

Capstone Project-3

Cardiovascular Risk Prediction

By-Prabir Debnath

Points of Discussion

1. The Problem Statement
2. Concept of ML Model
3. Summary of the Experience Set
4. Exploration and Pre-processing of Data
5. Building and Evaluation of Model-1
6. Building and Evaluation of Model-2
7. Building and Evaluation of Model-3
8. Final Conclusion

The Problem Statement

We are provided with a labeled dataset on the details of patients with or without cardiovascular disease. The task is to explore and analyze the data and to build a classification model for 10 Years Coronary Heart Disease prediction.

Concept of ML Model

The performance of a machine learning model depends on three factors:

i. Quality of Data

(cleaner experiences for better learning)

ii. Quantity of Data

(more experiences for better learning)

iii. Quality of Model

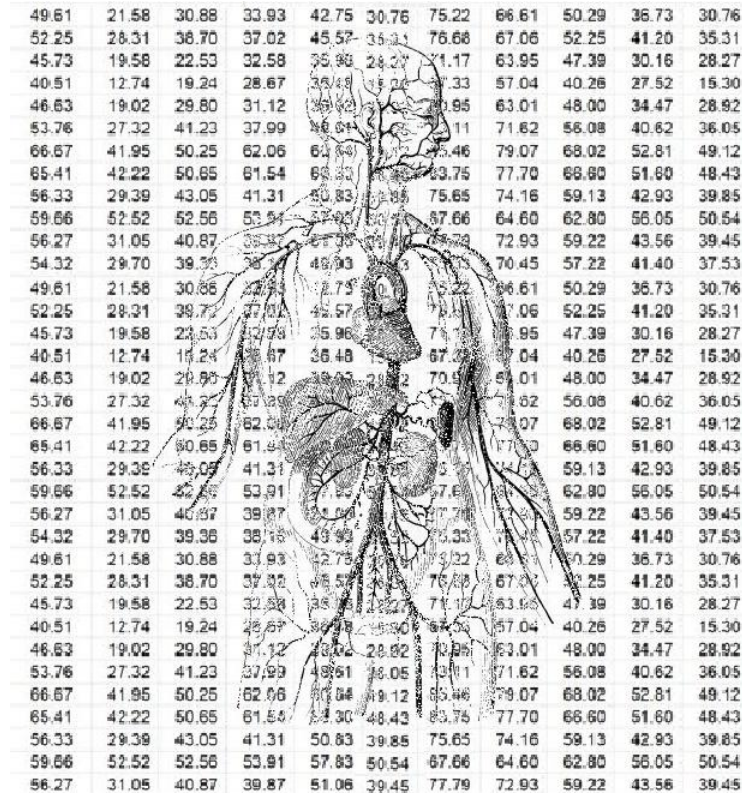
(right model and right hyperparameters for better learning)

Summary of the Experience Set

Here, the dataset has 3,390 rows, which means 3,390 experiences about patients with or without cardiovascular disease and

It has 17 columns, which means each experience is observed along 17 features or dimensions.

Let's Decode the Experiences !



The image displays a detailed anatomical illustration of a human torso, showing the internal organs and skeletal structure. Overlaid on this illustration is a grid of numerical data, organized into 10 rows and 10 columns. The numbers are arranged in a repeating pattern, with some variations in the middle rows. The grid is centered over the torso, with the numbers appearing to be part of a larger data set or a specific analysis related to the human body.

