



Ram Lal Anand College

(University of Delhi)

DATA MINING PROJECT

Department of Computer Science

Heart Disease Prediction

Name of Course : B.Sc. (H) Computer Science

Semester : 6th

Name of the Paper : Data Mining

Paper Code : 32347611

Name of the Student : Abhishek Kumar

Examination Roll No : 20058570051

About the Project:

Heart-Disease-Prediction:

Thus, preventing heart diseases has become more than necessary. Good data-driven systems for predicting heart diseases can improve the entire research and prevention process, making sure that more people can live healthy lives. This is where Machine Learning comes into play. Machine Learning helps in predicting the heart diseases, and the predictions made are quite accurate.

The project involved analysis of the heart disease patient dataset with proper data processing. Then, different models were trained and predictions are made with different algorithms KNN, Decision Tree, Random Forest, SVM, Logistic Regression etc This is the jupyter notebook code and dataset I've used for my Kaggle kernel 'Binary Classification with Sklearn and Keras'

I've used a variety of Machine Learning algorithms, implemented in Python, to predict the presence of heart disease in a patient. This is a classification problem, with input features as a variety of parameters, and the target variable as a binary variable, predicting whether heart disease is present or not.

Machine Learning algorithms used:

1. Logistic Regression (Scikit-learn)
2. Naive Bayes (Scikit-learn)
3. Support Vector Machine (Linear) (Scikit-learn)
4. K-Nearest Neighbours (Scikit-learn)
5. Decision Tree (Scikit-learn)
6. Random Forest (Scikit-learn)
7. XGBoost (Scikit-learn)
8. Artificial Neural Network with 1 Hidden layer (Keras)

Accuracy achieved: 95% (Random Forest)

Dataset used: <https://www.kaggle.com/ronitf/heart-disease-uci>

Heart Disease Prediction By Abhishek Kumar

▼ I. Importing essential libraries

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

%matplotlib inline

import os
print(os.listdir())

import warnings
warnings.filterwarnings('ignore')

['.ipynb_checkpoints', 'heart.csv', 'Heart_disease_prediction.ipynb', 'README.md']
```

▼ II. Importing and understanding our dataset

```
dataset = pd.read_csv("heart.csv")
```

▼ Verifying it as a 'dataframe' object in pandas

```
type(dataset)

pandas.core.frame.DataFrame
```

▼ Shape of dataset

```
dataset.shape

(303, 14)
```

▼ Printing out a few columns

```
dataset.head(5)
```

	age	sex	cp	trestbps	chol	fb	restecg	thalach	exang	oldpeak	slope	ca	th
0	63	1	3	145	233	1	0	150	0	2.3	0	0	0
1	37	1	2	130	250	0	1	187	0	3.5	0	0	0
2	41	0	1	130	204	0	0	172	0	1.4	2	0	0
3	56	1	1	120	236	0	1	178	0	0.8	2	0	0

dataset.sample(5)

	age	sex	cp	trestbps	chol	fb	restecg	thalach	exang	oldpeak	slope	ca	t
248	54	1	1	192	283	0	0	195	0	0.0	2	1	0
147	60	0	3	150	240	0	1	171	0	0.9	2	0	0
239	35	1	0	126	282	0	0	156	1	0.0	2	0	0
4	57	0	0	120	354	0	1	163	1	0.6	2	0	0
7	44	1	1	120	263	0	1	173	0	0.0	2	0	0

▼ Description

dataset.describe()

	age	sex	cp	trestbps	chol	fb	restecg
count	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000
mean	54.366337	0.683168	0.966997	131.623762	246.264026	0.148515	0.52805
std	9.082101	0.466011	1.032052	17.538143	51.830751	0.356198	0.52586
min	29.000000	0.000000	0.000000	94.000000	126.000000	0.000000	0.00000
25%	47.500000	0.000000	0.000000	120.000000	211.000000	0.000000	0.00000
50%	55.000000	1.000000	1.000000	130.000000	240.000000	0.000000	1.00000
75%	61.000000	1.000000	2.000000	140.000000	274.500000	0.000000	1.00000
max	77.000000	1.000000	3.000000	200.000000	564.000000	1.000000	2.00000

dataset.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 303 entries, 0 to 302
Data columns (total 14 columns):
age      303 non-null int64
sex      303 non-null int64
cp       303 non-null int64
trestbps 303 non-null int64
chol     303 non-null int64
fb       303 non-null int64
```

```

restecg      303 non-null int64
thalach     303 non-null int64
exang       303 non-null int64
oldpeak      303 non-null float64
slope        303 non-null int64
ca           303 non-null int64
thal         303 non-null int64
target       303 non-null int64
dtypes: float64(1), int64(13)
memory usage: 33.2 KB

