

# Machine Learning Project

Name: Jhansi Bhaskarla

Name: Surendra Pothuri

## Table of Contents

1) Import Packages

2) EDA

3) Preparing ML models

4) Models evaluation

5) Ensembling

6) Conclusion

## Packages Required

```
In [4]: 1 import pandas as pd
        2 import numpy as np
        3 import matplotlib.pyplot as plt
        4 import seaborn as sns
        5
        6
        7 sns.set(style="whitegrid")
```

In [28]:

```

1 from sklearn.preprocessing import StandardScaler
2 from sklearn.model_selection import train_test_split
3 from sklearn.metrics import confusion_matrix, accuracy_score, roc_auc_score
4 from sklearn.linear_model import LogisticRegression
5 from sklearn.naive_bayes import GaussianNB
6 from sklearn.naive_bayes import GaussianNB
7 from xgboost import XGBClassifier
8 from sklearn.ensemble import RandomForestClassifier
9 from sklearn.tree import DecisionTreeClassifier
10 from sklearn.neighbors import KNeighborsClassifier
11 from sklearn.naive_bayes import GaussianNB
12 from sklearn.svm import SVC
13 from mlxtend.classifier import StackingCVClassifier
14 from collections import Counter

```

In [8]:

```
1 df = pd.read_csv('heart.csv')
```

In [14]:

```
1 df.head()
```

Out[14]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
0	52	1	0	125	212	0	1	168	0	1.0	2	2	3	0
1	53	1	0	140	203	1	0	155	1	3.1	0	0	3	0
2	70	1	0	145	174	0	1	125	1	2.6	0	0	3	0
3	61	1	0	148	203	0	1	161	0	0.0	2	1	3	0
4	62	0	0	138	294	1	1	106	0	1.9	1	3	2	0

In [16]:

```
1 df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1025 entries, 0 to 1024
Data columns (total 14 columns):
 #   Column      Non-Null Count  Dtype
---  -
 0   age         1025 non-null   int64
 1   sex         1025 non-null   int64
 2   cp          1025 non-null   int64
 3   trestbps    1025 non-null   int64
 4   chol        1025 non-null   int64
 5   fbs         1025 non-null   int64
 6   restecg     1025 non-null   int64
 7   thalach     1025 non-null   int64
 8   exang       1025 non-null   int64
 9   oldpeak     1025 non-null   float64
10  slope       1025 non-null   int64
11  ca          1025 non-null   int64
12  thal        1025 non-null   int64
13  target      1025 non-null   int64
dtypes: float64(1), int64(13)
memory usage: 112.2 KB
```

In [9]:

```
1 print(df.isnull().sum())
2

age         0
sex         0
cp          0
trestbps    0
chol        0
fbs         0
restecg     0
thalach     0
exang       0
oldpeak     0
slope       0
ca          0
thal        0
target      0
dtype: int64
```

In [10]:

```
1 print(df.describe())
2
```

	age	sex	cp	trestbps	chol
\					
count	1025.000000	1025.000000	1025.000000	1025.000000	1025.000000
mean	54.434146	0.695610	0.942439	131.611707	246.000000
std	9.072290	0.460373	1.029641	17.516718	51.59251
min	29.000000	0.000000	0.000000	94.000000	126.000000
25%	48.000000	0.000000	0.000000	120.000000	211.000000
50%	56.000000	1.000000	1.000000	130.000000	240.000000
75%	61.000000	1.000000	2.000000	140.000000	275.000000
max	77.000000	1.000000	3.000000	200.000000	564.000000

