

```
import pandas as pd
```

```
df = pd.read_csv("train.csv")
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

Futrelle, Mrs. Jacques Heath (Lily May ...

Start coding or [generate](#) with AI.

```
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
```

```
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

```
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
```

## ✦ Initial Data Summary:

- Dataset contains information about Titanic passengers like Age, Fare, Sex, Pclass, and Survived.
- Using `.info()` and `.isnull().sum()` we found:
  - 177 missing values in Age
  - 687 missing values in Cabin
  - 2 missing values in Embarked
- `.describe()` shows Age ranges from 0.42 to 80 and Fare up to 512.
- Next, we will clean the data before analysis.

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

```
# Drop 'Cabin' column due to many missing values
df.drop('Cabin', axis=1, inplace=True)
```

```
# Drop rows where Embarked is missing
df.dropna(inplace=True)
```

⚠ C:\Users\admin\AppData\Local\Temp\ipykernel\_16408\457690679.py:2: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series. The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values is 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, inplace=True)

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

```
# 1. Fill missing Age values with median
df['Age'].fillna(df['Age'].median(), inplace=True)
```

```
# 2. Drop 'Cabin' column (too many missing values)
if 'Cabin' in df.columns:
    df.drop('Cabin', axis=1, inplace=True)
```

```
# 3. Drop rows with missing 'Embarked' values
df.dropna(subset=['Embarked'], inplace=True)
```

⚠ C:\Users\admin\AppData\Local\Temp\ipykernel\_16408\4007831684.py:2: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series. The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values is 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, inplace=True)

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

```
# Fill missing Age values with median (correct & safe method)
df.loc[:, 'Age'] = df['Age'].fillna(df['Age'].median())
```

```
# Drop 'Cabin' column if exists
if 'Cabin' in df.columns:
    df.drop('Cabin', axis=1, inplace=True)
```

```
# Drop rows where 'Embarked' is missing
df.dropna(subset=['Embarked'], inplace=True)
```

```
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
```

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

```
-----
NameError                                Traceback (most recent call last)
Cell In[10], line 1
----> 1 plt.figure(figsize=(10,6))
      2 sns.heatmap(df.corr(), annot=True, cmap='coolwarm')
      3 plt.title("Correlation Heatmap")

NameError: name 'plt' is not defined
```

```
import matplotlib.pyplot as plt
```

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

```
-----
NameError                                Traceback (most recent call last)
Cell In[12], line 2
----> 1 plt.figure(figsize=(10,6))
      2 sns.heatmap(df.corr(), annot=True, cmap='coolwarm')
      3 plt.title("Correlation Heatmap")
      4 plt.show()

NameError: name 'sns' is not defined

<Figure size 1000x600 with 0 Axes>
```

```
# Import all required libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

```
# Show plots inside the notebook
%matplotlib inline
```

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

```
# 2 Countplot - Gender vs Survival
plt.figure(figsize=(6,4))
sns.countplot(x='Survived', hue='Sex', data=df)
plt.title("Survival Count by Gender")
plt.xlabel("Survived (0 = No, 1 = Yes)")
plt.ylabel("Number of Passengers")
plt.show()
```

```
# 3 Boxplot - Age vs Survival
plt.figure(figsize=(6,4))
sns.boxplot(x='Survived', y='Age', data=df)
plt.title("Age Distribution by Survival")
plt.show()
```

```
# 4 Histogram - Age
plt.figure(figsize=(6,4))
```

```
df['Age'].hist(bins=20, edgecolor='black')
plt.title("Age Distribution")
plt.xlabel("Age")
plt.ylabel("Count")
plt.show()

# 5 Histogram - Fare
plt.figure(figsize=(6,4))
df['Fare'].hist(bins=20, edgecolor='black', color='orange')
plt.title("Fare Distribution")
plt.xlabel("Fare")
plt.ylabel("Count")
plt.show()
```

```
-----
ValueError                                Traceback (most recent call last)
Cell In[13], line 12
     10 # 1 Correlation Heatmap
     11 plt.figure(figsize=(10,6))
--> 12 sns.heatmap(df.corr(), annot=True, cmap='coolwarm')
     13 plt.title("Correlation Heatmap")
     14 plt.show()

File ~\anaconda3\Lib\site-packages\pandas\core\frame.py:11049, in DataFrame.corr(self, method, min_periods, numeric_only)
    11047 cols = data.columns
    11048 idx = cols.copy()
-> 11049 mat = data.to_numpy(dtype=float, na_value=np.nan, copy=False)
    11051 if method == "pearson":
    11052     correl = libalgos.nancorr(mat, minp=min_periods)

