

# Titanic Dataset Exploratory Data Analysis (EDA) Using Jupyter Notebook

I created a Jupyter Notebook file named **EDA\_Titanic.ipynb** in VS Code and wrote Python scripts to perform exploratory data analysis on the Titanic dataset (**train.csv**). I imported essential libraries such as Pandas, NumPy, Matplotlib, and Seaborn to process and visualize the data.

I began the analysis by loading the dataset and using functions like `head()`, `info()`, and `isnull().sum()` to understand the structure of the data and identify missing values. I handled null values by applying appropriate data-cleaning techniques: **missing Age values were replaced with the median**, and the **Embarked column was filled using the mode**. I also examined categorical value counts for features such as **Sex** and **Embarked** to understand their distributions.

After data cleaning, I performed a series of visualisations to explore patterns and trends within the dataset. These included a **Histogram** for passenger age distribution, a **Box Plot** to detect outliers, a **Scatter Plot** for Age vs Fare, a **Count Plot** for Gender vs Survival, and a **Bar Plot** to compare survival rates across passenger classes. These analyses helped uncover key insights about passenger demographics, fare variation, and survival patterns.

File Edit Selection View Go Run Terminal Help

EDATitanic.ipynb X Task1.py train.csv

Intership > EDA.Titanic.ipynb > #Dataset info

Generate + Code + Markdown | Run All | Restart | Clear All Outputs | Jupyter Variables | Outline ...

Python 3.13.5

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

[1] ✓ 1.3s

```
df = pd.read_csv("D:/New folder/Internship/train.csv")
df.head()
```

[1] ✓ 0.0s

	PassengerId	Survived	Passenger class	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

#Dataset info
df.info()

#Summary statistics
df.describe()

[1] ✓ 0.0s

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890

File Edit Selection View Go Run Terminal Help

EDATitanic.ipynb X Task1.py train.csv

Intership > EDA.Titanic.ipynb > #Dataset info

Generate + Code + Markdown | Run All | Restart | Clear All Outputs | Jupyter Variables | Outline ...

Python 3.13.5

#Summary statistics
df.describe()

[1] ✓ 0.0s

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 12 columns):
 # Column Non-Null Count Dtype
 --- --- --- --- --- --- --- --- --- --- --- --- ---

	PassengerId	Survived	Passenger class	Age	SibSp	Parch	Fare
count	891.000000	891.000000	891.000000	891.000000	891.000000	891.000000	891.000000
mean	446.000000	0.383838	2.308642	29.361582	0.523008	0.381594	32.204208
std	257.353842	0.486592	0.836071	13.019697	1.102743	0.806057	49.693429
min	1.000000	0.000000	0.420000	0.000000	0.000000	0.000000	0.000000
25%	223.500000	0.000000	2.000000	22.000000	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	35.000000	1.000000	0.000000	31.000000
max	891.000000	1.000000	3.000000	80.000000	6.000000	512.329200	512.329200

memory usage: 83.7+ KB

The screenshot shows a Jupyter Notebook interface with the following details:

