

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
```

## Read Data

```
In [2]: # We are reading our data
df = pd.read_excel("heart.xlsx")
```

```
In [3]: # First 5 rows of our data
df.head()
```

```
Out[3]:
```

	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	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

## Data Exploration

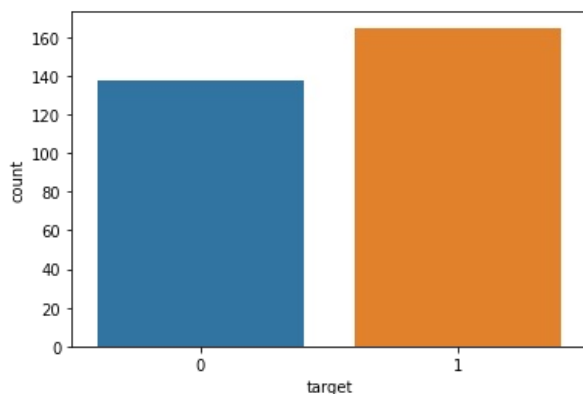
```
In [4]: df.target.value_counts()
```

```
Out[4]:
```

1	165
0	138

Name: target, dtype: int64

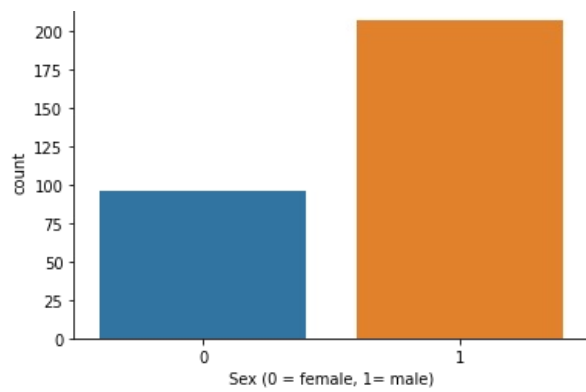
```
In [5]: sns.countplot(x="target", data=df)
plt.show()
```



```
In [6]: countNoDisease = len(df[df.target == 0])
countHaveDisease = len(df[df.target == 1])
print("Percentage of Patients Haven't Heart Disease: {:.2f}%".format((countNoDisease / (len(df.target))*100)))
print("Percentage of Patients Have Heart Disease: {:.2f}%".format((countHaveDisease / (len(df.target))*100)))
```

Percentage of Patients Haven't Heart Disease: 45.54%  
 Percentage of Patients Have Heart Disease: 54.46%

```
In [7]: sns.countplot(x='sex', data=df)
plt.xlabel("Sex (0 = female, 1= male)")
plt.show()
```



```
In [8]: countFemale = len(df[df.sex == 0])
countMale = len(df[df.sex == 1])
print("Percentage of Female Patients: {:.2f}%".format((countFemale / (len(df.sex))*100)))
print("Percentage of Male Patients: {:.2f}%".format((countMale / (len(df.sex))*100)))
```

Percentage of Female Patients: 31.68%  
Percentage of Male Patients: 68.32%

```
In [9]: df.groupby('target').mean()
```

```
Out[9]:
```

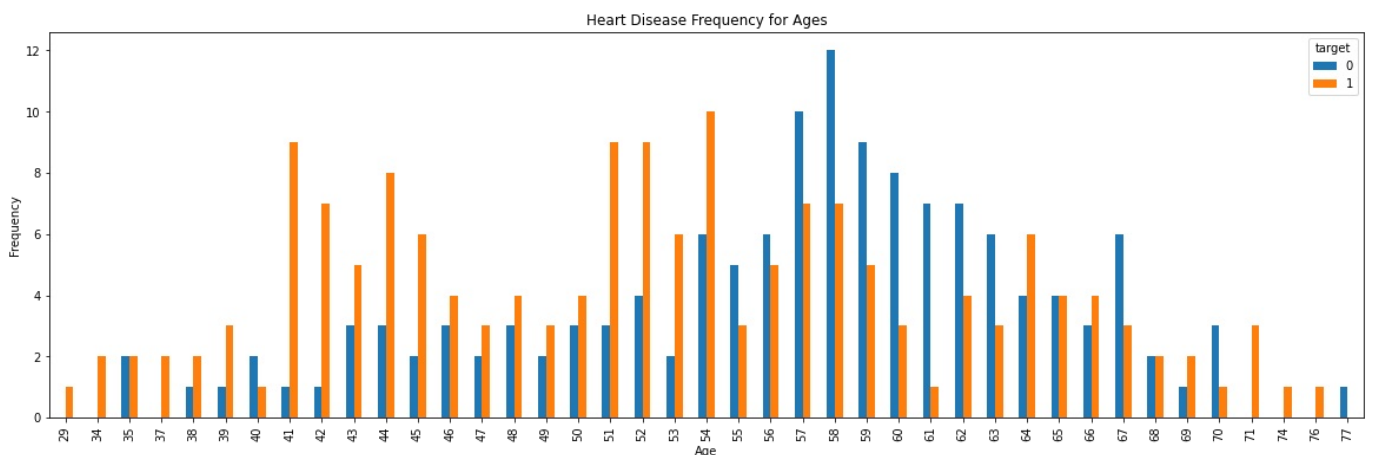
	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal
target													
0	56.601449	0.826087	0.478261	134.398551	251.086957	0.159420	0.449275	139.101449	0.550725	1.585507	1.166667	1.166667	2.543478
1	52.496970	0.563636	1.375758	129.303030	242.230303	0.139394	0.593939	158.466667	0.139394	0.583030	1.593939	0.363636	2.121212

