```
In [ ]:  NAME:HARINI KARTHIKA V
         TASK NO:2
         Prodigy InfoTech
```

```python
In [1]:  import pandas as pd

         # Load the Titanic dataset
         df = pd.read_csv('titanic.csv')

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

```
   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
```

```python
In [2]: # Get a summary of the dataframe
        print(df.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
None
```

```python
In [3]: # Check for missing values
        print(df.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
```

```
In [4]:   # Fill missing 'Age' values with the median age
          df['Age'].fillna(df['Age'].median(), inplace=True)

          # Fill missing 'Embarked' values with the mode (most common value)
          df['Embarked'].fillna(df['Embarked'].mode()[0], inplace=True)

          # Drop the 'Cabin' column because it has too many missing values
          df.drop(columns=['Cabin'], inplace=True)

          # Convert 'Sex' and 'Embarked' to categorical variables
          df['Sex'] = df['Sex'].astype('category')
          df['Embarked'] = df['Embarked'].astype('category')

          # Verify that there are no more missing values
          print(df.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
```

```python
In [7]: import matplotlib.pyplot as plt
        import seaborn as sns

        # Set the style for the plots
        sns.set(style="whitegrid")

        # Plot the survival rate by gender
        plt.figure(figsize=(8, 6))
        sns.countplot(x='Survived', hue='Sex', data=df)
        plt.title('Survival Count by Gender')
        plt.show()

        # Plot the distribution of ages
        plt.figure(figsize=(8, 6))
        sns.histplot(df['Age'], bins=30, kde=True)
        plt.title('Age Distribution')
        plt.show()

        # Plot the survival rate by age
        plt.figure(figsize=(8, 6))
        sns.histplot(data=df, x='Age', hue='Survived', multiple='stack', bins=30)
        plt.title('Survival Rate by Age')
        plt.show()

        # Plot the survival rate by passenger class
        plt.figure(figsize=(8, 6))
        sns.countplot(x='Pclass', hue='Survived', data=df)
        plt.title('Survival Rate by Passenger Class')
        plt.show()

