

HEART DISEASE PREDICTION BY NIKITA

I. IMPORTING THE ESSENTIAL LIBRARIES

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

```
%matplotlib inline
```

```
import warnings
warnings.filterwarnings('ignore')
```

II. IMPORTING AND UNDERSTANDING OUR DATASET

```
df = pd.read_excel("1645792390_ccp1_dataset.xlsx")
```

Shape of Dataset

df.shape

(303, 14)

Printing few Columns

df.head()

```
   age  sex  cp  trestbps  chol  fbs  restecg  thalach  exang  oldpeak  slope  ca  thal  target
0   63   1   3    145   233   1    0    150   0    2.3   0  0   1   1
1   37   1   2    130   250   0    1    167   0    3.5   0  0   2   1
2   41   0   1    120   204   0    0    172   0    1.4   2  0   2   1
3   56   1   1    120   236   0    1    170   0    0.8   2  0   2   1
4   57   0   0    120   354   0    1    163   1    0.6   2  0   2   1
```

Description

df.describe()

```
count    303.000000    303.000000    303.000000    303.000000    303.000000    303.000000    303.000000    303.000000    303.000000    303.000000    303.000000    303.000000    303.000000
mean      54.363377    0.683168    0.969997    131.623762    246.264026    0.148515    0.528053    149.646805    0.326733    1.039604    1.396340    0.729373    2.313531    0.544554
std       9.082101    0.460611    1.032052    17.538143    51.830751    0.356198    0.525860    22.905161    0.469794    1.161075    0.616226    1.022606    0.612277    0.498835
min       29.000000    0.000000    0.000000    94.000000    126.000000    0.000000    0.000000    71.000000    0.000000    0.000000    0.000000    0.000000    0.000000    0.000000
25%      47.500000    0.000000    0.000000    120.000000    211.000000    0.000000    0.000000    133.500000    0.000000    0.000000    1.000000    0.000000    2.000000    0.000000
50%      55.000000    1.000000    1.000000    130.000000    240.000000    0.000000    1.000000    153.000000    0.000000    0.800000    1.000000    0.000000    2.000000    1.000000
75%      61.000000    1.000000    2.000000    140.000000    274.500000    0.000000    1.000000    166.000000    1.000000    1.600000    2.000000    1.000000    3.000000    1.000000
max       77.000000    1.000000    3.000000    200.000000    564.000000    1.000000    2.000000    202.000000    1.000000    6.200000    2.000000    4.000000    3.000000    1.000000
```

df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 303 entries, 0 to 302
Data columns (total 14 columns):
 #   Column      Non-Null Count  Dtype
---  ---
 0   age         303 non-null    int64
 1   sex         303 non-null    int64
 2   cp         303 non-null    int64
 3   trestbps    303 non-null    int64
 4   chol       303 non-null    int64
 5   fbs        303 non-null    int64
 6   restecg     303 non-null    int64
 7   thalach     303 non-null    int64
 8   exang       303 non-null    int64
 9   oldpeak     303 non-null    float64
10  slope       303 non-null    int64
11  ca         303 non-null    int64
12  thal       303 non-null    int64
13  target      303 non-null    int64
dtypes: float64(1), int64(13)
memory usage: 33.3 KB
```

info = "age", "1: male, 0: female", "chest pain type, 1: typical angina, 2: atypical angina, 3: non-anginal pain, 4: asymptomatic", "resting blood pressure", " " sex

```
for i in range(len(info)):
    print(df.columns[i]+"-"+info[i])
```

age: age
sex: 1: male, 0: female
trestbps: resting blood pressure
chol: serum cholesterol in mg/dl
fbs: fasting blood sugar > 120 mg/dl
restecg: resting electrocardiographic results (values 0,1,2)
thalach: maximum heart rate achieved
exang: exercise induced angina
oldpeak: oldpeak = ST depression induced by exercise relative to rest
slope: the slope of the peak exercise ST segment
ca: number of major vessels (0-3) colored by fluoroscopy
thal: thal: 3 = normal; 6 = fixed defect; 7 = reversible defect

Analysing the target variable

df["target"].describe()

