

## ▼ *Titanic Survival*

### About this Project

the Titanic dataset and performing various data processing, analysis, visualization, and modeling steps. Here's a summary of what i've done in each section and a concluding statement:

1. **Importing Dependencies:** You've imported the necessary libraries like NumPy, Pandas, Matplotlib, Seaborn, and scikit-learn modules for data analysis, visualization, and modeling.
2. **Data Collection & Processing:** You've loaded the Titanic dataset using Pandas, checked the shape and information of the data, and identified missing values in columns like 'Age', 'Cabin', and 'Embarked'. You've dropped the 'Cabin' column, filled missing values in the 'Age' column with the mean value, and replaced missing values in the 'Embarked' column with the mode.
3. **Analysis and Visualization:** You've performed descriptive statistics on some numeric columns, counted the number of survivors and non-survivors, and created count plots and bar plots using Seaborn to visualize the distribution of features like 'Sex', 'Pclass', and their relationship with 'Survived'.
4. **Encoding Categories:** You've replaced categorical values like 'Sex' and 'Embarked' with numerical values (0 or 1) for model compatibility.
5. **Separating Features & Target:** You've separated the features ('X') and the target ('Y') columns for training and testing the model.
6. **Splitting the Data:** You've split the data into training and testing sets using the `train_test_split` function.
7. **Model Training:** You've trained a Logistic Regression model on the training data.
8. **Model Evaluation:** You've evaluated the model's accuracy on both the training and testing data using the `accuracy_score` metric.

**Conclusion:** The Logistic Regression model achieved an accuracy of approximately 80.76% on the training data and 78.21% on the testing data. This indicates that the model generalizes reasonably well to new, unseen data, although there might be room for improvement by trying out different models or performing more advanced feature engineering and hyperparameter tuning.

Remember that this is just one step in the data science pipeline, and further exploration and refinement could lead to even better results.

## ▼ Import the Dependencies

```
import numpy as np
import pandas as pd
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
```

## ▼ Data Collection & Processing

```
# load the data
titanic_data = pd.read_csv('/content/train.csv')
```

```
titanic_data.head()
```

	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	
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85	

```
# number of rows and columns
titanic_data.shape
```

```
(891, 12)
```

```
# info
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
```

```
#check the number of the missing values
titanic_data.isnull().sum()
```

```
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
```

## ▼ Handling the missing value

```
# drop the cabin cloumns the data
titanic_data = titanic_data.drop(columns='Cabin', axis=1)
```

```
# replacing the misssing values with mean values (age)
titanic_data['Age'].fillna(titanic_data['Age'].mean(), inplace=True)
```

```
#findind the mode value of column
print(titanic_data['Embarked'].mode())
```

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

```
print(titanic_data['Embarked'].mode()[0])
```

```
S
```

```
# replacing the missing value in mode
titanic_data['Embarked'].fillna(titanic_data['Embarked'].mode()[0], inplace=True)
```

```
titanic_data.isnull().sum()
```

```
PassengerId    0
Survived        0
```

```
Pclass      0
Name        0
Sex         0
Age         0
SibSp       0
Parch       0
Ticket      0
Fare        0
Embarked    0
dtype: int64
```

▼ Analysis

```
# statistical measures
titanic_data.describe()
```

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

```
# find the survived and not survived
titanic_data['Survived'].value_counts()
```

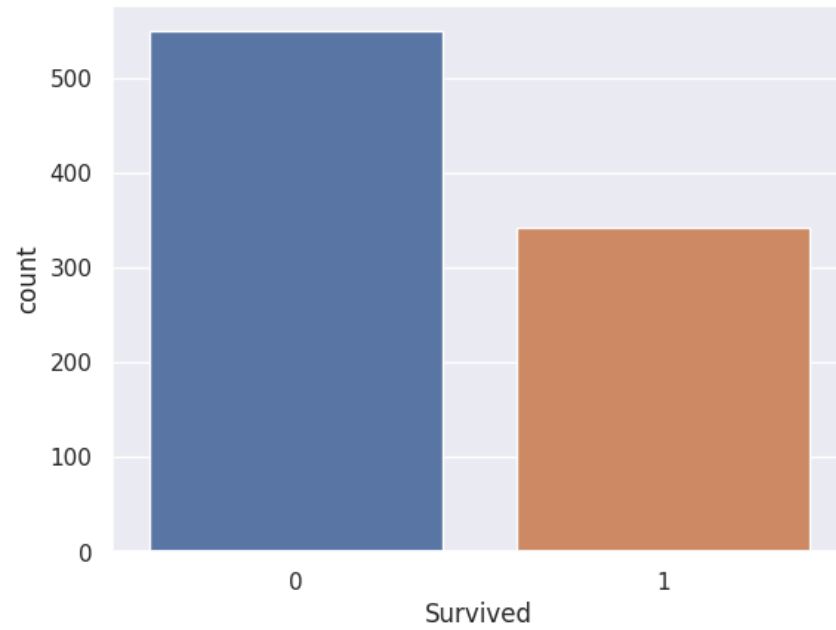
```
0    549
1    342
Name: Survived, dtype: int64
```

