

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

```
In [6]: # Load the data
titanic_data = pd.read_csv('train.csv')
```

```
In [7]: #print first five rows of data
titanic_data.head()
```

Out[7]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	NaN	S
1	2	1	1	Cummings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85	C
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	NaN	S
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	C123	S
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	NaN	S

```
In [8]: #print last five rows of the data
titanic_data.tail()
```

Out[8]:

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
886	887	0	2	Montvila, Rev. Juozas	male	27.0	0	0	211536	13.00	NaN	S
887	888	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.00	B42	S
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	NaN	1	2	W./C. 6607	23.45	NaN	S
889	890	1	1	Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.00	C148	C
890	891	0	3	Dooley, Mr. Patrick	male	32.0	0	0	370376	7.75	NaN	Q

```
In [9]: # finding no. of rows and columns
titanic_data.shape
```

Out[9]: (891, 12)

```
In [10]: # getting some information about data
titanic_data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 12 columns):
 #   Column        Non-Null Count  Dtype  
---  -
 0   PassengerId   891 non-null   int64  
 1   Survived      891 non-null   int64  
 2   Pclass        891 non-null   int64  
 3   Name          891 non-null   object  
 4   Sex           891 non-null   object  
 5   Age           714 non-null   float64 
 6   SibSp         891 non-null   int64  
 7   Parch         891 non-null   int64  
 8   Ticket        891 non-null   object  
 9   Fare          891 non-null   float64 
10   Cabin         204 non-null   object  
11   Embarked      889 non-null   object  
dtypes: float64(2), int64(5), object(5)
memory usage: 83.7+ KB
```

```
In [11]: # checking the missing values in each column
titanic_data.isnull().sum()
```

```
Out[11]: PassengerId    0
Survived              0
Pclass                0
Name                  0
Sex                   0
Age                  177
SibSp                 0
Parch                 0
Ticket                0
Fare                  0
Cabin                 687
Embarked              2
dtype: int64
```

```
In [13]: #drop the Cabin column from data set as it has too many missing values
titanic_data = titanic_data.drop(columns = 'Cabin', axis = 1 )
```

```
In [14]: titanic_data.isnull().sum()
```

```
Out[14]: PassengerId    0
Survived              0
Pclass                0
Name                  0
Sex                   0
Age                  177
SibSp                 0
Parch                 0
Ticket                0
Fare                  0
Embarked              2
dtype: int64
```

```
In [16]: # replacing value in the age column with the mean value
titanic_data['Age'].fillna(titanic_data['Age'].mean(), inplace = True)
```

```
In [17]: titanic_data.isnull().sum()
```

```
Out[17]: PassengerId    0
         Survived      0
         Pclass      0
         Name        0
         Sex         0
         Age         0
         SibSp       0
         Parch       0
         Ticket      0
         Fare        0
         Embarked    2
         dtype: int64
```

```
In [21]: # finding the mode value in Embarked column
         print(titanic_data['Embarked'].mode())
         print(titanic_data['Embarked'].mode()[0])
```

```
0    S
Name: Embarked, dtype: object
S
```

```
In [22]: # replacing the missing values in Embarked column with mode value
         titanic_data['Embarked'].fillna(titanic_data['Embarked'].mode()[0], inplace = True)
```

```
In [24]: titanic_data.isnull().sum()
```

```
Out[24]: PassengerId    0
         Survived      0
         Pclass      0
         Name        0
         Sex         0
         Age         0
         SibSp       0
         Parch       0
         Ticket      0
         Fare        0
         Embarked    0
         dtype: int64
```

```
In [25]: # now data analysis
         #getting some statical measure about data
         titanic_data.describe()
```

```
Out[25]:
```

