

```
In [12]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt

In [16]: df = pd.read_csv('train.csv')

In [18]: 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(4), object(6)
memory usage: 83.7+ KB

In [10]: print(df.describe(include="all"))

PassengerId  Survived  Pclass      Name  Sex  \
count      891.000000    891.000000    891.000000    891    2
unique      NaN         NaN         NaN      Braund, Mr. Owen Harris  male
freq        NaN         NaN         NaN      1          517
mean      44.000000    0.343601    2.308642    NaN      NaN
std       257.353842    0.475972    0.836071    NaN      NaN
min        0.000000    0.000000    1.000000    NaN      NaN
25%       223.000000    0.000000    2.000000    NaN      NaN
50%       446.000000    0.000000    3.000000    NaN      NaN
75%       665.000000    1.000000    3.000000    NaN      NaN
max      891.000000    1.000000    3.000000    NaN      NaN

PassengerId  Survived  Pclass      Name  Sex  \
count      714.000000    891.000000    891.000000    891    204
unique      NaN         NaN         NaN      851      147
top         NaN         NaN         NaN      347682    NaN      B56 B58
freq        NaN         NaN         NaN      7          4
mean      29.694118    0.523008    0.351354    NaN      32.294298
std       14.526497    1.105743    0.806057    NaN      49.693429
min        0.420000    0.000000    0.000000    NaN      0.000000
25%       20.125000    0.000000    0.000000    NaN      7.918400
50%       28.000000    0.000000    0.000000    NaN      14.454200
75%       36.000000    1.000000    0.000000    NaN      31.000000
max       80.000000    8.000000    6.000000    NaN      512.329200

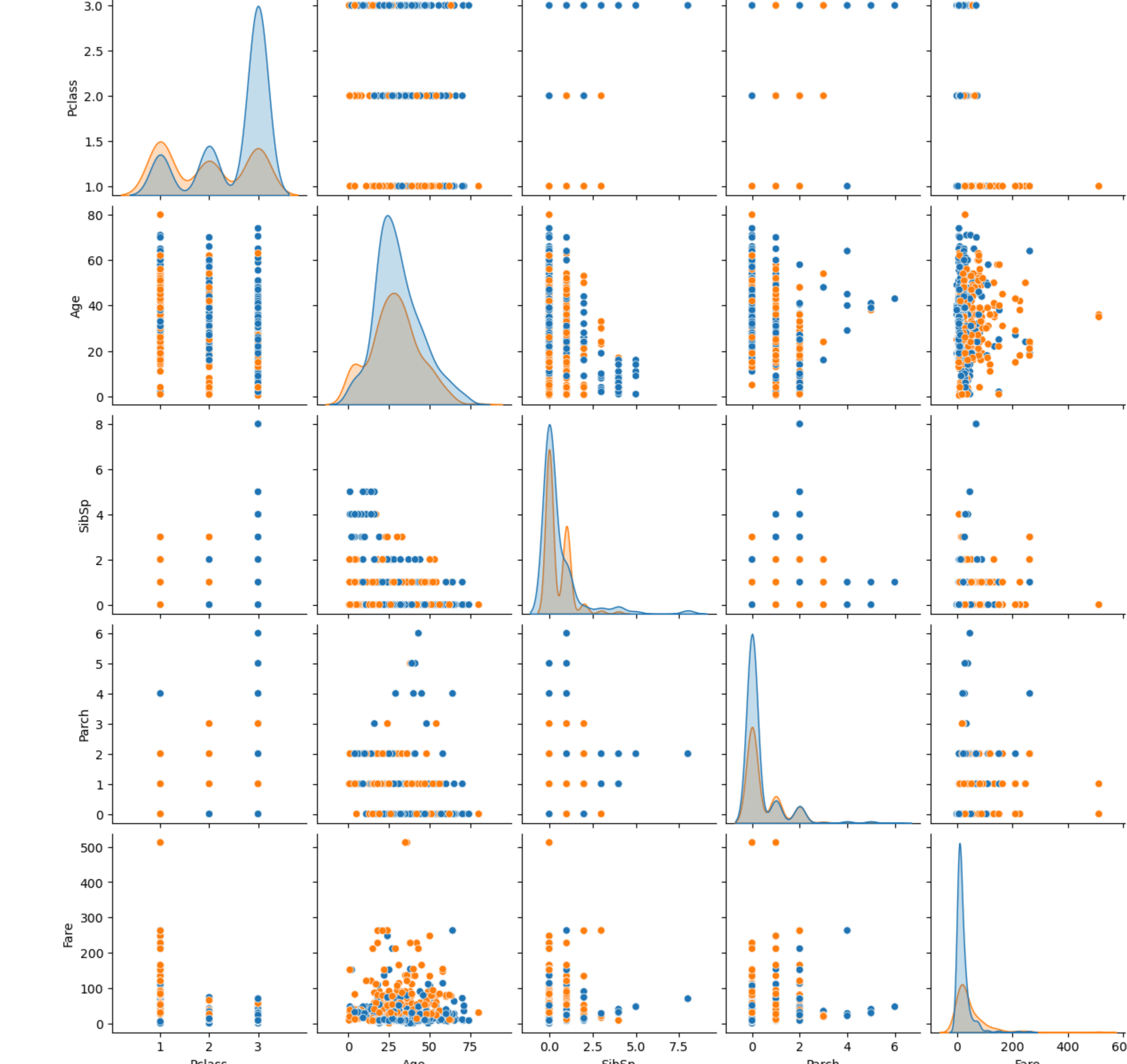
Embarked
count      889
unique      3
top         S
freq        687
mean      64.4
std       NaN
min       NaN
25%      NaN
50%      NaN
75%      NaN
max       NaN

In [10]: print("Missing Values:\n", df.isnull().sum())

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

for col in df.select_dtypes(include='object').columns: print(f"rValue counts for {col}:\n", df[col].value_counts())

In [16]: sns.pairplot(df[['Survived', 'Pclass', 'Age', 'SibSp', 'Parch', 'Fare']], hue='Survived')
plt.suptitle("Pairplot of Numeric Features (Colored by Survival)", y=1.02)
plt.show()
```



