

```

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

data = pd.read_csv('/content/train.csv')

data.head()

{"summary":{"\n  \"name\": \"data\", \n  \"rows\": 891, \n  \"fields\": [\n    {\n      \"column\": \"PassengerId\", \n      \"properties\": {\n        \"dtype\": \"number\", \n        \"std\": 257, \n        \"min\": 1, \n        \"max\": 891, \n        \"num_unique_values\": 891, \n        \"samples\": [\n          710, \n          440, \n          841\n        ], \n        \"semantic_type\": \"\", \n        \"description\": \"\" \n      }, \n      \"column\": \"Survived\", \n      \"properties\": {\n        \"dtype\": \"number\", \n        \"std\": 0, \n        \"min\": 0, \n        \"max\": 1, \n        \"num_unique_values\": 2, \n        \"samples\": [\n          1, \n          0\n        ], \n        \"semantic_type\": \"\", \n        \"description\": \"\" \n      }, \n      \"column\": \"Pclass\", \n      \"properties\": {\n        \"dtype\": \"number\", \n        \"std\": 0, \n        \"min\": 1, \n        \"max\": 3, \n        \"num_unique_values\": 3, \n        \"samples\": [\n          3, \n          1\n        ], \n        \"semantic_type\": \"\", \n        \"description\": \"\" \n      }, \n      \"column\": \"Name\", \n      \"properties\": {\n        \"dtype\": \"string\", \n        \"num_unique_values\": 891, \n        \"samples\": [\n          \"Moubarek, Master. Halim Gonios (\\\"William George\\\")\", \n          \"Kvillner, Mr. Johan Henrik Johannesson\"\n        ], \n        \"semantic_type\": \"\", \n        \"description\": \"\" \n      }, \n      \"column\": \"Sex\", \n      \"properties\": {\n        \"dtype\": \"category\", \n        \"num_unique_values\": 2, \n        \"samples\": [\n          \"female\", \n          \"male\"\n        ], \n        \"semantic_type\": \"\", \n        \"description\": \"\" \n      }, \n      \"column\": \"Age\", \n      \"properties\": {\n        \"dtype\": \"number\", \n        \"std\": 14.526497332334044, \n        \"min\": 0.42, \n        \"max\": 80.0, \n        \"num_unique_values\": 88, \n        \"samples\": [\n          0.75, \n          22.0\n        ], \n        \"semantic_type\": \"\", \n        \"description\": \"\" \n      }, \n      \"column\": \"SibSp\", \n      \"properties\": {\n        \"dtype\": \"number\", \n        \"std\": 1, \n        \"min\": 0, \n        \"max\": 8, \n        \"num_unique_values\": 7, \n        \"samples\": [\n          1, \n          0\n        ], \n        \"semantic_type\": \"\", \n        \"description\": \"\" \n      }, \n      \"column\": \"Parch\", \n      \"properties\": {\n        \"dtype\": \"number\", \n        \"std\": 0, \n        \"min\": 0, \n        \"max\": 6, \n        \"num_unique_values\": 7, \n        \"samples\": [\n          0, \n          1\n        ], \n        \"semantic_type\": \"\", \n        \"description\": \"\" \n      }