	PassengerId	Survived	Pclass	Age	SibSp	Parch	Fare
<b>count</b>	891.000000	891.000000	891.000000	891.000000	891.000000	891.000000	891.000000
<b>mean</b>	446.000000	0.383838	2.308642	29.699118	0.523008	0.381594	32.204208
<b>std</b>	257.353842	0.486592	0.836071	13.002015	1.102743	0.806057	49.693429
<b>min</b>	1.000000	0.000000	1.000000	0.420000	0.000000	0.000000	0.000000
<b>25%</b>	223.500000	0.000000	2.000000	22.000000	0.000000	0.000000	7.910400
<b>50%</b>	446.000000	0.000000	3.000000	29.699118	0.000000	0.000000	14.454200
<b>75%</b>	668.500000	1.000000	3.000000	35.000000	1.000000	0.000000	31.000000
<b>max</b>	891.000000	1.000000	3.000000	80.000000	8.000000	6.000000	512.329200

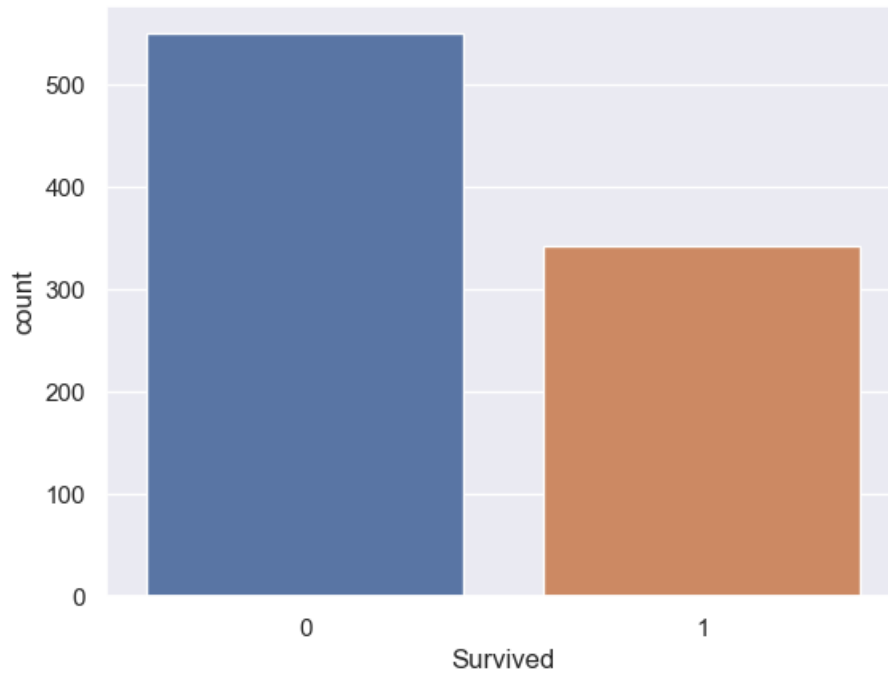
```
In [26]: # finding tha no. of people survived and not survived
         titanic_data['Survived'].value_counts()
```

```
Out[26]: Survived
0      549
1      342
Name: count, dtype: int64
```

```
In [48]: # data visualization
         sns.set()
```

```
In [51]: # making a count plot for survived column  
sns.countplot(x = 'Survived', data = titanic_data)
```

```
Out[51]: <Axes: xlabel='Survived', ylabel='count'>
```

