

DATA ANALYST INTERNSHIP

Task 5: Exploratory Data Analysis (EDA)

Objective: Extract insights using visual and statistical exploration.

Tools: Python (Pandas, Matplotlib, Seaborn)

Deliverables: Jupyter Notebook + PDF report of findings

Hints/Mini Guide:

- ✓ Use `describe()`, `.info()`, `.value_counts()`
- ✓ Use `sns. pairplot()`, `sns. heatmap()` for visualization
- ✓ Identify relationships and trends
- ✓ Plot histograms, boxplots, scatterplots
- ✓ Write observations for each visual
- ✓ Provide summary of findings

Upload a dataset in python

```
import pandas as pd
df=pd.read_csv("Titanic-Dataset.csv")
df
```

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
...
886	887	0	2	Montvila, Rev. Juozas	male	27.0	0	0	211536	13.0000	NaN	S
887	888	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.0000	B42	S
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	NaN	1	2	W./C. 6607	23.4500	NaN	S
889	890	1	1	Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.0000	C148	C
890	891	0	3	Dooley, Mr. Patrick	male	32.0	0	0	370376	7.7500	NaN	Q

891 rows × 12 columns

Check first five rows

```
df.head()
```

	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

Check last 5 rows

```
df.tail()
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
886	887	0	2	Montvila, Rev. Juozas	male	27.0	0	0	211536	13.00	NaN	S
887	888	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.00	B42	S
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	NaN	1	2	W./C. 6607	23.45	NaN	S
889	890	1	1	Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.00	C148	C
890	891	0	3	Dooley, Mr. Patrick	male	32.0	0	0	370376	7.75	NaN	Q

Check data information

```
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
```

Analyze data description

```
df.describe()
```

	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

Analyze data dimensions

```
df.shape
```

```
(891, 12)
```

counts how many times each unique value appears in a column

```
df.value_counts
```

```
<bound method DataFrame.value_counts of
0      1      0      3
1      2      1      1
2      3      1      3
3      4      1      1
4      5      0      3
..     ...     ...     ...
886    887      0      2
887    888      1      1
888    889      0      3
889    890      1      1
890    891      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
..     ...     ...     ...     ...
886  Montvila, Rev. Juozas    male  27.0      0
887  Graham, Miss. Margaret Edith  female  19.0      0
888  Johnston, Miss. Catherine Helen "Carrie"  female   NaN      1
889  Behr, Mr. Karl Howell    male  26.0      0
890  Dooley, Mr. Patrick    male  32.0      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
..     ...     ...     ...     ...
886    0      211536   13.0000   NaN      S
887    0      112053   30.0000   B42      S
888    2      W./C. 6607   23.4500   NaN      S
889    0      111369   30.0000  C148      C
890    0      370376    7.7500   NaN      Q

[891 rows x 12 columns]>
```

Check is there any null value present

```
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
```

Handle null values

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

df['Embarked'].fillna(df['Embarked'].mode()[0], inplace=True)
df.drop('Cabin', axis=1, inplace=True)
```