```

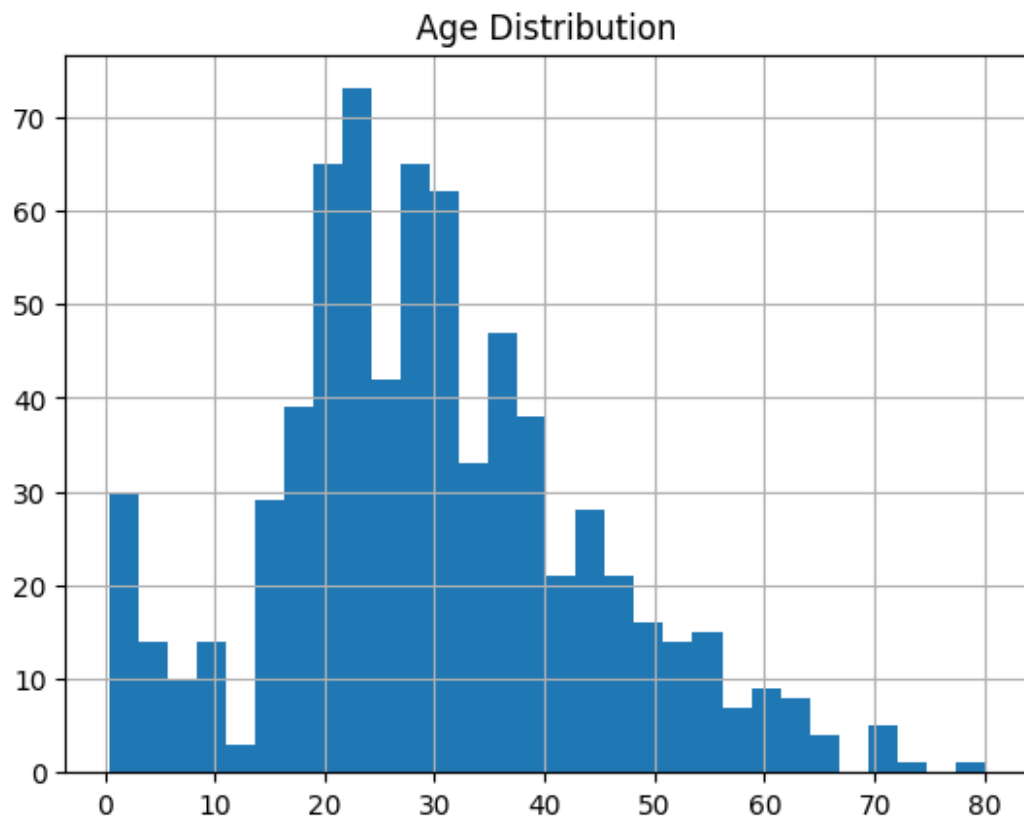


```
PassengerId      0
Survived          0
Pclass           0
Name             0
Sex              0
Age             177
SibSp            0
Parch            0
Ticket           0
Fare             0
Cabin           687
Embarked         2
dtype: int64

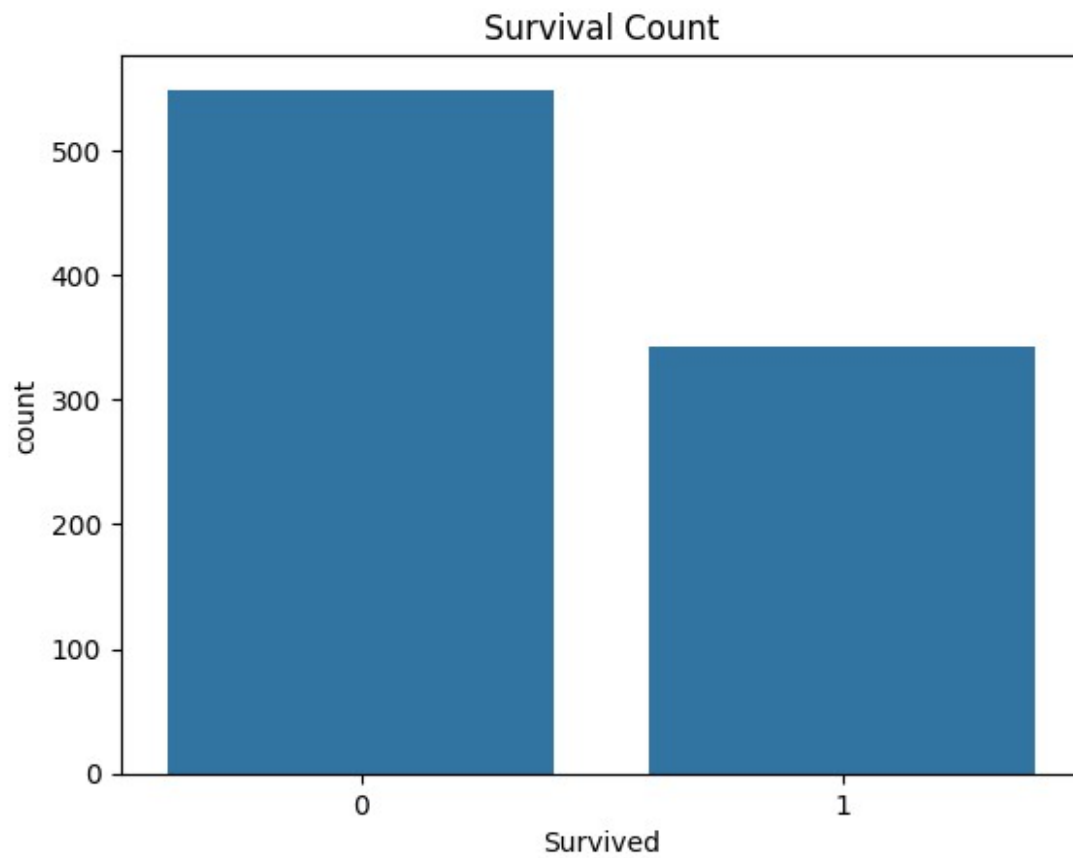
data.nunique()

PassengerId      891
Survived          2
Pclass           3
Name             891
Sex              2
Age             88
SibSp            7
Parch            7
Ticket          681
Fare            248
Cabin           147
Embarked         3
dtype: int64

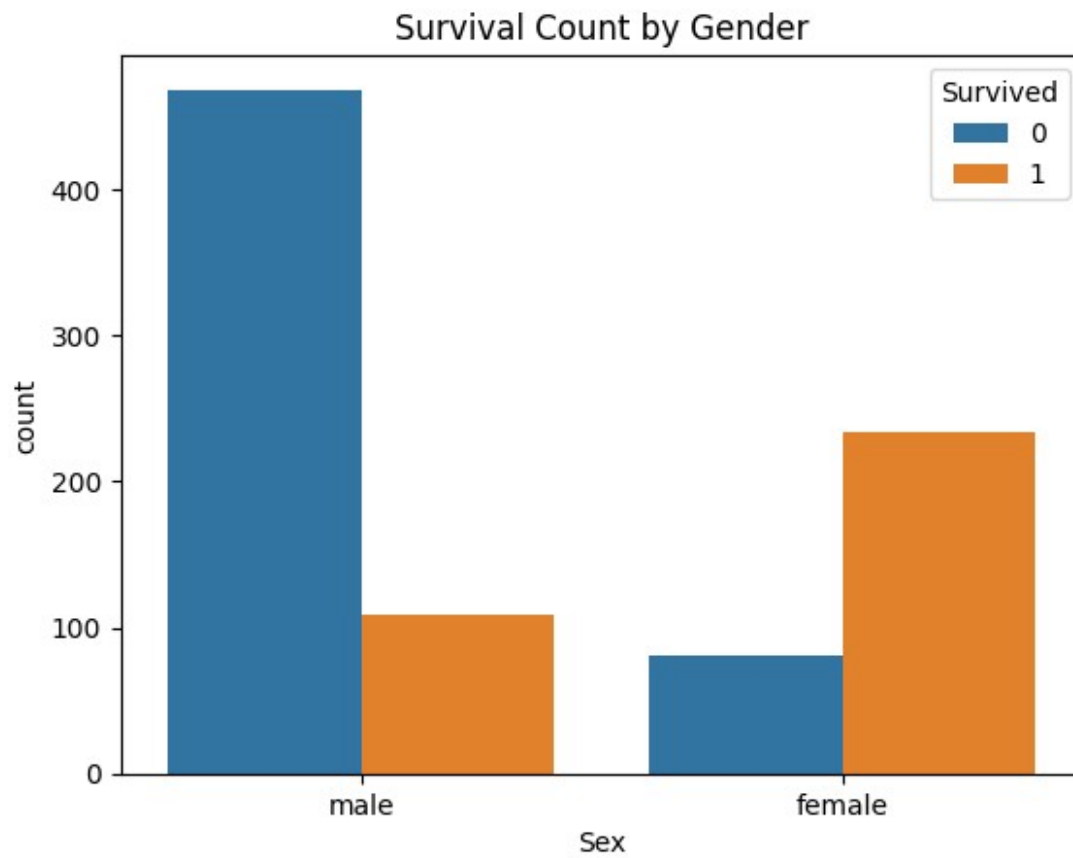
data['Age'].hist(bins=30)
plt.title('Age Distribution')
plt.show()
```