Observation: Higher fares and younger ages show slightly higher survival rates.

Pclass strongly separates survival chances.

```
In [18]: plt.figure(figsize=(10, 6))

Out[18]: <Figure size 1000x600 with 0 Axes>
<Figure size 1000x600 with 0 Axes>

In [20]: sns.heatmap(df.corr(), annot=True, cmap='coolwarm', fnc='.2f')

ValueError: could not convert string to float: 'Braund, Mr. Owen Harris'
Cell In[20], line 1                                Traceback (most recent call last)
----> 1 sns.heatmap(df.corr(), annot=True, cmap='coolwarm', fnc='.2f')

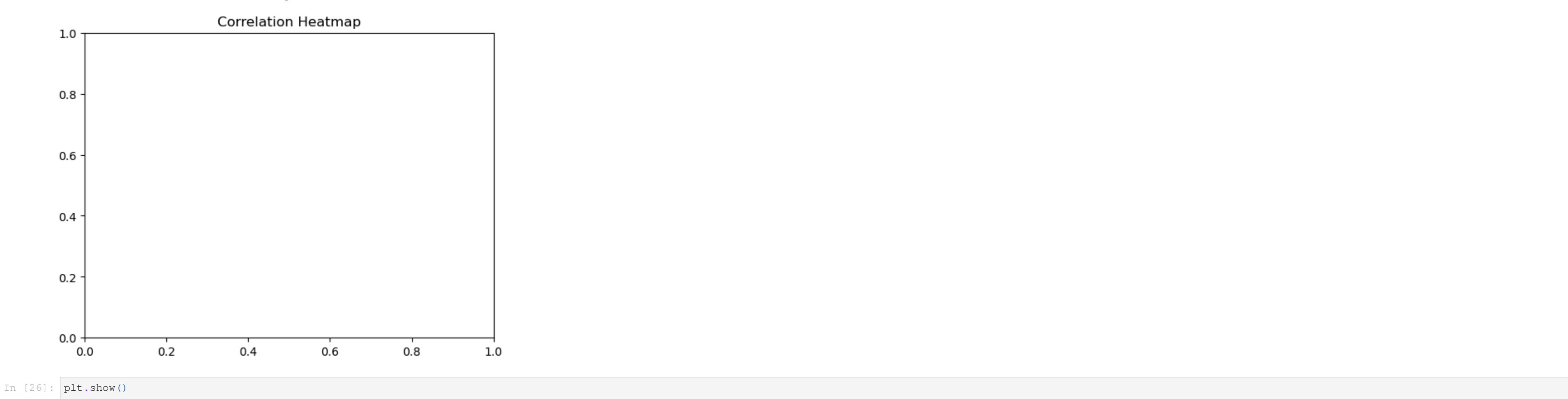
File C:\ProgramData\Anaconda3\Lib\site-packages\pandas\core\frame.py:11049, in DataFrame.corr(self, method, min_periods, numeric_only)
11047 cols = data.columns
11048 sds = cols.copy()
> 11049 mat = data.to_numpy(dtype=float, na_value=np.nan, copy=False)
11051 if method == 'pearson':
11052     corrval = libalgos.nancorr(mat, min_periods)

File C:\ProgramData\Anaconda3\Lib\site-packages\pandas\core\frame.py:1993, in DataFrame.to_numpy(self, dtype, copy, na_value)
1991 if dtype is not None:
1992     dtype = np.dtype(dtype)
-> 1993 result = self._get_notna_array(dtype=dtype, copy=copy, na_value=na_value)
1994 if result.dtype is not dtype:
1995     result = np.asarray(result, dtype=dtype)

File C:\ProgramData\Anaconda3\Lib\site-packages\pandas\core\internals\managers.py:1694, in BlockManager._interleave(self, dtype, copy, na_value)
1692     arr = self._interleave(dtype=dtype, na_value=na_value)
-> 1694 arr = self._interleave(dtype=dtype, na_value=na_value)
1695     # The underlying data was copied within _interleave, so no need
1696     # to further copy if copy=True or setting na_value
1698 if na_value is lib.no_default:

File C:\ProgramData\Anaconda3\Lib\site-packages\pandas\core\internals\managers.py:1753, in BlockManager._interleave(self, dtype, na_value)
1751     else:
1752         arr = blk.get_values(dtype)
-> 1753 result[i].indexer = arr
1754 itemmask[i].indexer = 1
1756 if not itemmask.all():

ValueError: could not convert string to float: 'Braund, Mr. Owen Harris'
```