        # Plot the survival rate by embarkation point
        plt.figure(figsize=(8, 6))
        sns.countplot(x='Embarked', hue='Survived', data=df)
        plt.title('Survival Rate by Embarkation Point')
        plt.show()
```
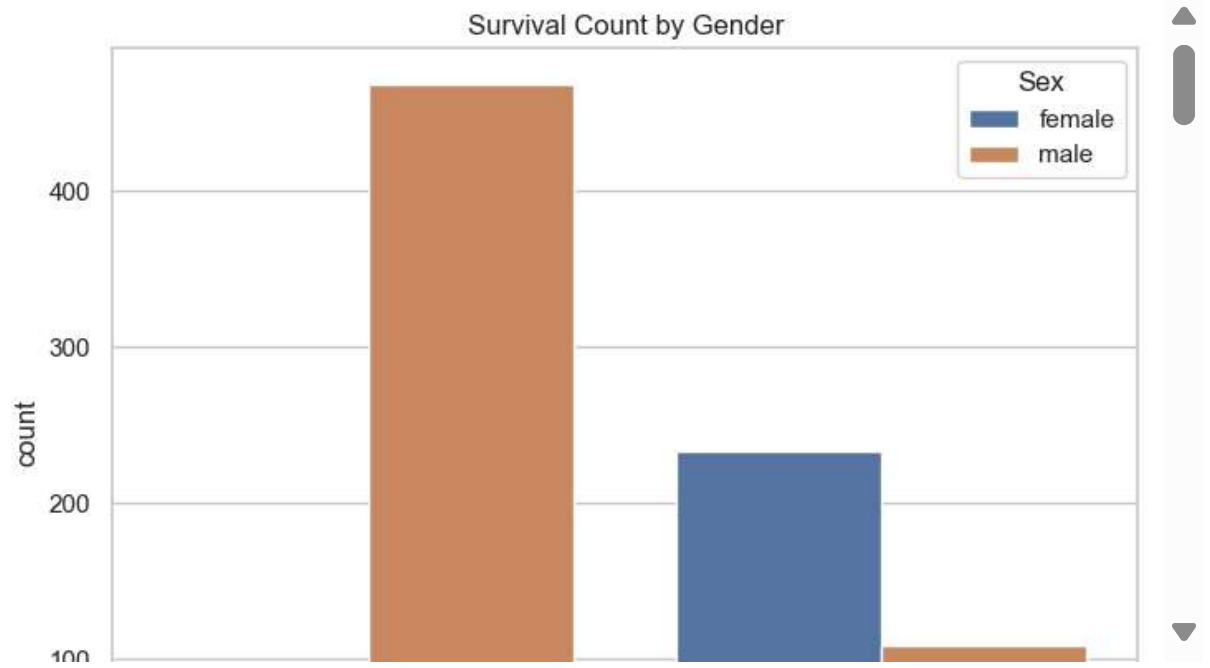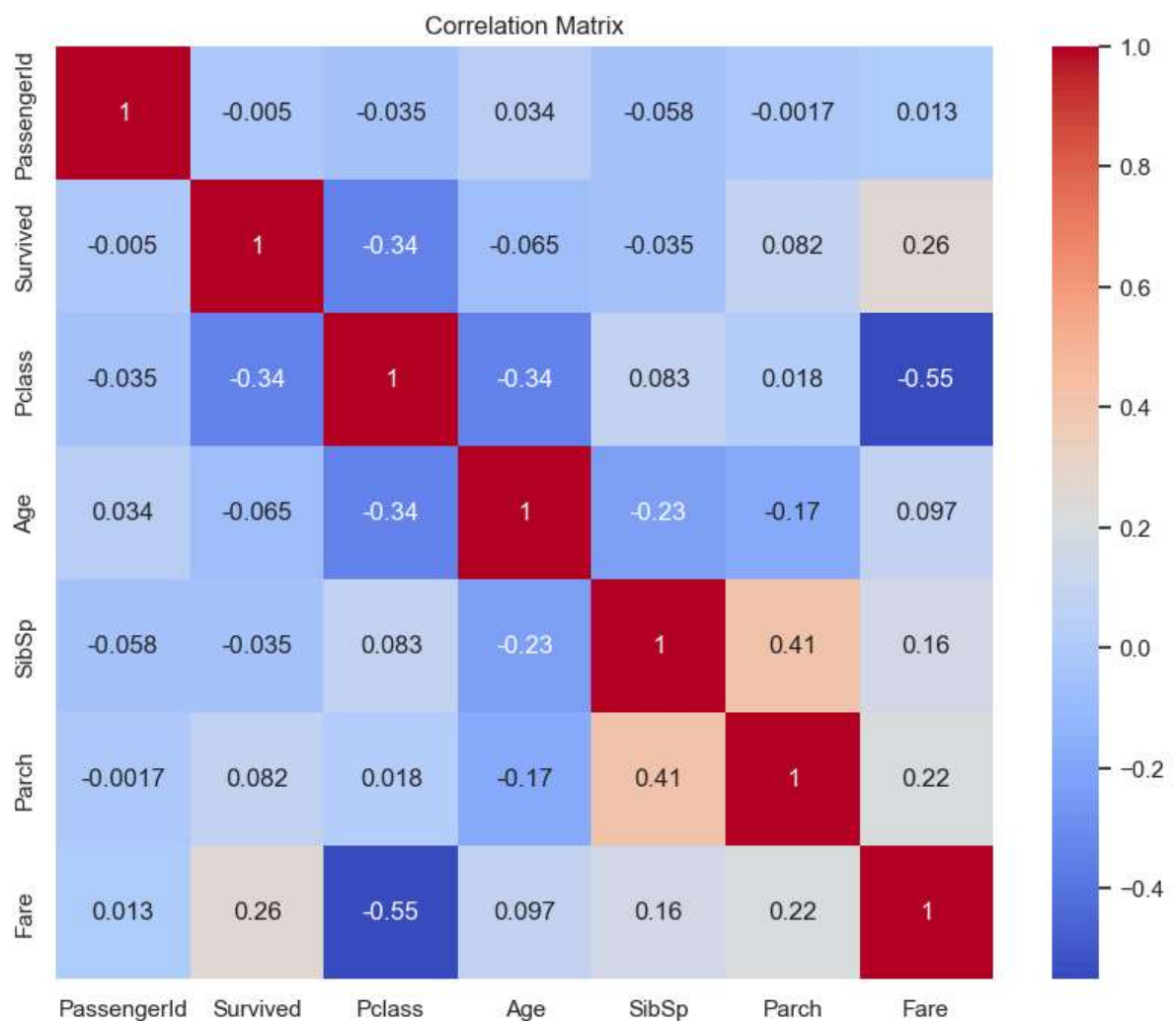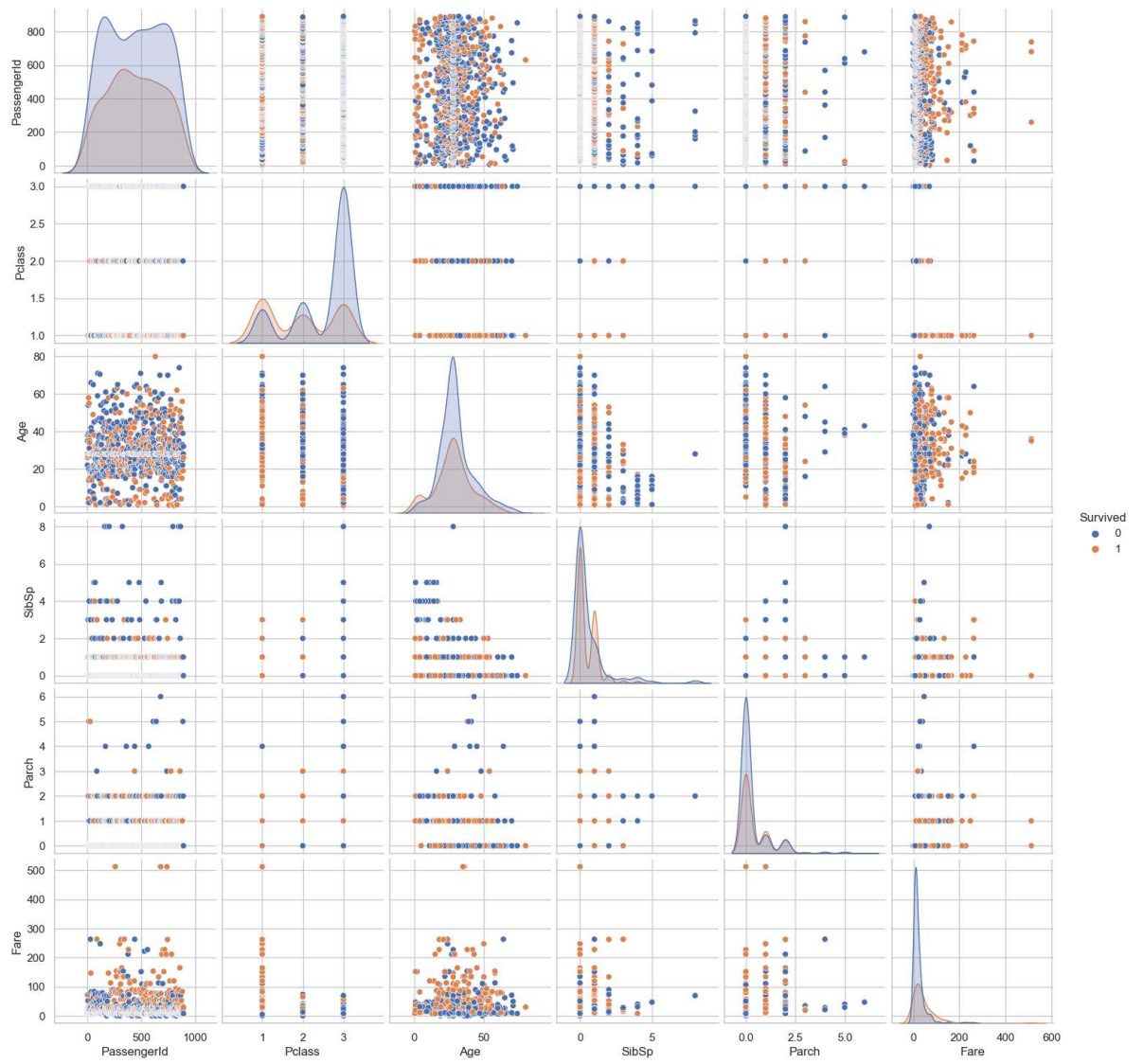
Survival Count by Gender

count

400

300

200

100

Sex
female
male

In [8]:
```python
# Plot the correlation matrix
plt.figure(figsize=(10, 8))
sns.heatmap(df.corr(), annot=True, cmap='coolwarm')
plt.title('Correlation Matrix')
plt.show()

# Pair plot to explore relationships between features
sns.pairplot(df, hue='Survived', diag_kind='kde')
plt.show()
```

C:\Users\Karthika\AppData\Local\Temp\ipykernel_9024\2245879067.py:3: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.
  sns.heatmap(df.corr(), annot=True, cmap='coolwarm')



Correlation Matrix

In [ ]: