# LAB NO 01

## ⌄ Task 1 => Dataset Cleaning

1. Load the Titanic dataset (train.csv).
2. Display the first 10 rows.
3. Check for missing values in each column.
4. Fill missing values in the "Age" column with the mean age.
5. Drop rows where the "Embarked" column is missing.

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

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

```
df.head(10)
```

| | 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 |
| **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 |
| **5** | 6 | 0 | 3 | Moran, Mr. James | male | NaN | 0 | 0 | 330877 | 8.4583 | NaN | Q |
| **6** | 7 | 0 | 1 | McCarthy, Mr. Timothy J | male | 54.0 | 0 | 0 | 17463 | 51.8625 | E46 | S |
| **7** | 8 | 0 | 3 | Palsson, Master. Gosta Leonard | male | 2.0 | 3 | 1 | 349909 | 21.0750 | NaN | S |

Next steps: Generate code with `df`    New interactive sheet

```
df.isnull().sum()
```

| | 0 |
|---|---|
| **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

```
df['Age'].mean()
```

```
np.float64(29.69911764705882)
```

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

```
/tmp/ipython-input-694922604.py:1: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained as
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[co

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

```
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 |
| **3** | 4 | 1 | 1 | Futrelle, Mrs. Jacques Heath | female | 35.0 | 1 | 0 | 113803 | 53.1000 | C123 | S |

Next steps: `Generate code with df`  `New interactive sheet`

```
df.isnull().sum()
```

| | 0 |
|---|---|
| **PassengerId** | 0 |
| **Survived** | 0 |
| **Pclass** | 0 |
| **Name** | 0 |
| **Sex** | 0 |
| **Age** | 0 |
| **SibSp** | 0 |
| **Parch** | 0 |
| **Ticket** | 0 |
| **Fare** | 0 |
| **Cabin** | 687 |
| **Embarked** | 2 |

**dtype:** int64

```
df.dropna(subset=['Embarked'])
```

| | PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 0 | 3 | Braund, Mr. Owen Harris | male | 22.000000 | 1 | 0 | A/5 21171 | 7.2500 | NaN | S |
| **1** | 2 | 1 | 1 | Cumings, Mrs. John Bradley (Florence Briggs Th... | female | 38.000000 | 1 | 0 | PC 17599 | 71.2833 | C85 | C |
| **2** | 3 | 1 | 3 | Heikkinen, Miss. Laina | female | 26.000000 | 0 | 0 | STON/O2. 3101282 | 7.9250 | NaN | S |
| **3** | 4 | 1 | 1 | Futrelle, Mrs. Jacques Heath (Lily May Peel) | female | 35.000000 | 1 | 0 | 113803 | 53.1000 | C123 | S |
| **4** | 5 | 0 | 3 | Allen, Mr. William Henry | male | 35.000000 | 0 | 0 | 373450 | 8.0500 | NaN | S |
| **...** | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| **886** | 887 | 0 | 2 | Montvila, Rev. Juozas | male | 27.000000 | 0 | 0 | 211536 | 13.0000 | NaN | S |
| **887** | 888 | 1 | 1 | Graham, Miss. Margaret Edith | female | 19.000000 | 0 | 0 | 112053 | 30.0000 | B42 | S |

```
df.shape
```

```
(891, 12)
```

## Task 2 ->Encoding Categorical Data

1. Convert the "Sex" column into numeric (0 = Male, 1 = Female).
2. Apply One-Hot Encoding on the "Embarked" column.

```
from sklearn.preprocessing import OneHotEncoder
```

```
df['Sex']=df['Sex'].map({'male':0, 'female':1})
```

```
pd.get_dummies(df,columns=['Embarked'],drop_first=True)
```

| | PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked_Q | Embarked_S |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 0 | 3 | Braund, Mr. Owen Harris | 0 | 22.000000 | 1 | 0 | A/5 21171 | 7.2500 | NaN | False | True |
| **1** | 2 | 1 | 1 | Cumings, Mrs. John Bradley (Florence Briggs Th... | 1 | 38.000000 | 1 | 0 | PC 17599 | 71.2833 | C85 | False | False |
| **2** | 3 | 1 | 3 | Heikkinen, Miss. Laina | 1 | 26.000000 | 0 | 0 | STON/O2. 3101282 | 7.9250 | NaN | False | True |
| **3** | 4 | 1 | 1 | Futrelle, Mrs. Jacques Heath (Lily May Peel) | 1 | 35.000000 | 1 | 0 | 113803 | 53.1000 | C123 | False | True |
| **4** | 5 | 0 | 3 | Allen, Mr. William Henry | 0 | 35.000000 | 0 | 0 | 373450 | 8.0500 | NaN | False | True |
| **...** | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |

## Task 3 -> Feature Scaling & Splitting

1. Select features: Age, Fare, Sex, Pclass.
2. Apply StandardScaler to normalize them.
3. Split data into 80% training and 20% testing.

```python
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
```

```python
x=df[['Age','Fare','Sex','Pclass']]
y=df['Survived']
```

```python
scaler = StandardScaler()
X_scaled = scaler.fit_transform(x)
```

```python
X_train, X_test, y_train, y_test = train_test_split(
    X_scaled, y, test_size=0.2, random_state=42
)
```