49.61	21.58	30.88	33.93	42.75	30.76	75.22	66.61	50.29	36.73	30.76
52.25	28.31	38.70	37.02	45.57	35.31	76.68	67.06	52.25	41.20	35.31
45.73	19.58	22.53	32.58	35.96	24.27	71.17	63.95	47.39	30.16	28.27
40.51	12.74	19.24	28.67	35.48	15.06	67.33	57.04	40.26	27.52	15.30
46.63	19.02	29.80	31.12	39.62	20.82	70.95	63.01	48.00	34.47	28.92
53.76	27.32	41.23	37.99	46.61	26.05	73.11	71.62	56.08	40.62	36.05
66.67	41.95	50.25	62.06	61.64	35.46	79.07	68.02	52.81	49.12	
65.41	42.22	50.65	61.54	63.33	33.75	77.70	68.80	51.80	48.43	
56.33	29.39	43.05	41.31	50.83	39.85	75.65	74.16	59.13	42.93	39.85
59.66	52.52	52.56	53.91	57.83	50.54	67.66	64.60	62.80	56.05	50.54
56.27	31.05	40.87	39.87	51.06	39.45	77.79	72.93	59.22	43.56	39.45
54.32	29.70	39.36	38.16	49.93	30.33	75.33	70.45	57.22	41.40	37.53
49.61	21.58	30.88	33.93	42.75	30.76	75.22	66.61	50.29	36.73	30.76
52.25	28.31	38.70	37.02	45.57	35.31	76.68	67.06	52.25	41.20	35.31
45.73	19.58	22.53	32.58	35.96	24.27	71.17	63.95	47.39	30.16	28.27
40.51	12.74	19.24	28.67	35.48	15.06	67.33	57.04	40.26	27.52	15.30
46.63	19.02	29.80	31.12	39.62	20.82	70.95	63.01	48.00	34.47	28.92
53.76	27.32	41.23	37.99	46.61	26.05	73.11	71.62	56.08	40.62	36.05
66.67	41.95	50.25	62.06	61.64	35.46	79.07	68.02	52.81	49.12	
65.41	42.22	50.65	61.54	63.33	33.75	77.70	68.80	51.80	48.43	
56.33	29.39	43.05	41.31	50.83	39.85	75.65	74.16	59.13	42.93	39.85
59.66	52.52	52.56	53.91	57.83	50.54	67.66	64.60	62.80	56.05	50.54
56.27	31.05	40.87	39.87	51.06	39.45	77.79	72.93	59.22	43.56	39.45
54.32	29.70	39.36	38.16	49.93	30.33	75.33	70.45	57.22	41.40	37.53
49.61	21.58	30.88	33.93	42.75	30.76	75.22	66.61	50.29	36.73	30.76
52.25	28.31	38.70	37.02	45.57	35.31	76.68	67.06	52.25	41.20	35.31
45.73	19.58	22.53	32.58	35.96	24.27	71.17	63.95	47.39	30.16	28.27
40.51	12.74	19.24	28.67	35.48	15.06	67.33	57.04	40.26	27.52	15.30
46.63	19.02	29.80	31.12	39.62	20.82	70.95	63.01	48.00	34.47	28.92
53.76	27.32	41.23	37.99	46.61	26.05	73.11	71.62	56.08	40.62	36.05
66.67	41.95	50.25	62.06	61.64	35.46	79.07	68.02	52.81	49.12	
65.41	42.22	50.65	61.54	63.33	33.75	77.70	68.80	51.80	48.43	
56.33	29.39	43.05	41.31	50.83	39.85	75.65	74.16	59.13	42.93	39.85
59.66	52.52	52.56	53.91	57.83	50.54	67.66	64.60	62.80	56.05	50.54
56.27	31.05	40.87	39.87	51.06	39.45	77.79	72.93	59.22	43.56	39.45

Exploration and Pre-processing of Data

We have done the exploration and pre-processing in seven steps to transform raw data into quality data for our ml model.

1. Connection with the Data
2. First Feelings of the Data
3. Deeper Understanding of the Data
4. Cleaning the Data
5. Treating Anomalies in the Data
6. Final Feature Selection from the Data
7. Preparation of Input and Output Data