```
In [10]: Age=pd.crosstab(df.age,df.target)
Age.head()
```

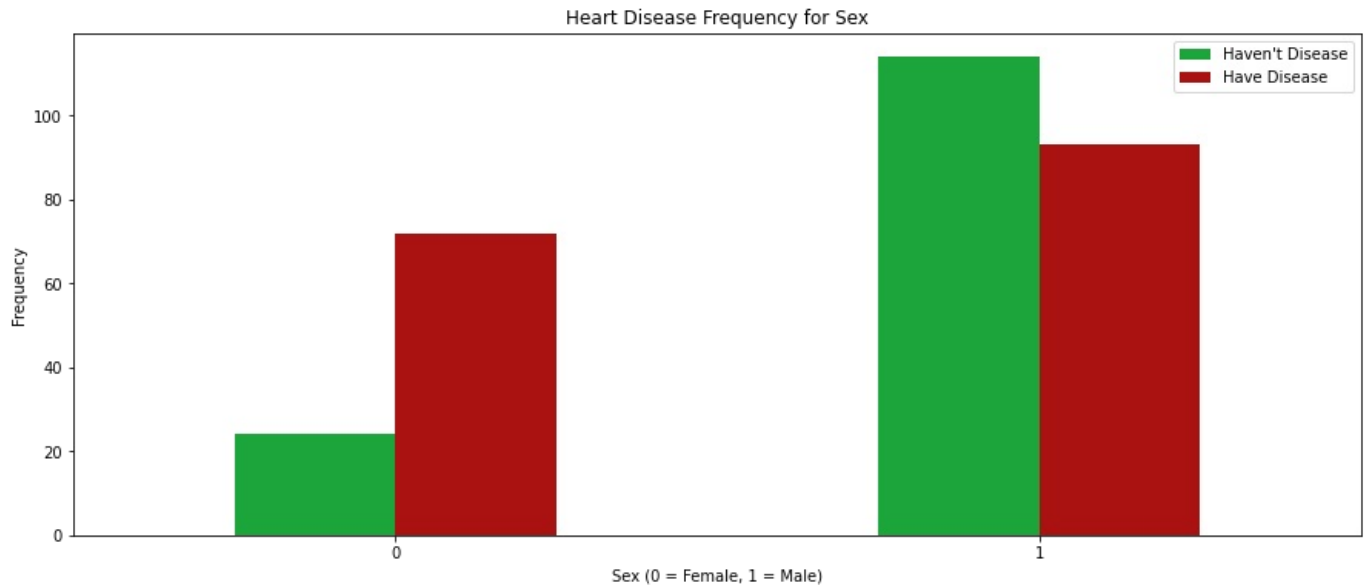
```
Out[10]:
```

target	0	1
age		
29	0	1
34	0	2
35	2	2
37	0	2
38	1	2

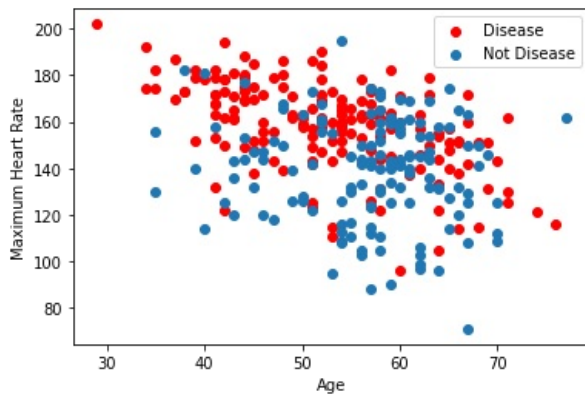
```
In [11]: Age.plot(kind="bar",figsize=(20,6))
plt.title('Heart Disease Frequency for Ages')
plt.xlabel('Age')
plt.ylabel('Frequency')
plt.savefig('heartDiseaseAndAges.png')
plt.show()
```



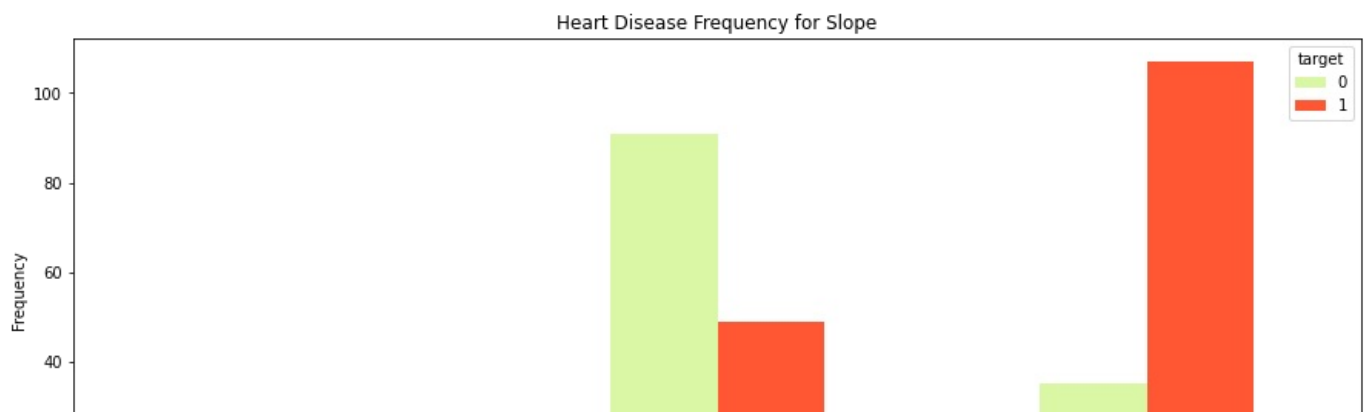
```
In [12]: pd.crosstab(df.sex,df.target).plot(kind="bar",figsize=(15,6),color=['#1CA53B','#AA1111' ])
plt.title('Heart Disease Frequency for Sex')
plt.xlabel('Sex (0 = Female, 1 = Male)')
plt.xticks(rotation=0)
plt.legend(["Haven't Disease", "Have Disease"])
plt.ylabel('Frequency')
plt.show()
```



```
In [13]: plt.scatter(x=df.age[df.target==1], y=df.thalach[(df.target==1)], c="red")
plt.scatter(x=df.age[df.target==0], y=df.thalach[(df.target==0)])
plt.legend(["Disease", "Not Disease"])
plt.xlabel("Age")
plt.ylabel("Maximum Heart Rate")
plt.show()
```

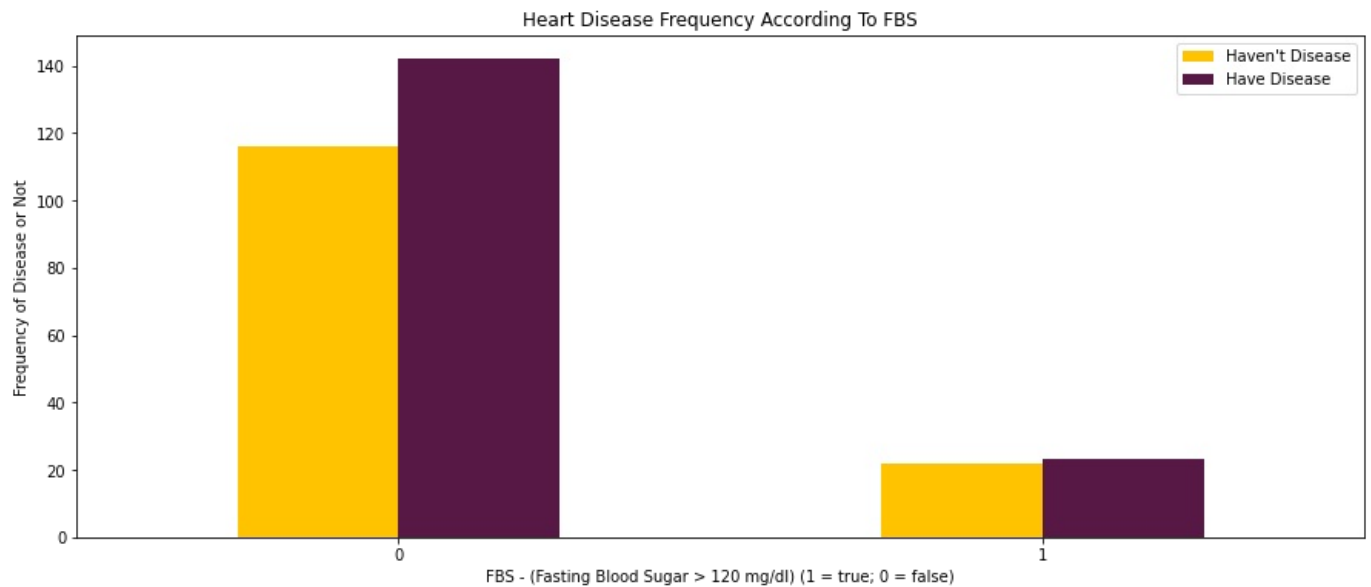


```
In [14]: pd.crosstab(df.slope,df.target).plot(kind="bar",figsize=(15,6),color=['#DAF7A6','#FF5733' ])
plt.title('Heart Disease Frequency for Slope')
plt.xlabel('The Slope of The Peak Exercise ST Segment ')
plt.xticks(rotation = 0)
plt.ylabel('Frequency')
plt.show()
```

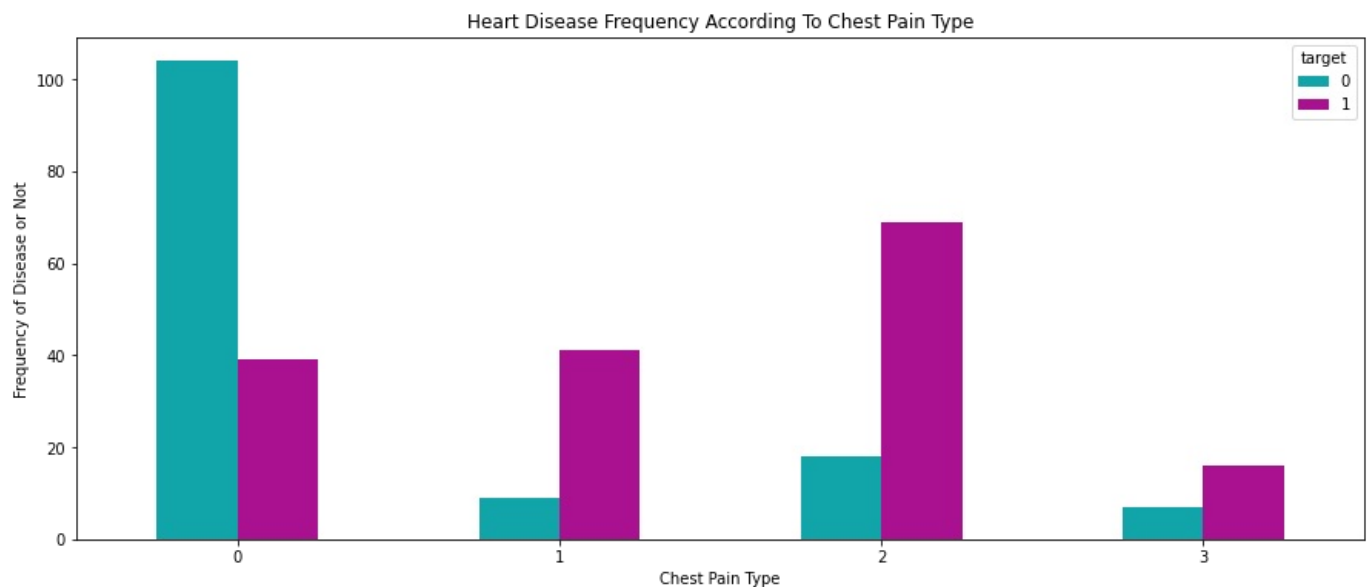




```
In [15]: pd.crosstab(df.fbs,df.target).plot(kind="bar",figsize=(15,6),color=['#FFC300','#581845' ])
plt.title('Heart Disease Frequency According To FBS')
plt.xlabel('FBS - (Fasting Blood Sugar > 120 mg/dl) (1 = true; 0 = false)')
plt.xticks(rotation = 0)
plt.legend(["Haven't Disease", "Have Disease"])
plt.ylabel('Frequency of Disease or Not')
plt.show()
```



```
In [16]: pd.crosstab(df.cp,df.target).plot(kind="bar",figsize=(15,6),color=['#11A5AA','#AA1190' ])
plt.title('Heart Disease Frequency According To Chest Pain Type')
plt.xlabel('Chest Pain Type')
plt.xticks(rotation = 0)
plt.ylabel('Frequency of Disease or Not')
plt.show()
```



## Creating Dummy Variables

Since 'cp', 'thal' and 'slope' are categorical variables we'll turn them into dummy variables.

