

Lab 1

To demonstrate various data pre-processing techniques for a given dataset.

Python code of "housing.csv" file.

```
import pandas as pd
```

```
# i) To load csv file into data frame
```

```
filename = "housing.csv"  
df = pd.read_csv(filename)
```

```
# ii) To display information of all columns
```

```
print("Columns")  
print(df.info())
```

```
# iii) Statistical information of all numerical  
print(df.describe())
```

```
# iv) Count of unique labels in 'ocean proximity'  
print(df['ocean proximity'].value_counts())
```

```
# v) Attributes or columns having missing  
value greater than zero
```

```
print(df.isnull().sum() [df.isnull().  
sum() > 0])
```


Python code to implement the following data preprocessing techniques for

①. Diabetes

Data preprocessing techniques:

1. Data cleaning: Handling missing values, Handling categorical data, handling outlier.
2. Data transformation: Min-max scaler / Normalization, Standard scaler

```
import pandas as pd
import numpy as np
from sklearn.preprocessing import MinMaxScaler, StandardScaler
from sklearn.preprocessing import LabelEncoder

data = pd.read_csv('/path/of/diabetes.csv')
```

1. Data cleaning

```
print("Missing values in each column:")
print(data.isnull().sum())

numerical_columns = ['bred', 'ex', 'HbA1c', 'cho', 'TG', 'HDL', 'LDL', 'VLDL', 'BMD']
for col in numerical_columns:
    data[col] = data[col].fillna((data[col].median()))

data.drop_duplicates(inplace=True)

print("Missing values after handling:")
print(data[numerical_columns].isnull().sum())
print("Shape after removing duplicates:")
print(data.shape)
```


Dataset of diabetes.csv

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

* In this dataset, there were no explicit missing values. This was checked when the dataframe was loaded to pandas `data.isnull().`

+ But missing values may exist in numerical columns

+ The numerical missing values were handled by filling it with median value of respective columns.

+ In categorical columns, no missing values were explicitly handled.

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

• The categorical columns were identified as Gender and Class, where Gender contains F & M values & Class contains N, P & Y

• These values were encoded using label encoder()

Gender: F \rightarrow 0, M \rightarrow 1
Class: N \rightarrow 0, P \rightarrow 1, Y \rightarrow 2

3. What is the difference between min max scaling & standardization? when would you use one of the other?

min-max:

Transform features to a fixed range [0,1].

$$\text{Formula: } x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

Preserves the shape of the dataset

Standardization:

Transform features to have a mean of 0 & standard deviation of 1.

Uses

min-max is used during uniform scaling [0,1]

Standardization is used for normally distributed data.

⇒ Adult income dataset.csv

1. Which column in the dataset has missing values? How did you verify them?

Columns with missing values:

workclass, occupation, native-country

2. Which categorical columns did you identify? How to encode them?

categorical columns: workclass, education, marital-status, occupation, relationships, race, gender, native-country, income.

use `labelEncoder()`

3. what is the difference between min max scaling and standardization?

→ min max scales data to a fixed range typically [0,1]

Used when need bounded values.

→ standardization transforms data to mean 0 & standard deviation 1.