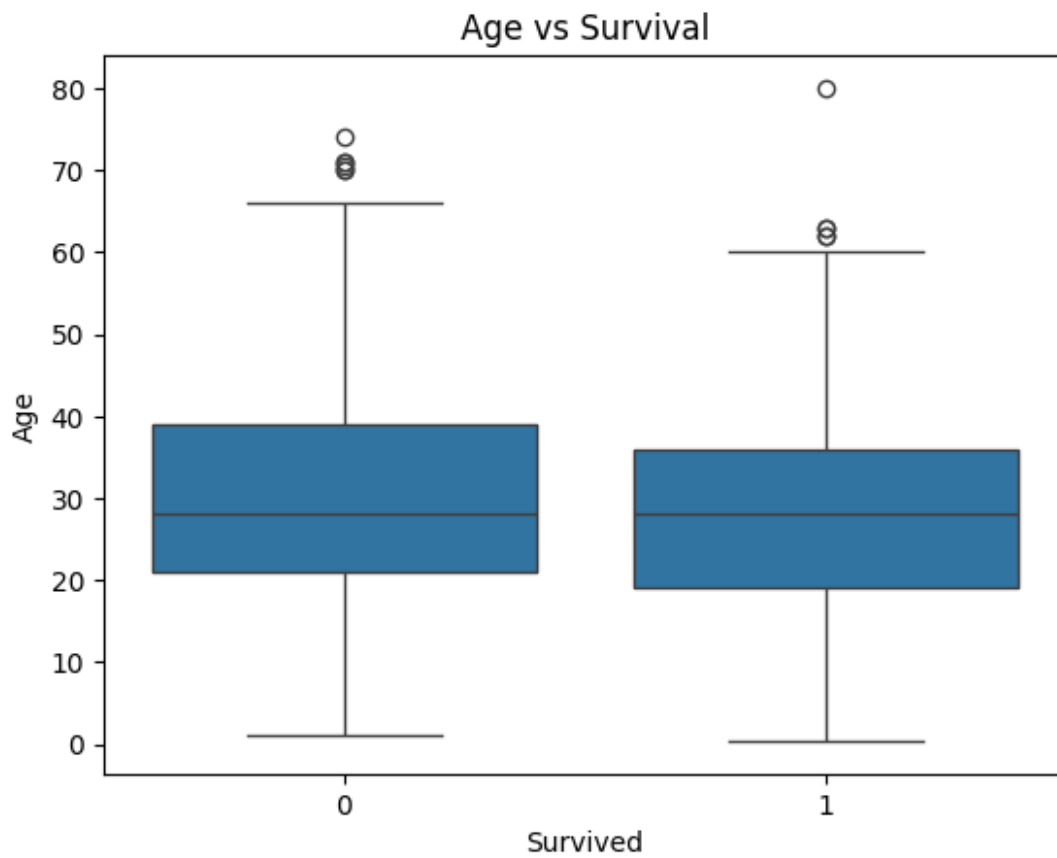
```
sns.countplot(x='Survived', data=data)
plt.title('Survival Count')
plt.show()
```