- File Bar:** File, Edit, Selection, View, Go, Run, Terminal, Help.
- Toolbar:** New folder, Python 3.13.5.
- Left Sidebar (EXPLORER):** OPEN EDITORS, NEW FOLDER, Internship, EDA.Titanic.ipynb, Raw\_data.xlsx, task1.py, train.csv, package1, tables, fmodule.py, 02\_01\_string.py, 02\_02\_striplining.py, 02\_03\_stringfunction..., 2module.py, 03\_practiceset\_01.py, 03\_practiceset\_02.py, 03\_practiceset\_03.py, 03\_practiceset\_04.py, 3module.py, 04\_intro\_methods.py, 04\_intro\_slicing.py, 04\_ifelse.py, 04\_practice\_set.py, 04\_tuples\_methods.py, 04\_tuples.py, 05\_1\_0\_range.py, 05\_bytearray.py, 05\_dictionary\_metho..., 05\_dictionary.py, 05\_set\_operations.py, 05\_sets\_methods.py, 05\_sets.py, 06\_1\_0\_ifelse\_nested..., 06\_if\_else\_els.py, 06\_inits.py, 06\_practiceset.py, 07\_1\_0\_forloop.py, 07\_1\_1\_nestedforloop..., 07\_2\_0whileloop.py, > OUTLINE, 27°C Partly cloudy.
- Code Cells:**
  - #Missing values  
df.isnull().sum()  
0.0s
  - #Categorical value counts  
df['Embarked'].value\_counts()  
0.0s
  - #Categorical value counts  
df['Sex'].value\_counts()  
0.0s
  - df['Age'].fillna(df['Age'].median(), inplace=True)  
0.0s
- Output Cells:**
  - #Categorical value counts  
df['Sex'].value\_counts()  
0.0s
  - df['Age'].fillna(df['Age'].median(), inplace=True)  
0.0s
  - C:\Users\iddmtn\Anaconda\local\Temp\ipykernel\_6892\1933487976.py:1: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method. The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always behaves as a copy.  
For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method([col: value], inplace=True)' or 'df[col].method(value)' instead, to perform the operation inplace on the original DataFrame or Series.  
df['Age'].fillna(df['Age'].median(), inplace=True)
  - df['Embarked'].fillna(df['Embarked'].mode()[0], inplace=True)  
#Missing values  
df.isnull().sum()  
0.0s
- Bottom Bar:** Search, Python 3.13.5, Cell 5 of 13, ENG IN, 16:57, 20-11-2025.

The screenshot shows a Jupyter Notebook interface with the following details:

- File Bar:** File, Edit, Selection, View, Go, Run, Terminal, Help.
- Toolbar:** New folder, Python 3.13.5.
- Left Sidebar (EXPLORER):** OPEN EDITORS, NEW FOLDER, Internship, EDA.Titanic.ipynb, Raw\_data.xlsx, task1.py, train.csv, package1, tables, fmodule.py, 02\_01\_string.py, 02\_02\_striplining.py, 02\_03\_stringfunction..., 2module.py, 03\_practiceset\_01.py, 03\_practiceset\_02.py, 03\_practiceset\_03.py, 03\_practiceset\_04.py, 3module.py, 04\_intro\_methods.py, 04\_intro\_slicing.py, 04\_ifelse.py, 04\_practice\_set.py, 04\_tuples\_methods.py, 04\_tuples.py, 05\_1\_0\_range.py, 05\_bytearray.py, 05\_dictionary\_metho..., 05\_dictionary.py, 05\_set\_operations.py, 05\_sets\_methods.py, 05\_sets.py, 06\_1\_0\_ifelse\_nested..., 06\_if\_else\_els.py, 06\_inits.py, 06\_practiceset.py, 07\_1\_0\_forloop.py, 07\_1\_1\_nestedforloop..., 07\_2\_0whileloop.py, > OUTLINE, 27°C Partly cloudy.
- Code Cells:**
  - #Categorical value counts  
df['Sex'].value\_counts()  
0.0s
  - df['Age'].fillna(df['Age'].median(), inplace=True)  
0.0s
  - C:\Users\iddmtn\Anaconda\local\Temp\ipykernel\_6892\1933487976.py:1: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method. The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always behaves as a copy.  
For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method([col: value], inplace=True)' or 'df[col].method(value)' instead, to perform the operation inplace on the original DataFrame or Series.  
df['Age'].fillna(df['Age'].median(), inplace=True)
  - df['Embarked'].fillna(df['Embarked'].mode()[0], inplace=True)  
#Missing values  
df.isnull().sum()  
0.0s
- Bottom Bar:** Search, Python 3.13.5, Cell 6 of 13, ENG IN, 16:57, 20-11-2025.

