

10/3/25

ML-Lab-2



PAGE :

DATE : / /

housing dataset :-

1)

import pandas as pd

# Load the csv file into a DataFrame

housing\_df = pd.read\_csv('/content/housing(1).csv')  
print(housing\_df.columns)# display column information  
print("Column Information:")  
print(housing\_df.info())# Show statistical summary of numerical columns  
print("In Statistical Summary:")  
print(housing\_df.describe())# Count unique values in the 'Ocean Proximity' column  
print("In Unique values in 'Ocean proximity':")  
print(housing\_df['Ocean proximity'].value\_counts())# Identify columns with missing values  
print("In Columns with missing values:")  
print(housing\_df.isnull().sum())

o/p: Column Information:

&lt;class 'pandas.core.frame.DataFrame'&gt;

RangeIndex: 20640 entries, 0 to 20639

Data columns (total 10 columns):

| # | column          | non-Null count | Dtype   |
|---|-----------------|----------------|---------|
| 0 | longitude       | 20640 non-null | float64 |
| : | :               | :              | :       |
| 9 | Ocean Proximity | 20640 non-null | object  |





### Statistical Summary :

|       | Longitude    | Latitude | housing-median-age | median-house-value |
|-------|--------------|----------|--------------------|--------------------|
| count | 20640.000000 |          |                    | 20640.000000       |
| mean  |              |          |                    |                    |
| std   |              |          |                    |                    |
| min   |              |          |                    |                    |
| 25%   |              |          |                    |                    |
| 50%   |              |          |                    |                    |
| 75%   |              |          |                    |                    |
| Max   | 114.310000   |          |                    | 500001.000000      |

### unique values in 'Ocean Proximity':

#### Ocean Proximity

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

Name : count, dtype : int64

### columns with missing values :

|                    |     |
|--------------------|-----|
| longitude          | 0   |
| latitude           | 0   |
| housing-median-age | 0   |
| total-rooms        | 0   |
| total-bedrooms     | 207 |
| population         | 0   |
| households         | 0   |
| median-income      | 0   |
| median-house-value | 6   |
| ocean-proximity    | 0   |
| dtype : int64      | 0   |



## Diabetes and Adult income Datasets :

### Data Preprocessing

Write python code to implement the following data preprocessing techniques for diabetes & Adult income data sets

1. Data cleaning: Handling Missing values, Handling categorical data, Handling outliers
2. Data Transformation: Min-Max scales/ Normalization, standard scales

adult.csv & dataset of diabetes.csv

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

- Adult Income Dataset (adult.csv)  
No columns had missing values

- Diabetes dataset (Data set of Diabetes.csv)  
No columns had missing values

- Handling method : Since no missing values were found, no imputation was performed





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

\* Adult Income Dataset (adult.csv)

- categorical columns:

- workclass
- education
- marital-status
- occupation
- relationship
- race
- gender
- native-country
- income

Encoding method: one-hot encoding was applied using `pd.get_dummies(df, drop_first=True)`, which converted these categorical columns into numerical format.

- Diabetes Dataset (Dataset of diabetes.csv)

- categories columns:
  - gender
  - ss

• encoding method: one-hot encoding was applied `pd.get_dummies(df, drop_first=True)` to ~~the~~ transformation.



them into numeric format.

3) what is the difference b/w Min-Max scaling & standardization? when would you use one over the other.

|            |   |                                    |
|------------|---|------------------------------------|
| feature    | Min-Max Scaling                                 | Standardization                    |
| Definition | scales values b/w a fixed range (usually 0 & 1) | centers data around mean(0) with a |

|         |  |  |
|---------|--|--|
| Formula | $x_{scaled} = \frac{x_{max} - x_{min}}{x_{max} - x_{min}} \cdot (1 - 0)$ | $\left[ \frac{x - \mu}{SD} \right]$<br>Standardized = $\frac{x - \mu}{\sigma}$ |
|---------|--|--|

|        |   |                                  |
|--------|---|----------------------------------|
| effect | Preserves original data distribution but scales it within a limited range | changes data have zero variance. |
|--------|---|----------------------------------|

|              |  |  |
|--------------|--|--|
| when to use? | when preserving relationship & original range of data is important | when data follows a normal or distribution or contains outliers. |
|--------------|--|--|

4) when the one over the other

- Min Max scaling is useful when  
- you need data to be within a fixed range. The dataset does not contain extreme outliers.





- Standardization (Z-score normalization) is useful when

- The dataset has varying units and a gaussian (normal) distribution

- There are significant outliers, as standardization is less sensitive to them

Sc 1 p. 03