In [2]:
```python
#Q.4)Perform EDA on a Different Data Set.
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
df = pd.read_csv("titanic.csv")
df.head()
```

Out[2]:

| | 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 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 [3]:
```python
# Dataset shape
df.shape

# Column information
df.info()

# Statistical summary
df.describe()
```

```
<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
```

Out[3]:

| | PassengerId | Survived | Pclass | Age | SibSp | Parch | Fare |
|---|---|---|---|---|---|---|---|
| count | 891.000000 | 891.000000 | 891.000000 | 714.000000 | 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 | 14.526497 | 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 | 20.125000 | 0.000000 | 0.000000 | 7.910400 |
| 50% | 446.000000 | 0.000000 | 3.000000 | 28.000000 | 0.000000 | 0.000000 | 14.454200 |
| 75% | 668.500000 | 1.000000 | 3.000000 | 38.000000 | 1.000000 | 0.000000 | 31.000000 |
| max | 891.000000 | 1.000000 | 3.000000 | 80.000000 | 8.000000 | 6.000000 | 512.329200 |

In [4]:
```python
df.isnull().sum()
```

```
Out[4]:  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
```
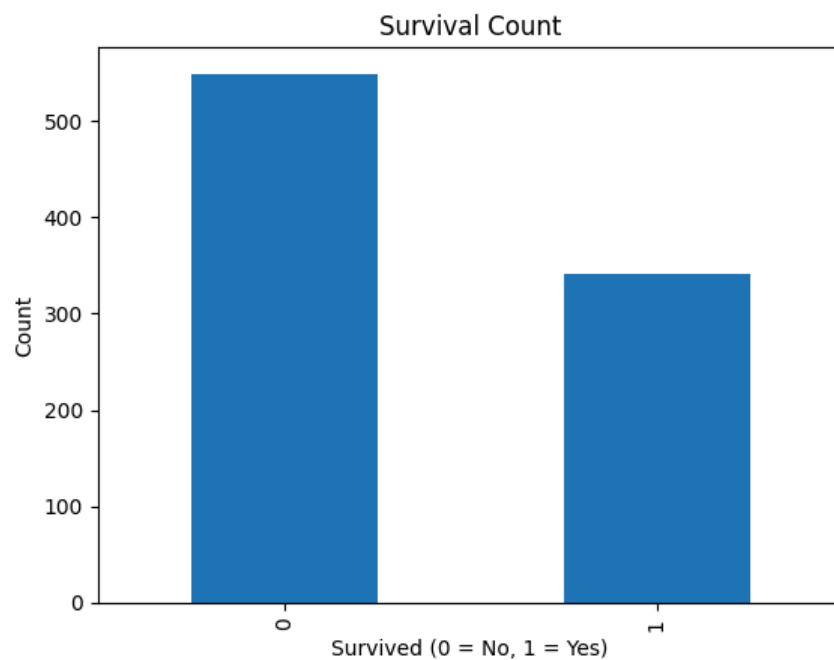
```python
In [6]:  # Fill missing Age values safely
         df['Age'] = df['Age'].fillna(df['Age'].median())

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

         # Drop Cabin column
         df = df.drop(columns=['Cabin'])
```
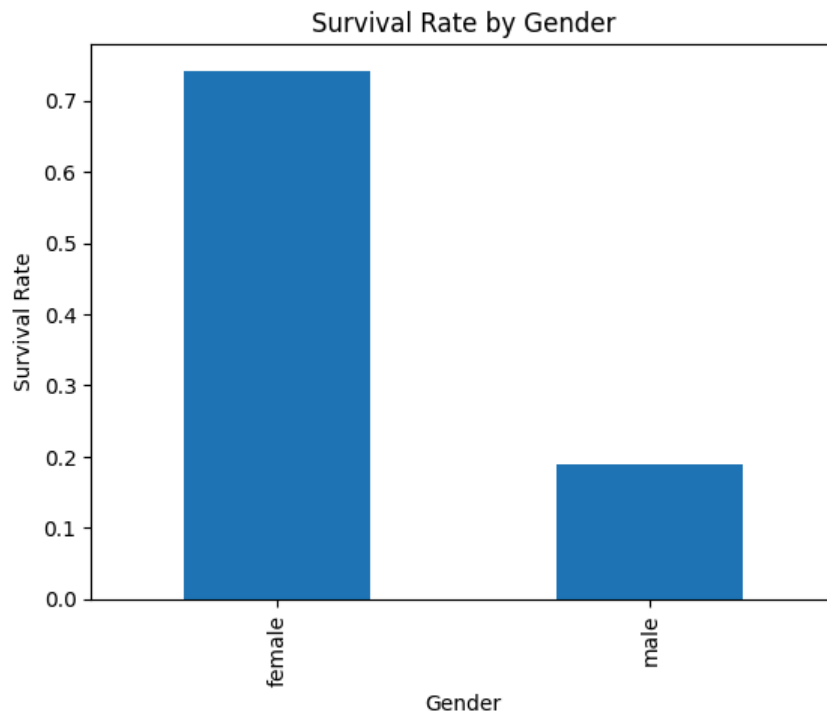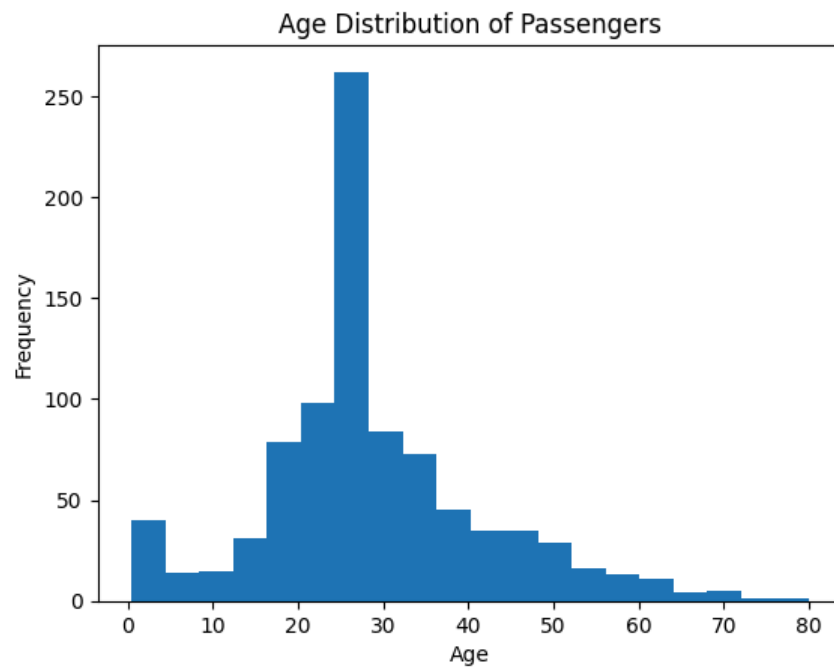
```python
In [8]:  plt.figure()
         df['Survived'].value_counts().plot(kind='bar')
         plt.title("Survival Count")
         plt.xlabel("Survived (0 = No, 1 = Yes)")
         plt.ylabel("Count")
         plt.show()
```
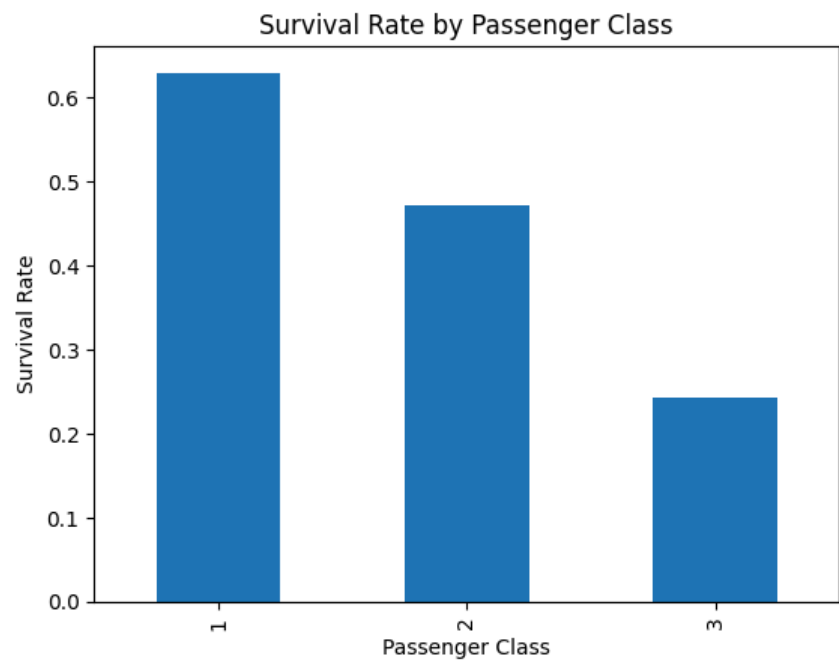


```python
In [9]:  plt.figure()
         df.groupby('Sex')['Survived'].mean().plot(kind='bar')
         plt.title("Survival Rate by Gender")
         plt.xlabel("Gender")
         plt.ylabel("Survival Rate")
         plt.show()
```