```
count    303.000000
mean      0.544554
std       0.498835
min       0.000000
25%       0.000000
50%       1.000000
75%       1.000000
max       1.000000
Name: target, dtype: float64
```

df["target"].unique()

array([1, 0], dtype=int64)

Checking Correlation between columns

```
print(df.corr()["target"].abs().sort_values(ascending=False))
```

```
target    1.000000
exang     0.436757
exang     0.433798
oldpeak   0.430696
thalach   0.421741
ca        0.391724
slope     0.345877
thal      0.344029
sex       0.280937
age       0.225439
trestbps  0.144931
restecg   0.137238
chol      0.085239
fbs       0.028048
Name: target, dtype: float64
```

#'fbs' is weakely correlated

Eploratory Data Analytics (EDA)

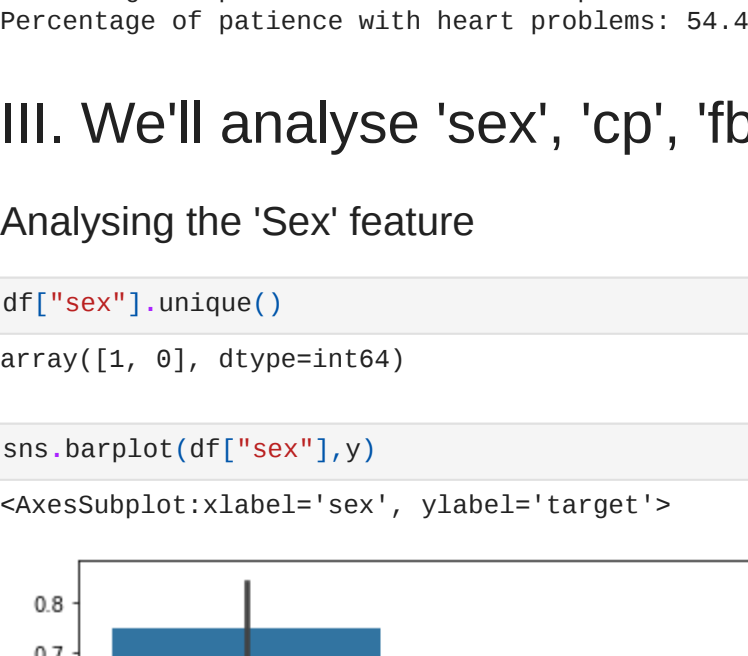
Analysing the target Variable

y = df["target"]

sns.countplot(y)

