


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

# Load the Titanic dataset
titanic_df = pd.read_csv('https://raw.githubusercontent.com/datasciencedojo/datasets/master/titanic.csv')

# Display the first few rows of the dataset
print(titanic_df.head())

# Check the data types and missing values
print(titanic_df.info())

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



	PassengerId	Survived	Pclass	\		
0	1	0	3			
1	2	1	1			
2	3	1	3			
3	4	1	1			
4	5	0	3			

	Name	Sex	Age	SibSp	\	
0	Braund, Mr. Owen Harris	male	22.0	1		
1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1		
2	Heikkinen, Miss. Laina	female	26.0	0		
3	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1		
4	Allen, Mr. William Henry	male	35.0	0		

	Parch	Ticket	Fare	Cabin	Embarked
0	0	A/5 21171	7.2500	NaN	S
1	0	PC 17599	71.2833	C85	C
2	0	STON/O2. 3101282	7.9250	NaN	S
3	0	113803	53.1000	C123	S
4	0	373450	8.0500	NaN	S

<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
None

	PassengerId	Survived	Pclass	Age	SibSp	\	
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

```
# Drop unnecessary columns
titanic_df = titanic_df.drop(['PassengerId', 'Name', 'Ticket', 'Cabin'], axis=1)

# Fill missing values in the Age column with the median age
titanic_df['Age'].fillna(titanic_df['Age'].median(), inplace=True)

# Fill missing values in the Embarked column with the mode
mode_embarked = titanic_df['Embarked'].mode()[0]
titanic_df['Embarked'].fillna(mode_embarked, inplace=True)

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

# Check for any remaining missing values
print(titanic_df.isnull().sum())
```

```
Survived      0
Pclass        0
Age            0
SibSp         0
Parch         0
Fare          0
Sex_male      0
Embarked_Q    0
Embarked_S    0
dtype: int64
```

```
# Visualize the distribution of Age
plt.figure(figsize=(10, 6))
sns.histplot(titanic_df['Age'], bins=20, kde=True)
plt.title('Distribution of Age')
plt.xlabel('Age')
plt.ylabel('Count')
plt.show()
```

```
# Explore the survival rate by gender
plt.figure(figsize=(8, 5))
sns.countplot(x='Survived', hue='Sex_male', data=titanic_df)
plt.title('Survival Count by Gender')
plt.xlabel('Survived')
plt.ylabel('Count')
plt.xticks([0, 1], ['No', 'Yes'])
plt.legend(['Female', 'Male'])
plt.show()
```

```
# Explore the survival rate by passenger class
plt.figure(figsize=(8, 5))
sns.countplot(x='Survived', hue='Pclass', data=titanic_df)
plt.title('Survival Count by Passenger Class')
plt.xlabel('Survived')
plt.ylabel('Count')
plt.xticks([0, 1], ['No', 'Yes'])
plt.legend(title='Passenger Class')
plt.show()
```

```
# Explore the relationship between fare and survival
plt.figure(figsize=(10, 6))
sns.boxplot(x='Survived', y='Fare', data=titanic_df)
plt.title('Survival by Fare')
plt.xlabel('Survived')
plt.ylabel('Fare')
plt.xticks([0, 1], ['No', 'Yes'])
plt.show()
```