## Survival Rate by Gender
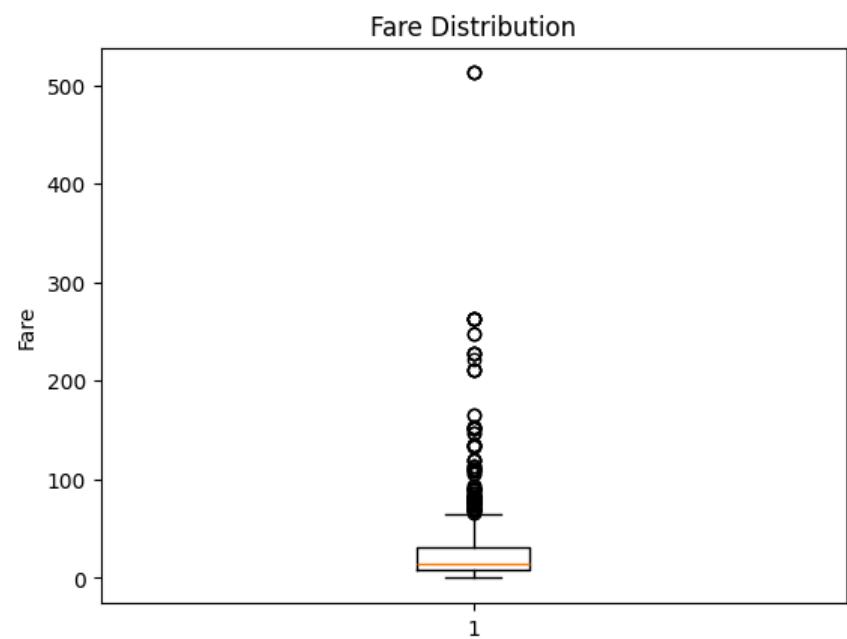


```
In [10]:  plt.figure()
          plt.hist(df['Age'], bins=20)
          plt.title("Age Distribution of Passengers")
          plt.xlabel("Age")
          plt.ylabel("Frequency")
          plt.show()
```

## Age Distribution of Passengers



```
In [11]:  plt.figure()
          df.groupby('Pclass')['Survived'].mean().plot(kind='bar')
          plt.title("Survival Rate by Passenger Class")
          plt.xlabel("Passenger Class")
          plt.ylabel("Survival Rate")
          plt.show()
```

## Survival Rate by Passenger Class



```
In [12]: plt.figure()
         plt.boxplot(df['Fare'])
         plt.title("Fare Distribution")
         plt.ylabel("Fare")
         plt.show()
```

## Fare Distribution



```
In [13]: df[['Age', 'Fare', 'Survived', 'Pclass']].corr()
```

Out[13]:

|          | Age       | Fare      | Survived  | Pclass    |
|----------|-----------|-----------|-----------|-----------|
| **Age**  | 1.000000  | 0.096688  | -0.064910 | -0.339898 |
| **Fare** | 0.096688  | 1.000000  | 0.257307  | -0.549500 |
| **Survived** | -0.064910 | 0.257307 | 1.000000 | -0.338481 |
| **Pclass** | -0.339898 | -0.549500 | -0.338481 | 1.000000 |

```
In [ ]:
```