```
target_temp = df.target.value_counts()
print(target_temp)
```

```
1    165
0     128
Name: target, dtype: int64
```



```
print(Percentage of patience without heart problems: "%str(round(target_temp[0]/100*303,2))")
print(Percentage of patience with heart problems: "%str(round(target_temp[1]/100*303,2))")
```

Percentage of patience without heart problems: 45.54
Percentage of patience with heart problems: 54.46

III. We'll analyse 'sex', 'cp', 'fbs', 'restecg', 'exang', 'slope', 'ca' and 'thal' features

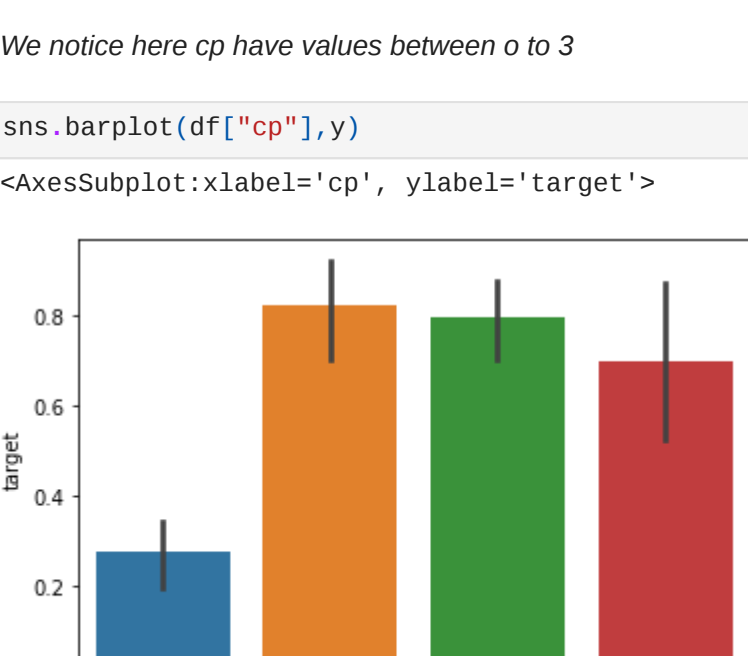
Analysing the 'Sex' feature

df["sex"].unique()

array([1, 0], dtype=int64)

sns.barplot(df["sex"],y)

<AxesSubplot: xlabel='sex', ylabel='target'>



We notice that Female have more chances of getting heart problems

Analysing the chest pain feature*

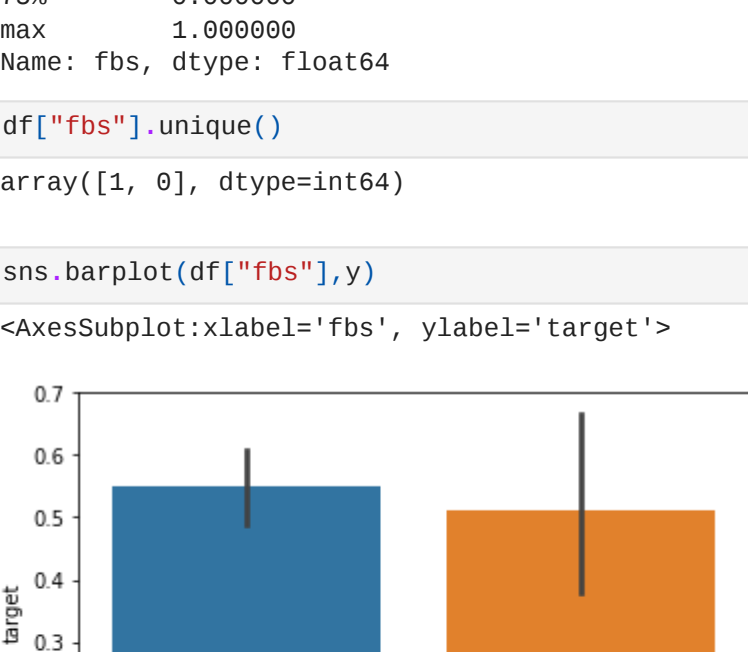
df["cp"].unique()

array([3, 2, 1, 0], dtype=int64)

We notice here cp have values between 0 to 3

sns.barplot(df["cp"],y)

