

Results Screen shots:

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

%matplotlib inline

import os

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

II. Importing and understanding our dataset

```
In [12]: dataset = pd.read_csv("heart.csv")
```

Verifying it as a 'dataframe' object in pandas

```
In [13]: type(dataset)
```

```
Out[13]: pandas.core.frame.DataFrame
```

Shape of dataset

```
In [14]: dataset.shape
```

```
Out[14]: (303, 14)
```

Printing out a few columns

```
In [15]: dataset.head(7)
```

```
Out[15]:
```

	age	sex	cp	trestbps	chol	fb	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	187	0	3.5	0	0	2	1
2	41	0	1	130	204	0	0	172	0	1.4	2	0	2	1
3	56	1	1	120	236	0	1	178	0	0.8	2	0	2	1
4	57	0	0	120	354	0	1	163	1	0.6	2	0	2	1
5	57	1	0	140	192	0	1	148	0	0.4	1	0	1	1
6	56	0	1	140	294	0	0	153	0	1.3	1	0	2	1

```
In [16]: dataset.sample(5)
```

```
Out[16]:
```

	age	sex	cp	trestbps	chol	fb	restecg	thalach	exang	oldpeak	slope	ca	thal	target
23	61	1	2	150	243	1	1	137	1	1.0	1	0	2	1
184	50	1	0	150	243	0	0	128	0	2.6	1	0	3	0
207	60	0	0	150	258	0	0	157	0	2.6	1	2	3	0
212	39	1	0	118	219	0	1	140	0	1.2	1	0	3	0
36	54	0	2	135	304	1	1	170	0	0.0	2	0	2	1

Description

```
: dataset.describe()
```

	age	sex	cp	trestbps	chol	fb	restecg	thalach	exang	oldpeak	slope	ca
count	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.
mean	54.366337	0.683168	0.966997	131.623762	246.264026	0.148515	0.528053	149.646865	0.326733	1.039604	1.399340	0.729373
std	9.082101	0.466011	1.032052	17.538143	51.830751	0.356198	0.525860	22.905161	0.469794	1.161075	0.616226	1.022606
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
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
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
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
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

```
: dataset.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)
```

In [19]: *###Luckily, we have no missing values*

Let's understand our columns better:

```
In [20]: info = ["age","1: male, 0: female","chest pain type, 1: typical angina, 2: atypical angina, 3: non-angina pain, 4: asymptomatic"]
for i in range(len(info)):
    print(dataset.columns[i]+":\t\t"+info[i])
age:                               age
sex:                             1: male, 0: female
cp:                            chest pain type, 1: typical angina, 2: atypical angina, 3: non-anginal pain, 4: asymptomatic
c:                                 
trestbps:                         resting blood pressure
chol:                            serum cholestral 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 flourosopy
thal:                            thal: 3 = normal; 6 = fixed defect; 7 = reversable defect
```

Analysing the 'target' variable

```
In [21]: dataset["target"].describe()
```

```
Out[21]: 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
```

```
In [22]: dataset["target"].unique()
```

```
Out[22]: array([1, 0])
```

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

Checking correlation between columns

```
In [23]: 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
```

```
fbs      0.028046
Name: target, dtype: float64
```

```
In [24]: #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:

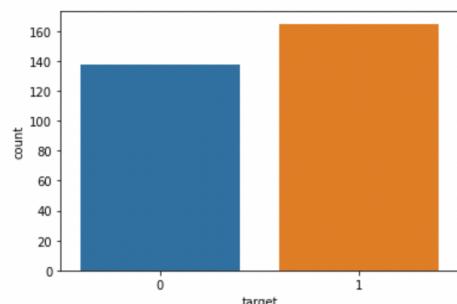
```
In [25]: y = dataset["target"]
```

```
sns.countplot(y)
```

```
target_temp = dataset.target.value_counts()
```

```
print(target_temp)
```

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



```
In [26]: 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

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

Analysing the 'Sex' feature

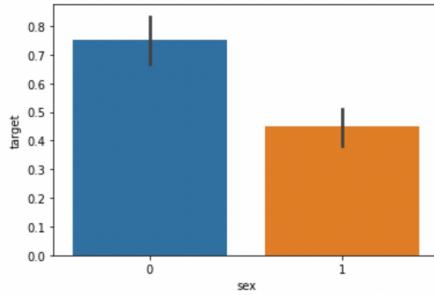
```
In [27]: dataset["sex"].unique()
```

Out[27]: array([1, 0])

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

```
In [28]: sns.barplot(dataset["sex"],y)
```

