



```
In [10]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

## SETUP NOTEBOOK

```
In [11]: df = pd.read_csv("titanic_train_dataset.csv")
df.head()
```

```
Out[11]:
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2 3101282
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450

## INITIAL EXPLORATION

```
In [12]: df.shape
```

```
Out[12]: (891, 12)
```

```
In [13]: 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

```

```
In [14]: df.describe()
```

```
Out[14]:
```

	PassengerId	Survived	Pclass	Age	SibSp	Parch
<b>count</b>	891.000000	891.000000	891.000000	714.000000	891.000000	891.000000
<b>mean</b>	446.000000	0.383838	2.308642	29.699118	0.523008	0.381594
<b>std</b>	257.353842	0.486592	0.836071	14.526497	1.102743	0.806057
<b>min</b>	1.000000	0.000000	1.000000	0.420000	0.000000	0.000000
<b>25%</b>	223.500000	0.000000	2.000000	20.125000	0.000000	0.000000
<b>50%</b>	446.000000	0.000000	3.000000	28.000000	0.000000	0.000000
<b>75%</b>	668.500000	1.000000	3.000000	38.000000	1.000000	0.000000
<b>max</b>	891.000000	1.000000	3.000000	80.000000	8.000000	6.000000

```
In [15]: df.isnull().sum()
```

```
Out[15]: 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 [16]: df.duplicated().sum()
```

```
Out[16]: np.int64(0)
```

```
In [21]: # Safely fill missing Age values
df['Age'] = df['Age'].fillna(df['Age'].median())

# Safely drop Cabin column if it exists
if 'Cabin' in df.columns:
    df.drop(columns='Cabin', inplace=True)

# Fill missing Embarked values with mode
df['Embarked'] = df['Embarked'].fillna(df['Embarked'].mode()[0])
```

```
In [22]: df.isnull().sum()
```

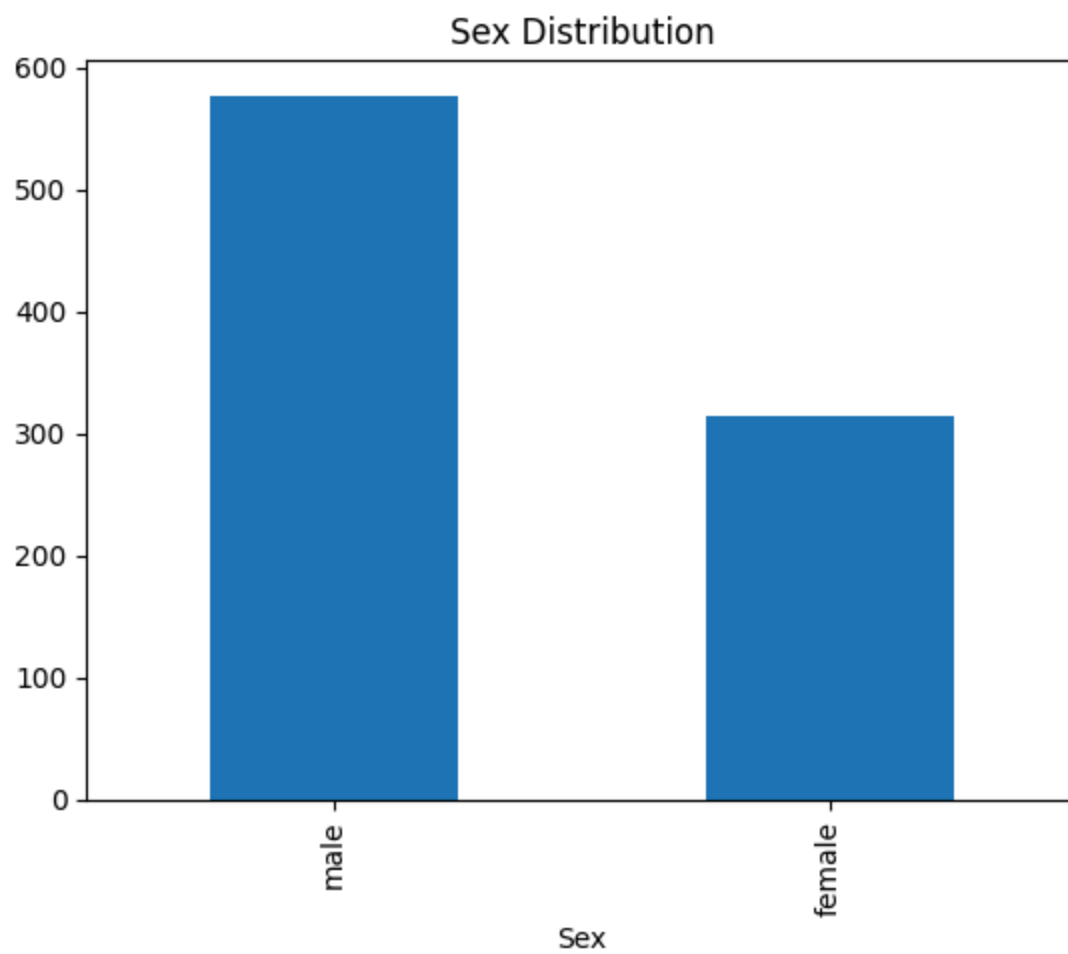
```
Out[22]: PassengerId    0
Survived              0
Pclass               0
Name                 0
Sex                  0
Age                  0
SibSp                0
Parch                0
Ticket              0
Fare                 0
Embarked             0
dtype: int64
```