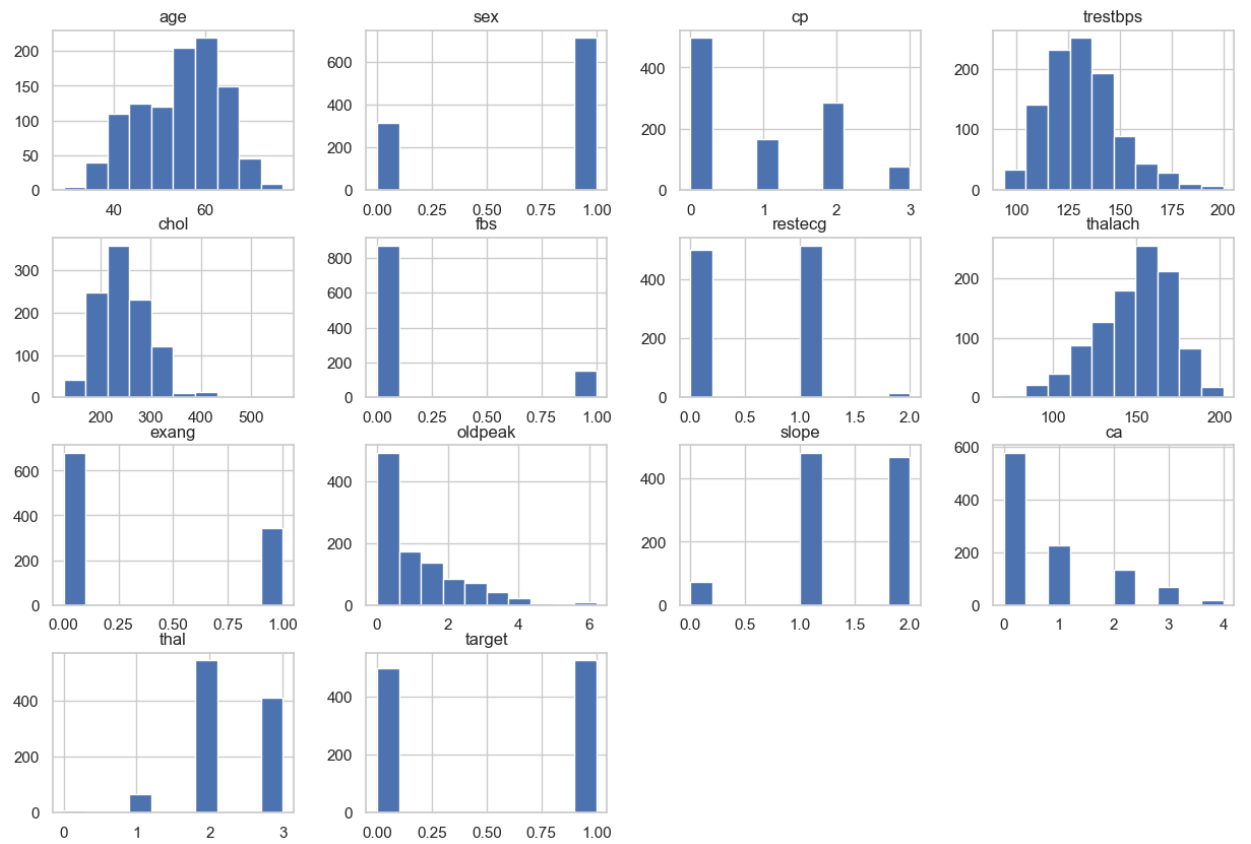
	fbs	restecg	thalach	exang	oldpea
k \					
count	1025.000000	1025.000000	1025.000000	1025.000000	1025.000000
0					
mean	0.149268	0.529756	149.114146	0.336585	1.07151
2					
std	0.356527	0.527878	23.005724	0.472772	1.17505
3					
min	0.000000	0.000000	71.000000	0.000000	0.000000
0					
25%	0.000000	0.000000	132.000000	0.000000	0.000000
0					
50%	0.000000	1.000000	152.000000	0.000000	0.800000
0					
75%	0.000000	1.000000	166.000000	1.000000	1.800000
0					
max	1.000000	2.000000	202.000000	1.000000	6.200000
0					