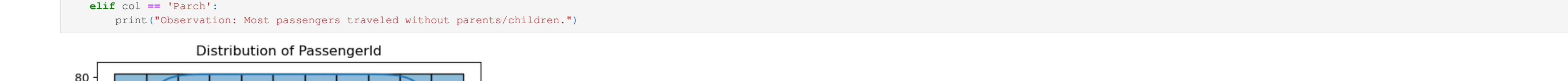
Observation: Fare and Pclass are negatively correlated (-0.55).
Survival is positively correlated with Fare (0.26) and negatively with Pclass (-0.34).

4. Histograms

```
In [28]: for col in df.select_dtypes(include=np.number).columns:
    col in [28], line 1
    for col in df.select_dtypes(include=np.number).columns:
        SyntaxError: incomplete input
```

```
In [30]: for col in df.select_dtypes(include=np.number).columns:
    plt.figure(figsize=(6, 4))
    sns.histplot(df[col].dropna(), kde=True)
    plt.title(f"Distribution of {col}")
    plt.show()

# Add observation for each:
if col == 'Age':
    print("Observation: Age distribution is right-skewed with most passengers in 20-40 age range.")
elif col == 'Fare':
    print("Observation: Fare is heavily right-skewed; a few passengers paid very high fares.")
elif col == 'SibSp':
    print("Observation: Most passengers traveled without siblings/spouses.")
elif col == 'Parch':
    print("Observation: Most passengers traveled without parents/children.")
```



Observation: Age distribution is right-skewed with most passengers in 20-40 age range.



Observation: Most passengers traveled without siblings/spouses.



Observation: Most passengers traveled without parents/children.



Observation: Fare is heavily right-skewed; a few passengers paid very high fares.

5. Boxplots

```
In [32]: for col in ['Age', 'Fare']:
    plt.figure(figsize=(6, 4))
    sns.boxplot(df[col])
    plt.title(f"Boxplot of {col}")
    plt.show()

if col == 'Age':
    print("Observation: Few outliers in age (very young and very old).")
elif col == 'Fare':
    print("Observation: Significant outliers in Fare - some tickets were extremely expensive.")
```



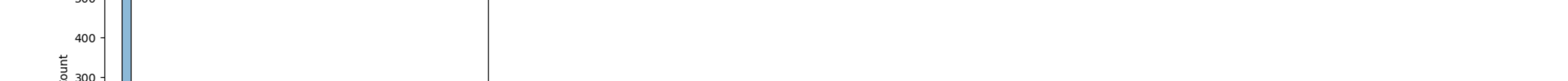
Observation: Few outliers in age (very young and very old).



Observation: Significant outliers in Fare - some tickets were extremely expensive.

6. Scatterplots

```
In [34]: plt.figure(figsize=(6, 4))
sns.scatterplot(x='Age', y='Fare', hue='Survived', data=df)
plt.title("Age vs Fare (Colored by Survival)")
plt.show()
```



Observation: Higher fares generally linked to higher survival; survival more common among younger passengers.

7. Summary of Findings

```
In [38]: print("""
Summary:
1. Dataset has 891 rows, 12 columns, with missing data in Age (177), Cabin (687), Embarked (2).
2. Most passengers were male (577/891).
3. Survival rate is higher among females, higher-class passengers, and those paying higher fares.
4. Age distribution is concentrated between 20-40 years.
5. Fare distribution is right-skewed with extreme outliers.
6. Pclass is strongly related to Fare and survival chance.
7. Missing cabin data suggests many passengers traveled in shared or lower-class accommodations.
""")
```

Summary:
1. Dataset has 891 rows, 12 columns, with missing data in Age (177), Cabin (687), Embarked (2).
2. Most passengers were male (577/891).
3. Survival rate is higher among females, higher-class passengers, and those paying higher fares.
4. Age distribution is concentrated between 20-40 years.