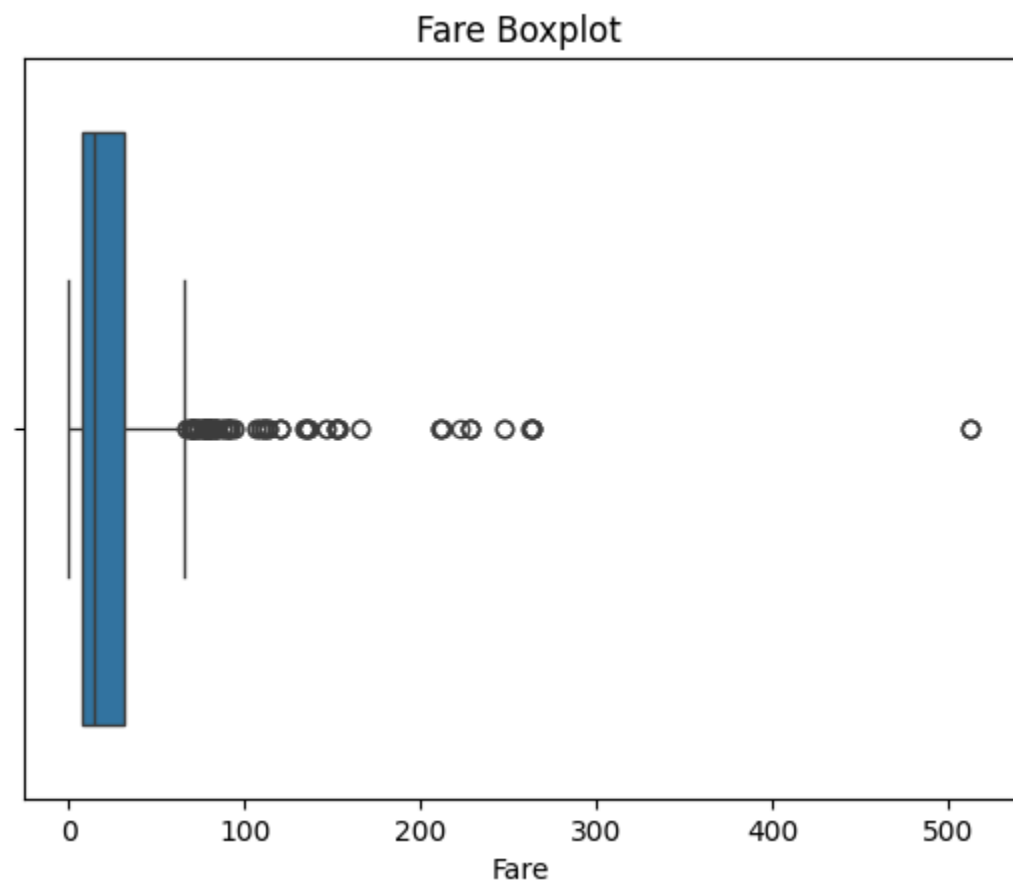
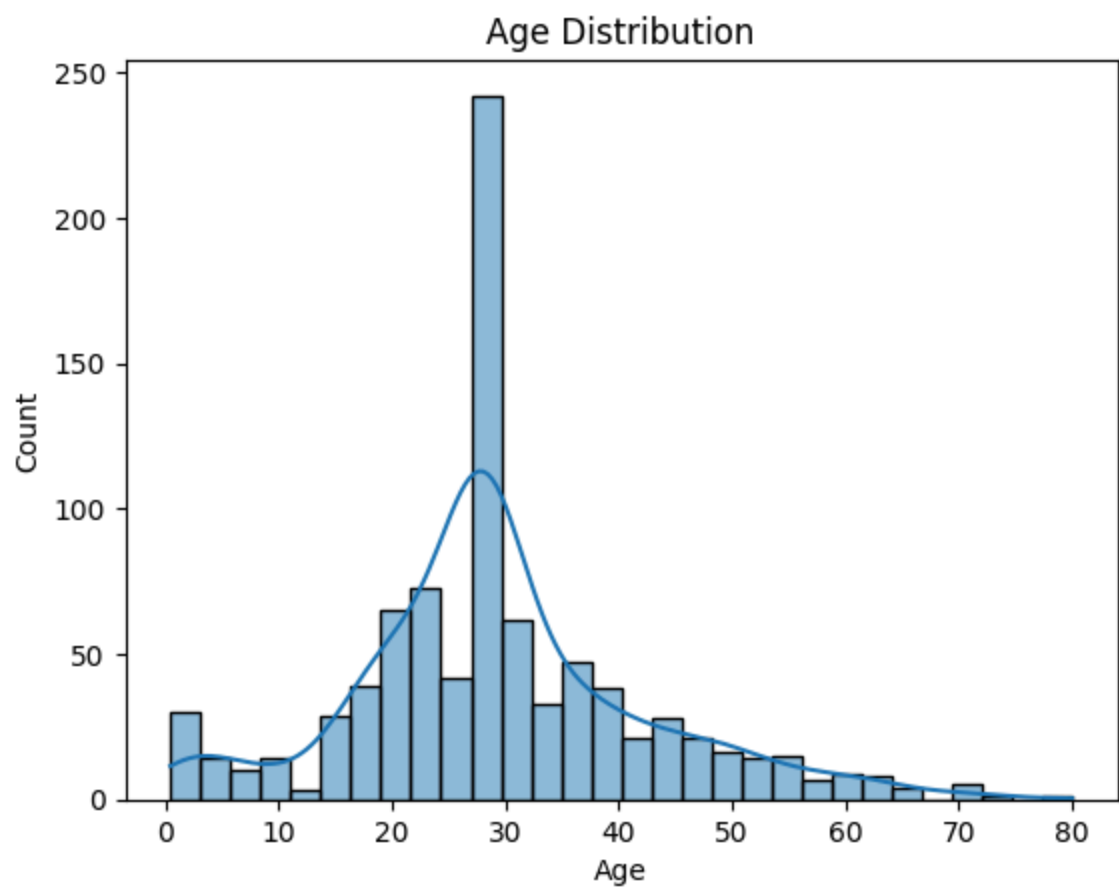
## UNIVARIATE ANALYSIS (SINGLE FEATURE )

```
In [23]: # Categorical
df['Sex'].value_counts().plot(kind='bar', title='Sex Distribution')
plt.show()

# Numerical
sns.histplot(df['Age'].dropna(), kde=True)
plt.title('Age Distribution')
plt.show()

sns.boxplot(x=df['Fare'])
plt.title('Fare Boxplot')
plt.show()
```

