

Date = 3/03/25

Lab-1

Date \_\_\_\_\_  
Page \_\_\_\_\_

### Data Processing Lab

1.) Load CSV file

Import pandas as pd

⇒ `file = pd.read_csv("housing.csv")`

2.) Information of all col's

⇒ `file1 = file.info()`

3.) Statistical ~~column~~ info of all numerical columns

⇒ `file.describe()`

4.)

Count of unique labels, values for  
"ocean\_proximity"

⇒ `file['ocean_proximity'].value_counts()`

5.) To count attribute value having missing  
values count 70

⇒ `missing_value = df.isnull().sum()`

~~print~~ count = `missing_value[missing_value`  
70]

`print(count)`

3-3-25

After Execution,

Q.3.)

Q.1.) Which columns in the dataset had missing values? How did you handle them?

Ans → ~~Both the~~ No column in both the dataset had any missing values.

Ans →

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

For Adult Dataset,

Ans → Income column was categorical. Handled using ordinal encoder with  $< 50K$  &  $750K$

Work class column was handled using one-hot encoding.

→ Diabetes dataset:

~~No column had~~  
Gender and class were categorical.  
Gender → ordinal encoding [4 M, 4 F]

Class → one-hot encoding

Q.3) What is the difference between Min-Max scaling and standardization. When would you use one over the other?

Ans →

Min Max Scaling

$$X_i = \frac{X_i - \min}{\max - \min}$$

Range between  $[0, 1]$

ideal for distance-based Model but can distort data with extreme outliers.

Standardization:

$$X_i = \frac{X_i - \text{mean}}{\text{std}}$$

range =  $[-1, 1]$

Works well for normally distributed data but sensitive to outliers.

~~Q.4~~  
~~3/2/25~~