


```
from google.colab import files
uploaded = files.upload()
```

 Choose Files train.csv


- train.csv(text/csv) - 61194 bytes, last modified: 6/9/2025 - 100% done

Saving train.csv to train.csv

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

%matplotlib inline
```

```
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

Putrelle Mrs. Jacques Heath


Next steps:

[Generate code with df](#)

[View recommended plots](#)


[New interactive sheet](#)

```
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()
```

What can I help you build?


 



	0
<b>PassengerId</b>	0
<b>Survived</b>	0
<b>Pclass</b>	0
<b>Name</b>	0
<b>Sex</b>	0
<b>Age</b>	177
<b>SibSp</b>	0
<b>Parch</b>	0
<b>Ticket</b>	0
<b>Fare</b>	0
<b>Cabin</b>	687
<b>Embarked</b>	2

df.head()


```
df['Sex'].value_counts()
df['Pclass'].value_counts()
df['Embarked'].value_counts()
```



	count
<b>Embarked</b>	
<b>S</b>	644
<b>C</b>	168
<b>Q</b>	77

df.head()

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

 <ipython-input-8-1326d7c8ea59>:1: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment. The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting

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)

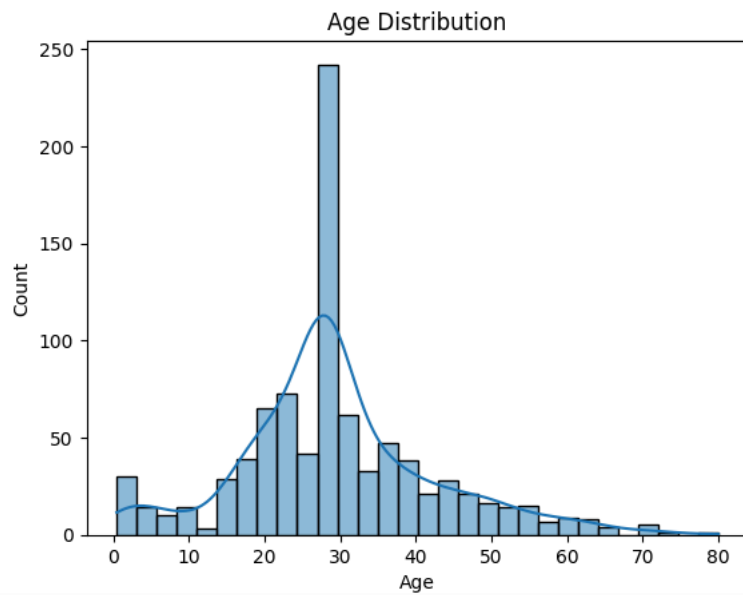
```
df['Age'].fillna(df['Age'].median(), inplace=True)
```

<ipython-input-8-1326d7c8ea59>:2: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment. The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting

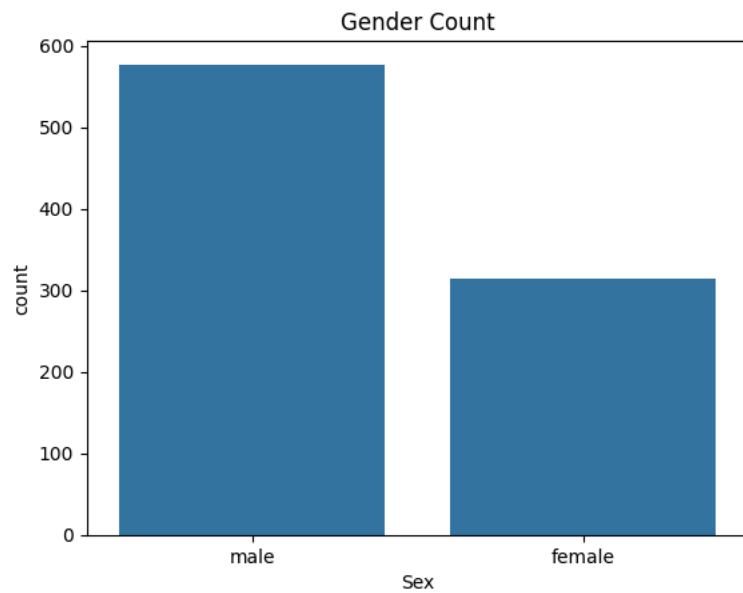
For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col].method(value, inplace=True)