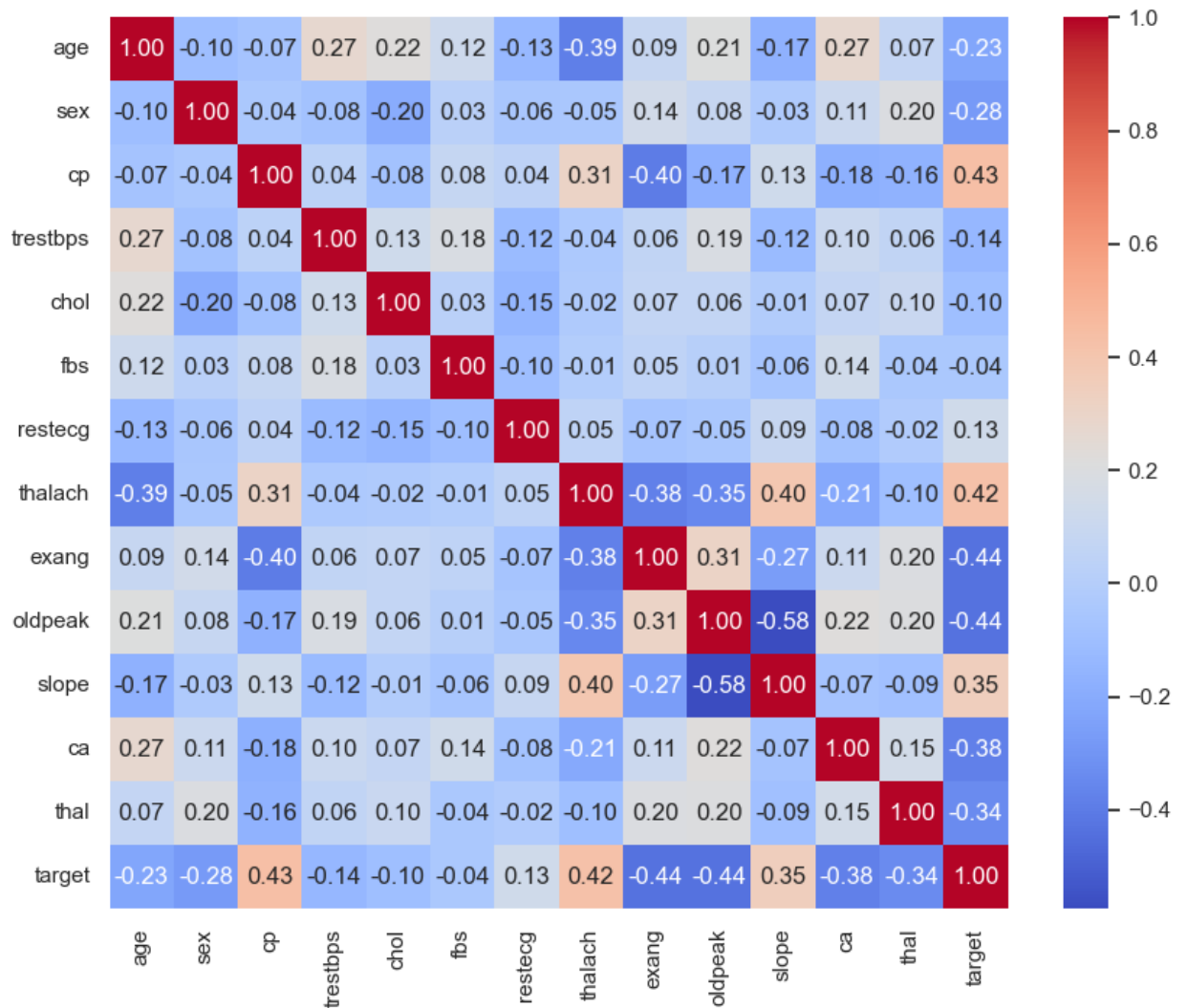
	slope	ca	thal	target
count	1025.000000	1025.000000	1025.000000	1025.000000
mean	1.385366	0.754146	2.323902	0.513171
std	0.617755	1.030798	0.620660	0.500070
min	0.000000	0.000000	0.000000	0.000000
25%	1.000000	0.000000	2.000000	0.000000
50%	1.000000	0.000000	2.000000	1.000000
75%	2.000000	1.000000	3.000000	1.000000
max	2.000000	4.000000	3.000000	1.000000

