

EXPLORATORY DATA ANALYSIS (EDA) - TITANIC DATASET

```
[1] import pandas as pd  
import seaborn as sns  
import matplotlib.pyplot as plt  
  
| df = pd.read_csv('/content/drive/My Drive/titanic/train.csv')
```

Import pandas, seaborn, and matplotlib.pyplot to handle data and create visualizations. We load the Titanic dataset from Google Drive using pd.read_csv() and store it in a DataFrame called df.

```
[9] print("Jumlah data duplikat:", df.duplicated().sum())
→ Jumlah data duplikat: 0
```

Use `df.duplicated().sum()` to check for duplicate rows, and the result shows there are none in the dataset.

```
[10] print("Missing value sebelum ditangani:")
      print(df.isna().sum())
```

→ Missing value sebelum ditangani:

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	

By using `df.isna().sum()`, we identify the number of missing values in each column, revealing that the Age column has 177 missing entries, Cabin has 687, and Embarked has 2, which indicates the need for proper data cleaning before analysis.

```
[10] print("Missing value sebelum ditangani:")
      print(df.isna().sum())
```

→ Missing value sebelum ditangani:

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	

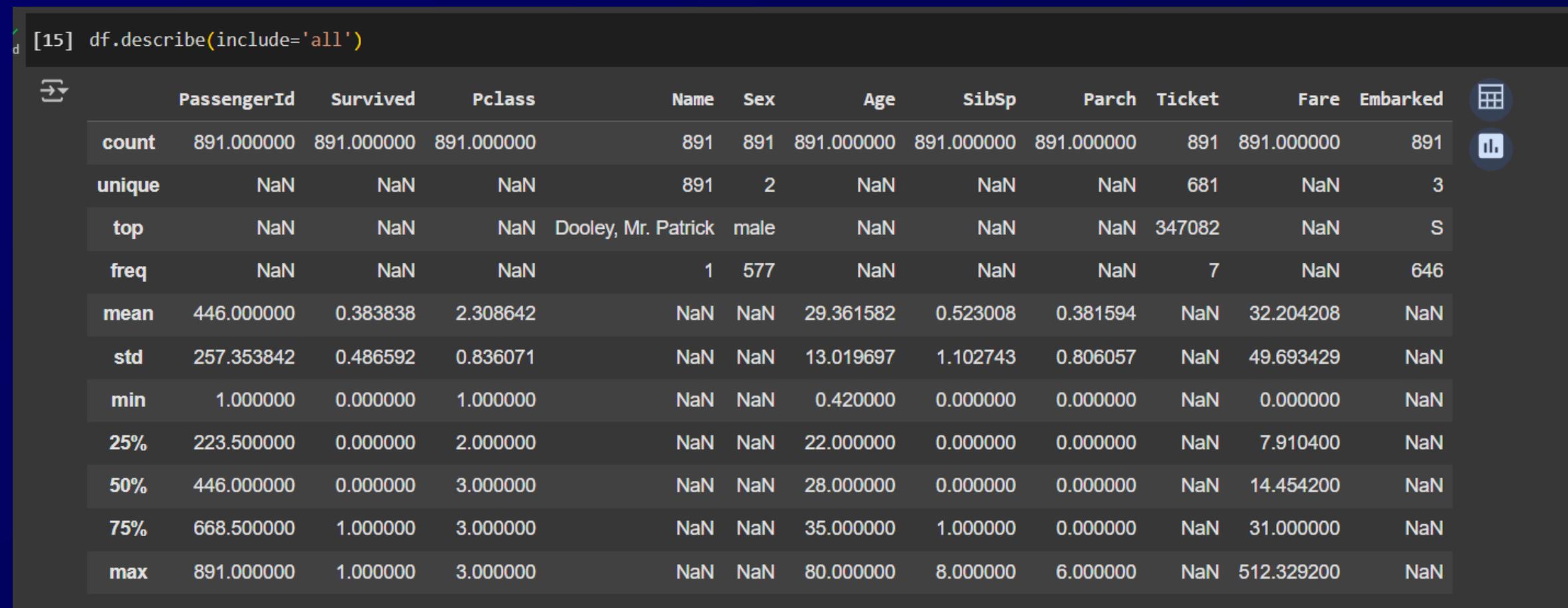
By using `df.isna().sum()`, we identify the number of missing values in each column, revealing that the Age column has 177 missing entries, Cabin has 687, and Embarked has 2, which indicates the need for proper data cleaning before analysis.

```
df['Age'].fillna(df['Age'].median(), inplace=True)
ipython-input-11-63d4fb902a4f:1: FutureWarning: A value is trying to be set on a
The behavior will change in pandas 3.0. This inplace method will never work because
For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method
df['Age'].fillna(df['Age'].median(), inplace=True)

[12] df.drop('Cabin', axis=1, inplace=True)

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

To clean the dataset, we replaced missing values in the Age column with the median value to avoid distortion by outliers, dropped the Cabin column entirely due to a large number of missing entries, and filled the missing values in the Embarked column with its most frequent value (mode), ensuring the dataset is complete and ready for analysis.



	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Embarked
count	891.000000	891.000000	891.000000	891	891	891.000000	891.000000	891.000000	891	891.000000	891
unique	Nan	Nan	Nan	891	2	Nan	Nan	Nan	681	Nan	3
top	Nan	Nan	Nan	Dooley, Mr. Patrick	male	Nan	Nan	Nan	347082	Nan	S
freq	Nan	Nan	Nan	1	577	Nan	Nan	Nan	7	Nan	646
mean	446.000000	0.383838	2.308642	Nan	Nan	29.361582	0.523008	0.381594	Nan	32.204208	Nan
std	257.353842	0.486592	0.836071	Nan	Nan	13.019697	1.102743	0.806057	Nan	49.693429	Nan
min	1.000000	0.000000	1.000000	Nan	Nan	0.420000	0.000000	0.000000	Nan	0.000000	Nan
25%	223.500000	0.000000	2.000000	Nan	Nan	22.000000	0.000000	0.000000	Nan	7.910400	Nan
50%	446.000000	0.000000	3.000000	Nan	Nan	28.000000	0.000000	0.000000	Nan	14.454200	Nan
75%	668.500000	1.000000	3.000000	Nan	Nan	35.000000	1.000000	0.000000	Nan	31.000000	Nan
max	891.000000	1.000000	3.000000	Nan	Nan	80.000000	8.000000	6.000000	Nan	512.329200	Nan

Using `df.describe(include='all')`, we generate a summary of both numerical and categorical columns, showing statistics like mean, median, standard deviation for numeric data (Age, Fare, etc.) and frequency distribution for categorical data (Sex, Embarked), which helps us understand the overall distribution, spread, and common values in the dataset.

[Home](#)[About](#)[Content](#)[Others](#)[Page 10](#)

THANK YOU