Out[28]: <AxesSubplot:xlabel='sex', ylabel='target'>



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

Analysing the 'Chest Pain Type' feature

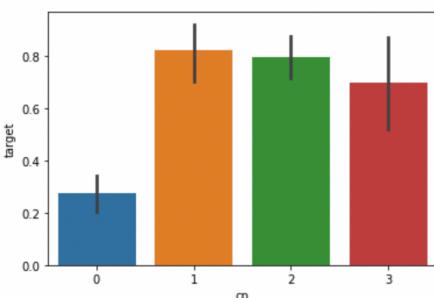
```
In [29]: dataset["cp"].unique()
```

Out[29]: array([3, 2, 1, 0])

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

```
In [30]: sns.barplot(dataset["cp"],y)
```

Out[30]: <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

```
In [31]: dataset["fbs"].describe()
```

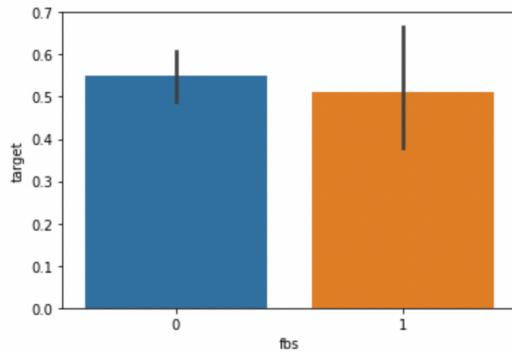
```
Out[31]: count    303.000000
mean      0.148515
std       0.356198
min       0.000000
25%       0.000000
50%       0.000000
75%       0.000000
max       1.000000
Name: fbs, dtype: float64
```

```
In [32]: dataset["fbs"].unique()
```

```
Out[32]: array([1, 0])
```

```
In [33]: sns.barplot(dataset["fbs"],y)
```

```
Out[33]: <AxesSubplot:xlabel='fbs', ylabel='target'>
```



Nothing extraordinary here

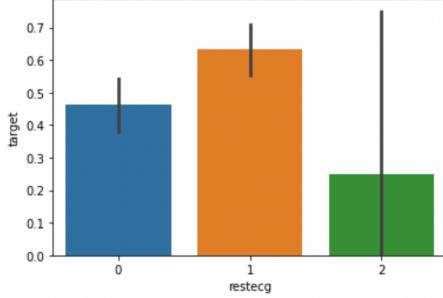
Analysing the restecg feature

```
In [34]: dataset["restecg"].unique()
```

```
Out[34]: array([0, 1, 2])
```

```
In [35]: sns.barplot(dataset["restecg"],y)
```

```
Out[35]: <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'

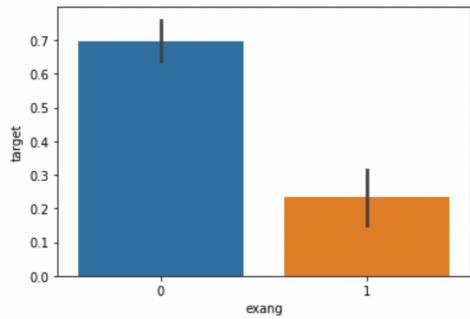
Analysing the 'exang' feature

```
In [36]: dataset["exang"].unique()
```

```
Out[36]: array([0, 1])
```

```
In [37]: sns.barplot(dataset["exang"],y)
```

```
Out[37]: <AxesSubplot:xlabel='exang', ylabel='target'>
```



People with `exang=1` i.e. Exercise induced angina are much less likely to have heart problems 7

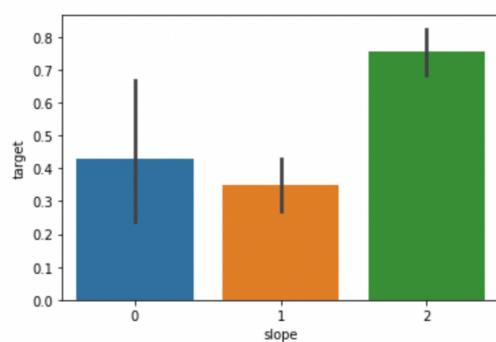
Analysing the Slope feature

```
In [38]: dataset["slope"].unique()
```

```
Out[38]: array([0, 2, 1])
```

```
In [39]: sns.barplot(dataset["slope"],y)
```

```
Out[39]: <AxesSubplot:xlabel='slope', ylabel='target'>
```



Analysing the 'ca' feature

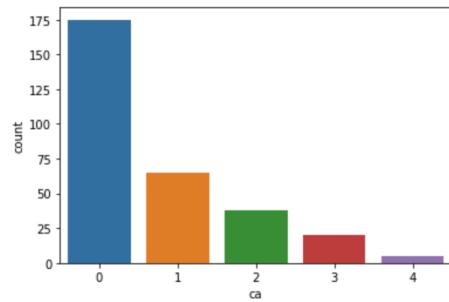
```
In [40]: #number of major vessels (0-3) colored by flourosopy
```