## EDA

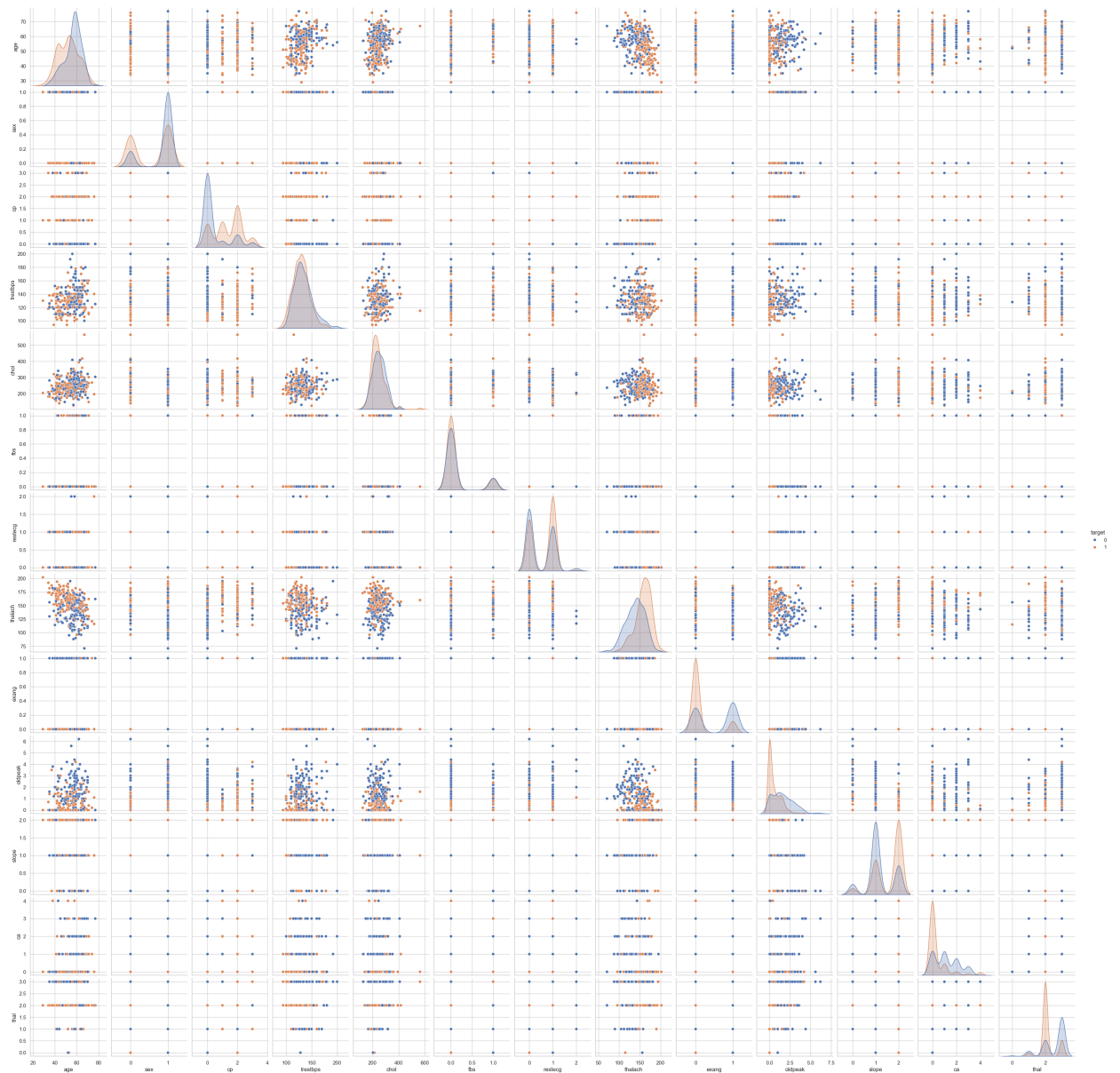
```
In [11]: 1 df.hist(figsize=(15, 10))  
         2 plt.show()  
         3
```



```
In [12]: 1 plt.figure(figsize=(10, 8))
2          sns.heatmap(df.corr(), annot=True, fmt=".2f", cmap="coolwarm")
3          plt.show()
4
```



```
In [13]: 1 sns.pairplot(df, hue='target')
          2 plt.show()
          3
```



## Model prepration

```
In [22]: 1 y = df["target"]
          2 X = df.drop('target',axis=1)
          3 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
```

```
In [29]: 1 print(y_test.unique())  
         2 Counter(y_train)
```

[1 0]

Out[29]: Counter({1: 419, 0: 401})

```
In [30]: 1 scaler = StandardScaler()  
         2 X_train = scaler.fit_transform(X_train)  
         3 X_test = scaler.transform(X_test)
```

## ML models

**Logistic Regression**

**Naive Bayes**

**Random Forest Classifier**

**Extreme Gradient Boost**

**K-Nearest Neighbour**

**Decision Tree**

**Support Vector Machine**



```
In [31]: 1 m1 = 'Logistic Regression'
2 lr = LogisticRegression()
3 model = lr.fit(X_train, y_train)
4 lr_predict = lr.predict(X_test)
5 lr_conf_matrix = confusion_matrix(y_test, lr_predict)
6 lr_acc_score = accuracy_score(y_test, lr_predict)
7 print("confussion matrix")
8 print(lr_conf_matrix)
9 print("\n")
10 print("Accuracy of Logistic Regression:",lr_acc_score*100,'\n')
11 print(classification_report(y_test,lr_predict))
```

confussion matrix

```
[[ 77  21]
 [  7 100]]
```

Accuracy of Logistic Regression: 86.34146341463415

	precision	recall	f1-score	support
0	0.92	0.79	0.85	98
1	0.83	0.93	0.88	107
accuracy			0.86	205
macro avg	0.87	0.86	0.86	205
weighted avg	0.87	0.86	0.86	205

```
In [32]: 1 m2 = 'Naive Bayes'
2 nb = GaussianNB()
3 nb.fit(X_train,y_train)
4 nbpred = nb.predict(X_test)
5 nb_conf_matrix = confusion_matrix(y_test, nbpred)
6 nb_acc_score = accuracy_score(y_test, nbpred)
7 print("confussion matrix")
8 print(nb_conf_matrix)
9 print("\n")
10 print("Accuracy of Naive Bayes model:",nb_acc_score*100,'\n')
11 print(classification_report(y_test,nbpred))
```

```
confussion matrix
[[79 19]
 [11 96]]
```

