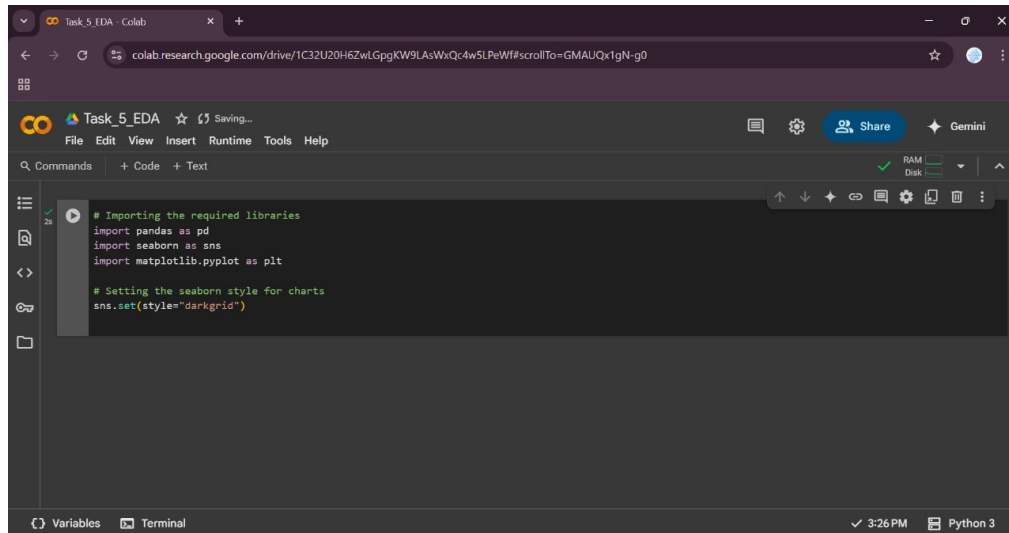# Task 5 – Exploratory Data Analysis (EDA)

## Tool Used: Google Colab (Python, Pandas, Seaborn, Matplotlib)

1. **Import Libraries**

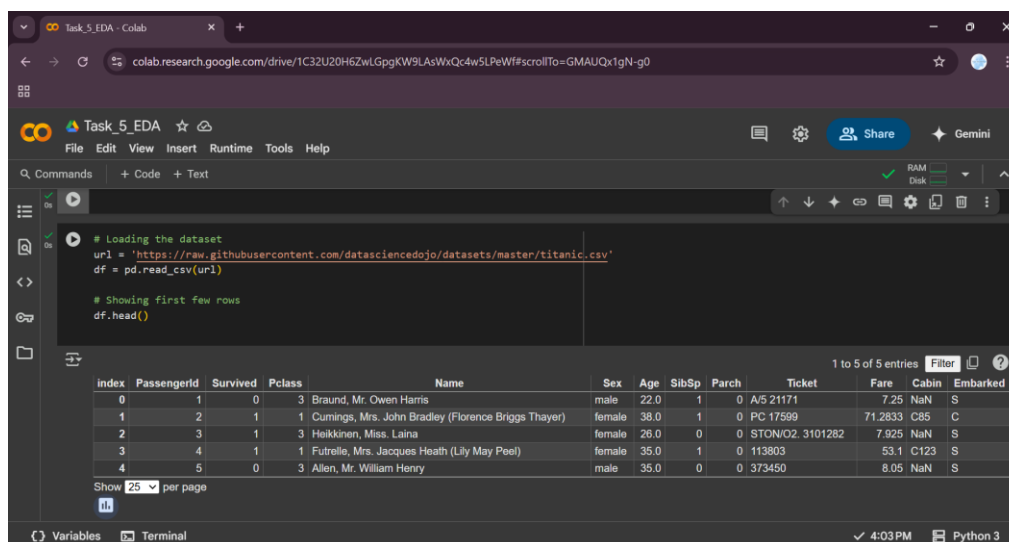

2. **df.head()**



**Observation:**
Shows the first 5 rows of the dataset. Columns like Survived, Pclass, Name, Sex, Age, Fare are available.