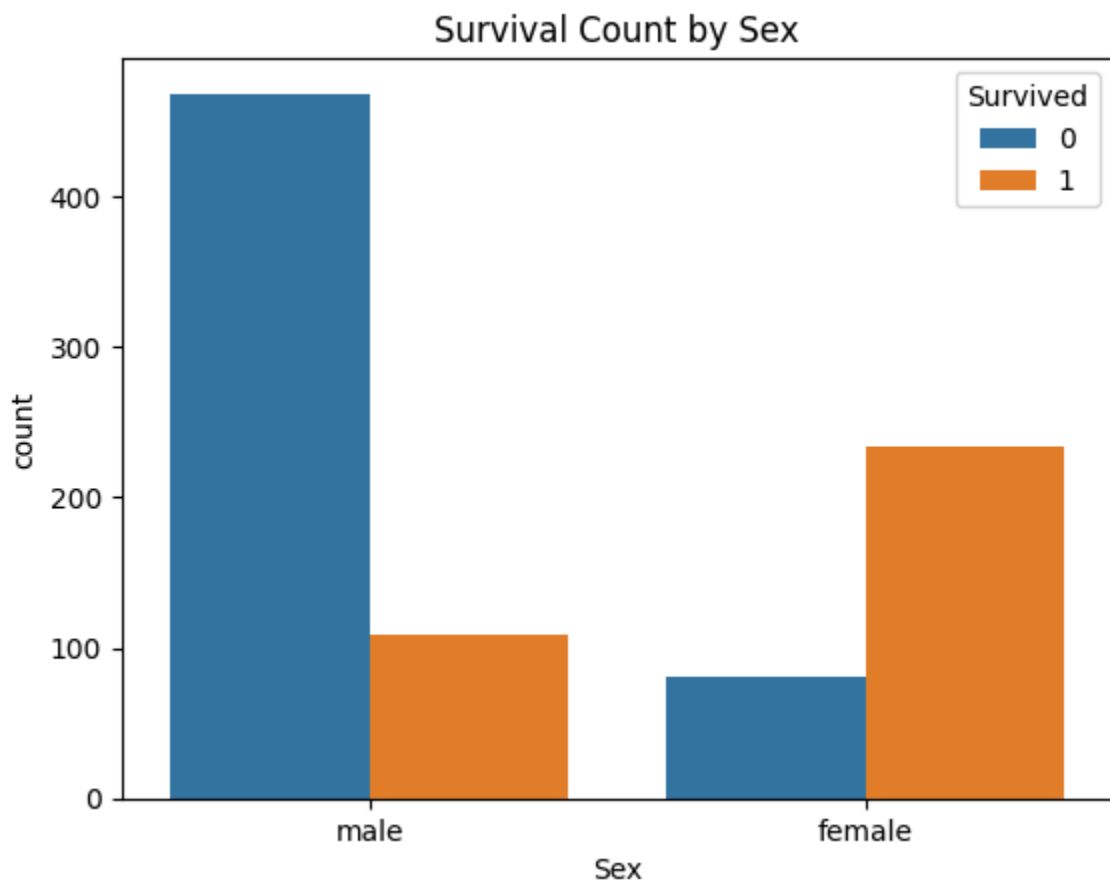


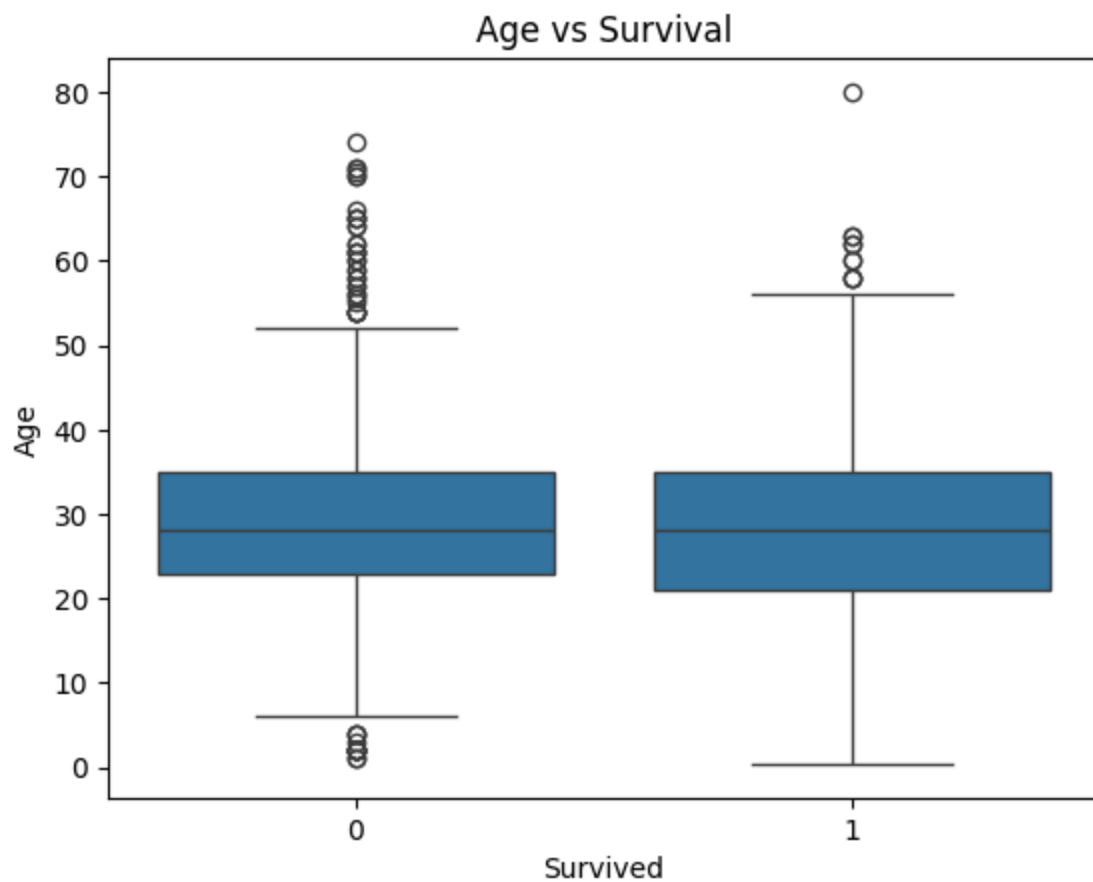


# BIVARIATE ANALYSIS

```
In [24]: # Survived vs Sex
sns.countplot(x='Sex', hue='Survived', data=df)
plt.title('Survival Count by Sex')
plt.show()

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

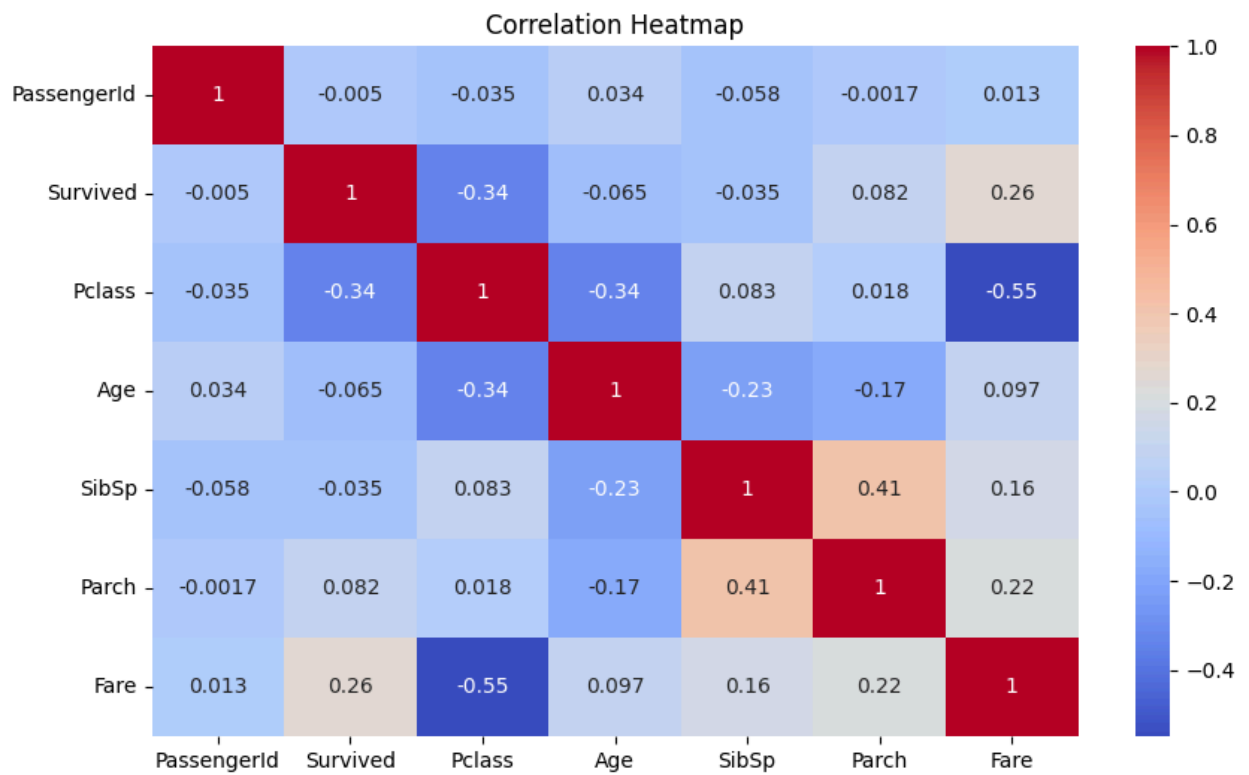




## MULTIVARIATE ANALYSIS

```
In [26]: # Select only numeric columns for correlation
numeric_df = df.select_dtypes(include=['number'])

# Plot heatmap
plt.figure(figsize=(10,6))
sns.heatmap(numeric_df.corr(), annot=True, cmap='coolwarm')
plt.title('Correlation Heatmap')
plt.show()
```



## Summary of Insights:

- Females had a higher survival rate than males.
- Passengers in 1st class had better chances of survival.
- Fare and Age had weak correlation with survival.
- Most passengers were in the age range 20-40.

In [ ]: