

5/9/25

Lab - 1

i) `import pandas as pd`
`df = pd.read_csv('housing.csv')`

ii) `df.info()`

iii) `df.describe()`

iv) `df["Ocean Proximity"].value_counts()`

v) `mis_val = df.isnull().sum()`
`val = mis_val[mis_val > 0]`
`print(val)`

Diabetes

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

`df = pd.read_csv('/content/dataset of Diabetes - csv')`

`print(df.head())`

#missing values

`print(df.isnull().sum())`

#Impute

`nc = df.select_dtypes(include=['float64', 'int64']).`
`columns`

`imputer = SimpleImputer(strategy='mean')`
`df[num_columns] = imputer.fit_transform(df[num_columns])`
`cat = df.select_dtypes(include=['object']).columns`
`imputer_cat = SimpleImputer(strategy='most_frequent')`
`df[cat] = imputer_cat.fit_transform(df[cat].values)`

Handling categorical data

`label_encoder = LabelEncoder()`
`df['gender'] = label_encoder.fit_transform(df['gender'])`
`df['class'] = label_encoder.fit_transform(df['class'])`

Handling outliers

`Q1 = df[num_columns].quantile(0.25)`
`Q3 = df[num_columns].quantile(0.75)`
`IQR = Q3 - Q1`

`df_clean = df[~((df[num_columns] < (Q1 - 1.5 * IQR)) |`
`(df[num_columns] > (Q3 + 1.5 * IQR))).any`
`(axis=1)]`

Data Transformation

Apply Min Max or standard scalar

`scaler_choice = 'minmax'`

`if scaler_choice == 'minmax':`
`scaler = MinMaxScaler()`

else:

`scaler = StandardScaler()`

`df_scaled = pd.DataFrame(scaler.fit_transform(df_clean`
`[num_columns]), columns=num_columns)`

```
df_final = pd.concat([df_clean[cat_columns],  
df_scaled], axis=1)
```

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

① A No columns had missing values.
We handled them by writing the code:
`df.isnull().sum()`

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

① A Gender and Class
We encoded them by using label encoding

③ What is the difference between MinMax Scaling & Standardization? When would you use one over the other?

① A Min-Max Scaling :-
Transforms values to fixed range (usually 0 to 1) and useful if there is a fixed range for features bounded range, and also used when no outliers are present

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

Standardization (z-score normalization) :-

Transforms data to have zero mean and unit variance.
outliers exist.

Gaussian distribution, such as in regression generally used in the normal distribution

Formula:
$$X' = \frac{X - \mu}{\sigma}$$

22/11/25