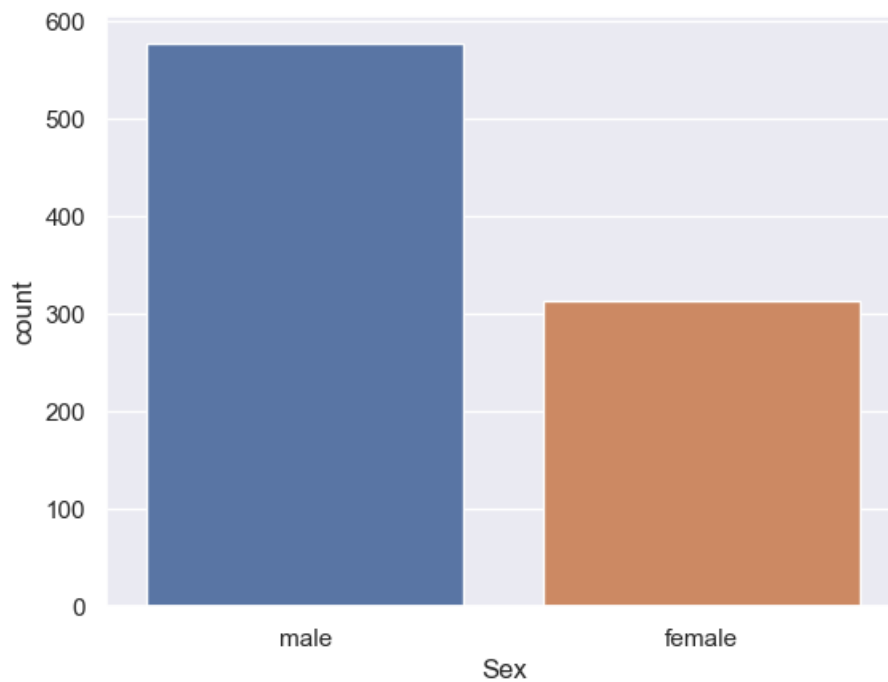


```
In [52]: titanic_data['Sex'].value_counts()
```

```
Out[52]: Sex  
male      577  
female    314  
Name: count, dtype: int64
```

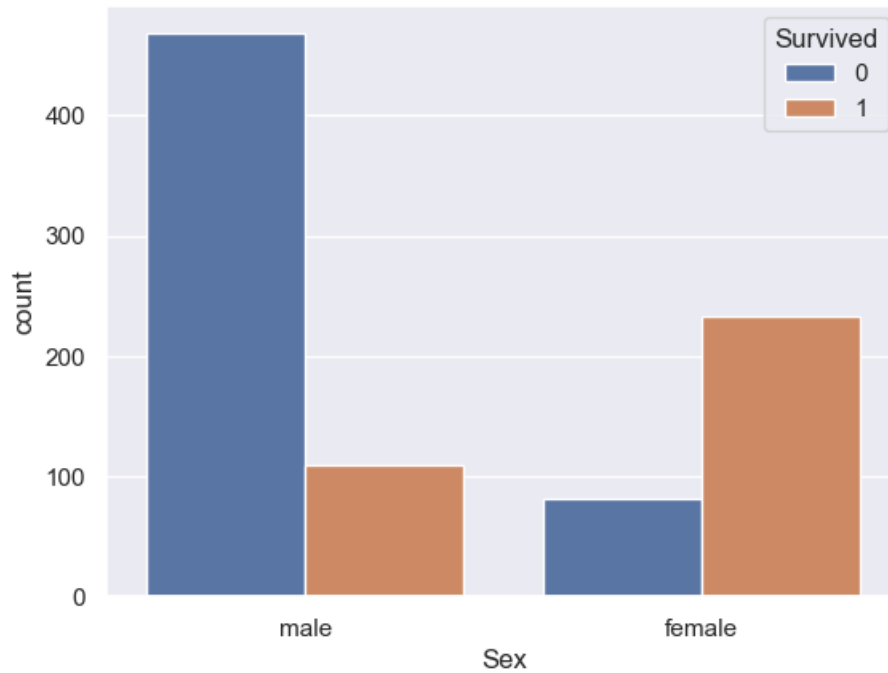
```
In [53]: sns.countplot(x = 'Sex', data = titanic_data)
```

```
Out[53]: <Axes: xlabel='Sex', ylabel='count'>
```