▼ Visualization

```
sns.set()
```

```
# count plot
sns.countplot(x='Survived', data=titanic_data)
```

```
<Axes: xlabel='Survived', ylabel='count'>
```

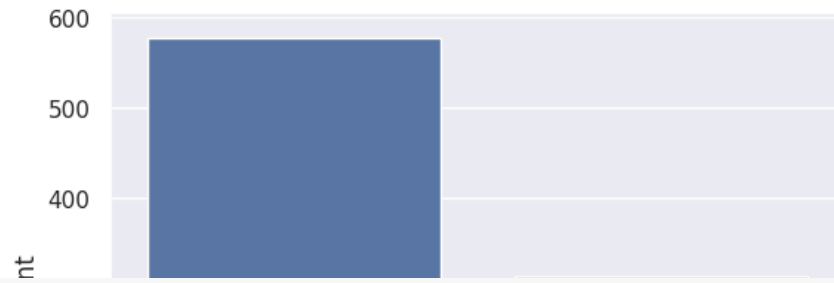


```
titanic_data['Sex'].value_counts()
```

```
male      577
female    314
Name: Sex, dtype: int64
```

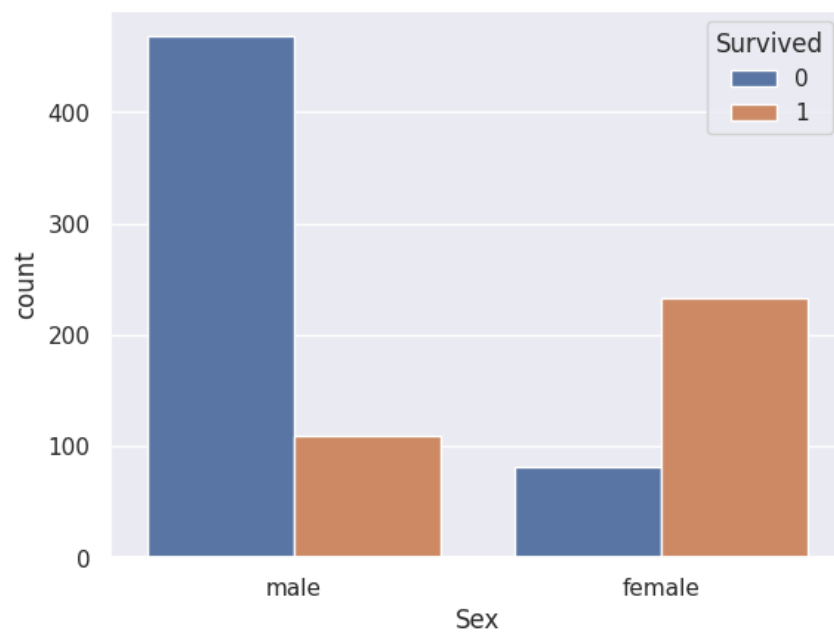
```
sns.countplot(x='Sex', data=titanic_data)
```

```
<Axes: xlabel='Sex', ylabel='count'>
```