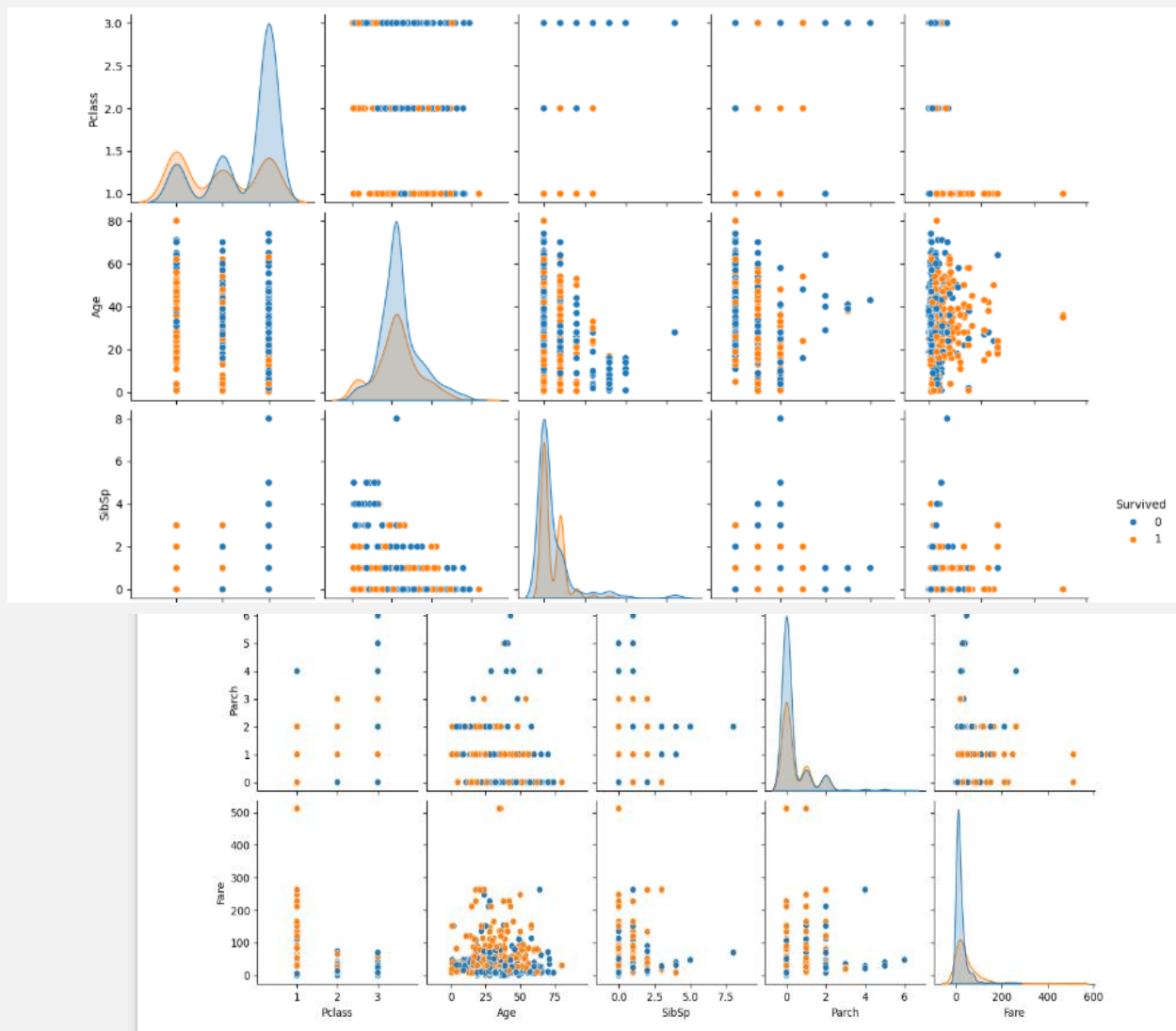
EDA

```
import seaborn as sns
import matplotlib.pyplot as plt

sns.pairplot(df[['Survived', 'Pclass', 'Age', 'SibSp', 'Parch', 'Fare']], hue='Survived')
plt.show()
```

Insights

- ✓ Survivors usually had lower Pclass (higher class).
- ✓ Survivors had higher Fare.
- ✓ Age is spread similarly for survived and non-survived.

