

EX: 01
22/09/25

DATA PREPROCESSING AND CLEANING

Sim:

To load the titanic dataset and convert it into a DataFrame and explore it

Program code

```
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.preprocessing import LabelEncoder,
                                 StandardScaler

titanic = sns.load_dataset('titanic')
print(titanic.head())
print("Data info:")
print(titanic.info())
print("Missing values per column:")
print(titanic.isnull().sum())

titanic['age'] = titanic['age'].fillna(method
                                         = .ffill)

titanic['deck'] = titanic['deck'].cat.add_
                  categories(['unknown'])

titanic.Bin[missing - deck indices, 'deck']
           = 'unknown'
```

Output

	survived	pclass	age	sibsp	parch	fare	class
0	0	3	22.0	1	0	7.2500	Third
1	1	1	38.0	1	0	71.2813	First
2	1	3	26.0	0	0	2.9250	Third
3	0	1	35.0	1	0	53.1000	First
4		3	38.0	0	0	8.0800	Third

Missing values per column

survived : 0
pclass : 0
sex : 0
age : 177
sibsp : 0
parch : 0

Duplicates removed : 56

de = Label Encoder ()

titanic['Sex'] = encoder.fit_transform(titanic['Sex'])

titanic['Age'] = age

scaler = StandardScaler()

titanic['fare'] = scaled.fit_transform(titanic['fare'])
(titanic['fare'] = scaled)

sns.pairplot(titanic, vars = ['pclass', 'sex', 'embarked',
'age', 'sibsp'])

plt.suptitle('pairplot of selected features
y=1.02')

corr_features = ['pclass', 'age', 'sibsp',
'parch', 'fare']

plt.figure(figsize=(8, 6))

sns.heatmap(corr, annot=True, cmap='coolwarm',
font="-.2f")

plt.title('correlation Heatmap')

plt.show()

Thus the data preprocessing and cleaning
using titanic dataset has been executed successfully