Cleaning the Data

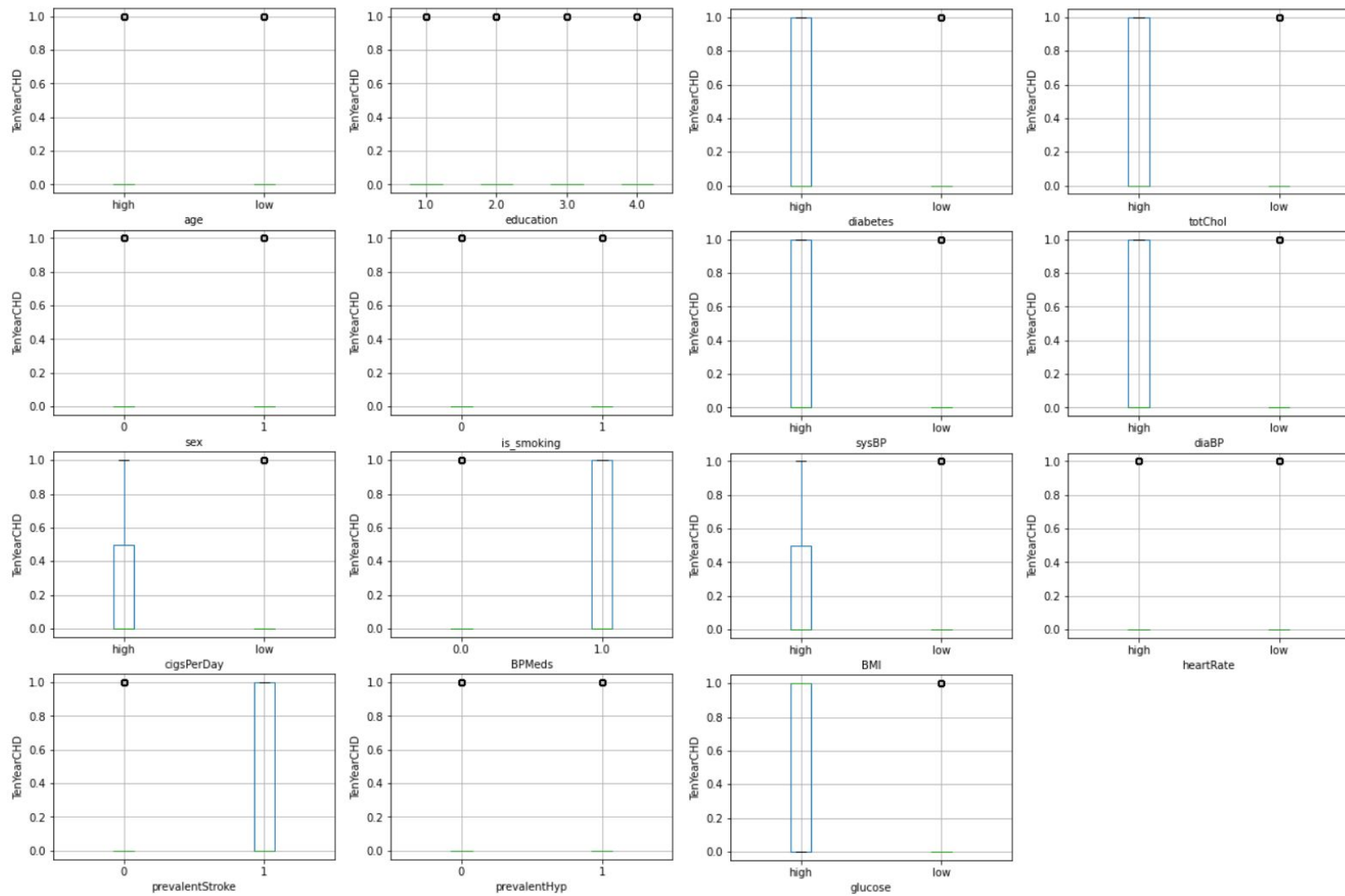
- > We have handled all null values in 'cigsPerDay', 'totChol', 'BMI', 'heartRate', 'glucose', 'education', and 'BPMeds' columns with imputation. Thus there is no loss of data.
- > We have encoded 'sex' column with two categories: F: 0, M: 1
- > We have encoded 'is_smoking' column with two categories: NO:0, YES:1

Treating Anomalies in the Data

In our dataset, for most of the features class 1 targets are outliers. Thus we need more experience with class 1 targets to bring a balance in prediction.

In simple words, with the available experience set, our model will be expert in predicting the features for which there will be no heart disease.