```

Luckily, we have no missing values

▼ Let's understand our columns better:

```

info = ["age","1: male, 0: female","chest pain type, 1: typical angina, 2: atypical angina

for i in range(len(info)):
    print(dataset.columns[i]+":\t\t\t"+info[i])

age:                               age
sex:                             1: male, 0: female
cp:                            chest pain type, 1: typical angina, 2: atypical angina, 3: no
                                resting blood pressure
trestbps:                         serum cholestorol in mg/dl
chol:                            fasting blood sugar > 120 mg/dl
fbs:                             resting electrocardiographic results (values 0,1,2)
restecg:                           maximum heart rate achieved
thalach:                          exercise induced angina
exang:                            oldpeak = ST depression induced by exercise relative
oldpeak:                           the slope of the peak exercise ST segment
slope:                            number of major vessels (0-3) colored by flourosopy
ca:                               thal: 3 = normal; 6 = fixed defect; 7 = reversable defect
thal:

```



▼ Analysing the 'target' variable

```
dataset["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

```
dataset["target"].unique()  
  
array([1, 0])
```

Clearly, this is a classification problem, with the target variable having values '0' and '1'

▼ Checking correlation between columns

```
print(dataset.corr()["target"].abs().sort_values(ascending=False))  
  
target      1.000000  
exang      0.436757  
cp        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     0.137230  
chol       0.085239  
fbs        0.028046  
Name: target, dtype: float64
```

#This shows that most columns are moderately correlated with target, but 'fbs' is very weakly correlated.

▼ Exploratory Data Analysis (EDA)

▼ First, analysing the target variable:

```
y = dataset["target"]  
  
sns.countplot(y)  
  
target_temp = dataset.target.value_counts()  
  
print(target_temp)
```

```
1    165
0    138
Name: target, dtype: int64
```



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

```
#Alternatively,
# print("Percentage of patients with heart problems: "+str(y.where(y==1).count()*100/303))
# print("Percentage of patients with heart problems: "+str(y.where(y==0).count()*100/303))

# #Or,
# countNoDisease = len(df[df.target == 0])
# countHaveDisease = len(df[df.target == 1])

Percentage of patients without heart problems: 45.54
Percentage of patients with heart problems: 54.46
```

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

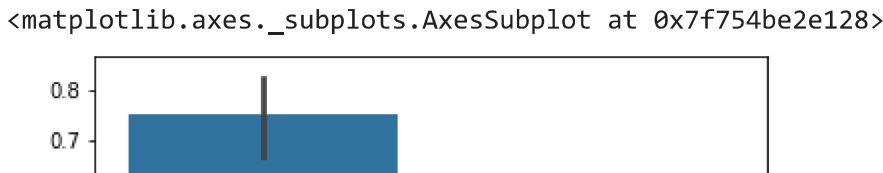
▼ Analysing the 'Sex' feature

```
dataset["sex"].unique()
```

```
array([1, 0])
```

▼ We notice, that as expected, the 'sex' feature has 2 unique features

```
sns.barplot(dataset["sex"],y)
```



We notice, that females are more likely to have heart problems than males



▼ Analysing the 'Chest Pain Type' feature

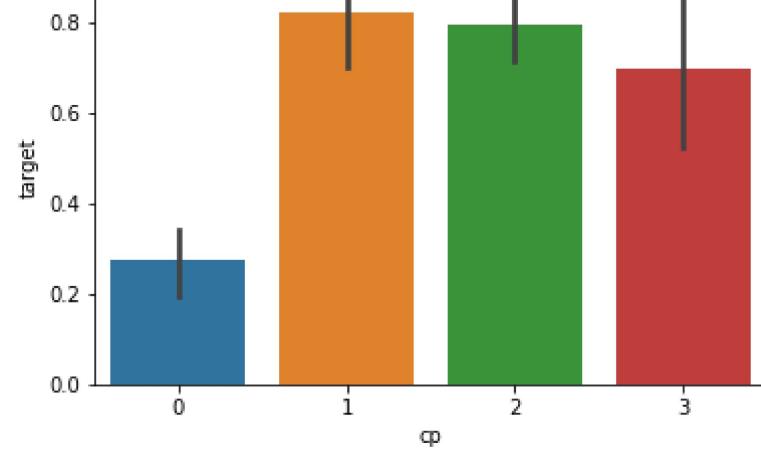