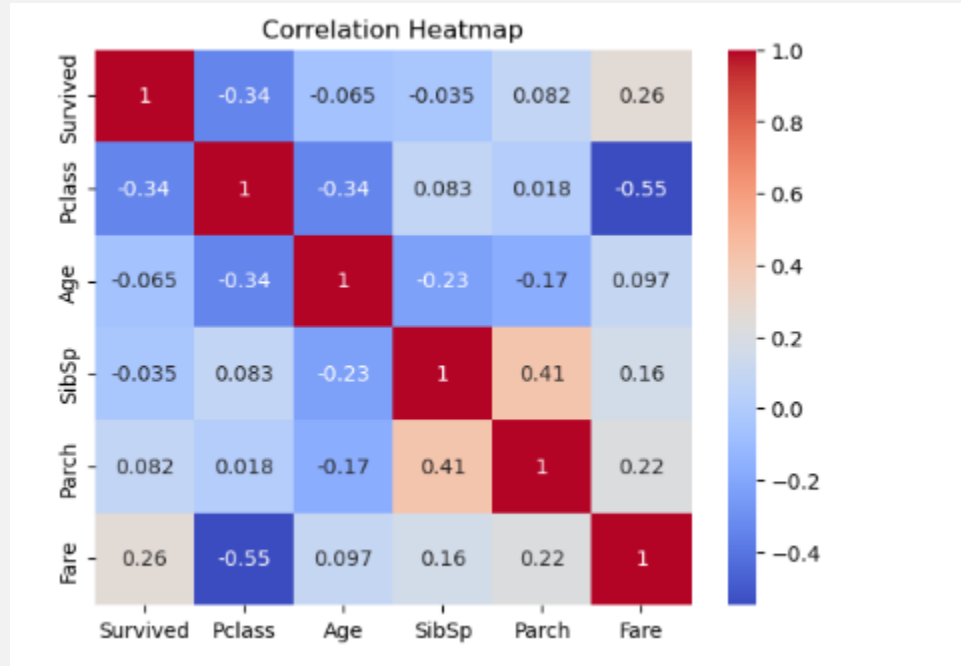


```

: # Correlation matrix
corr = df[['Survived', 'Pclass', 'Age', 'SibSp', 'Parch', 'Fare']].corr()

# Plot heatmap
sns.heatmap(corr, annot=True, cmap='coolwarm')
plt.title("Correlation Heatmap")
plt.show()

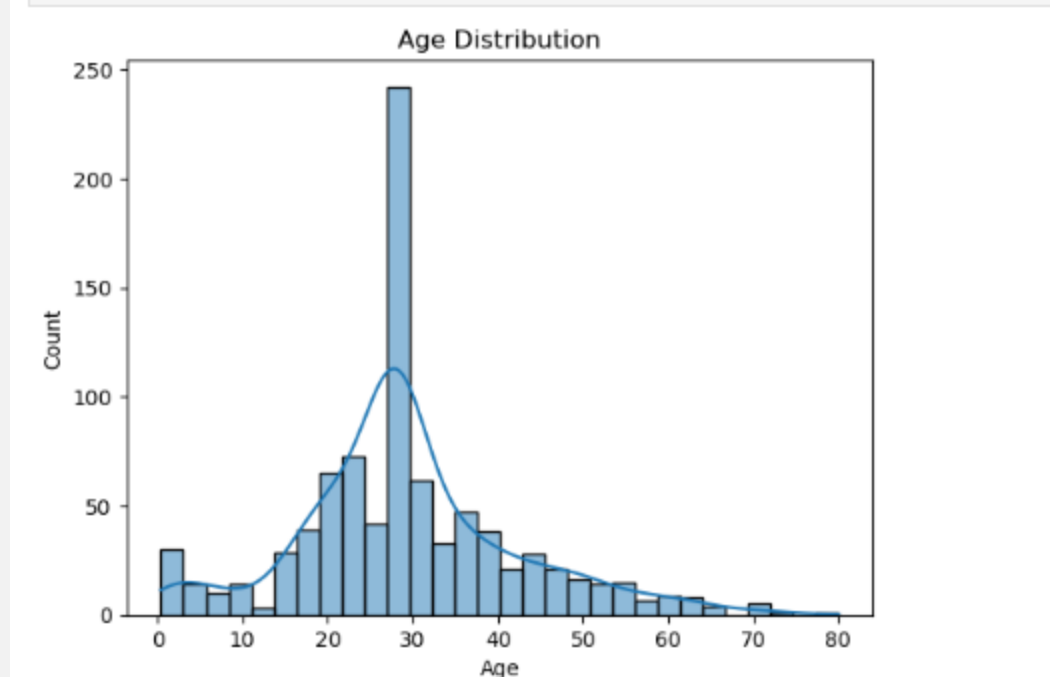
```