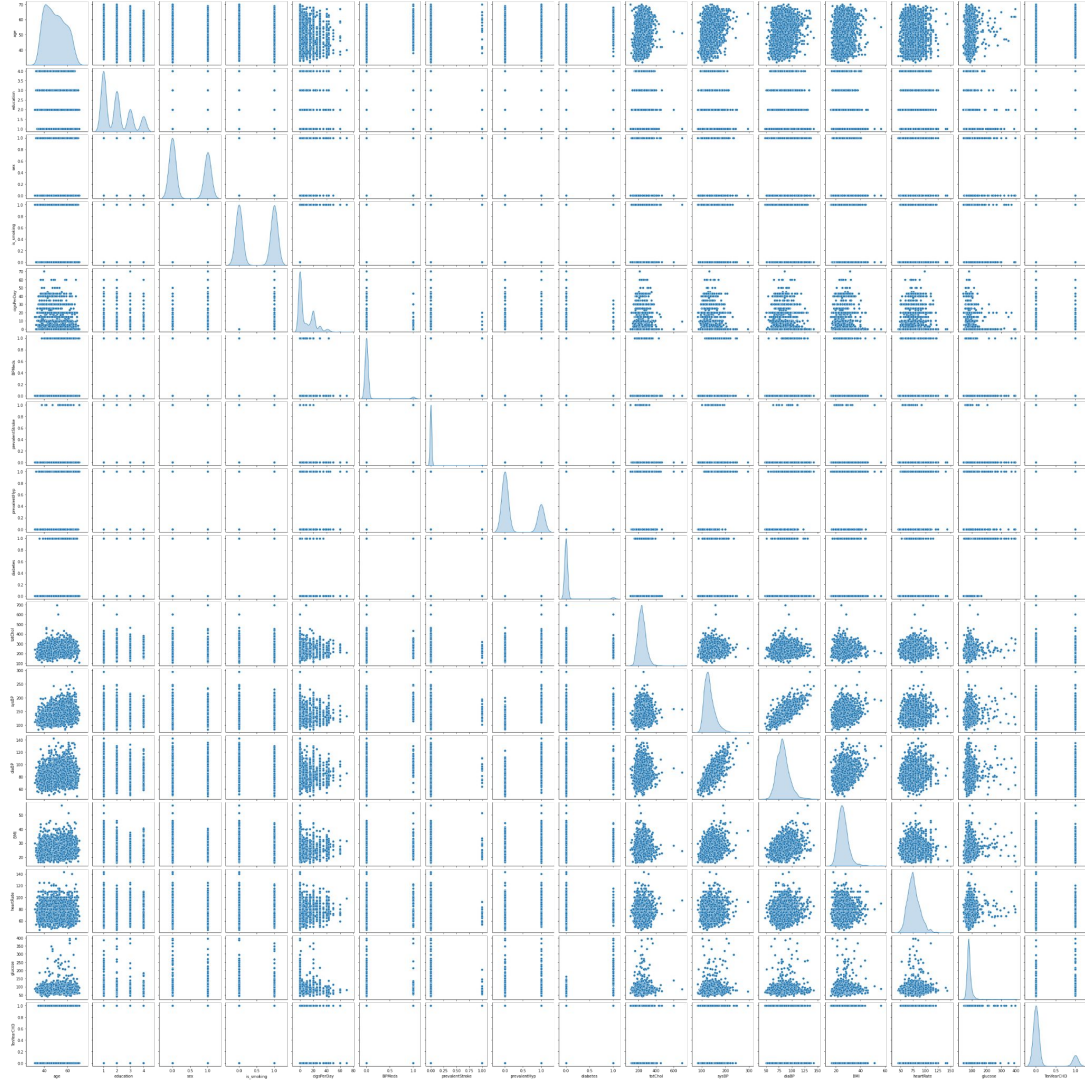
Let's check the boxplot!



Overall Feature Understanding

Here, the distribution of 'cigsPerDay', 'totChol', 'sysBP', 'diaBP', 'BMI', 'heartRate' and 'glucose' are positively skewed. Thus we have done log transformation on these features to normalize their distribution.

Let's check the pairplot!



Looking for Truly Independent Features

We have removed 'sysBP', 'diaBP', 'BMI', 'Age', 'heartRate', 'totChol', 'is_smoking' and 'Diabetes' in sequence from our dataset to bring all the VIF values below 10. Thus all our input variables became truly independent.

Let's check the heatmap!

Correlation between Variables



	variables	VIF
0	education	3.655103
1	sex	1.957236
2	cigsPerDay	1.733063
3	BPMeds	1.117896
4	prevalentStroke	1.024059
5	prevalentHyp	1.561774
6	glucose	4.497470

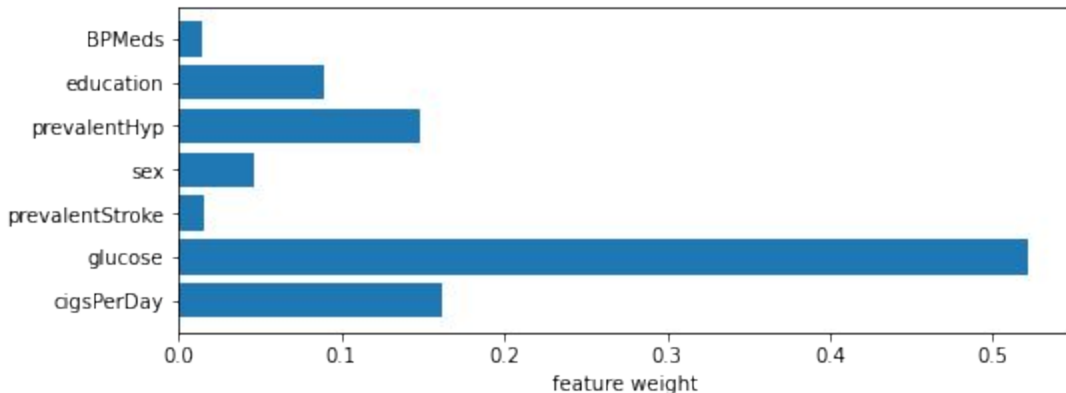
Let's Train the Models !