```
dataset["cp"].unique()
```

```
array([3, 2, 1, 0])
```

▼ As expected, the CP feature has values from 0 to 3

```
sns.barplot(dataset["cp"],y)
```



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

```
dataset["fbs"].describe()
```

count	303.000000
mean	0.148515
std	0.356198
min	0.000000
25%	0.000000
50%	0.000000
75%	0.000000

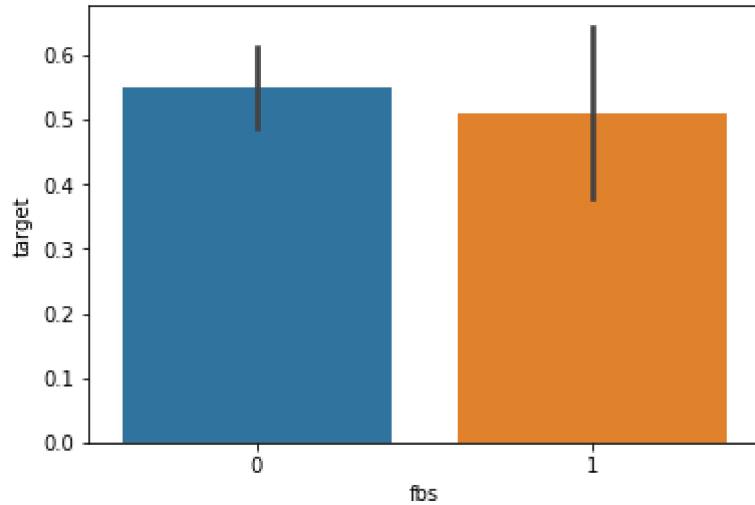
```
max      1.000000  
Name: fbs, dtype: float64
```

```
dataset["fbs"].unique()
```

```
array([1, 0])
```

```
sns.barplot(dataset["fbs"],y)
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f754bdda90>
```



Nothing extraordinary here

▼ Analysing the restecg feature

```
dataset["restecg"].unique()
```

```
array([0, 1, 2])
```

```
sns.barplot(dataset["restecg"],y)
```

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



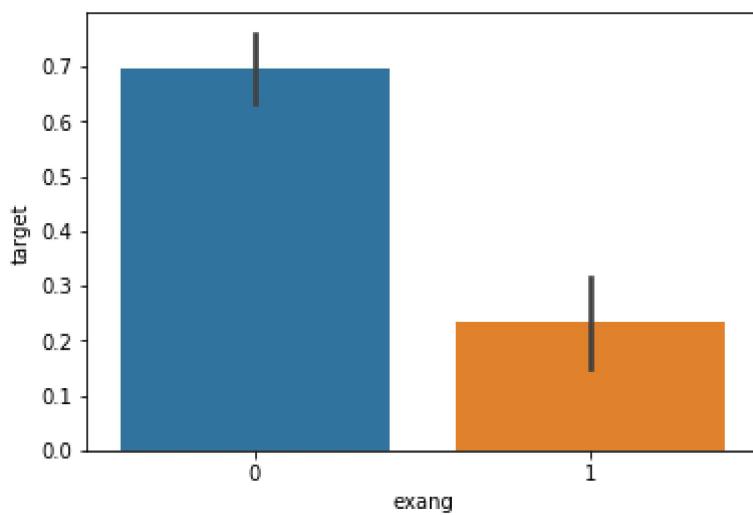
```
dataset["exang"].unique()
```

```
array([0, 1])
```



```
sns.barplot(dataset["exang"],y)
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f754bc9eda0>
```



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

▼ Analysing the Slope feature