File ~\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._mgr.as_array(dtype=dtype, copy=copy, na_value=na_value)
    1994 if result.dtype is not dtype:
    1995     result = np.asarray(result, dtype=dtype)

File ~\anaconda3\Lib\site-packages\pandas\core\internals\managers.py:1694, in BlockManager.as_array(self, dtype, copy, na_value)
    1692     arr.flags.writeable = False
    1693 else:
-> 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 ~\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[rl.indexer] = arr
    1754     itemmask[rl.indexer] = 1
    1756 if not itemmask.all():

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

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

```
%matplotlib inline
```

```
df.head()
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Embarked
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	S
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	S
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	S

```
df = pd.read_csv("train.csv")

# Heatmap
plt.figure(figsize=(10,6))
sns.heatmap(df.corr(numeric_only=True), annot=True, cmap='coolwarm')
plt.title("Correlation Heatmap")
plt.show()

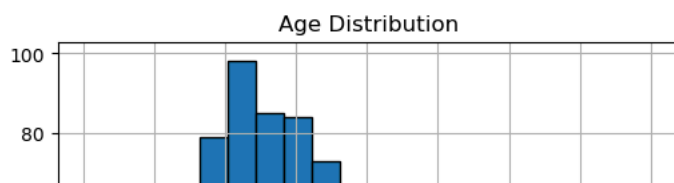
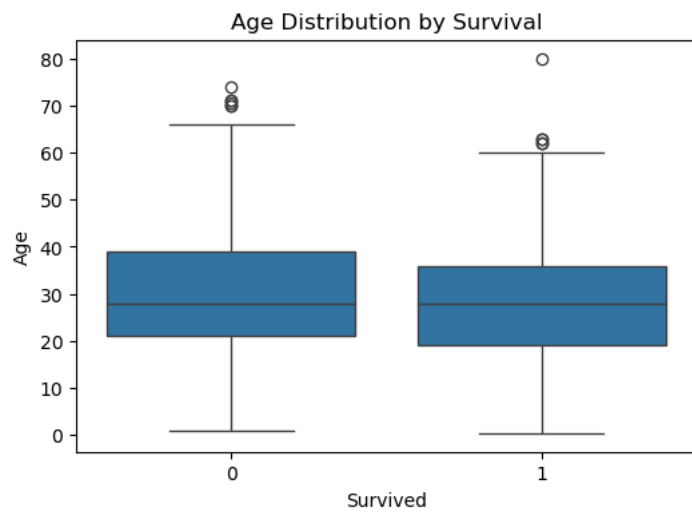
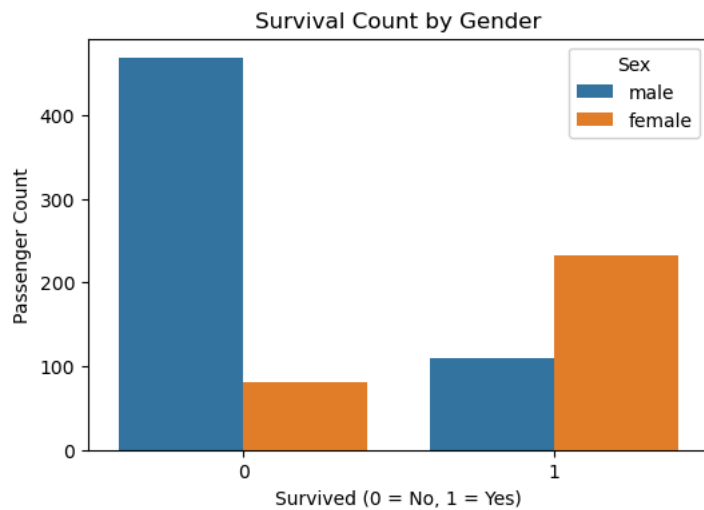
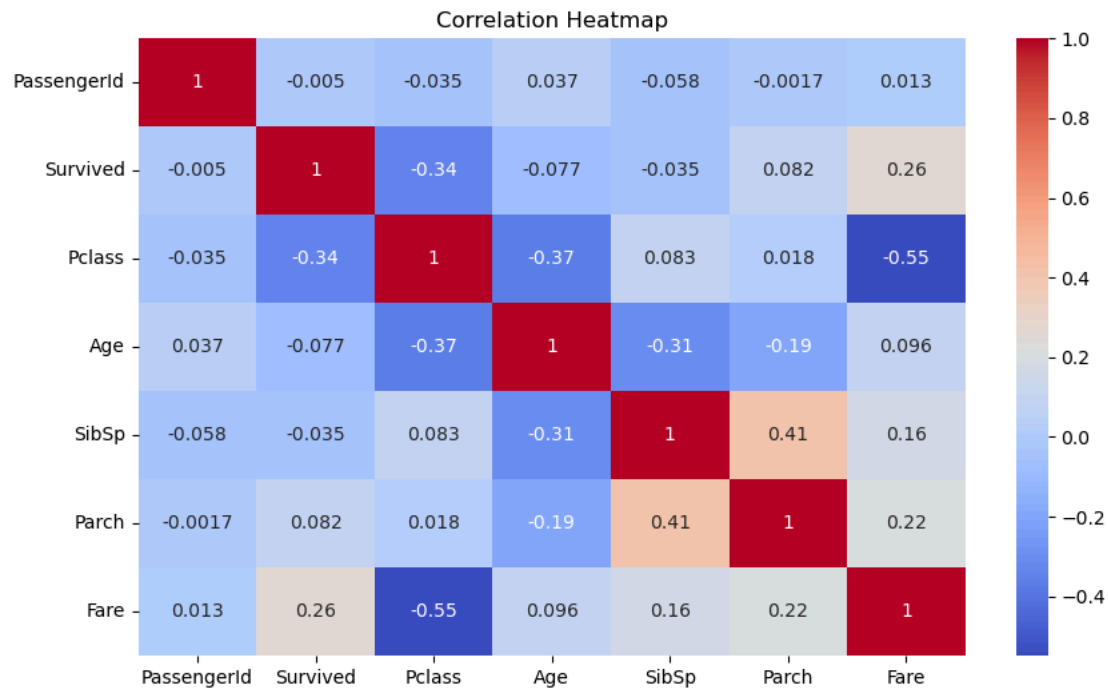
# Countplot: Gender vs Survival
plt.figure(figsize=(6,4))
sns.countplot(x='Survived', hue='Sex', data=df)
plt.title("Survival Count by Gender")
plt.xlabel("Survived (0 = No, 1 = Yes)")
plt.ylabel("Passenger Count")
plt.show()