Building and Evaluation of Model-1

Final Random Forest Model:

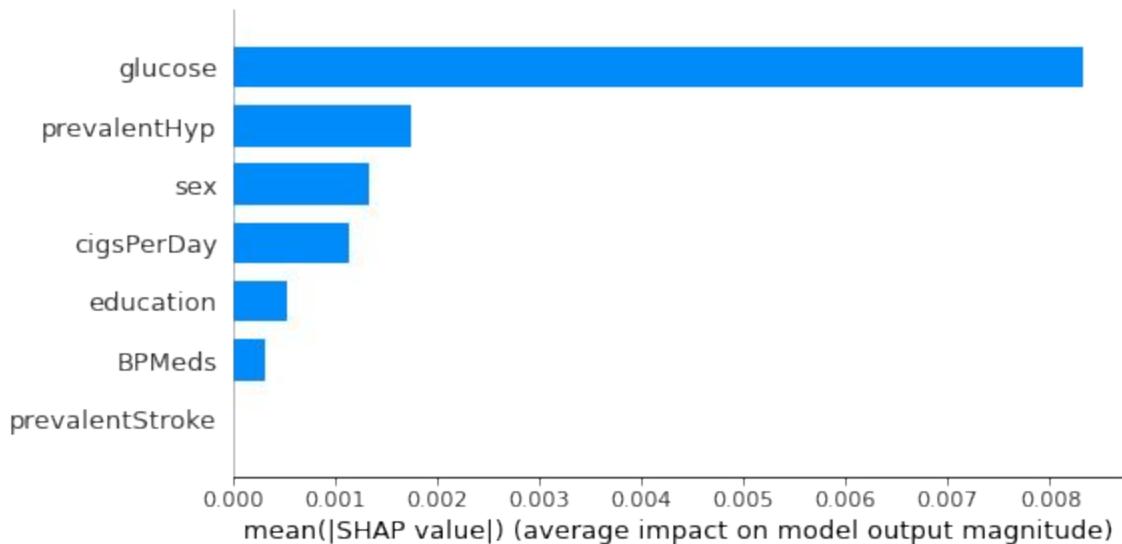
After cross validation and hyperparameter tuning, the best parameters are {'max_depth': 25, 'max_features': 'auto', 'min_samples_leaf': 5, 'min_samples_split': 15, 'n_estimators': 100} (test accuracy is 84% and variance in prediction is 2%)



Building and Evaluation of Model-2

Final KNN Model:

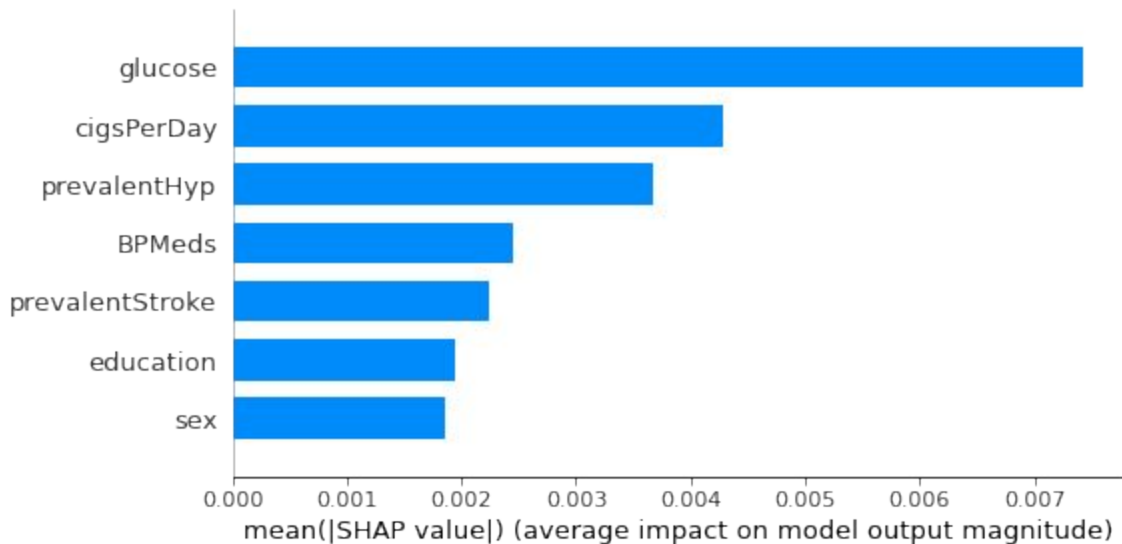
After cross validation and hyperparameter tuning, the best parameters are {'leaf_size': 30, 'n_neighbors': 19} (test accuracy is 84% and variance in prediction is 1%)