```
dataset["slope"].unique()
```

```
array([0, 2, 1])
```

```
sns.barplot(dataset["slope"],y)
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f754bc75710>
```



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



▼ Analysing the 'ca' feature



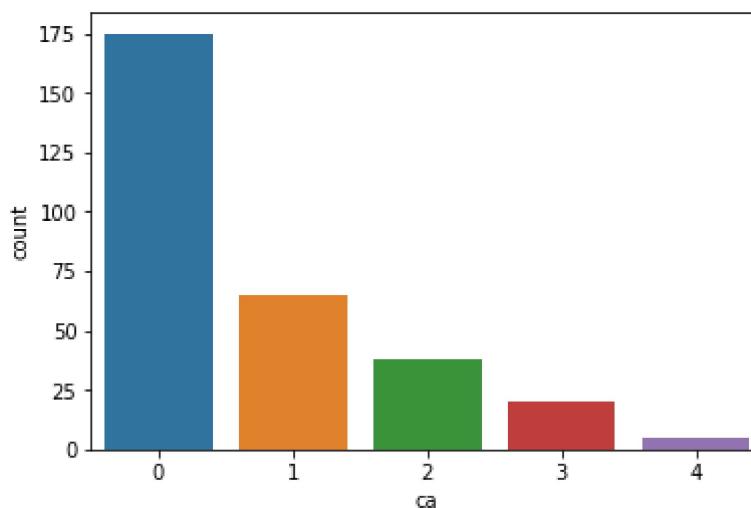
```
#number of major vessels (0-3) colored by flourosopy
```

```
dataset["ca"].unique()
```

```
array([0, 2, 1, 3, 4])
```

```
sns.countplot(dataset["ca"])
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f754bd13940>
```



```
sns.barplot(dataset["ca"],y)
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f754bc34c88>
```

- ▼ ca=4 has astonishingly large number of heart patients



Analysing the 'thal' feature



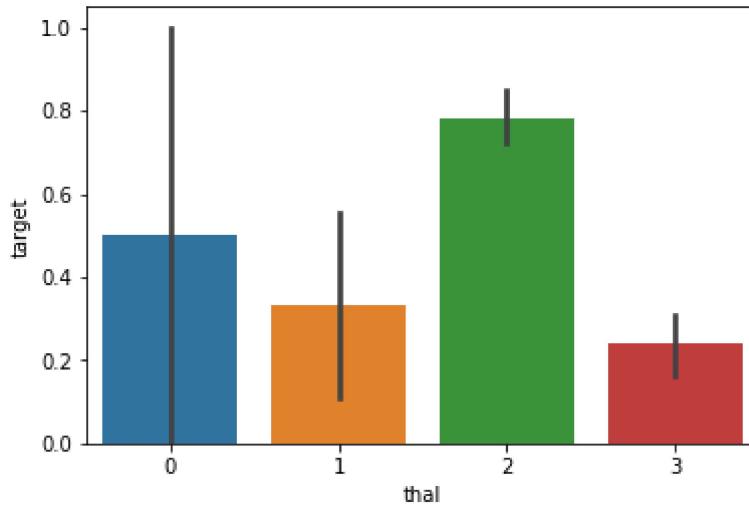
dataset["thal"].unique()

array([1, 2, 3, 0])



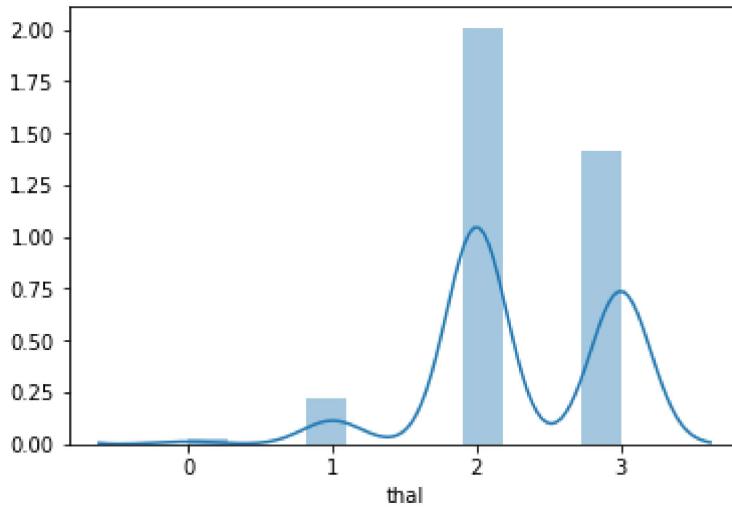
sns.barplot(dataset["thal"],y)

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f754bb89128>
```



sns.distplot(dataset["thal"])