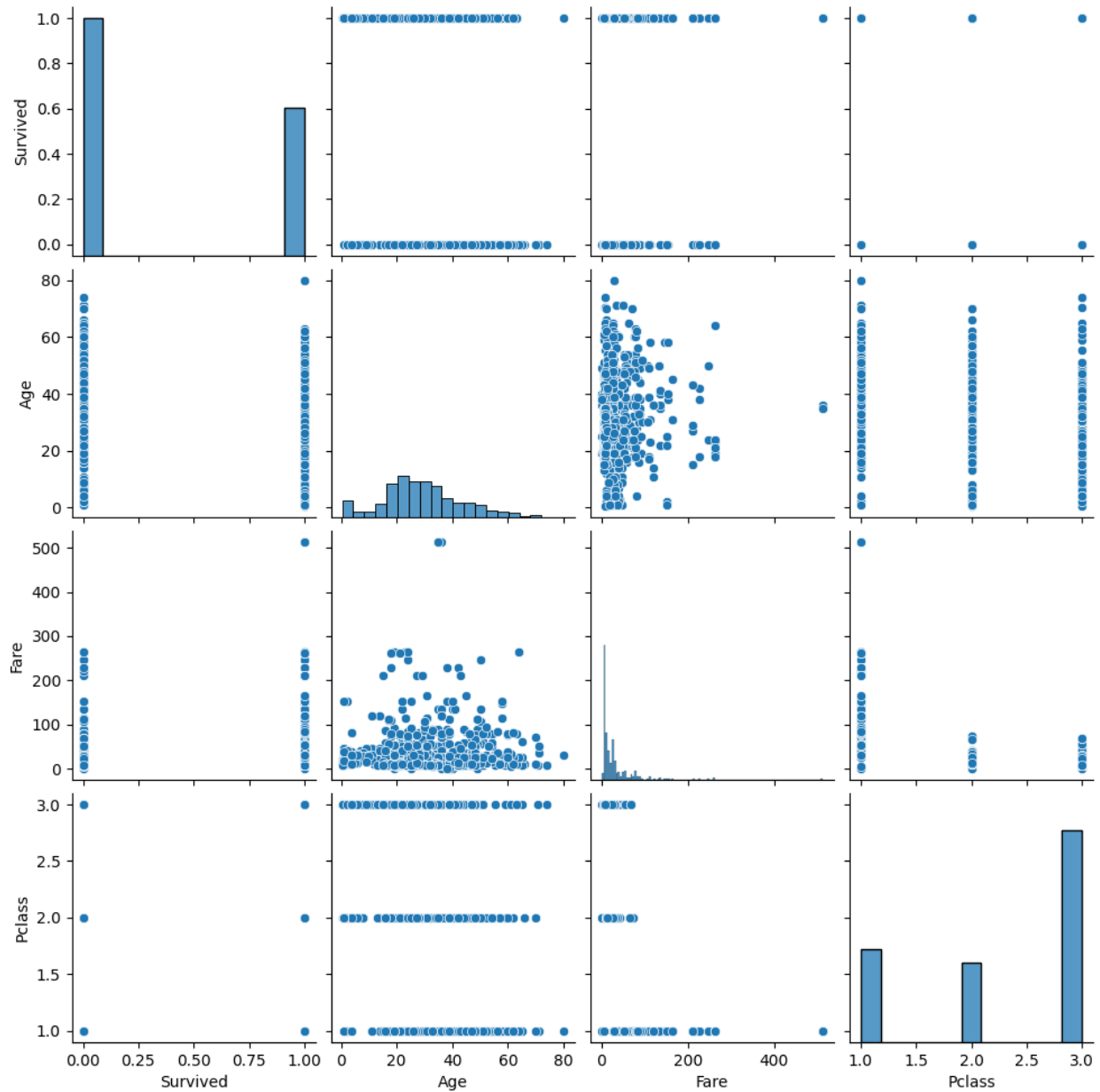
```
sns.countplot(x='Sex', hue='Survived', data=data)
plt.title('Survival Count by Gender')
plt.show()
```



```
sns.boxplot(x='Survived', y='Age', data=data)
plt.title('Age vs Survival')
plt.show()
```