3. **df.info()**



**Observation:**
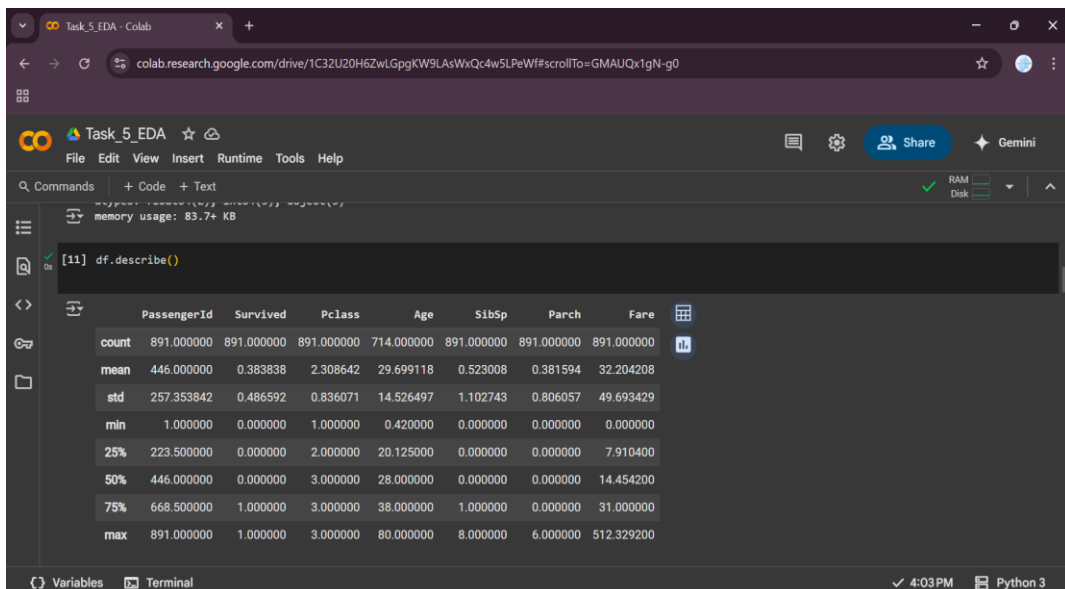Shows that Age and Cabin columns have missing values. Data types are mostly correct.

4. **df.describe()**



**Observation:**
Summary statistics show how Fare has a high range, and Age has outliers and a wide distribution.

5. **df['Survived'].value_counts()**



**Observation:**

More people didn't survive (0) than survived (1). Helps understand the class imbalance.

6. **Age Histogram**



**Observation:**

Most passengers are between 20 and 40 years old. Very few children or seniors.

