

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

```
# Load Titanic dataset
df = pd.read_csv("/kaggle/input/titanic/train.csv")

# Display the first five rows
df.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	S
1	2	1	1	Cumings, 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 Mav Peel)	female	35.0	1	0	113803	53.1000	C123	S

```
# Checking for missing values
missing_values = df.isnull().sum()
print(missing_values)
```

```
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
```

```
# Summary statistics
print(df.describe())

# Check the dimensions of the dataset
print("Dataset dimensions:", df.shape)
```

```
count    891.000000    891.000000    891.000000    714.000000    891.000000 \
mean      446.000000     0.383838     2.308642    29.699118     0.523008
std       257.353842     0.486592     0.836071    14.526497     1.102743
min         1.000000     0.000000     1.000000     0.420000     0.000000
25%       223.500000     0.000000     2.000000    20.125000     0.000000
50%       446.000000     0.000000     3.000000    28.000000     0.000000
75%       668.500000     1.000000     3.000000    38.000000     1.000000
max        891.000000     1.000000     3.000000    80.000000     8.000000

      Parch      Fare
count    891.000000    891.000000
mean       0.381594    32.204208
std        0.806057    49.693429
min         0.000000     0.000000
25%         0.000000     7.910400
50%         0.000000    14.454200
75%         0.000000    31.000000
max         6.000000   512.329200
Dataset dimensions: (891, 12)
```

```
print(df.dtypes)
```

```
PassengerId    int64
Survived        int64
Pclass          int64
Name            object
Sex             object
Age            float64
SibSp           int64
Parch           int64
```

```
Ticket      object
Fare        float64
Cabin       object
Embarked    object
dtype: object
```

```
# Convert 'Survived' and 'Pclass' to categorical type
df['Survived'] = df['Survived'].astype('category')
df['Pclass'] = df['Pclass'].astype('category')
```

```
# Convert 'Age' and 'Fare' to float
df['Age'] = df['Age'].astype(float)
df['Fare'] = df['Fare'].astype(float)
```

```
df['Age'] = df['Age'].fillna(df['Age'].median())
```

```
df['Embarked'] = df['Embarked'].fillna(df['Embarked'].mode()[0])
```

```
df.drop(columns=['Cabin'], inplace=True)
```

```
# Convert categorical variables into dummy/indicator variables
df = pd.get_dummies(df, columns=['Sex', 'Embarked'], drop_first=True)
```

```
# Display first five rows after encoding
df.head()
```

	PassengerId	Survived	Pclass	Name	Age	SibSp	Parch	Ticket	Fare	Sex_male	Embarked_Q	Embarked_S
0	1	0	3	Braund, Mr. Owen Harris	22.0	1	0	A/5 21171	7.2500	True	False	True
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	38.0	1	0	PC 17599	71.2833	False	False	False
2	3	1	3	Heikkinen, Miss. Laina	26.0	0	0	STON/O2. 3101282	7.9250	False	False	True

```
df.to_csv("/kaggle/working/cleaned_titanic.csv", index=False)
```

```
df['Survived'] = df['Survived'].astype(int) # Ensure it's numeric
sns.countplot(x='Survived', data=df)
plt.show()
```



