

E23CSEU0055_Lab01_CSET301(P)

Name - Arihant Gupta - E23CSEU0055

Batch - EB02

Date - 31.07.2025

Ipynb File - [cset301-aiml-sem5/jul29/one.ipynb at main · GuptajiRocks/cset301-aiml-sem5](#)

1. Load the Dataset

```
a. 1 import pandas as pd
   2 from sklearn.preprocessing import LabelEncoder, StandardScaler
   3 import matplotlib.pyplot as plt
   4
   5 df =
   pd.read_csv("https://raw.githubusercontent.com/datasciencedojo/datasets/master/titanic.csv")
```

2. Identify and Handle Missing Values

```
a. 1 # Code
   2 print(df.isnull().sum())
   3
   4 #Output
   5 PassengerId      0
   6 Survived        0
   7 Pclass          0
   8 Name            0
   9 Sex             0
  10 Age            177
  11 SibSp           0
  12 Parch           0
  13 Ticket          0
  14 Fare            0
  15 Cabin          687
  16 Embarked        2
  17 dtype: int64
```

Checking if any null values are present. We can see columns “Age”, “Cabin”, “Embarked” are having null values. Age is the only numeric value, so we will impute the NaN values with the mean of the age column.

```
1 df["Age"] = df["Age"].fillna(df["Age"].mean())
2
3 #Age NaN values have been imputed with mean now.
```

Since the column Cabin has very high NaN values, is categorical in nature with every entry unique and not affecting the result in any way whatsoever. We will drop the cabin column and create a new DataFrame without it.

```
b. 1 # Code
2   df_n = df.drop("Cabin", axis=1)
3   print(df_n.isnull().sum())
4
5   #Output
6   PassengerId      0
7   Survived         0
8   Pclass           0
9   Name             0
10  Sex              0
11  Age              0
12  SibSp            0
13  Parch            0
14  Ticket           0
15  Fare             0
16  Embarked         2
17  dtype: int64
```

3. Deal with Duplicate Data

```
a. 1 # Code
2   print(df.duplicated().sum())
3
4   #Output
5   0
```

4. Converting Categorical Column to Numerical

a. We can see Age has no null values, Cabin has been dropped. Now comes Embarked column. Embarked acts as the class label and has three unique values namely “E”, “C” and “S”. To fill in the NaN value, I will impute a value “Z” as it would not mean anything but will remove the NaN. We’ll apply Label Encoding to the column and convert it to numerical value for further analysis.

```
b. 1 #Code
2   print(df_n[df_n["Embarked"].isnull()])
3
4   #Output
5   PassengerId  Survived  Pclass  Name \
6   61           62        1      1      Icard, Miss. Amelie
7   829          830        1      1  Stone, Mrs. George Nelson (Martha Evelyn)
```

8								
9		Sex	Age	SibSp	Parch	Ticket	Fare	Embarked
10	61	female	38.0	0	0	113572	80.0	NaN
11	829	female	62.0	0	0	113572	80.0	NaN

```
C. 1 # Fillna with "Z"
2 df_n["Embarked"] = df_n["Embarked"].fillna("Z")
3
4 le = LabelEncoder()
5 df_n["Embarked"] = le.fit_transform(df_n["Embarked"])
6 print(df_n["Embarked"].unique())
7
8 #Output
9 [2 0 1 3]
```

The unique numerical values of Embarked column suggest that the categorical data has been converted to numerical form.

5. Normalize appropriate numerical features.

```
a. 1 stdC = StandardScaler()
2 df_n["Age"] = stdC.fit_transform(df_n[["Age"]])
```

Age has been normalized using Z-score normalization - (StandardScaler() from Scikit-Learn)

6. Apply sorting and filtering logic

a. I will sort the values based on Age.

```
i. 1 # Sorting
2 print(df_n[["Name", "Sex", "Age"]].sort_values(by="Age"))
3
4 #Output
5 Name      Sex      Age
6 803      Thomas, Master. Assad Alexander    male -2.253155
7 755      Hamalainen, Master. Viljo    male -2.233917
8 644      Baclini, Miss. Eugenie    female -2.227761
9 469      Baclini, Miss. Helene Barbara    female -2.227761
10 831      Richards, Master. George Sibley    male -2.221604
11 ..      ...      ...
12 116      Connors, Mr. Patrick    male 3.139805
13 96      Goldschmidt, Mr. George B    male 3.178283
14 493      Artagaveytia, Mr. Ramon    male 3.178283
15 851      Svensson, Mr. Johan    male 3.409146
16 630      Barkworth, Mr. Algernon Henry Wilson    male 3.870872
17
18 [891 rows x 3 columns]
```

7. Engineer new column based on Logic.

- a. Filtering out the people who have Fare more than the average fare. And creating a new column of Status, wherein people who have fare higher than average are rich, and others are poor.

b.

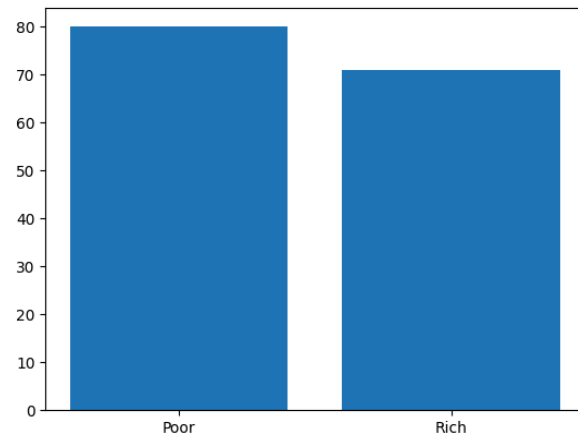
```
1 #Code
2 mufare = df_n["Fare"].mean()
3 df_n["Status"] = df["Fare"].apply(lambda x : "Rich" if x > mufare else "Poor")
4 print(df_n[["Fare", "Status"]])
5
6 #Output
7 Fare Status
8 0      7.2500  Poor
9 1     71.2833  Rich
10 2      7.9250  Poor
11 3     53.1000  Rich
12 4      8.0500  Poor
13 ..      ...    ...
14 886  13.0000  Poor
15 887  30.0000  Poor
16 888  23.4500  Poor
17 889  30.0000  Poor
18 890   7.7500  Poor
19
20 [891 rows x 2 columns]
```

8. Visualization

- a. Creating a Bar chart comparing the Age and Status of members in the data.
- b. We are comparing status with Age that is before normalization to get a numerical understanding.

c.

```
1 #Code
2 print(df["Age"])
3 plt.bar(df_n["Status"], height=df["Age"])
4 plt.show()
5
```



We can see that Poor people have Age range from 0 to 80, and Rich People have Age range from 0 to 70.