```
In [41]: dataset["ca"].unique()
```

```
Out[41]: array([0, 2, 1, 3, 4])
```

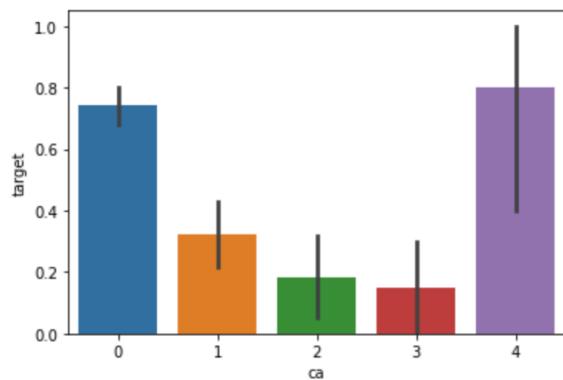
```
In [42]: sns.countplot(dataset["ca"])
```

```
Out[42]: <AxesSubplot:xlabel='ca', ylabel='count'>
```



```
In [43]: sns.barplot(dataset["ca"],y)
```

```
Out[43]: <AxesSubplot:xlabel='ca', ylabel='target'>
```



ca=4 has astonishingly large number of heart patients

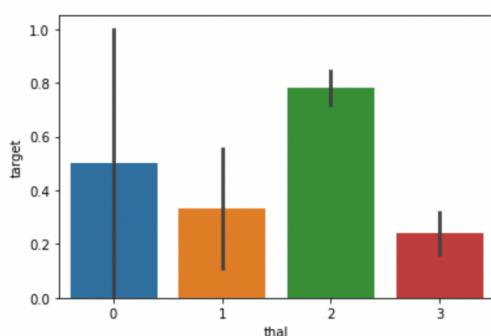
Analysing the 'thal' feature

```
In [45]: dataset["thal"].unique()
```

```
Out[45]: array([1, 2, 3, 0])
```

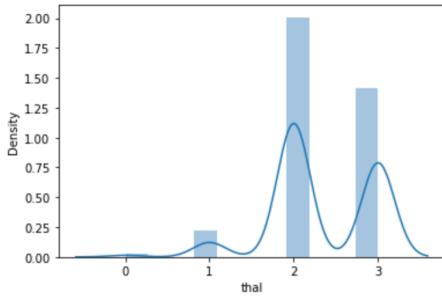
```
In [46]: sns.barplot(dataset["thal"],y)
```

```
Out[46]: <AxesSubplot:xlabel='thal', ylabel='target'>
```



```
In [47]: sns.distplot(dataset["thal"])

Out[47]: <AxesSubplot:xlabel='thal', ylabel='Density'>
```



IV. Train Test split

```
In [48]: 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_state=0)

In [49]: X_train.shape
Out[49]: (242, 13)

In [50]: X_test.shape
Out[50]: (61, 13)
```

```
In [51]: Y_train.shape
Out[51]: (242,)
```

```
In [52]: Y_test.shape
Out[52]: (61,)
```

V. Model Fitting

```
In [53]: from sklearn.metrics import accuracy_score
```

SVM

```
In [54]: from sklearn import svm

sv = svm.SVC(kernel='linear')

sv.fit(X_train, Y_train)

Y_pred_svm = sv.predict(X_test)

In [55]: Y_pred_svm.shape
Out[55]: (61,)
```

```
In [56]: print(Y_pred_svm)
[0 1 1 0 0 1 0 0 0 1 1 0 1 1 1 0 1 0 1 1 0 0 1 1 0 0 1 1 1 0 1 1 1 0
1 0 0 1 1 0 0 1 1 1 0 1 1 1 1 0 1 1 1 1]
```

```
In [57]: 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 %
```

Decision Tree

```
In [58]: 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)
```

```
In [59]: print(Y_pred_dt.shape)  
(61,)
```

```
In [60]: 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

```
In [ ]: 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)
```

```
In [ ]: Y_pred_rf.shape
```

```
In [ ]: score_rf = round(accuracy_score(Y_pred_rf,Y_test)*100,2)  
print("The accuracy score achieved using Random Forest is: "+str(score_rf)+" %")
```

XGBoost

```
In [107]: 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)

In [108]: Y_pred_xgb.shape
Out[108]: (61,)

In [109]: 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 %

In [110]: scores = [score_svm,score_dt,score_rf,score_xgb]
algorithms = ["Support Vector Machine","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 Support Vector Machine is: 81.97 %
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 %

In [111]: sns.set(rc={'figure.figsize':(15,8)})
plt.xlabel("Algorithms")
plt.ylabel("Accuracy score")

sns.barplot(algorithms,scores)

Out[111]: <AxesSubplot:xlabel='Algorithms', ylabel='Accuracy score'>
```

```
In [111]: sns.set(rc={'figure.figsize':(15,8)})
plt.xlabel("Algorithms")
plt.ylabel("Accuracy score")

sns.barplot(algorithms,scores)

Out[111]: <AxesSubplot:xlabel='Algorithms', ylabel='Accuracy score'>
```