```
In [17]: a = pd.get_dummies(df['cp'], prefix = "cp")
```

```
a = pd.get_dummies(df['cp'], prefix = "cp",
b = pd.get_dummies(df['thal'], prefix = "thal")
c = pd.get_dummies(df['slope'], prefix = "slope")
```

```
In [18]: frames = [df, a, b, c]
df = pd.concat(frames, axis = 1)
df.head()
```

```
Out[18]:
```

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	...	cp_1	cp_2	cp_3	thal_0	thal_1	thal_2	thal_3	slope_0	slope_1
0	63	1	3	145	233	1	0	150	0	2.3	...	0	0	1	0	1	0	0	1	0
1	37	1	2	130	250	0	1	187	0	3.5	...	0	1	0	0	0	1	0	1	0
2	41	0	1	130	204	0	0	172	0	1.4	...	1	0	0	0	0	1	0	0	0
3	56	1	1	120	236	0	1	178	0	0.8	...	1	0	0	0	0	1	0	0	0
4	57	0	0	120	354	0	1	163	1	0.6	...	0	0	0	0	0	1	0	0	0

5 rows × 25 columns

```
In [19]: df = df.drop(columns = ['cp', 'thal', 'slope'])
df.head()
```

```
Out[19]:
```

	age	sex	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	ca	...	cp_1	cp_2	cp_3	thal_0	thal_1	thal_2	thal_3	slope_0	slope_1
0	63	1	145	233	1	0	150	0	2.3	0	...	0	0	1	0	1	0	0	1	0
1	37	1	130	250	0	1	187	0	3.5	0	...	0	1	0	0	0	1	0	1	0
2	41	0	130	204	0	0	172	0	1.4	0	...	1	0	0	0	0	1	0	0	0
3	56	1	120	236	0	1	178	0	0.8	0	...	1	0	0	0	0	1	0	0	0
4	57	0	120	354	0	1	163	1	0.6	0	...	0	0	0	0	0	1	0	0	0

5 rows × 22 columns

## Creating Model for Logistic Regression

We can use sklearn library or we can write functions ourselves. Let's them both. Firstly we will write our functions after that we'll use sklearn library to calculate score.

```
In [20]: y = df.target.values
x_data = df.drop(['target'], axis = 1)
```

## Normalize Data

```
In [21]: # Normalize
x = (x_data - np.min(x_data)) / (np.max(x_data) - np.min(x_data)).values
```

```
In [22]: x_train, x_test, y_train, y_test = train_test_split(x,y,test_size = 0.2,random_state=0)
```

```
In [23]: #transpose matrices
x_train = x_train.T
y_train = y_train.T
x_test = x_test.T
y_test = y_test.T
```

Let's say weight = 0.01 and bias = 0.0

```
In [24]: #initialize
def initialize(dimension):

    weight = np.full((dimension,1),0.01)
    bias = 0.0
    return weight,bias
```

## Siamoid Function

```
In [25]: def sigmoid(z):
        y_head = 1/(1+ np.exp(-z))
        return y_head
```

## Gradient Descent

```
In [26]: def forwardBackward(weight,bias,x_train,y_train):
        # Forward

        y_head = sigmoid(np.dot(weight.T,x_train) + bias)
        loss = -(y_train*np.log(y_head) + (1-y_train)*np.log(1-y_head))
        cost = np.sum(loss) / x_train.shape[1]

        # Backward
        derivative_weight = np.dot(x_train,((y_head-y_train).T))/x_train.shape[1]
        derivative_bias = np.sum(y_head-y_train)/x_train.shape[1]
        gradients = {"Derivative Weight" : derivative_weight, "Derivative Bias" : derivative_bias}

        return cost,gradients
```

```
In [27]: def update(weight,bias,x_train,y_train,learningRate,iteration) :
        costList = []
        index = []