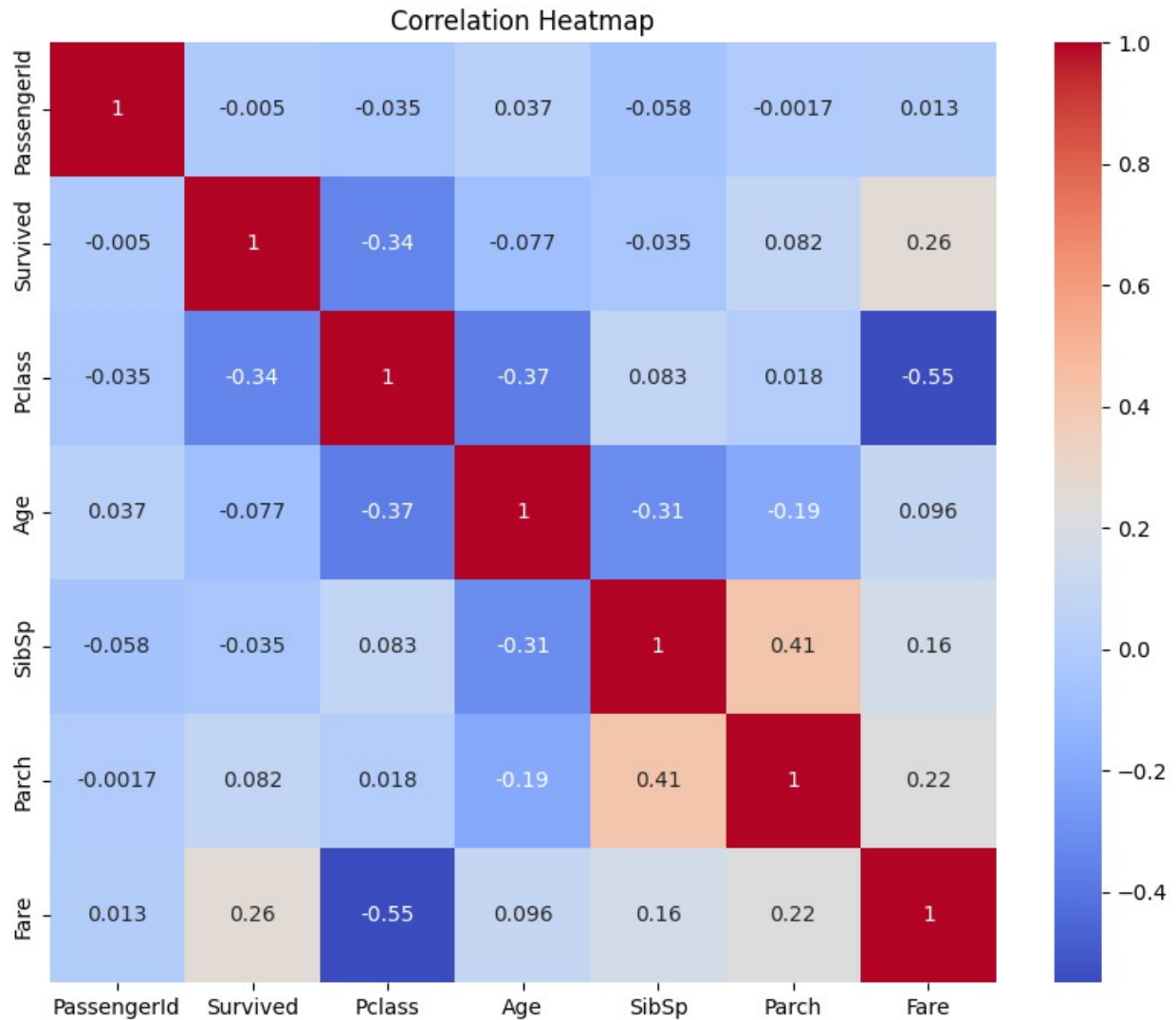


```
sns.pairplot(data[['Survived', 'Age', 'Fare', 'Pclass']])  
plt.show()
```



```
plt.figure(figsize=(10,8))
numeric_data = data.select_dtypes(include=['int64', 'float64']) #
Select only numeric columns
sns.heatmap(numeric_data.corr(), annot=True, cmap='coolwarm')
plt.title('Correlation Heatmap')
plt.show()
```





```
data['Age'].fillna(data['Age'].median(), inplace=True)
data['Embarked'].fillna(data['Embarked'].mode()[0], inplace=True)
```