<AxesSubplot: xlabel='cp', ylabel='target'>



We notice, that chest pain of '0', i.e. the ones with typical angina are much less likely to have heart problems

Analysing the fbs feature

df["fbs"].describe()

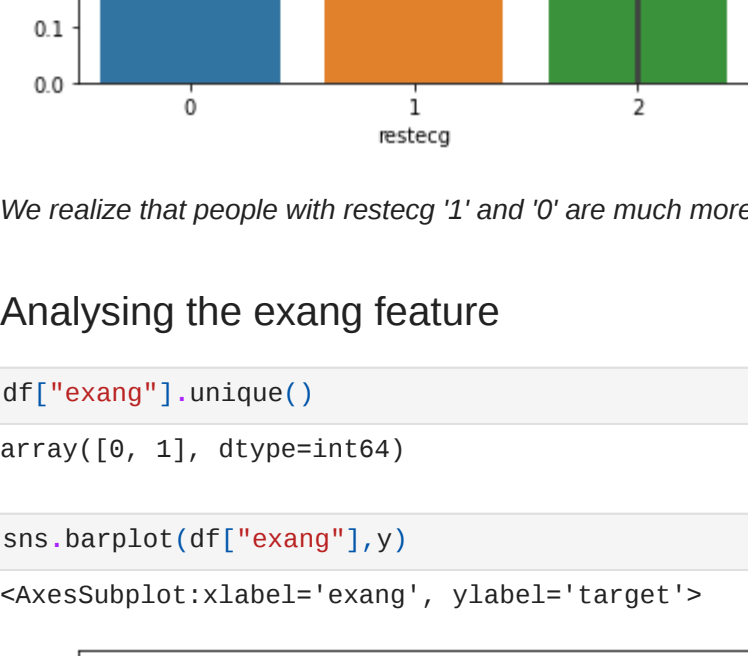
```
count    303.000000
mean      0.148515
std       0.356198
min       0.000000
25%       0.000000
50%       0.000000
75%       0.000000
max       1.000000
Name: fbs, dtype: float64
```

df["fbs"].unique()

array([1, 0], dtype=int64)

sns.barplot(df["fbs"],y)

<AxesSubplot: xlabel='fbs', ylabel='target'>



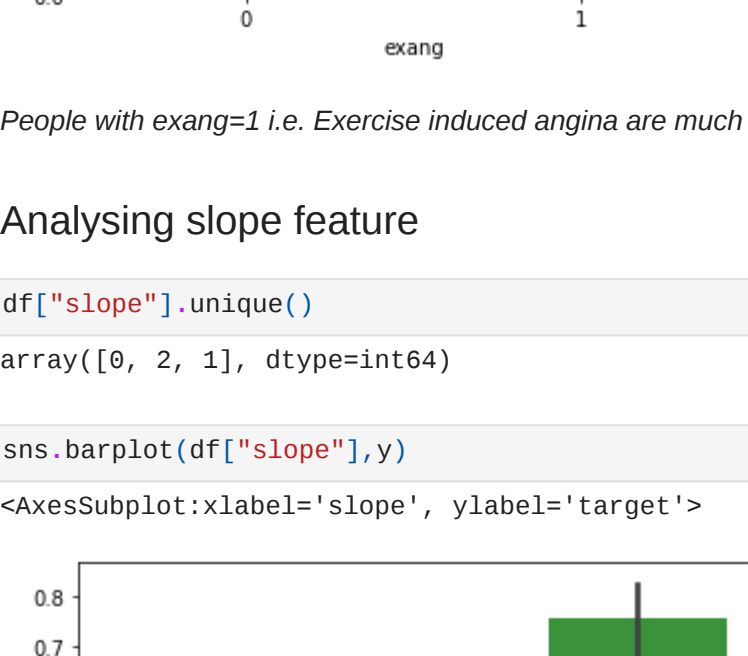
Analysing the restecg feature

df["restecg"].unique()

array([0, 1, 2], dtype=int64)

sns.barplot(df["restecg"],y)

<AxesSubplot: xlabel='restecg', ylabel='target'>



We realize that people with restecg '1' and '0' are much more likely to have a heart disease than with restecg '2'

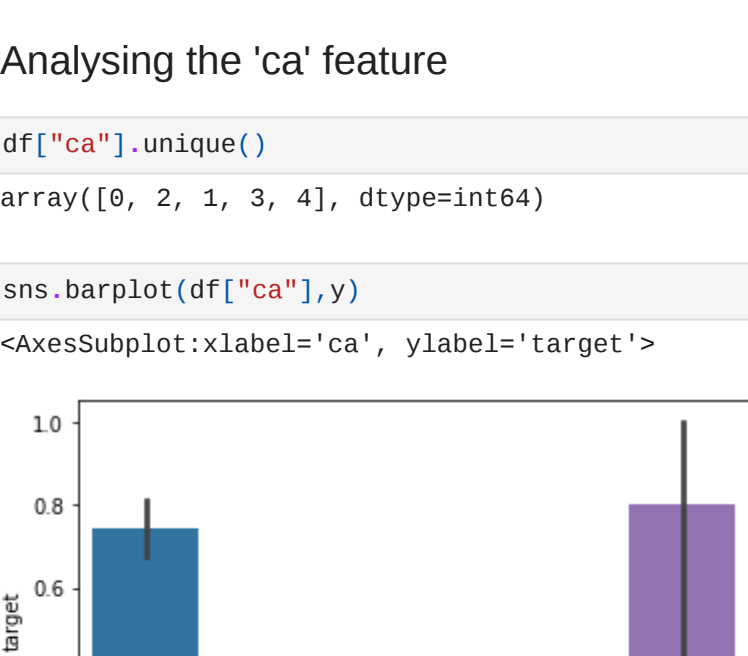
Analysing the exang feature

df["exang"].unique()

array([0, 1], dtype=int64)