File Edit Selection View Go Run Terminal Help

OPEN EDITORS EDA.Titanic.ipynb Task1.py train.csv

Intership EDA\_Titanic.ipynb #Categorical value counts

Generate + Code + Markdown Run All Restart Clear All Outputs Jupyter Variables Outline Python 3.13.5

```
#histogram (passanger age distribution)
sns.histplot(df['Age'], kde=True)
plt.title("Age Distribution of Passengers")
plt.xlabel("Age")
plt.ylabel("Count")
plt.show()
```

045

```
#Boxplot((To detect outliers))
sns.boxplot(x=df['Fare'])
plt.title("Fare Boxplot")
plt.xlabel("Fare")
plt.show()
```

015

File Edit Selection View Go Run Terminal Help

OPEN EDITORS EDA.Titanic.ipynb Task1.py train.csv

Intership EDA\_Titanic.ipynb #Categorical value counts

Generate + Code + Markdown Run All Restart Clear All Outputs Jupyter Variables Outline Python 3.13.5

```
#Boxplot((To detect outliers))
sns.boxplot(x=df['Fare'])
plt.title("Fare Boxplot")
plt.xlabel("Fare")
plt.show()
```

015

```
# Scatterplot( Age vs Fare)
sns.scatterplot(x='Age', y='Fare', data=df)
plt.title("Age vs Fare")
plt.xlabel("Age")
```

025

File Edit Selection View Go Run Terminal Help

EXPLORER OPEN EDITORS EDA.Titanic.ipynb Task1.py train.csv

Intership EDA.Titanic.ipynb #Categorical value counts

Generate + Code + Markdown Run All Restart Clear All Outputs Jupyter Variables Outline Python 3.13.5

```
# Scatterplot( Age vs Fare)
sns.scatterplot(x='Age', y='Fare', data=df)
plt.title("Age vs Fare")
plt.xlabel("Age")
plt.ylabel("Fare")
plt.show()
```

02s

#Countplot - Gender vs Survival

Spaces: 4 CRLF ENG IN 18:57 Cell 6 of 13 20-11-2025

27°C Partly cloudy

File Edit Selection View Go Run Terminal Help

EXPLORER OPEN EDITORS EDA.Titanic.ipynb Task1.py train.csv

Intership EDA.Titanic.ipynb #Categorical value counts

Generate + Code + Markdown Run All Restart Clear All Outputs Jupyter Variables Outline Python 3.13.5

```
#Countplot - Gender vs Survival
sns.countplot(x='Sex', hue='Survived', data=df)
plt.title("Survival Count by Gender")
plt.xlabel("Gender")
plt.ylabel("Count")
plt.show()
```

01s

#Barplot - Survival Rate by Class

Spaces: 4 CRLF ENG IN 18:58 Cell 6 of 13 20-11-2025

27°C Partly cloudy

The screenshot shows a Jupyter Notebook interface with the following details:

- File Bar:** File, Edit, Selection, View, Go, Run, Terminal, Help.
- Toolbar:** New folder, Python 3.13.5.
- Left Sidebar (Explorer):** Shows a tree view of files and folders. The current folder is 'Internship'. Files include 'EDA.Titanic.ipynb', 'Raw\_data.xlsx', 'train.csv', 'Task1.ipynb', and many practice files starting with '01.' through '07.'. A weather icon indicates it's 27°C and partly cloudy.
- Top Bar:** EDA.Titanic.ipynb, Task1.ipynb, train.csv.
- Bottom Bar:** Spaces: 4, CRLF, Cell 6 of 13, ENG IN, 18:58, 20-11-2025.
- Code Cell:** Contains Python code to generate a bar plot:

```
#Barplot - Survival Rate by Class
sns.barplot(x='Passenger class', y='Survived', data=df)
plt.title("Survival Rate by Passenger Class")
plt.xlabel("Class")
plt.ylabel("Survival Rate")
plt.show()
```
- Output Cell:** Displays the generated bar plot titled "Survival Rate by Passenger Class". The x-axis is labeled "Class" with categories 1, 2, and 3. The y-axis is labeled "Survival Rate" ranging from 0.0 to 0.7. The survival rates are approximately 0.63 for Class 1, 0.46 for Class 2, and 0.24 for Class 3.