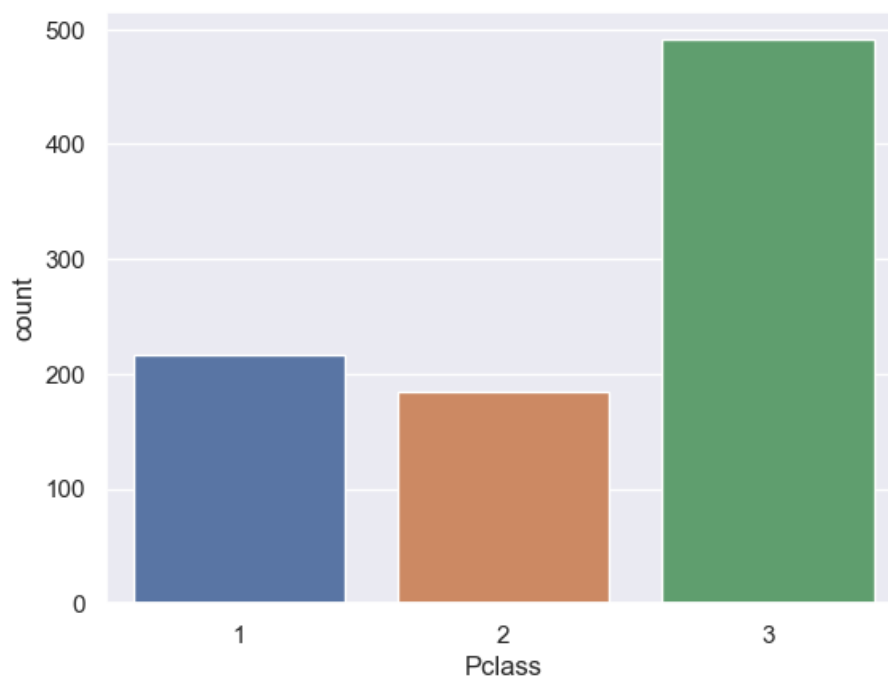
```
In [60]: # no. of survivors gender wise  
sns.countplot(x = 'Sex', hue = 'Survived', data = titanic_data)
```

Out[60]: <Axes: xlabel='Sex', ylabel='count'>



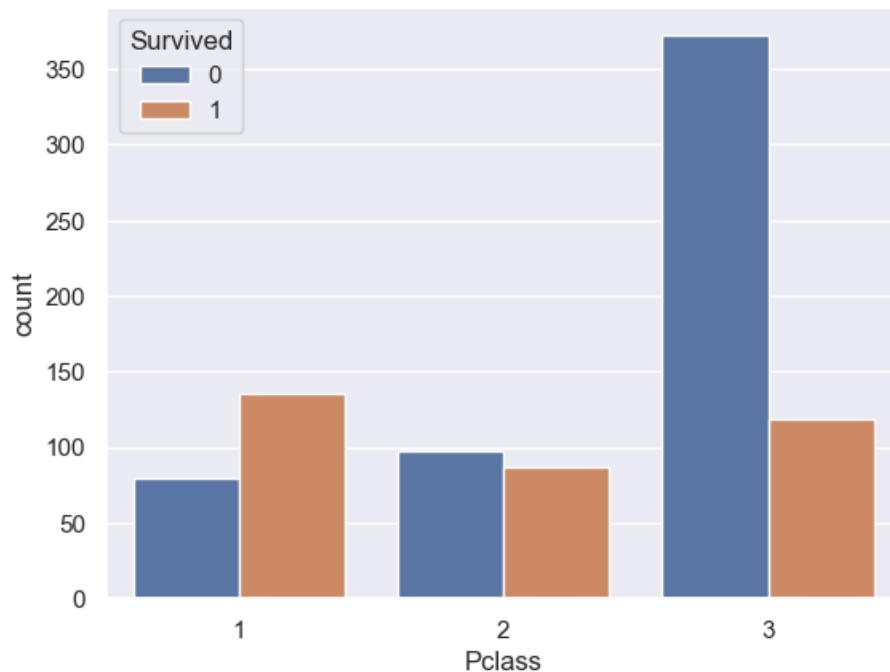
```
In [61]: # making a countplot for Pclass  
sns.countplot(x = 'Pclass', data = titanic_data)
```