sns.barplot(df["exang"],y)

<AxesSubplot: xlabel='exang', ylabel='target'>



People with exang=1 i.e. Exercise induced angina are much less likely to have heart problems

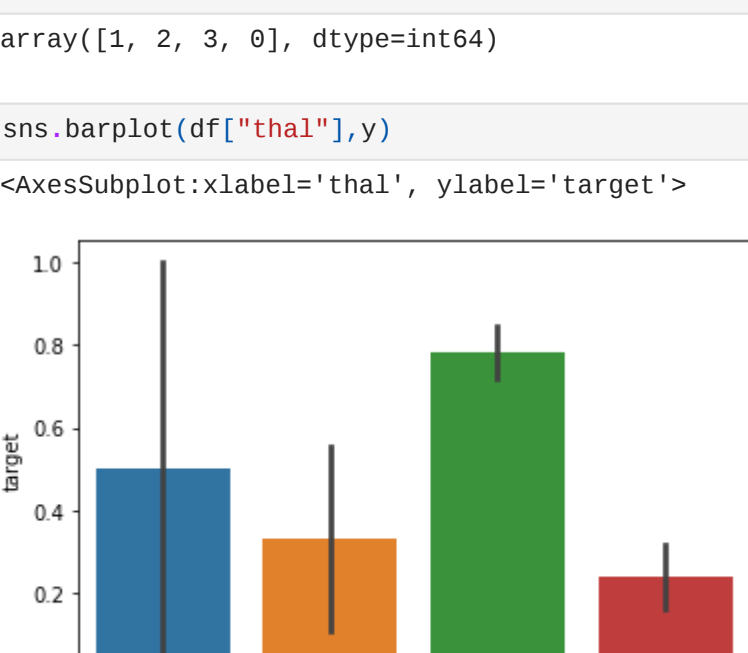
Analysing slope feature

df["slope"].unique()

array([0, 2, 1], dtype=int64)

sns.barplot(df["slope"],y)

<AxesSubplot: xlabel='slope', ylabel='target'>



We observe, that Slope '2' causes heart pain much more than Slope '0' and '1'

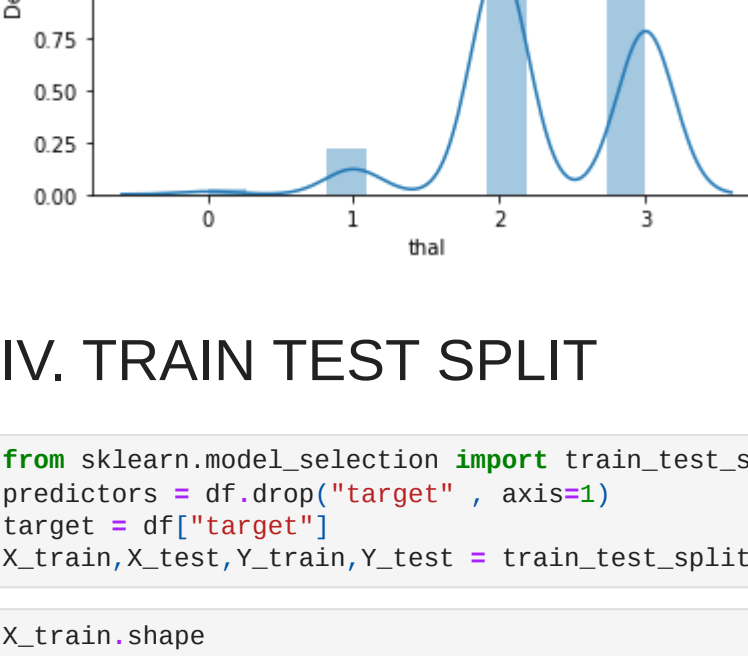
Analysing the 'ca' feature

df["ca"].unique()

array([0, 2, 1, 3, 4], dtype=int64)

sns.barplot(df["ca"],y)

<AxesSubplot: xlabel='ca', ylabel='target'>



ca=4 has astonishingly large number of heart patients

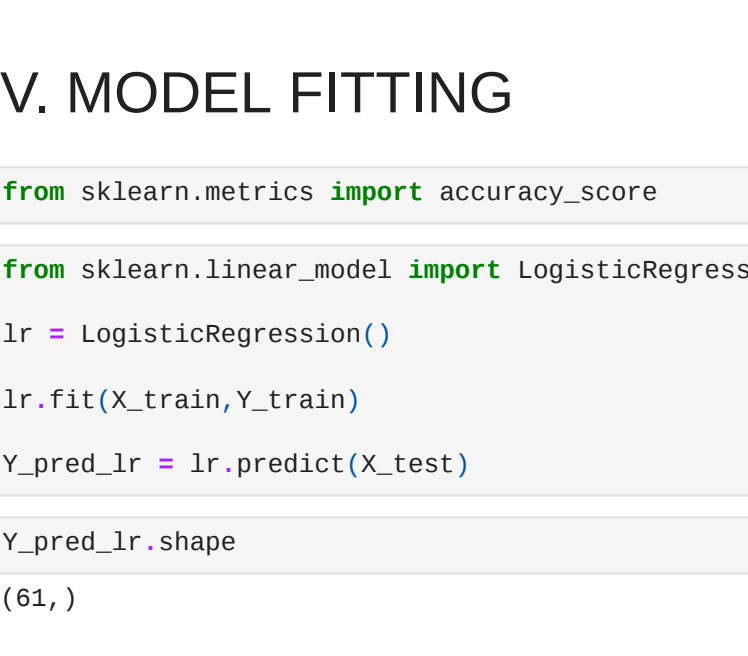
Analysing the 'thal' feature

df["thal"].unique()

array([1, 2, 3, 0], dtype=int64)

sns.barplot(df["thal"],y)

<AxesSubplot: xlabel='thal', ylabel='target'>



sns.distplot(df["thal"],y)

<AxesSubplot: xlabel='thal', ylabel='Density'>



IV. TRAIN TEST SPLIT