Insights

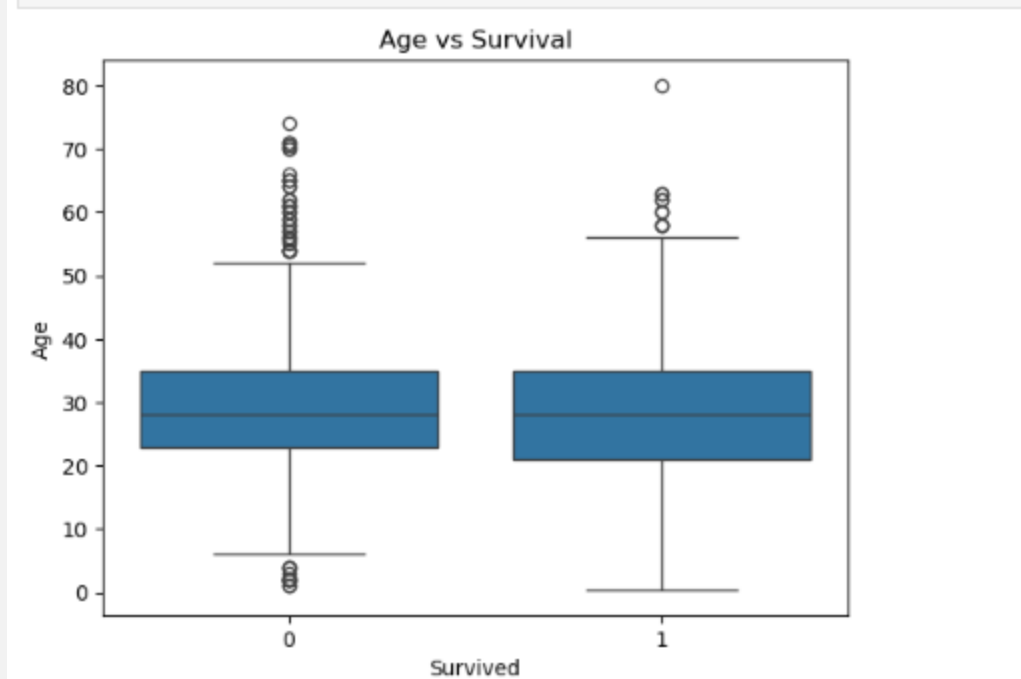
- ✓ Pclass and Fare are negatively correlated (lower Pclass = higher Fare).
- ✓ Fare positively correlates with Survived.

```
: # Plot histogram of Age
sns.histplot(df['Age'], kde=True)
plt.title('Age Distribution')
plt.show()
```



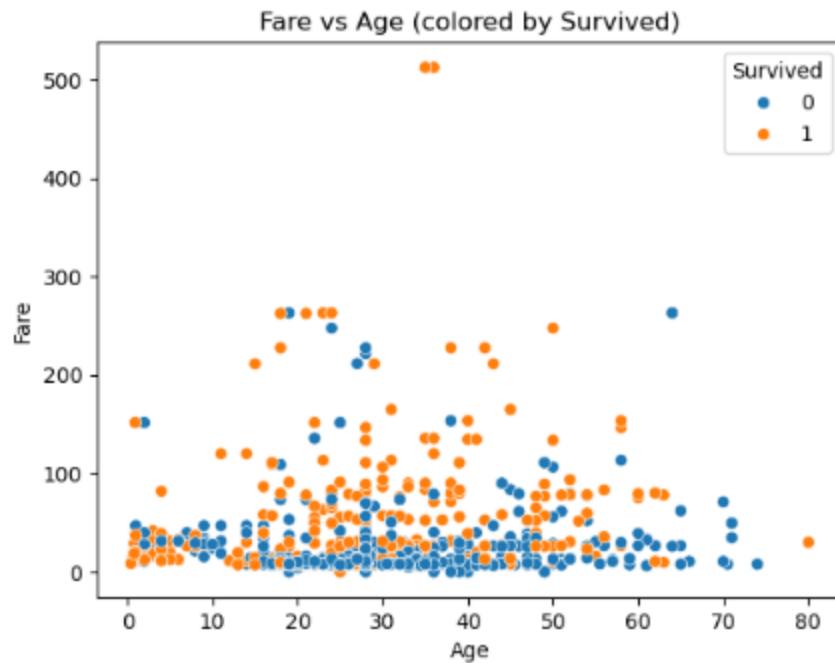
Insights Most passengers are between 20 to 40 years old.


