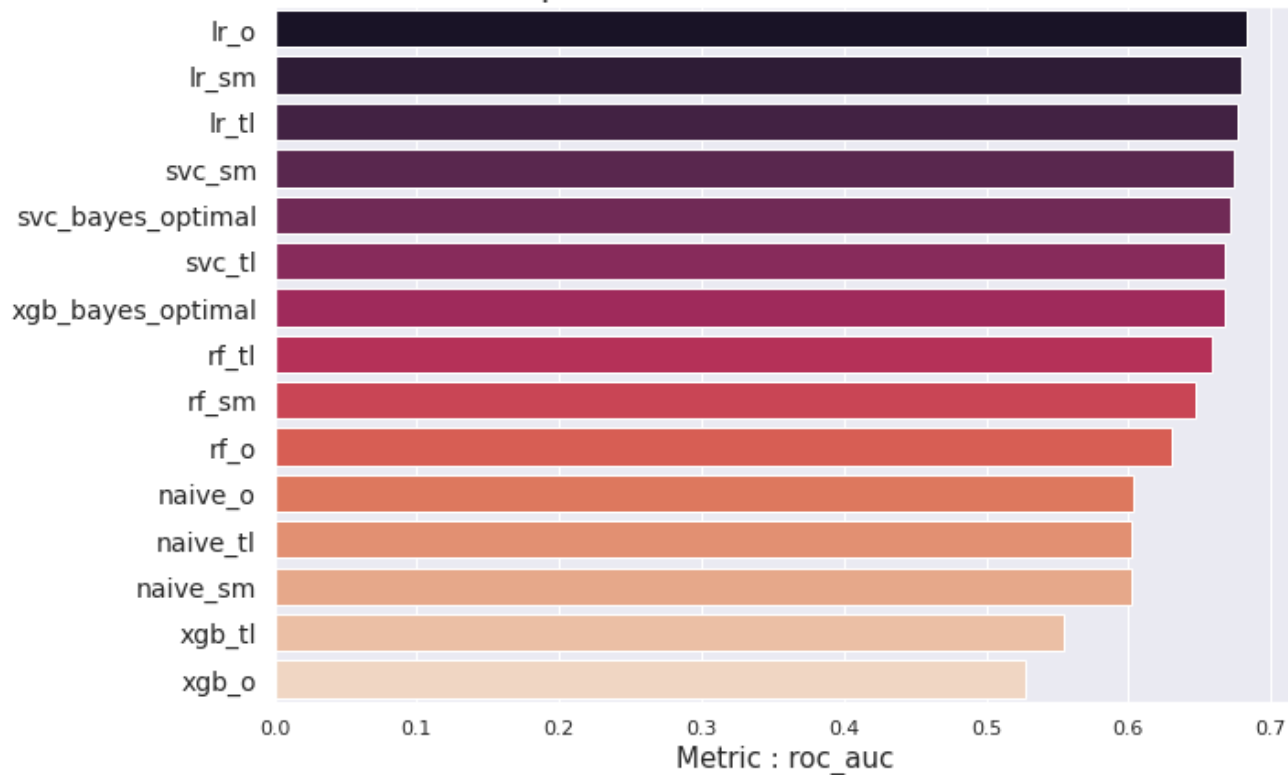
y_tl:
Undersampled with
tomeklins but still
imbalanced train dataset

y_smote :
Oversampled with
SMOTE, so balanced
train dataset

y_test:
Imbalanced test dataset

ML Models' performance

Comparison of classification Models



Abbreviations

Models

lr :

Logistic Regression

svc:

Support vector machine

rf :

Random forest

xgb :

XGBoost

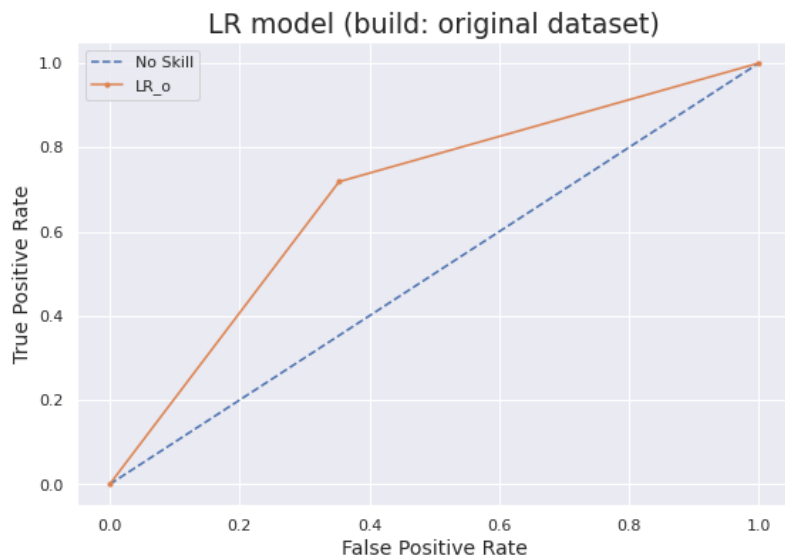
Dataset

o : original

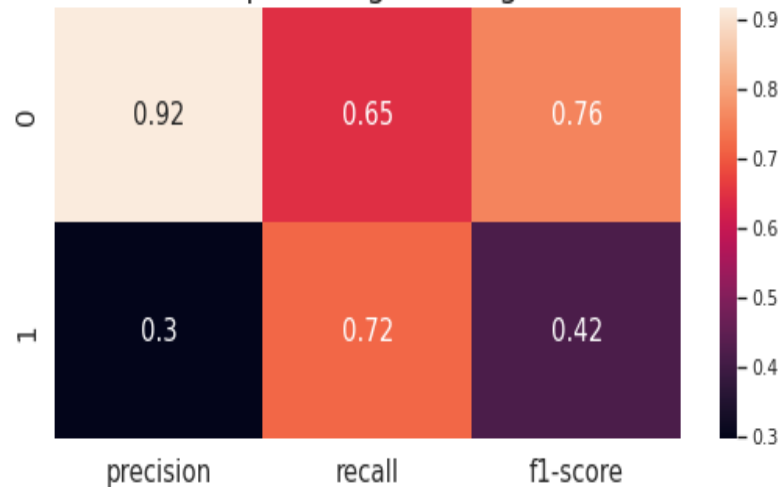
tl : under-sampled

sm : over-sampled

Best model : Logistic Regression model



Classification report: Logistic Regression model

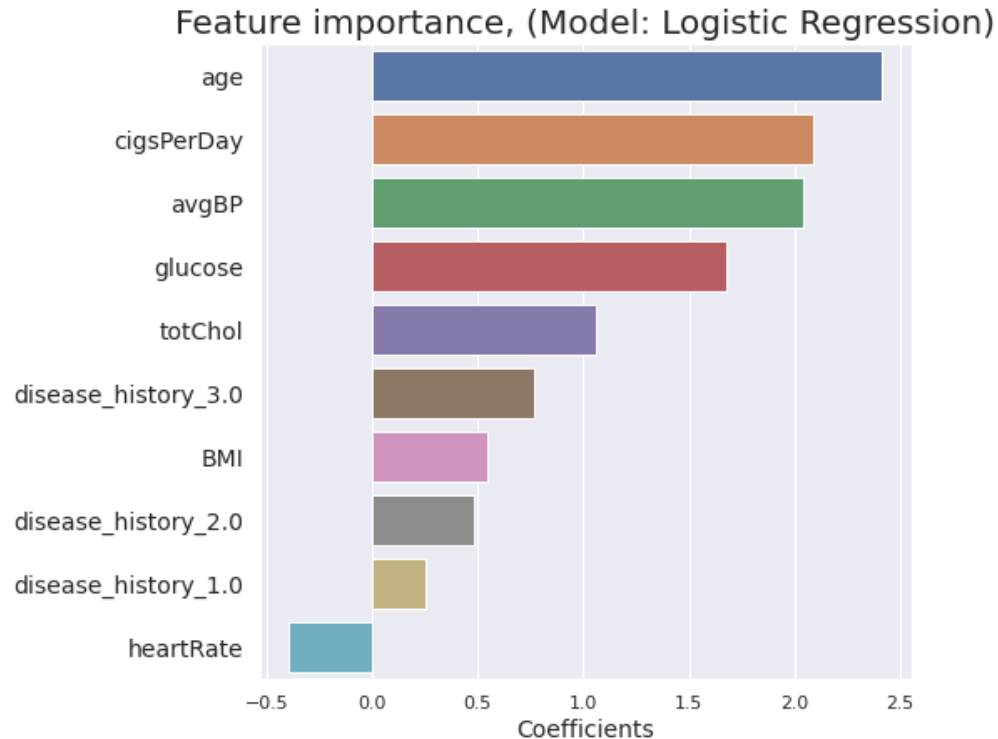


Properties of the model:

The model has the best precision and the best auc_roc score.

The model is really fast to react, takes less than 1 second to give out the results.

Feature Importance



Top 3 features

The top 3 features are +vely correlated with the dependent variable.

The top risk factor is age followed by cigarettes per day and blood pressure.

Conclusions

- Selected features for model building : age, cigsPerDay, totChol, avgBP, BMI, heartRate, glucose, disease_history.
- Three train datasets were used for the analysis.
- Built naive bayes, logistic regression, support vector machine, random forest and XGBoost model using the 3 train datasets.
- On comparing all the models, Logistic regression model trained with original dataset gave best performance.
- We found out that, Logistic Regression model is performing the best with the minimum roc_auc score of 0.68, accuracy of 0.66, recall of 0.69 and precision of 0.263.
- Top 3 features that are helpful in predicting the class in decreasing order of their importance are age, cigs_per_day, avgBP.

Thank you!