        #for each iteration, update weight and bias values
        for i in range(iteration):
            cost,gradients = forwardBackward(weight,bias,x_train,y_train)
            weight = weight - learningRate * gradients["Derivative Weight"]
            bias = bias - learningRate * gradients["Derivative Bias"]

            costList.append(cost)
            index.append(i)

        parameters = {"weight": weight,"bias": bias}

        print("iteration:",iteration)
        print("cost:",cost)

        plt.plot(index,costList)
        plt.xlabel("Number of Iteration")
        plt.ylabel("Cost")
        plt.show()

        return parameters, gradients
```

```
In [28]: def predict(weight,bias,x_test):
        z = np.dot(weight.T,x_test) + bias
        y_head = sigmoid(z)

        y_prediction = np.zeros((1,x_test.shape[1]))

        for i in range(y_head.shape[1]):
            if y_head[0,i] <= 0.5:
                y_prediction[0,i] = 0
            else:
                y_prediction[0,i] = 1
        return y_prediction
```

```
In [29]: def logistic_regression(x_train,y_train,x_test,y_test,learningRate,iteration):
        dimension = x_train.shape[0]
        weight,bias = initialize(dimension)

        parameters, gradients = update(weight,bias,x_train,y_train,learningRate,iteration)

        y_prediction = predict(parameters["weight"],parameters["bias"],x_test)

        print("Manuel Test Accuracy: {:.2f}%".format((100 - np.mean(np.abs(y_prediction - y_test))*100)/100*100))
```

logistic\_regression(x\_train,y\_train,x\_test,y\_test,1,100)

## Sklearn Logistic Regression

```
In [30]: lr = LogisticRegression()
lr.fit(x_train.T,y_train.T)
print("Test Accuracy {:.2f}%".format(lr.score(x_test.T,y_test.T)*100))
```

Test Accuracy 86.89%

## K-Nearest Neighbour (KNN) Classification

```
In [31]: # KNN Model
from sklearn.neighbors import KNeighborsClassifier
knn = KNeighborsClassifier(n_neighbors = 2) # n_neighbors means k
knn.fit(x_train.T, y_train.T)
prediction = knn.predict(x_test.T)

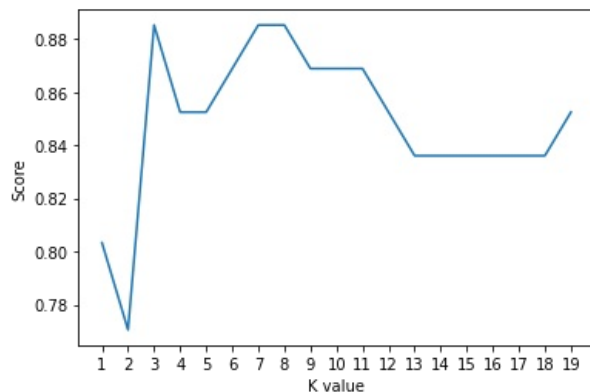
print("{} NN Score: {:.2f}%".format(2, knn.score(x_test.T, y_test.T)*100))
```

2 NN Score: 77.05%

```
In [32]: # try to find best k value
scoreList = []
for i in range(1,20):
    knn2 = KNeighborsClassifier(n_neighbors = i) # n_neighbors means k
    knn2.fit(x_train.T, y_train.T)
    scoreList.append(knn2.score(x_test.T, y_test.T))

plt.plot(range(1,20), scoreList)
plt.xticks(np.arange(1,20,1))
plt.xlabel("K value")
plt.ylabel("Score")
plt.show()

print("Maximum KNN Score is {:.2f}%".format((max(scoreList))*100))
```



Maximum KNN Score is 88.52%

## Support Vector Machine (SVM) Algorithm

Now we will use SVM algorithm.

```
In [33]: from sklearn.svm import SVC
svm = SVC(random_state = 1)
svm.fit(x_train.T, y_train.T)
```

```
Out[33]: SVC
SVC(random_state=1)
```

```
In [34]: print("Test Accuracy of SVM Algorithm: {:.2f}%".format(svm.score(x_test.T,y_test.T)*100))
```

Test Accuracy of SVM Algorithm: 88.52%

## Naive Bayes Algorithm

```
In [35]: from sklearn.naive_bayes import GaussianNB
nb = GaussianNB()
nb.fit(x_train.T, y_train.T)
print("Accuracy of Naive Bayes: {:.2f}%".format(nb.score(x_test.T,y_test.T)*100))
```

Accuracy of Naive Bayes: 86.89%

## Decision Tree Algorithm

