

10 / 03 / 25

ML-LAB 13 - 01

⇒ following. CSU

import panelas as pd

```
housing_df = pd.read_csv("content/housing.csv")
```

print ("Information of all columns: ")

```
print(housing_df.info())
```

print("In statistical information of numerical columns")

```
print("\\nCount of unique labels for 'Ocean Proximity'"  
      column : ")
```

column .  
in housing-df. columns:  
ocean-pronimisry

if 'ocean\_proximity' in df['ocean\_proximity'].value\_counts():  
 print('housing - ', df['ocean\_proximity'].value\_counts())

else:  
    print(f"Column '{ocean\_proximity}' not found  
        in the DataFrame. Available columns: {housing.columns}")

~~in the following columns, to list + 14" of columns, to list + 14" of columns with missing values;~~

~~if columns. to list + ( ) y' 19~~

~~df = columns, to list + 19')~~  
~~grouping columns with missing values;~~

~~pointing in columns with misspellings~~

~~Prin~~  
missing-values = housing if .isnull().sum()>0  
Dnint(missing-values > 75)

## → Handling Missing data,

- Data processing techniques for  
Diabetes and Adult income datasets.

```
import pandas as pd
```

```
import numpy as np
```

```
from sklearn.impute import SimpleImputer
```

```
from sklearn.preprocessing import MinMaxScaler,  
StandardScaler
```

```
from supy import stats
```

```
diabetes_file_path = '/content/dataset_of/Diabetes.csv'
```

```
diabetes_df = pd.read_csv(diabetes_file_path)
```

```
numerical_cols = diabetes_df.select_dtypes(include=[  
    np.number]).columns
```

```
categorical_cols = diabetes_df.select_dtypes(exclude=[  
    np.number]).columns
```

```
imputer_numerical = SimpleImputer(strategy='mean')
```

```
diabetes_df[numerical_cols] = imputer_numerical.fit_  
transform(diabetes_df[numerical_cols])
```

```
imputer_categorical = SimpleImputer(strategy='most  
frequent')
```

```
diabetes_df[categorical_cols] = imputer_categorical.  
fit_transform(diabetes_df[categorical_cols])
```

=> observations on Data processing techniques  
for the Diabetes & Adult income datasets

1. which columns in the Data set had missing values?  
How did you handle them?

### Diabetes dataset:

→ columns with missing value: Urea, Cr, HbA1c,  
chol, TG, HDL, LDL, VLDL, BMI.

→ For numeric columns for missing ~~dat~~ values,  
we applied mean imputation. The missing values  
were replaced with the mean of each respective  
column.

→ For categorical columns missing value were  
handled using most frequent value (mode)  
imputation.

### Adult income Dataset:

→ ~~columns~~: In this dataset, missing values  
were represented by ? In the og capaset

→ we replaced all occurrences of ? with NaN  
to make the missing values identifiable

→ After this, we applied mode imputation  
for categorical columns to replace missing values

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

Diabetes:

- columns: gender, race, class
- For encoding categorical variables, one-hot encoding or label encoding could be applied to convert into binary columns.

Adult Income Dataset:

- columns: workclass, education, marital-status, occupation, relationship, race, sex, native-country
- one-hot encoding was applied to convert the categorical columns into multiple binary columns. For instance, workclass would be converted into ~~one~~ multiple binary column
- Using pd.get\_dummies() with drop first. This allowed us to avoid multicollinearity by removing the first column from each categorical variable set.

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

→ Min-Max Scaling:

① normalization transforms data to a specific range usually b/w 0 & 1

- using the formula

$$x_{\text{scaled}} = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

when to use: Min-Max scaling is preferred when the dataset contains features with varying units or scales & you want to scale them to a fixed range.

Standardization:

② Z-Score normalization transforms data to have a mean of 0 & a standard deviation of 1. It is calculated using the formula:

$$x_{\text{scaled}} = \frac{x - \mu}{\sigma}$$

$\mu$  - mean &  $\sigma$  - standard deviation.

→ Standardisation is typically used when the data is assumed to follow a Gaussian distribution, especially when working with algorithms that assume normally distributed data regression, logistic regression.

Schall  
10/3/25