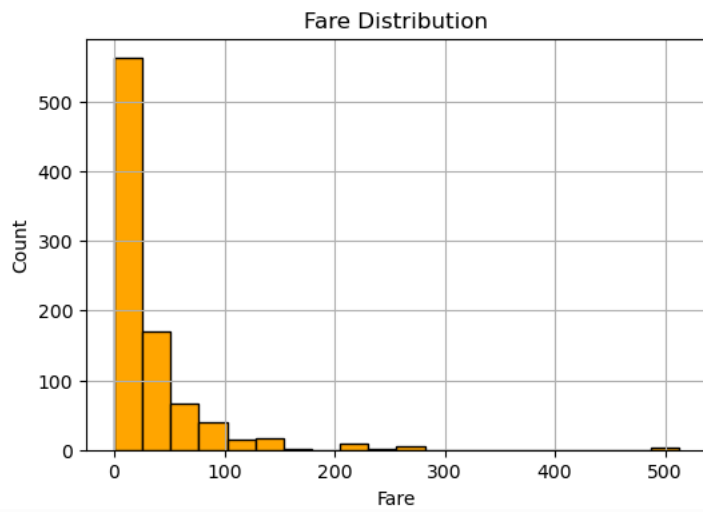
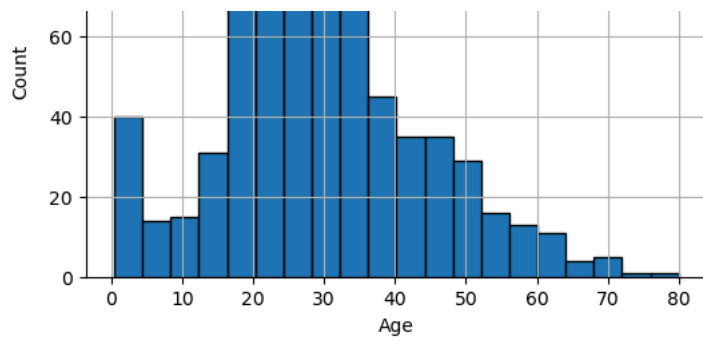
# Boxplot: Age vs Survival
plt.figure(figsize=(6,4))
sns.boxplot(x='Survived', y='Age', data=df)
plt.title("Age Distribution by Survival")
plt.show()

# Histogram - Age
plt.figure(figsize=(6,4))
df['Age'].hist(bins=20, edgecolor='black')
plt.title("Age Distribution")
plt.xlabel("Age")
plt.ylabel("Count")
plt.show()

# Histogram - Fare
plt.figure(figsize=(6,4))
df['Fare'].hist(bins=20, edgecolor='black', color='orange')
plt.title("Fare Distribution")
plt.xlabel("Fare")
plt.ylabel("Count")
plt.show()
```

<Figure size 1000x600 with 0 Axes>





✓ ☒ Final Summary