

## Lab-1

Write python code for fall considering "housing.csv"

① To load .csv file into dataframe

```
import pandas as pd
```

```
filename = "/content/housing.csv"
```

```
df = pd.read_csv(filename)
```

② To display info of all columns

```
print("Dataset Info:")
```

```
print(df.info())
```

Dataset Info:

#	Column	Non-Null count	Dtype
0	Longitude	20640	float64
1	Latitude	20640	— " —
2	housing-median-age	20640	— " —
3	total-rooms	20640	— " —
4	total-bedrooms	20640	— " —
5	population	— " —	— " —
6	households	— " —	— " —
7	median-income	— " —	— " —
8	median-house-value	— " —	— " —
9	ocean-proximity	— " —	object

③ To display statistical info of all numerical columns

```
print("In Statistical summary of Numerical columns:")
```

```
print(df.describe())
```

④ To display count of unique labels for "ocean proximity" columns

```
print("\n Unique Value counts for 'Ocean Proximity':")
```

```
print(df["ocean-proximity"].value_counts())
```

⑤ To display which columns in a dataset have missing values  
count greater than zero

$\text{missing-values} = \text{df.isnull}().\text{sum}()$   
 $\text{missing-columns} = \text{missing-values}[\text{missing-values} > 0]$   
 $\text{print}(\text{"In Columns with missing values: "})$   
 $\text{print}(\text{missing-columns})$

o/p:-

Unique Value Counts for 'Ocean Proximity':

ocean-proximity	
<1H OCEAN	9136
INLAND	6551
NEAR OCEAN	2658
NEAR BAY	2290
ISLAND	5

Columns with Missing Values:

total-bedrooms: 207

Questions:-

1) Which columns in the datasets had missing values? How did you handle them?

Ans: Missing value columns:

Adult income dataset  $\rightarrow$  Age, Salary

Diabetes dataset  $\rightarrow$  Glucose, BMI

Handling approach:

Adult income dataset  $\rightarrow$  for age  $\rightarrow$  used median since its less sensitive to outliers  
 $\rightarrow$  for salary  $\rightarrow$  used mean as salaries typically follow normal distribution

Diabetes dataset  $\rightarrow$  Glucose  $\rightarrow$  used median since glucose levels may have outliers  
 $\rightarrow$  BMI  $\rightarrow$  used mean assuming normal distribution.



2) Which categorical columns did you identify in the dataset?  
How did you encode them?

Ans: Adult Income dataset:

Categorical columns: Gender  $\rightarrow$  original encoding

City  $\rightarrow$  One-hot encoding

Diabetes dataset:

Categorical columns: Gender  $\rightarrow$  original encoding

Outcome  $\rightarrow$  one-hot encoding

3) What is the difference b/w Min-Max scaling and standardization?  
When would you use one over another?

Ans: Min Max Scaling:

$$\rightarrow X' = \frac{X - X_{\min}}{X_{\max} - X_{\min}}$$

$\rightarrow$  scales values b/w 0 & 1

$\rightarrow$  sensitive to outliers

Standardization:

$$\rightarrow X' = \frac{X - \mu}{\sigma}$$

$\rightarrow$  Transforms data to have mean = 0 and variance = 1

$\rightarrow$  less affected by outliers

When data is not normally distributed and has known bounds, min-max scaling is used.

When data follows a normal distribution, standardization is used

Shreya B  
10/3/25