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f754baf1fd0>
```



- ▼ IV. Train Test split

```
from sklearn.model_selection import train_test_split

predictors = dataset.drop("target",axis=1)
target = dataset["target"]

X_train,X_test,Y_train,Y_test = train_test_split(predictors,target,test_size=0.20,random_s

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
```

▼ Logistic Regression

```
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 %

▼ 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 %
```

▼ 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 %
```

▼ K Nearest Neighbors

```
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 %
```

▼ 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

#print(max_accuracy)
#print(best_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 %
```

▼ 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: 95.08 %
```

▼ 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)

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: 85.25 %
```

▼ Neural Network

```
from keras.models import Sequential
from keras.layers import Dense

    Using TensorFlow backend.

# https://stats.stackexchange.com/a/136542 helped a lot in avoiding overfitting

model = Sequential()
model.add(Dense(11,activation='relu',input_dim=13))
model.add(Dense(1,activation='sigmoid'))

model.compile(loss='binary_crossentropy',optimizer='adam',metrics=['accuracy'])

model.fit(X_train,Y_train,epochs=300)

Epoch 1/300
242/242 [=====] - 3s 10ms/step - loss: 4.7817 - acc: 0.60
Epoch 2/300
242/242 [=====] - 0s 123us/step - loss: 2.6633 - acc: 0.60
Epoch 3/300
242/242 [=====] - 0s 124us/step - loss: 2.3371 - acc: 0.60
Epoch 4/300
242/242 [=====] - 0s 127us/step - loss: 2.2938 - acc: 0.60
Epoch 5/300
242/242 [=====] - 0s 127us/step - loss: 2.0712 - acc: 0.60
Epoch 6/300
242/242 [=====] - 0s 125us/step - loss: 2.0645 - acc: 0.60
Epoch 7/300
242/242 [=====] - 0s 124us/step - loss: 2.0023 - acc: 0.60
Epoch 8/300
242/242 [=====] - 0s 127us/step - loss: 1.9928 - acc: 0.60
Epoch 9/300
242/242 [=====] - 0s 129us/step - loss: 1.9670 - acc: 0.60
Epoch 10/300
242/242 [=====] - 0s 123us/step - loss: 1.9316 - acc: 0.60
Epoch 11/300
242/242 [=====] - 0s 123us/step - loss: 1.9223 - acc: 0.60
Epoch 12/300
242/242 [=====] - 0s 124us/step - loss: 1.8955 - acc: 0.60
Epoch 13/300
242/242 [=====] - 0s 124us/step - loss: 1.8573 - acc: 0.60
Epoch 14/300
242/242 [=====] - 0s 125us/step - loss: 1.8268 - acc: 0.60
Epoch 15/300
242/242 [=====] - 0s 122us/step - loss: 1.7915 - acc: 0.60
Epoch 16/300
242/242 [=====] - 0s 121us/step - loss: 1.7727 - acc: 0.60
Epoch 17/300
242/242 [=====] - 0s 122us/step - loss: 1.7275 - acc: 0.60
Epoch 18/300
242/242 [=====] - 0s 122us/step - loss: 1.7143 - acc: 0.70
```

```
242/242 [=====] - 0s 124us/step - loss: 1.6761 - acc: 0.6
Epoch 20/300
242/242 [=====] - 0s 136us/step - loss: 1.6591 - acc: 0.6
Epoch 21/300
242/242 [=====] - 0s 128us/step - loss: 1.6486 - acc: 0.7
Epoch 22/300
242/242 [=====] - 0s 130us/step - loss: 1.6004 - acc: 0.6
Epoch 23/300
242/242 [=====] - 0s 125us/step - loss: 1.5765 - acc: 0.7
Epoch 24/300
242/242 [=====] - 0s 128us/step - loss: 1.5428 - acc: 0.7
Epoch 25/300
242/242 [=====] - 0s 125us/step - loss: 1.5179 - acc: 0.6
Epoch 26/300
242/242 [=====] - 0s 124us/step - loss: 1.5084 - acc: 0.7
Epoch 27/300
242/242 [=====] - 0s 124us/step - loss: 1.4619 - acc: 0.6
Epoch 28/300
```