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



Insights Young children had a slightly higher survival rate.

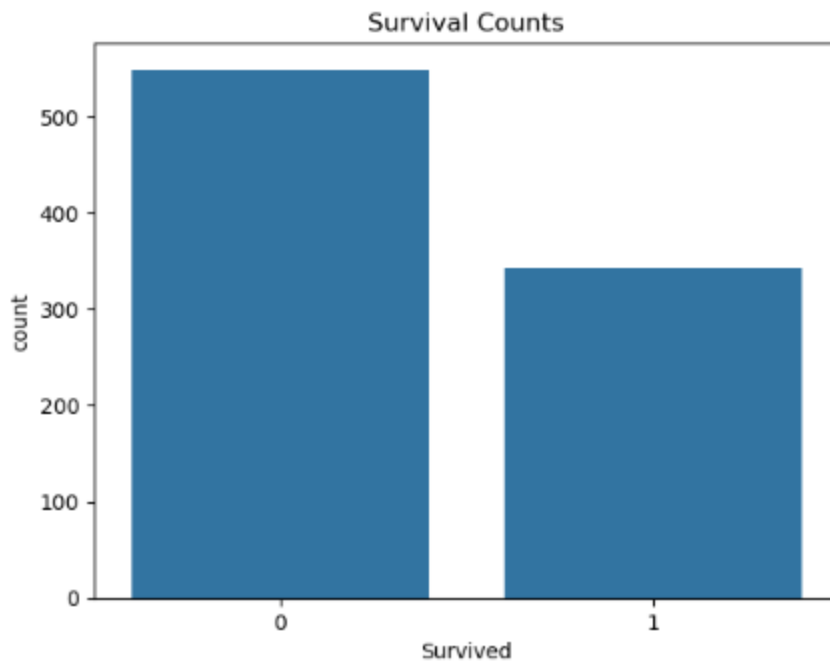
```
# Scatterplot Fare vs Age
sns.scatterplot(x='Age', y='Fare', hue='Survived', data=df)
plt.title('Fare vs Age (colored by Survived)')
plt.show()
```