Building and Evaluation of Model-3

Final SVC Model:

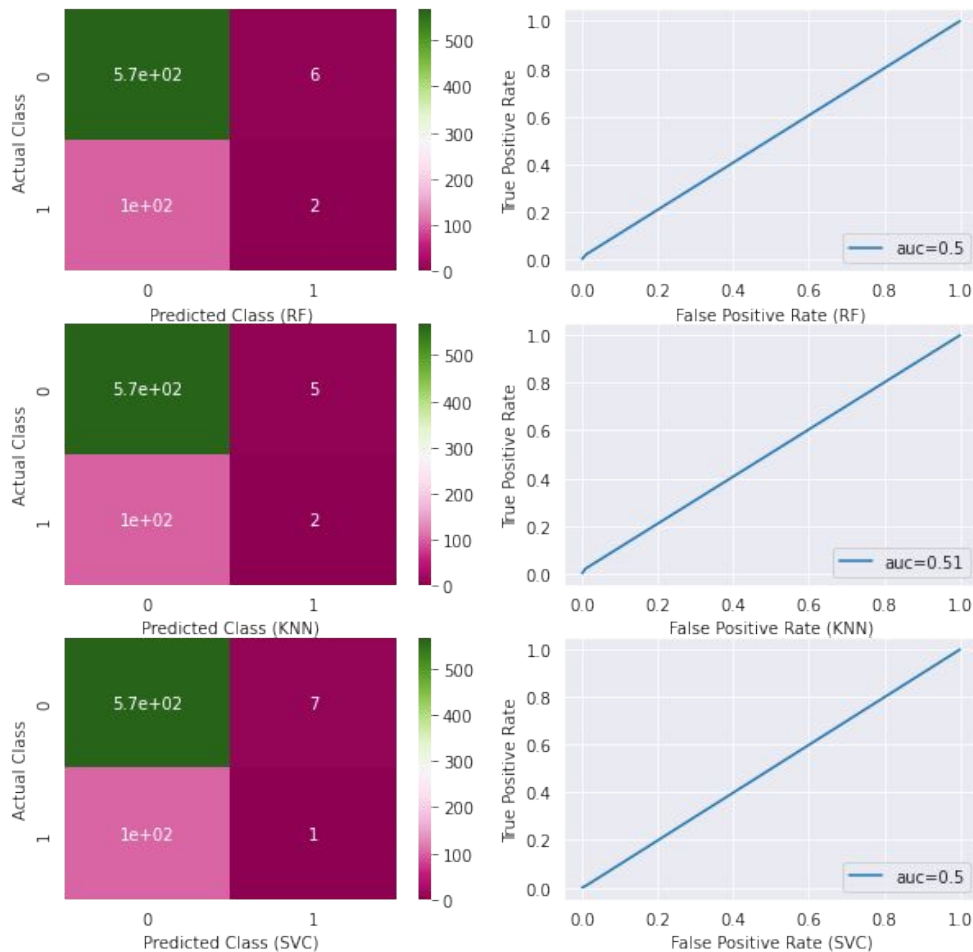
After cross validation and hyperparameter tuning, the best parameters are {'C': 6, 'gamma': 0.1} (test accuracy is 84% and variance in prediction is 2%)



RF- Train_Accuracy: 0.86, Test_Accuracy: 0.84, Test_F1: 0.47

KNN- Train_Accuracy: 0.85, Test_Accuracy: 0.84, Test_F1: 0.48

SVC- Train_Accuracy: 0.86, Test_Accuracy: 0.84, Test_F1: 0.46



Final Conclusion

On the basis of the performance study of our three models, we are selecting **KNN** classifier (*the best warrior*) for predicting 10 Years Coronary heart disease, as it has low variance in prediction, good f1_score and good ROC_AUC score among all three models



Thank you !