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

```
sns.histplot(df['Age'], bins=30, kde=True)
plt.title('Age Distribution')
plt.show()
```



```
sns.countplot(x='Sex', data=df)
plt.title('Gender Count')
plt.show()
```

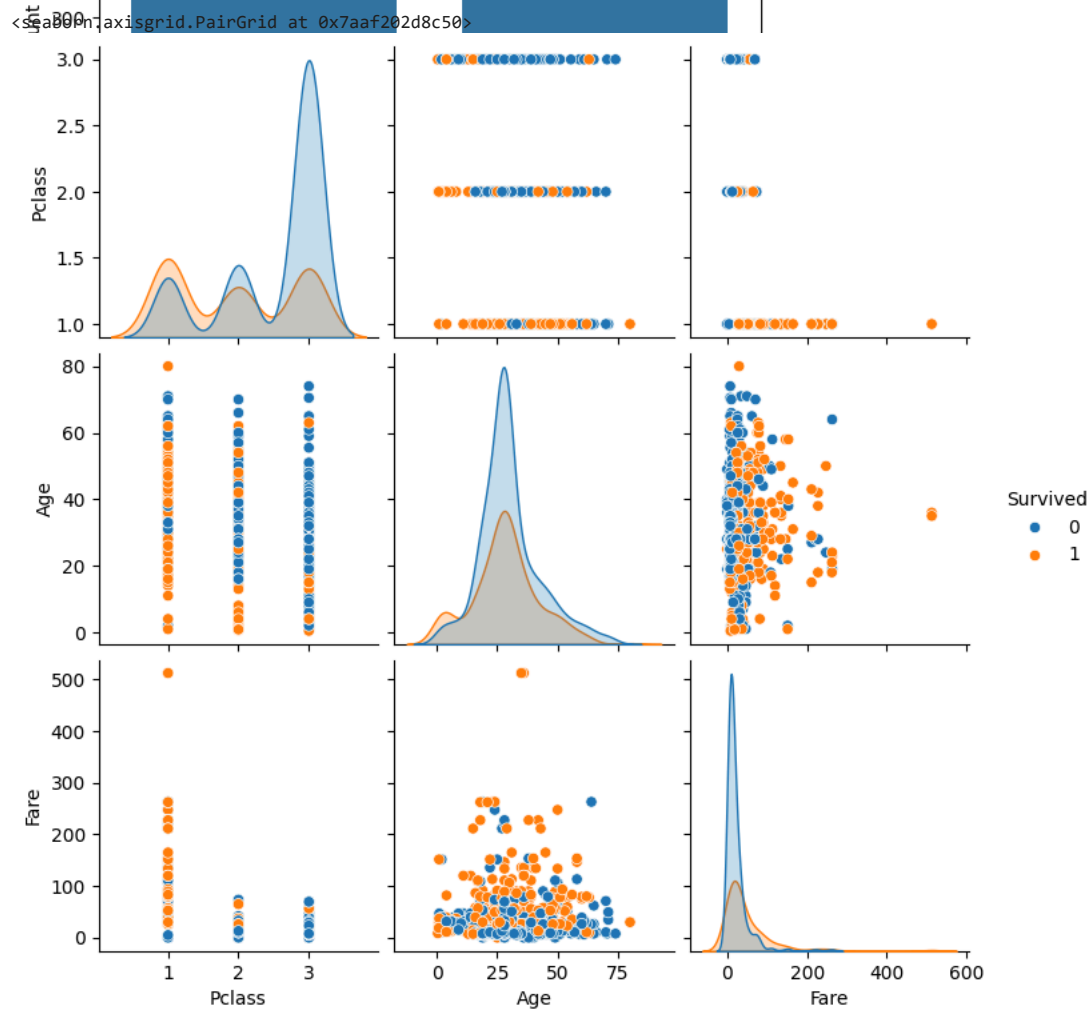


```
sns.countplot(x='Survived', data=df)
plt.title('Survival Count')
plt.show()
```



Survival Count

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



### Summary of Insights:

- Most passengers were between 20–40 years old.
- More male passengers than female.
- Fewer people survived than died.
- Women and 1st class passengers had a better survival rate.
- Fare had a positive correlation with survival.