Insights

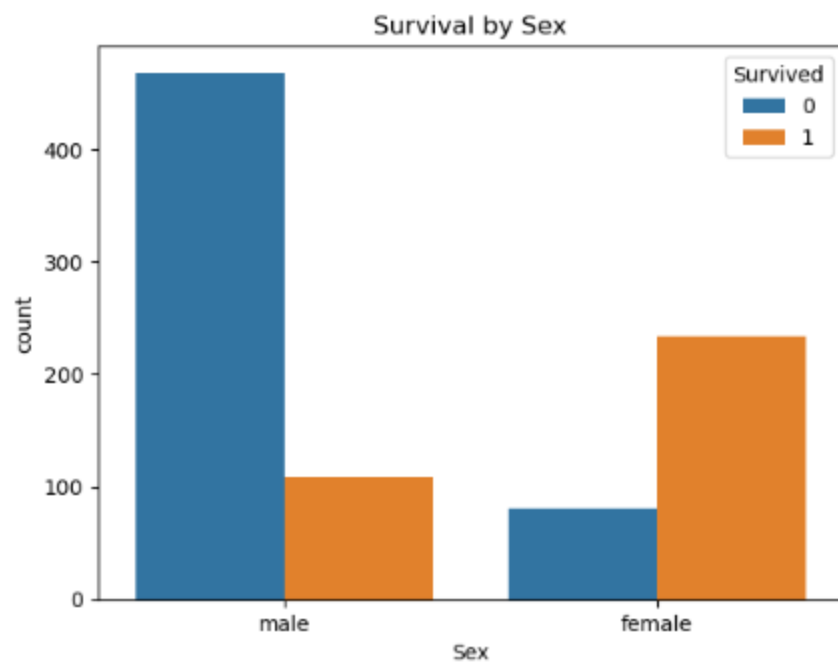
- ✓ Passengers paying higher fare had better survival.
- ✓ Most younger, low-fare passengers did not survive.

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



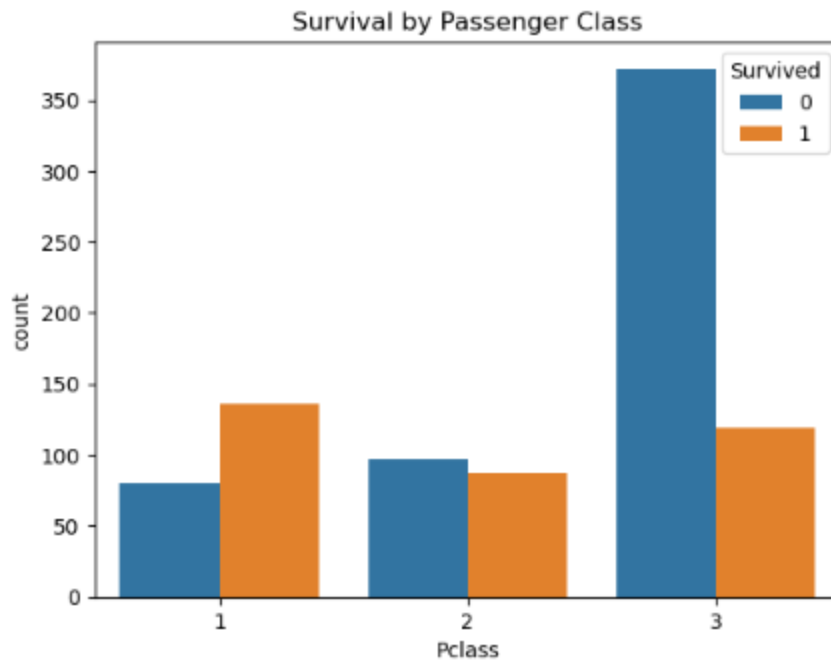
- ✓ More people died (0) than survived (1).
- ✓ Survival rate is about ~38% overall.

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



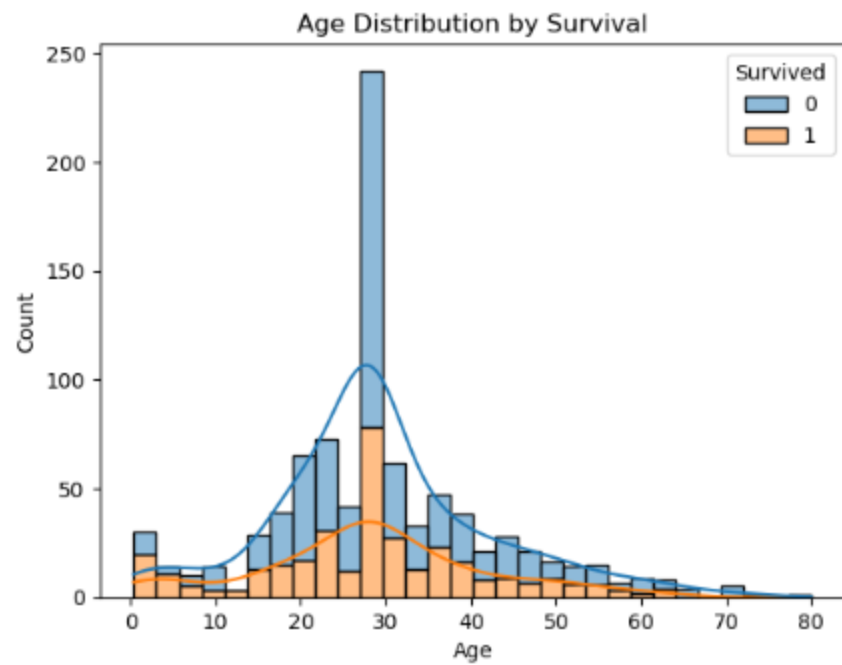
- ✓ Females had a much higher survival rate than males.
- ✓ Titanic applied "women and children first" during evacuation.

```
: sns.countplot(x='Pclass', hue='Survived', data=df)
plt.title('Survival by Passenger Class')
plt.show()
```



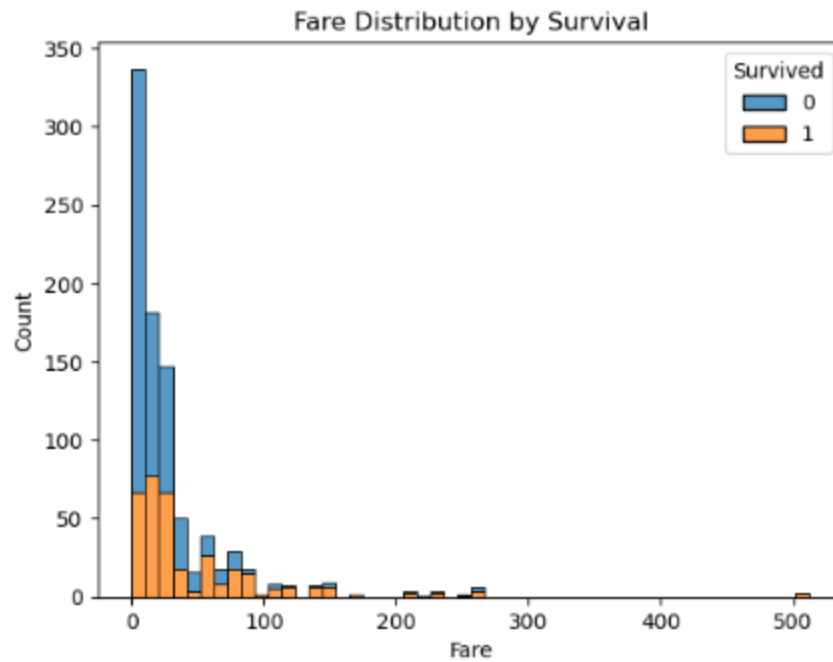
- ✓ First Class passengers survived much more than 2nd and 3rd class.
- ✓ Wealth was strongly linked to survival chance.

```
sns.histplot(data=df, x='Age', hue='Survived', multiple='stack', kde=True)  
plt.title('Age Distribution by Survival')  
plt.show()
```



- ✓ Young children had better survival rates.
- ✓ Many passengers in their 20s and 30s died.

```
: sns.histplot(data=df, x='Fare', hue='Survived', multiple='stack', bins=50)  
plt.title('Fare Distribution by Survival')  
plt.show()
```



- ✓ **Passengers who paid higher fares survived more often.**
- ✓ **Most of the deaths happened among passengers who paid lower fares (cheap tickets).**

```
: sns.violinplot(x='Survived', y='Age', data=df)  
plt.title('Age Distribution by Survival')  
plt.show()
```



- ✓ **Survivors had a more even spread across ages.**
- ✓ **Non-survivors had many more adults (young and middle-aged).**