```
from sklearn.model_selection import train_test_split
predictors = df.drop("target", axis=1)
X_train,X_test,Y_train,Y_test = train_test_split(predictors,target,test_size=0.20,random_state=0)
```

X_train.shape

(242, 13)

X_test.shape

(61, 13)

Y_train.shape

(242,)

Y_test.shape

(61,)

V. MODEL FITTING

```
from sklearn.metrics import accuracy_score
```

```
from sklearn.linear_model import LogisticRegression
```

```
lr = LogisticRegression()
```

```
lr.fit(X_train,Y_train)
```

```
Y_pred_lr = lr.predict(X_test)
```

Y_pred_lr.shape

(61,)

```
score_lr = round(accuracy_score(Y_pred_lr,Y_test)*100,2)
```

```
print("The accuracy score achieved using Logistic Regression is: "+str(score_lr)+" %")
```

The accuracy score achieved using Logistic Regression is: 85.25 %

(A). Naive Bayes

```
from sklearn.naive_bayes import GaussianNB
```

```
nb = GaussianNB()
```

```
nb.fit(X_train,Y_train)
```

```
Y_pred_nb = nb.predict(X_test)
```

Y_pred_nb.shape

(61,)

```
score_nb = round(accuracy_score(Y_pred_nb,Y_test)*100,2)
```

```
print("The accuracy score achieved using Naive Bayes is: "+str(score_nb)+" %")
```

The accuracy score achieved using Naive Bayes is: 85.25 %

(B). SVM

```
from sklearn import svm
```

```
sv = svm.SVC(kernel="linear")
```

```
sv.fit(X_train,Y_train)
```

```
Y_pred_svm = sv.predict(X_test)
```

Y_pred_svm.shape

(61,)

```
score_svm = round(accuracy_score(Y_pred_svm,Y_test)*100,2)
```

```
print("The accuracy score achieved using Linear SVM is: "+str(score_svm)+" %")
```

The accuracy score achieved using Linear SVM is: 81.97 %

(C). K Nearest Neighbours

```
from sklearn.neighbors import KNeighborsClassifier
```

```
knn = KNeighborsClassifier(n_neighbors=7)
```

```
knn.fit(X_train,Y_train)
```

```
Y_pred_knn=knn.predict(X_test)
```

Y_pred_knn.shape

(61,)

