# TITANIC DATASET ANALYSIS AND PREPROCESSING USING PANDAS AND SIMPLE IMPUTER.

5/8/27

Aim :

To read titanic dataset from csv, handle missing values using simple imputer, analyze key passanger features.

procedure / Alg:

Step 1: load titanic.csv into a data frame

Step 2: Explore dataset shap, info and summary statistics.

step 3: Use simple Imputer to fill missing age.

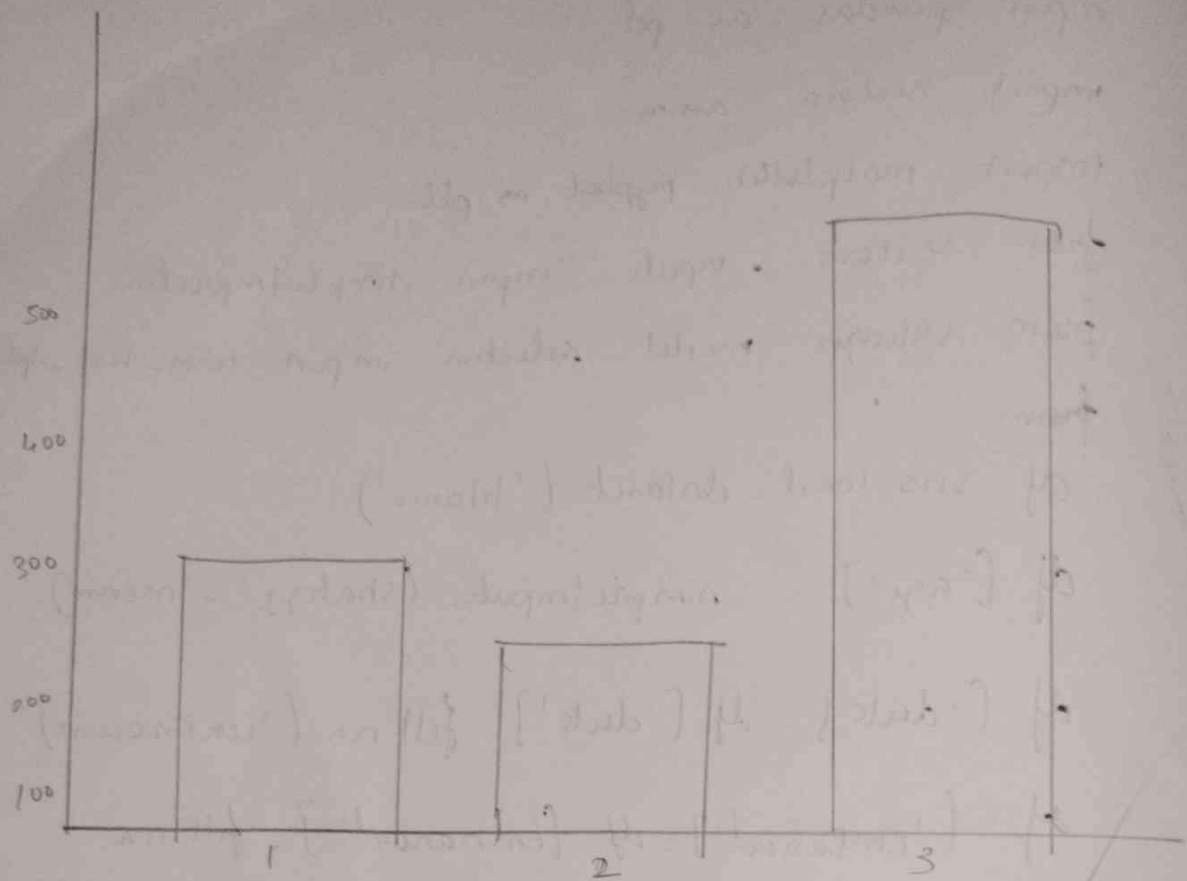Step 4: fill missing Cabin with "unknown" and embarked with mode

Step 5: visualize passanger class distribution with countplot.

Step 6: Identify top oldest.

Step 7: Split training and testing splits.

## program:

```python
import pandas as pd.
import seaborn as sns
import matplotlib.pyplot as plt.
from seaborn.inputs import simpleImputer.
from sklearn model_selection import train_test_split
from .

df = sns.load_dataset('titanic')

df['age'] = simpleImputer(strategy='mean')

df['duck'] = df['duck'].fillna('unknown')

df['embanad'] = df['embanad'].fillna
                    (df['embanad'].mode(1[0])

sns.countplot(x='pelans', data=df) plt
    file ('passanger dans distribution')

plt.show()


print("females who survied:"; df((df_sex ==
"female") & (df.survived==1)].index.list().
```

```
print ("male passengers who paid fare >
            100:", df [df.sex == 'male'] &
        (df.fare >100)].index.tolist()

print ("passengers embarced df 'c' and in
        class 2:", df [(df.embarked == 'c') &
        (df.pclass == 2)].
```

Result:

The program successfully identifies passengers with zero false and efficiently splits the datasets into 80% training and 20% testing sets.