```
Y_pred_nn = model.predict(X_test)
```

```
Y_pred_nn.shape
```

```
(61, 1)
```

```
rounded = [round(x[0]) for x in Y_pred_nn]
```

```
Y_pred_nn = rounded
```

```
score_nn = round(accuracy_score(Y_pred_nn,Y_test)*100,2)
```

```
print("The accuracy score achieved using Neural Network is: "+str(score_nn)+" %")
```

```
#Note: Accuracy of 85% can be achieved on the test set, by setting epochs=2000, and number
```

```
The accuracy score achieved using Neural Network is: 80.33 %
```

▼ VI. Output final score

```
scores = [score_lr,score_nb,score_svm,score_knn,score_dt,score_rf,score_xgb,score_nn]
algorithms = ["Logistic Regression","Naive Bayes","Support Vector Machine","K-Nearest Neig
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: 95.08 %

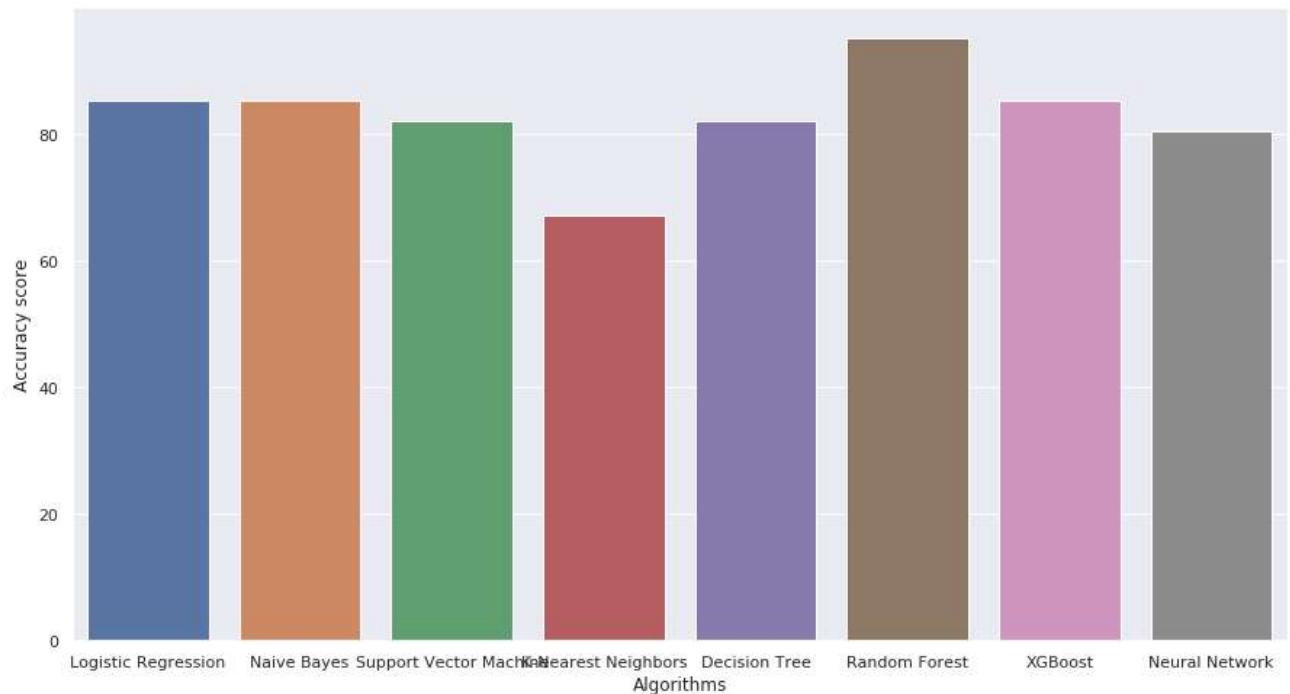
The accuracy score achieved using XGBoost is: 85.25 %

The accuracy score achieved using Neural Network is: 80.33 %

```
sns.set(rc={'figure.figsize':(15,8)})  
plt.xlabel("Algorithms")  
plt.ylabel("Accuracy score")
```

```
sns.barplot(algorithms,scores)
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7f74ea800eb8>
```



Hey Abhishek there random forest has good result as compare to other algorithms