```
In [36]: from sklearn.tree import DecisionTreeClassifier
dtc = DecisionTreeClassifier()
dtc.fit(x_train.T, y_train.T)
print("Decision Tree Test Accuracy {:.2f}%".format(dtc.score(x_test.T, y_test.T)*100))
```

Decision Tree Test Accuracy 77.05%

## Random Forest Classification

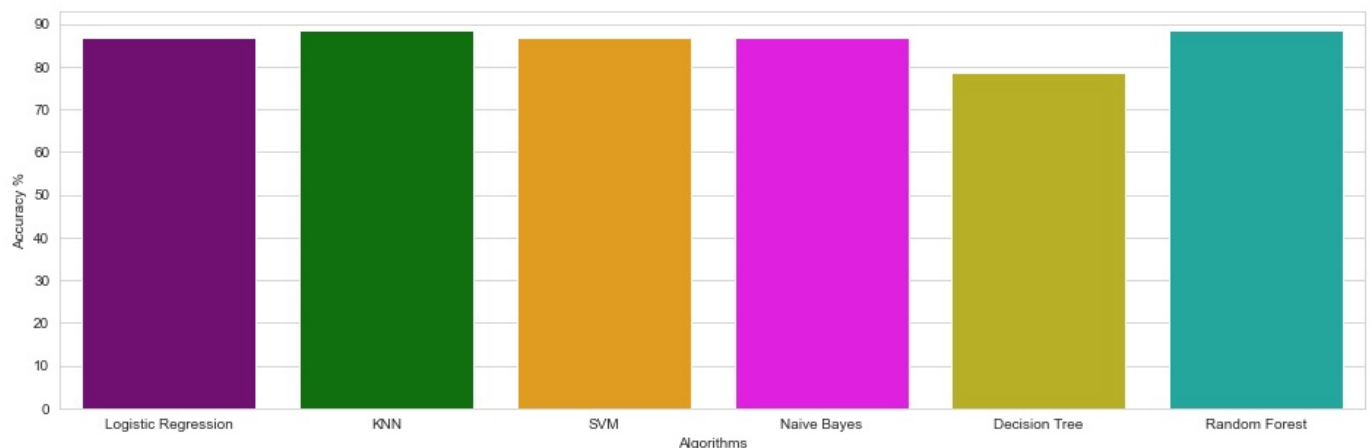
```
In [37]: # Random Forest Classification
from sklearn.ensemble import RandomForestClassifier
rf = RandomForestClassifier(n_estimators = 1000, random_state = 1)
rf.fit(x_train.T, y_train.T)
print("Random Forest Algorithm Accuracy Score : {:.2f}%".format(rf.score(x_test.T,y_test.T)*100))
```

Random Forest Algorithm Accuracy Score : 88.52%

## Comparing Models

```
In [38]: methods = ["Logistic Regression", "KNN", "SVM", "Naive Bayes", "Decision Tree", "Random Forest"]
accuracy = [86.89, 88.52, 86.89, 86.89, 78.69, 88.52]
colors = ["purple", "green", "orange", "magenta", "#CFC60E", "#0FBBAE"]

sns.set_style("whitegrid")
plt.figure(figsize=(16,5))
plt.yticks(np.arange(0,100,10))
plt.ylabel("Accuracy %")
plt.xlabel("Algorithms")
sns.barplot(x=methods, y=accuracy, palette=colors)
plt.show()
```





# Confusion Matrix

```
In [39]: # Predicted values
y_head_lr = lr.predict(x_test.T)
knn3 = KNeighborsClassifier(n_neighbors = 3)
knn3.fit(x_train.T, y_train.T)
y_head_knn = knn3.predict(x_test.T)
y_head_svm = svm.predict(x_test.T)
y_head_nb = nb.predict(x_test.T)
y_head_dtc = dtc.predict(x_test.T)
y_head_rf = rf.predict(x_test.T)
```

```
In [40]: from sklearn.metrics import confusion_matrix

cm_lr = confusion_matrix(y_test,y_head_lr)
cm_knn = confusion_matrix(y_test,y_head_knn)
cm_svm = confusion_matrix(y_test,y_head_svm)
cm_nb = confusion_matrix(y_test,y_head_nb)
cm_dtc = confusion_matrix(y_test,y_head_dtc)
cm_rf = confusion_matrix(y_test,y_head_rf)
```

```
In [41]: plt.figure(figsize=(24,12))

plt.suptitle("Confusion Matrixes",fontsize=24)
plt.subplots_adjust(wspace = 0.4, hspace= 0.4)

plt.subplot(2,3,1)
plt.title("Logistic Regression Confusion Matrix")
sns.heatmap(cm_lr,annot=True,cmap="Blues",fmt="d",cbar=False)

plt.subplot(2,3,2)
plt.title("K Nearest Neighbors Confusion Matrix")
sns.heatmap(cm_knn,annot=True,cmap="Blues",fmt="d",cbar=False)

plt.subplot(2,3,3)
plt.title("Support Vector Machine Confusion Matrix")
sns.heatmap(cm_svm,annot=True,cmap="Blues",fmt="d",cbar=False)

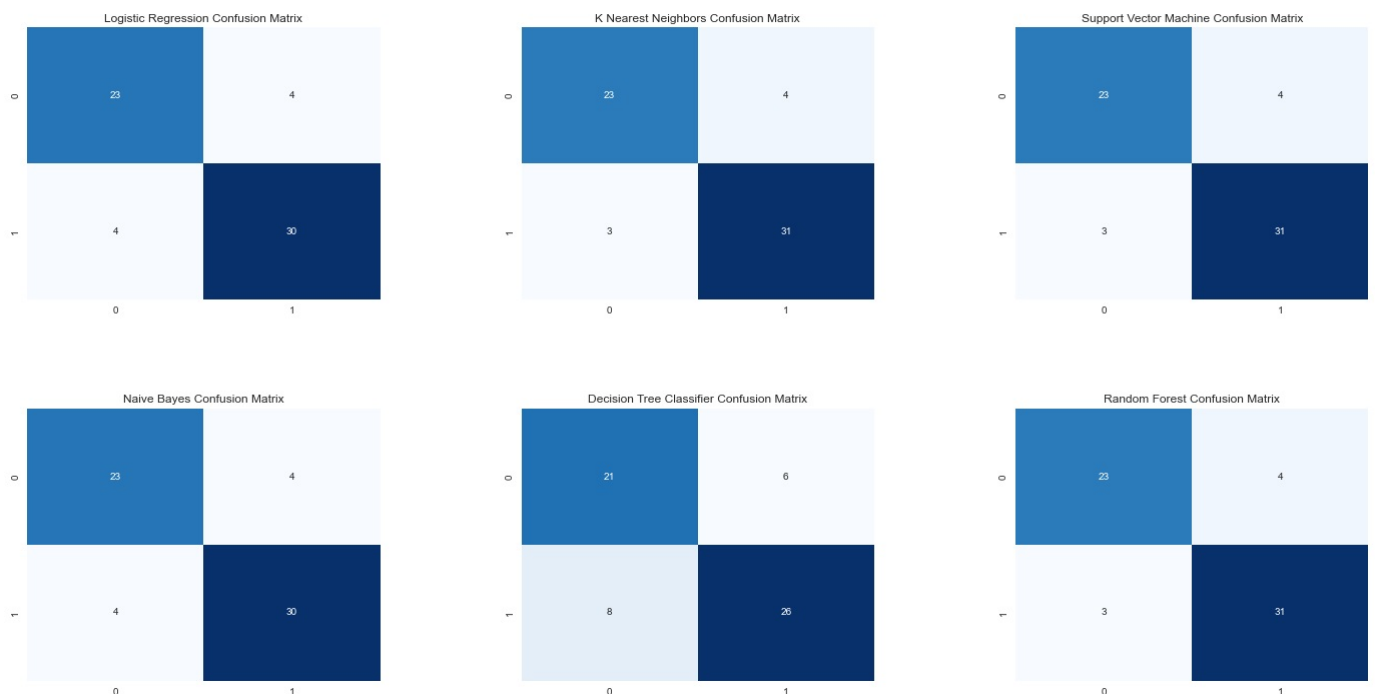
plt.subplot(2,3,4)
plt.title("Naive Bayes Confusion Matrix")
sns.heatmap(cm_nb,annot=True,cmap="Blues",fmt="d",cbar=False)

plt.subplot(2,3,5)
plt.title("Decision Tree Classifier Confusion Matrix")
sns.heatmap(cm_dtc,annot=True,cmap="Blues",fmt="d",cbar=False)

plt.subplot(2,3,6)
plt.title("Random Forest Confusion Matrix")
sns.heatmap(cm_rf,annot=True,cmap="Blues",fmt="d",cbar=False)

plt.show()
```

Confusion Matrixes



In [ ]:

Loading [MathJax]/jax/output/CommonHTML/fonts/TeX/fontdata.js