Accuracy of Naive Bayes model: 85.36585365853658

	precision	recall	f1-score	support
0	0.88	0.81	0.84	98
1	0.83	0.90	0.86	107
accuracy			0.85	205
macro avg	0.86	0.85	0.85	205
weighted avg	0.86	0.85	0.85	205

```
In [33]: 1 m3 = 'Random Forest Classifier'
2 rf = RandomForestClassifier(n_estimators=20, random_state=12,max_c
3 rf.fit(X_train,y_train)
4 rf_predicted = rf.predict(X_test)
5 rf_conf_matrix = confusion_matrix(y_test, rf_predicted)
6 rf_acc_score = accuracy_score(y_test, rf_predicted)
7 print("confussion matrix")
8 print(rf_conf_matrix)
9 print("\n")
10 print("Accuracy of Random Forest:",rf_acc_score*100,'\n')
11 print(classification_report(y_test,rf_predicted))
```

confussion matrix

```
[[ 88  10]
 [  3 104]]
```

Accuracy of Random Forest: 93.65853658536587

	precision	recall	f1-score	support
0	0.97	0.90	0.93	98
1	0.91	0.97	0.94	107
accuracy			0.94	205
macro avg	0.94	0.93	0.94	205
weighted avg	0.94	0.94	0.94	205

```
In [34]: 1 m4 = 'Extreme Gradient Boost'
2 xgb = XGBClassifier(learning_rate=0.01, n_estimators=25, max_depth
3                  reg_lambda=2, booster='dart', colsample_byleve
4 xgb.fit(X_train, y_train)
5 xgb_predicted = xgb.predict(X_test)
6 xgb_conf_matrix = confusion_matrix(y_test, xgb_predicted)
7 xgb_acc_score = accuracy_score(y_test, xgb_predicted)
8 print("confussion matrix")
9 print(xgb_conf_matrix)
10 print("\n")
11 print("Accuracy of Extreme Gradient Boost:",xgb_acc_score*100,'\n')
12 print(classification_report(y_test,xgb_predicted))
```

```
confussion matrix
[[ 84  14]
 [  2 105]]
```

Accuracy of Extreme Gradient Boost: 92.19512195121952

	precision	recall	f1-score	support
0	0.98	0.86	0.91	98
1	0.88	0.98	0.93	107
accuracy			0.92	205
macro avg	0.93	0.92	0.92	205
weighted avg	0.93	0.92	0.92	205

```
In [35]: 1 m5 = 'K-NeighborsClassifier'
2 knn = KNeighborsClassifier(n_neighbors=10)
3 knn.fit(X_train, y_train)
4 knn_predicted = knn.predict(X_test)
5 knn_conf_matrix = confusion_matrix(y_test, knn_predicted)
6 knn_acc_score = accuracy_score(y_test, knn_predicted)
7 print("confussion matrix")
8 print(knn_conf_matrix)
9 print("\n")
10 print("Accuracy of K-NeighborsClassifier:", knn_acc_score*100, '\n')
11 print(classification_report(y_test, knn_predicted))
```

```
confussion matrix
[[84 14]
 [11 96]]
```

Accuracy of K-NeighborsClassifier: 87.8048780487805

	precision	recall	f1-score	support
0	0.88	0.86	0.87	98
1	0.87	0.90	0.88	107
accuracy			0.88	205
macro avg	0.88	0.88	0.88	205
weighted avg	0.88	0.88	0.88	205

```
In [36]: 1 m6 = 'DecisionTreeClassifier'
2 dt = DecisionTreeClassifier(criterion = 'entropy', random_state=0, n
3 dt.fit(X_train, y_train)
4 dt_predicted = dt.predict(X_test)
5 dt_conf_matrix = confusion_matrix(y_test, dt_predicted)
6 dt_acc_score = accuracy_score(y_test, dt_predicted)
7 print("confussion matrix")
8 print(dt_conf_matrix)
9 print("\n")
10 print("Accuracy of DecisionTreeClassifier:", dt_acc_score*100, '\n')
11 print(classification_report(y_test, dt_predicted))
```