Out[61]: <Axes: xlabel='Pclass', ylabel='count'>



```
In [63]: # find the no of people survived based on Pclass
sns.countplot(x = 'Pclass', hue = 'Survived', data = titanic_data)
```

```
Out[63]: <Axes: xlabel='Pclass', ylabel='count'>
```



```
In [72]: # encoding the categorical column
titanic_data['Sex'].value_counts()
```

```
Out[72]: Sex
male      577
female    314
Name: count, dtype: int64
```

```
In [73]: titanic_data['Embarked'].value_counts()
```

```
Out[73]: Embarked
S      646
C      168
Q       77
Name: count, dtype: int64
```

```
In [80]: # converting categorical column
titanic_data.replace({'Sex':{'male':0, 'female':1}, 'Embarked':{'S':0, 'C':1, 'Q':2}}, inplace = True)
```

```
In [81]: titanic_data.head()
```

```
Out[81]:
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Embarked
0	1	0	3	Braund, Mr. Owen Harris	0	22.0	1	0	A/5 21171	7.2500	0
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...)	1	38.0	1	0	PC 17599	71.2833	1
2	3	1	3	Heikkinen, Miss. Laina	1	26.0	0	0	STON/O2. 3101282	7.9250	0
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	1	35.0	1	0	113803	53.1000	0
4	5	0	3	Allen, Mr. William Henry	0	35.0	0	0	373450	8.0500	0

```
In [83]: # seperating features and target
x = titanic_data.drop(columns = ['PassengerId', 'Name', 'Ticket', 'Survived'])
y = titanic_data['Survived']
```

In [84]: `print(x)`

	Pclass	Sex	Age	SibSp	Parch	Fare	Embarked
0	3	0	22.000000	1	0	7.2500	0
1	1	1	38.000000	1	0	71.2833	1
2	3	1	26.000000	0	0	7.9250	0
3	1	1	35.000000	1	0	53.1000	0
4	3	0	35.000000	0	0	8.0500	0
..	...	...	...	...	...	...	...
886	2	0	27.000000	0	0	13.0000	0
887	1	1	19.000000	0	0	30.0000	0
888	3	1	29.699118	1	2	23.4500	0
889	1	0	26.000000	0	0	30.0000	1
890	3	0	32.000000	0	0	7.7500	2

[891 rows x 7 columns]

In [85]: `print(y)`

```
0    0
1    1
2    1
3    1
4    0
..
886  0
887  1
888  0
889  1
890  0
```

Name: Survived, Length: 891, dtype: int64

In [89]: *#splitting data into training data and test data*

```
x_train, x_test, y_train, y_test = train_test_split(x,y, test_size = 0.2, random_state = 2)
```

In [90]: `print(x.shape, x_train.shape, x_test.shape)`

(891, 7) (712, 7) (179, 7)

In [94]: *# model training - we are using LogisticRegression model*

```
model = LogisticRegression()
```

In [96]: *#training the logistic regression with training data*

```
model.fit(x_train, y_train)
```

C:\Users\91993\anaconda3\Lib\site-packages\sklearn\linear\_model\\_logistic.py:460: ConvergenceWarning: lbfgs failed to converge (status=1):  
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max\_iter) or scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html> (<https://scikit-learn.org/stable/modules/preprocessing.html>)

Please also refer to the documentation for alternative solver options:

[https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression) ([https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression))

```
n_iter_i = _check_optimize_result(
```

Out[96]:

```
LogisticRegression
LogisticRegression()
```

In [99]: *# model evaluation*

*# accuracy on training data*

```
x_train_prediction = model.predict(x_train)
```

In [100]: `print(x_train_prediction)`