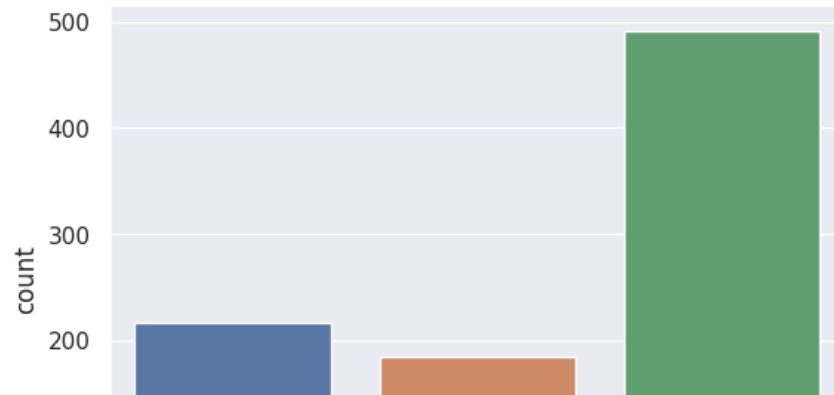
```
sns.countplot(x='Sex', hue='Survived', data=titanic_data)
```

```
<Axes: xlabel='Sex', ylabel='count'>
```



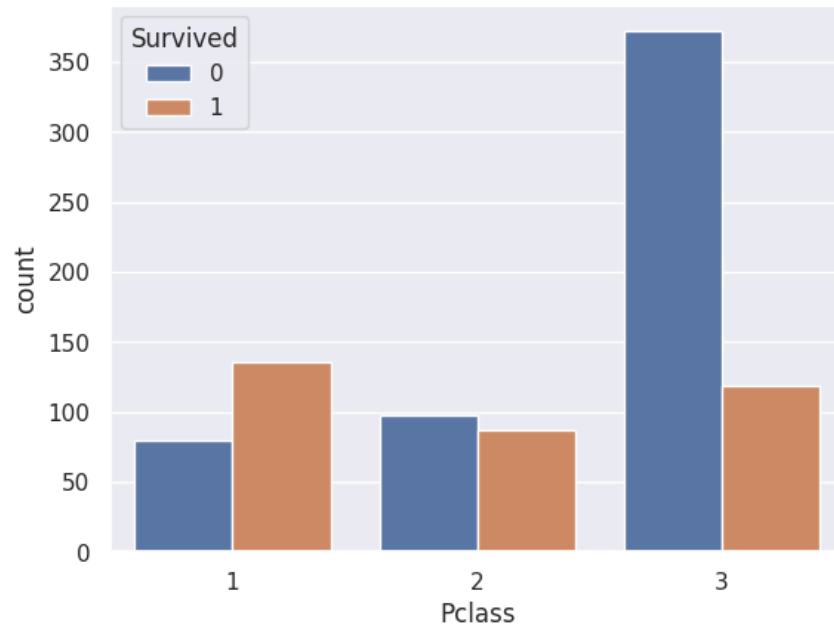
```
sns.countplot(x='Pclass', data=titanic_data)
```

```
<Axes: xlabel='Pclass', ylabel='count'>
```



```
sns.countplot(x='Pclass', hue='Survived', data=titanic_data)
```

```
<Axes: xlabel='Pclass', ylabel='count'>
```



## ▼ Encoding the Categories

```
#  
titanic_data['Sex'].value_counts()
```

```
male      577
female    314
Name: Sex, dtype: int64
```

```
titanic_data['Embarked'].value_counts()
```

```
S      646
C      168
Q       77
Name: Embarked, dtype: int64
```

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

```
titanic_data.head()
```

	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

## ▾ Separating Features & Target

```
X = titanic_data.drop(columns = ['PassengerId','Name','Ticket','Survived'],axis=1)
Y = titanic_data['Survived']
```

```
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]
```

```
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
```

## ▼ *\*Splitting the data into test Data \**

```
X_train, X_test, Y_train, Y_test = train_test_split(X,Y, test_size=0.2, random_state=2)
```

```
print(X.shape, X_train.shape, X_test.shape)
```

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

## ▼ **Model Training**

### **Logistic Regression**

```
model = LogisticRegression()
```

```
# training Logistic model
model.fit(X_train, Y_train)
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_logistic.py:458: ConvergenceWarning: lbfgs
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.
```

## ▼ Model Evalution

```
        n_iter_i = check_optimize_result(
# accuracy Score
X_train_prediction = model.predict(X_train)
```

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

```
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
```

```
# for the test data
X_test_prediction = model.predict(X_test)
```

```
print(X_test_prediction)
```

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

```
test_data_accuracy = accuracy_score(Y_test, X_test_prediction)
print('Accuracy score of test data : ', test_data_accuracy)
```

➤ Accuracy score of test data : 0.7821229050279329

+ Code

+ Text

✓ 0s completed at 12:21 PM