confussion matrix

```
[[95  3]
 [ 8 99]]
```

Accuracy of DecisionTreeClassifier: 94.6341463414634

	precision	recall	f1-score	support
0	0.92	0.97	0.95	98
1	0.97	0.93	0.95	107
accuracy			0.95	205
macro avg	0.95	0.95	0.95	205
weighted avg	0.95	0.95	0.95	205

```
In [37]: 1 m7 = 'Support Vector Classifier'
2 svc = SVC(kernel='rbf', C=2)
3 svc.fit(X_train, y_train)
4 svc_predicted = svc.predict(X_test)
5 svc_conf_matrix = confusion_matrix(y_test, svc_predicted)
6 svc_acc_score = accuracy_score(y_test, svc_predicted)
7 print("confussion matrix")
8 print(svc_conf_matrix)
9 print("\n")
10 print("Accuracy of Support Vector Classifier:", svc_acc_score*100, '%')
11 print(classification_report(y_test, svc_predicted))
```

confussion matrix

```
[[ 94   4]
 [  0 107]]
```

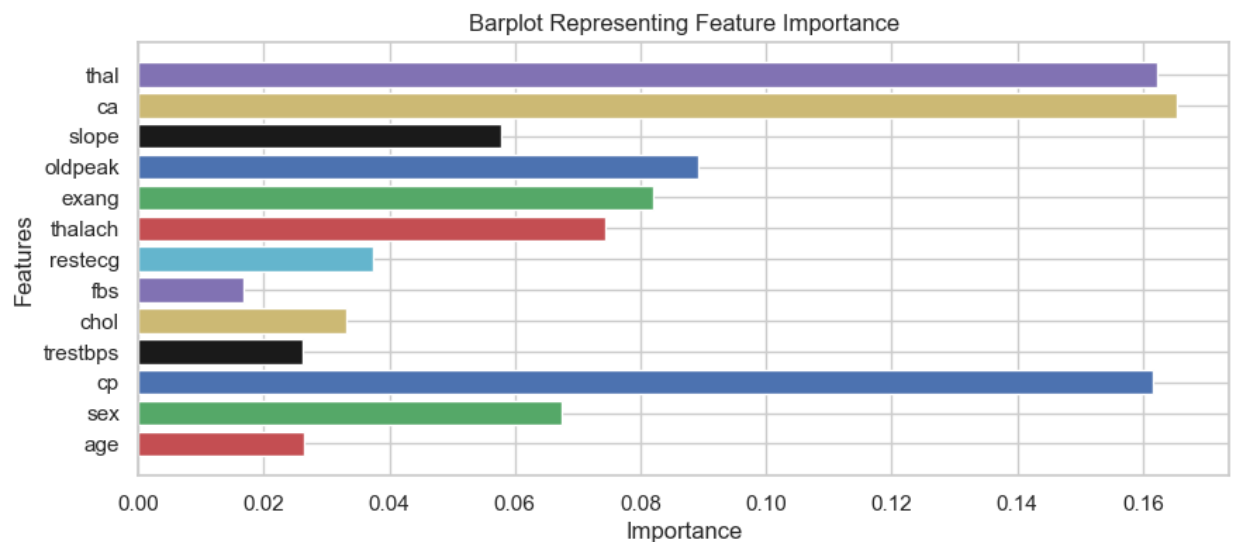
Accuracy of Support Vector Classifier: 98.04878048780488

	precision	recall	f1-score	support
0	1.00	0.96	0.98	98
1	0.96	1.00	0.98	107
accuracy			0.98	205
macro avg	0.98	0.98	0.98	205
weighted avg	0.98	0.98	0.98	205

```

In [39]: 1 imp_feature = pd.DataFrame({
2         'Feature': ['age', 'sex', 'cp', 'trestbps', 'chol', 'fbs', 're
3         'Importance': xgb.feature_importances_
4     })
5
6     # Plotting
7     plt.figure(figsize=(10, 4))
8     plt.title("Barplot Representing Feature Importance")
9     plt.xlabel("Importance")
10    plt.ylabel("Features")
11    plt.barh(imp_feature['Feature'], imp_feature['Importance'], color=
12    plt.show()