```
[0 1 0 0 0 0 0 1 0 0 0 1 0 0 1 0 1 0 0 0 0 0 1 0 0 1 0 0 1 0 1 1 0 0 1 0 1
0 0 0 0 0 0 1 1 0 0 1 0 1 0 1 0 0 0 0 0 0 0 1 0 1 0 0 1 1 0 0 1 1 0 1 0 0 1
0 0 0 0 0 0 1 0 0 0 1 0 0 0 1 0 1 0 0 1 0 0 0 1 1 1 0 1 0 0 0 0 0 1 0 0 0
1 1 0 0 1 0 0 1 0 0 1 0 0 1 0 1 0 1 0 1 1 1 1 1 1 0 0 1 1 1 0 0 1 0 0
0 0 0 0 1 0 1 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 1 1 0 0 1 0 1 0 1 1 1
0 0 0 1 0 0 0 1 0 0 1 0 0 0 1 1 0 1 0 0 0 0 0 0 1 1 0 1 1 1 1 0 0 0 0 0 0 0
0 1 0 0 1 1 1 0 0 1 0 1 1 1 0 0 1 0 0 0 1 0 0 0 1 0 0 0 1 0 1 0 1 0 0 0
0 0 0 0 0 0 1 0 1 0 0 1 0 0 1 0 1 0 1 1 0 0 0 0 1 0 1 0 0 1 0 0 0 1 0 0 0
0 1 1 0 0 0 0 0 0 1 0 1 0 0 0 0 0 1 1 1 0 0 0 1 0 1 0 0 0 0 0 0 1 1 0 1 1
0 1 1 1 0 0 0 0 0 0 0 0 0 0 1 0 0 1 1 1 0 1 0 0 0 0 1 1 0 0 0 1 0 1 1 1 0 0
0 0 1 0 0 0 1 1 0 0 1 0 0 0 0 1 0 0 0 0 0 1 0 0 0 0 1 0 1 1 1 0 1 1 0 0 0
0 1 0 1 0 0 1 1 0 0 0 0 1 0 0 0 0 1 1 0 1 0 1 0 0 0 0 0 1 0 0 0 0 1 1 0 0
1 0 1 0 0 1 0 0 0 0 0 0 0 0 1 0 0 1 1 0 0 0 1 1 0 1 0 0 1 0 0 0 1 1 0 1 0
0 0 0 0 1 0 0 1 0 1 1 0 0 1 0 0 1 0 0 0 1 0 1 1 0 0 1 1 0 1 0 1 1 1 0 1 0
0 1 0 0 1 0 0 1 0 0 0 0 1 1 0 0 1 0 1 0 0 0 0 0 1 1 1 0 0 1 1 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 1 1 0 0 0 0 0 0 0 0 0 0 0 1 0 1 0 0 0 0
0 0 1 0 0 0 0 0 1 0 1 0 1 0 0 0 1 0 1 1 1 0 0 0 1 0 1 0 0 0 1 1 1 0 0 1 1
0 0 0 1 0 1 0 0 0 0 0 1 1 0 1 1 1 0 0 0 1 0 0 0 0 1 0 0 0 1 0 0 1 0 0 0 0
1 0 0 1 0 1 0 0 0 1 1 1 1 1 0 0 1 1 0 1 1 1 0 0 0 1 1 0 0 1 0 0 0 0 0 0
0 0 0 1 1 0 0 1 0]
```

In [106]: `training_data_accuracy = accuracy_score(y_train, x_train_prediction)`  
`print('Accuracy score of training data', training_data_accuracy)`

Accuracy score of training data 0.8075842696629213

In [107]: `# accuracy on test data`  
`x_test_prediction = model.predict(x_test)`  
`print(x_test_prediction)`

```
[0 0 1 0 0 0 0 0 0 0 0 1 1 0 0 1 0 0 1 0 1 1 0 1 0 1 1 0 0 0 0 0 0 0 1 1
0 0 0 0 0 1 0 0 1 1 0 0 1 0 0 0 0 0 0 1 0 0 0 1 0 0 0 1 0 1 0 0 0 1 0 1 0
1 0 0 0 1 0 1 0 0 0 1 1 0 0 1 0 0 0 0 0 0 1 0 1 0 0 1 0 1 1 0 1 1 0 0 0 0
0 0 0 1 1 0 1 0 0 1 0 0 0 0 0 0 1 0 0 0 0 1 1 0 0 0 0 0 0 1 1 1 1 0 1 0 0
0 1 0 0 0 0 1 0 0 1 1 0 1 0 0 0 1 1 0 0 1 0 0 1 1 1 0 0 0 0 0]
```

In [109]: `test_data_accuracy = accuracy_score(y_test, x_test_prediction)`  
`print('Accuracy score of training data', test_data_accuracy)`

Accuracy score of training data 0.7821229050279329

In [ ]: