

Data Pre-processing Techniques

(i) `import pandas as pd.`  
`df = pd.read_csv("housing.csv")`  
`print("Data loaded into DataFrame")`

(ii) `print("In Information of all columns:")`  
`print(df.info())`

(iii) `print("In Statistical information:")`  
`print(df.describe())`

(iv) `print("In Count of unique labels for Ocean Proximity column:")`  
`print(df['Ocean Proximity'].value_counts())`

(v) `print("In Columns with missing values:")`  
`missing_value = df.isnull().sum()`  
`cm = missing_value[missing_value > 0]`  
`print(cm)`

① Diabetes Dataset: Column like Glucosol, Blood Pressure and BMI had missing values handled by imputing mean or median

Adult income: Column like occupation and native country had missing values handled by mode or `dropna()`

② Dialysis Dataset: The Outcome column is categorical, encoded using Label Encoding.

Adult Income Dataset: Columns like workclass, education were categorical encoded using one-hot encoding.

③ Min-Max Scaling scales the data to a fixed range (0 to 1) and is used when data is bounded the model is sensitive to the scale.

Standardization scales the data to have a mean of 0 and a standard deviation of 1 and is used when the data is normally distributed or when the model assumes a normal distribution.

Yash 23.03.24