```

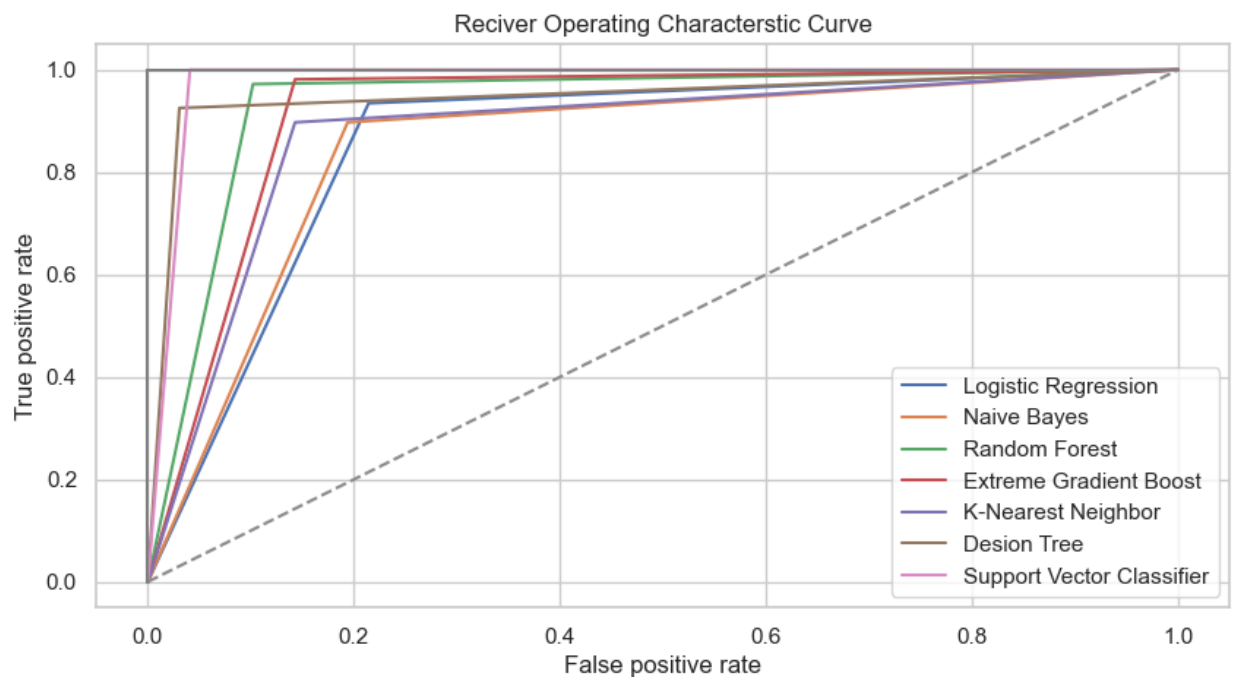




```

In [40]: 1 lr_false_positive_rate,lr_true_positive_rate,lr_threshold = roc_curve(
2 nb_false_positive_rate,nb_true_positive_rate,nb_threshold = roc_curve(
3 rf_false_positive_rate,rf_true_positive_rate,rf_threshold = roc_curve(
4 xgb_false_positive_rate,xgb_true_positive_rate,xgb_threshold = roc_curve(
5 knn_false_positive_rate,knn_true_positive_rate,knn_threshold = roc_curve(
6 dt_false_positive_rate,dt_true_positive_rate,dt_threshold = roc_curve(
7 svc_false_positive_rate,svc_true_positive_rate,svc_threshold = roc_curve(
8
9
10 sns.set_style('whitegrid')
11 plt.figure(figsize=(10,5))
12 plt.title('Reciver Operating Characterstic Curve')
13 plt.plot(lr_false_positive_rate,lr_true_positive_rate,label='Logistic Regression')
14 plt.plot(nb_false_positive_rate,nb_true_positive_rate,label='Naive Bayes')
15 plt.plot(rf_false_positive_rate,rf_true_positive_rate,label='Random Forest')
16 plt.plot(xgb_false_positive_rate,xgb_true_positive_rate,label='Extreme Gradient Boost')
17 plt.plot(knn_false_positive_rate,knn_true_positive_rate,label='K-Nearest Neighbor')
18 plt.plot(dt_false_positive_rate,dt_true_positive_rate,label='Decision Tree')
19 plt.plot(svc_false_positive_rate,svc_true_positive_rate,label='Support Vector Classifier')
20 plt.plot([0,1],ls='--')
21 plt.plot([0,0],[1,0],c='.5')
22 plt.plot([1,1],c='.5')
23 plt.ylabel('True positive rate')
24 plt.xlabel('False positive rate')
25 plt.legend()
26 plt.show()

```



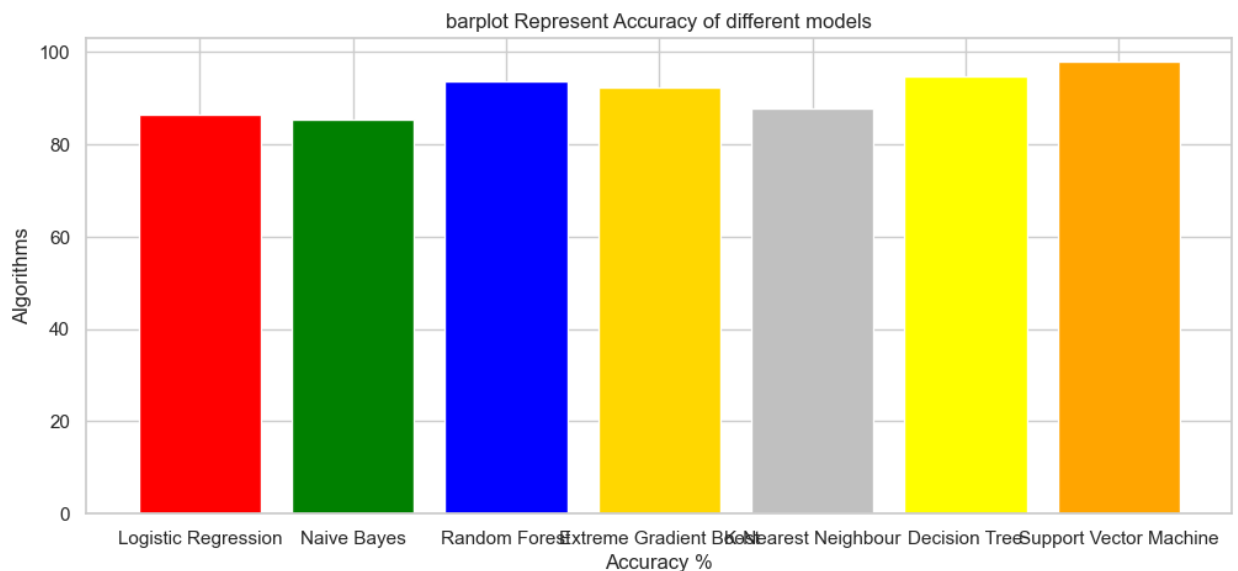
## Model Evaluation

```
In [41]: 1 model_ev = pd.DataFrame({'Model': ['Logistic Regression', 'Naive Bayes',
2                                     'K-Nearest Neighbour', 'Decision Tree', 'Support Vector Machine'],
3                                     'nb_acc_score': 86.341463, 'rf_acc_score': 93.658537, 'xgb_acc_score': 92.195122,
4                                     'knn_acc_score': 87.804878, 'dt_acc_score': 94.634146, 'svm_acc_score': 98.048780})
        model_ev
```

Out[41]:

	Model	Accuracy
0	Logistic Regression	86.341463
1	Naive Bayes	85.365854
2	Random Forest	93.658537
3	Extreme Gradient Boost	92.195122
4	K-Nearest Neighbour	87.804878
5	Decision Tree	94.634146
6	Support Vector Machine	98.048780

```
In [42]: 1 colors = ['red', 'green', 'blue', 'gold', 'silver', 'yellow', 'orange',]
2 plt.figure(figsize=(12,5))
3 plt.title("barplot Represent Accuracy of different models")
4 plt.xlabel("Accuracy %")
5 plt.ylabel("Algorithms")
6 plt.bar(model_ev['Model'],model_ev['Accuracy'],color = colors)
7 plt.show()
```



## Ensembling

In order to increase the accuracy of the model we use ensembling. Here we use stacking technique.

```
In [43]: 1 scv=StackingCVClassifier(classifiers=[xgb,knn,svc],meta_classifier
2 scv.fit(X_train,y_train)
3 scv_predicted = scv.predict(X_test)
4 scv_conf_matrix = confusion_matrix(y_test, scv_predicted)
5 scv_acc_score = accuracy_score(y_test, scv_predicted)
6 print("confussion matrix")
7 print(scv_conf_matrix)
8 print("\n")
9 print("Accuracy of StackingCVClassifier:",scv_acc_score*100,'\n')
10 print(classification_report(y_test,scv_predicted))
```

```
confussion matrix
[[ 94   4]
 [  0 107]]
```

Accuracy of StackingCVClassifier: 98.04878048780488

	precision	recall	f1-score	support
0	1.00	0.96	0.98	98
1	0.96	1.00	0.98	107
accuracy			0.98	205
macro avg	0.98	0.98	0.98	205
weighted avg	0.98	0.98	0.98	205

## Conclusion

1) The Support Vector Machine (SVM) and the StackingCVClassifier demonstrated the highest accuracy, each achieving an impressive 98.05%. The Decision Tree model also performed notably well, with an accuracy of 94.63%.

2) Exercise induced angina,Chest pain is major symptoms of heart attack.

3) Ensembling technique increase the accuracy of the model.

Type *Markdown* and LaTeX:  $\alpha^2$

In [ ]:

1