<ipython-input-13-f7b6d87cfc76>:1: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method.

The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always behaves as a copy.

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col].method(value) instead, to perform the operation inplace on the original object.

```
data['Age'].fillna(data['Age'].median(), inplace=True)
```

```
<ipython-input-13-f7b6d87cfc76>:2: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method.
```

The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always behaves as a copy.

For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col].method(value) instead, to perform the operation inplace on the original object.

```
data['Embarked'].fillna(data['Embarked'].mode()[0], inplace=True)
```

□ Main Insights: Age:

Most passengers were between 20 to 40 years old.

Very few passengers were very young (children) or very old (elderly).

Fare:

Most passengers paid a fare between \$0 and \$100.

A few outliers paid very high fares (above \$500), mostly from 1st class.

Gender and Survival:

Females had a much higher survival rate compared to males.

Most males did not survive, while a large proportion of females did.

Passenger Class and Survival:

First-class passengers had the highest survival rate.

Third-class passengers had the lowest survival rate.

Correlation Analysis:

Fare and Pclass had a strong negative correlation (higher class → higher fare).

Survival had a positive relationship with higher fare and higher passenger class.

Age had a weak or moderate relationship with survival — young children had slightly better survival rates.

Missing Values:

Columns like Age and Cabin had missing values.

Age missing values could be handled later using imputation if needed.

□ Overall Conclusion: Survival was strongly influenced by gender, passenger class, and ticket fare.

Being female, young, and/or from a higher class increased the chances of survival.

Rich passengers (paying higher fares) had much better survival rates compared to those who paid less.

□ End of Summary