```
score_knn = round(accuracy_score(Y_pred_knn,Y_test)*100,2)
```

```
print("The accuracy score achieved using KNN is: "+str(score_knn)+" %")
```

The accuracy score achieved using KNN is: 67.21 %

(D). Decision Tree

```
from sklearn.tree import DecisionTreeClassifier
```

```
max_accuracy = 0
```

```
for x in range(200):
    dt = DecisionTreeClassifier(random_state=x)
    dt.fit(X_train,Y_train)
    Y_pred_dt = dt.predict(X_test)
    current_accuracy = round(accuracy_score(Y_pred_dt,Y_test)*100,2)
    if(current_accuracy>max_accuracy):
        max_accuracy = current_accuracy
        best_x = x
```

```
dt = DecisionTreeClassifier(random_state=best_x)
dt.fit(X_train,Y_train)
Y_pred_dt = dt.predict(X_test)
```

print(Y_pred_dt.shape)

(61,)

```
score_dt = round(accuracy_score(Y_pred_dt,Y_test)*100,2)
```

```
print("The accuracy score achieved using Decision Tree is: "+str(score_dt)+" %")
```

The accuracy score achieved using Decision Tree is: 81.97 %

(E). Random Forest

```
from sklearn.ensemble import RandomForestClassifier
```

```
max_accuracy = 0
```

```
for x in range(2000):
    rf = RandomForestClassifier(random_state=x)
    rf.fit(X_train,Y_train)
    Y_pred_rf = rf.predict(X_test)
    current_accuracy = round(accuracy_score(Y_pred_rf,Y_test)*100,2)
    if(current_accuracy>max_accuracy):
        max_accuracy = current_accuracy
        best_x = x
```

```
#print(max_accuracy)
#print(best_x)
```

```
rf = RandomForestClassifier(random_state=best_x)
rf.fit(X_train,Y_train)
Y_pred_rf = rf.predict(X_test)
```

Y_pred_rf.shape

(61,)

```
score_rf = round(accuracy_score(Y_pred_rf,Y_test)*100,2)
```

```
print("The accuracy score achieved using Decision Tree is: "+str(score_rf)+" %")
```

The accuracy score achieved using Decision Tree is: 90.16 %

XGBoost

```
import xgboost as xgb
```

```
xgb_model = xgb.XGBClassifier(objective="binary:logistic", random_state=42)
```

```
xgb_model.fit(X_train,Y_train)
```

```
Y_pred_xgb = xgb_model.predict(X_test)
```

[11:52:22] WARNING: D:\bld\xgboost-split_164519815404\work\src\learner.cc:1115: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.

Y_pred_xgb.shape

(61,)

```
score_xgb = round(accuracy_score(Y_pred_xgb,Y_test)*100,2)
```

```
print("The accuracy score achieved using XGBoost is: "+str(score_xgb)+" %")
```

The accuracy score achieved using XGBoost is: 78.69 %

```
scores = [score_lr,score_nb,score_svm,score_knn,score_dt,score_rf,score_xgb]
algorithms = ["Logistic Regression","Naive Bayes","Support Vector Machine","K-Nearest Neighbors","Decision Tree","Random Forest","XGBoost"]
```

```
for i in range(len(algorithms)):
    print("The accuracy score achieved using "+algorithms[i]+" is: "+str(scores[i])+" %")
```

The accuracy score achieved using Logistic Regression is: 85.25 %
The accuracy score achieved using Naive Bayes is: 85.25 %
The accuracy score achieved using Support Vector Machine is: 81.97 %
The accuracy score achieved using K-Nearest Neighbors is: 67.21 %
The accuracy score achieved using Decision Tree is: 81.97 %
The accuracy score achieved using Random Forest is: 90.16 %
The accuracy score achieved using XGBoost is: 78.69 %

```
sns.set(rc={"figure.figsize":(15,8)})
```

```
plt.xlabel("Algorithms")
```

```
plt.ylabel("Accuracy score")
```

sns.barplot(algorithms,scores)

<AxesSubplot: xlabel='Algorithms', ylabel='Accuracy score'>

Here, Random Forest is giving Best accuracy

In []