7. **Fare Boxplot**



**Observation:**
There are a few very high fare values (outliers), indicating some rich passengers.

8. **Correlation Heatmap**

**Observation:**
Survived has slight positive correlation with Fare and negative correlation with Pclass.


9. **Pairplot**



**Observation:**
Shows relationship between variables. Survival might relate to Fare and Pclass.

**Final Summary:**

-The dataset contains 891 passengers with 12 columns.
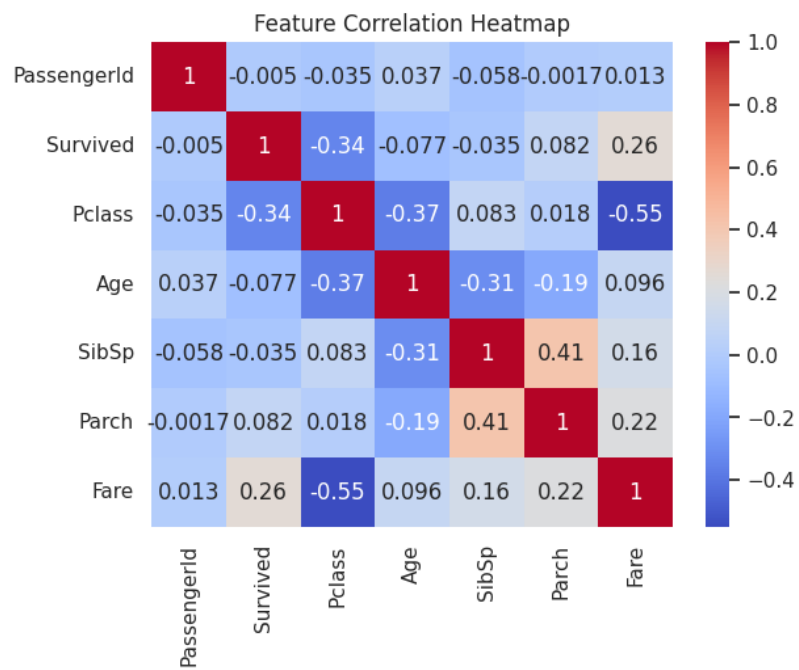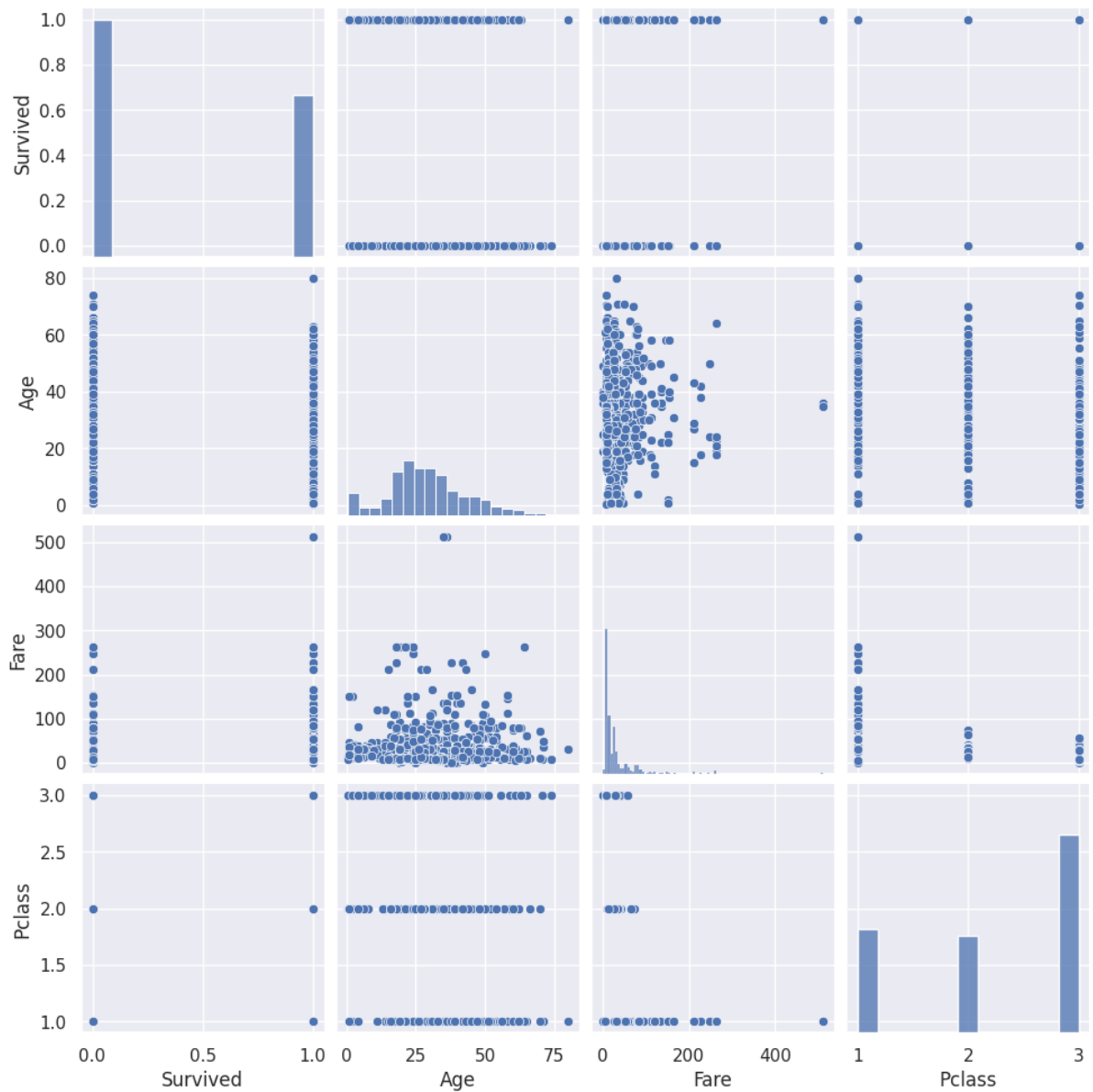
- Missing values found in 'Age', 'Cabin', and 'Embarked'.

- Most passengers are in the age group 20–40.

- Survival rate is low; more people did not survive.

- Fare has large outliers; some paid extremely high amounts.

- Correlation shows survival is affected by Fare and Pclass.

- Females had higher survival rates than males.