

import pandas as pd.

- i) df = pd.read_csv("housing.csv")
- ii) df.info()
- iii) df.describe()
- iv) df["Ocean Proximity"].value_counts()
- v) miss_val = df.isnull().sum()
val = miss_val[miss_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("Dataset of Diabetes.csv")

print(df, headers)

missing values

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

impute,

nc = df.select_dtypes(include=[float64, 'int64'])
imputer = SimpleImputer(strategy='mean').fit(nc)

$df[nc] = imputer$ $fit_transform(df[nc])$
 $cat_c = df.select_dtype(include = ['object'])$ $element$
 $imputer_cat = SimpleImputer(strategy = 'most_frequent')$
 $df[cat_c] = imputer_cat.fit_transform(df[cat_c])$

Handling categorical data.

$label_encoder = LabelEncoder()$
 $df['gender'] = label_encoder.fit_transform(df['gender'])$
 $df['CLASS'] = label_encoder.fit_transform(df['CLASS'])$

$Q1 = df[num_column].quantile(0.25)$
 $Q3 = df[num_column].quantile(0.75)$

$IQR = Q3 - Q1$

$df_clean = df[nc] (df[nc] \leq (Q1 + 1.5 * IQR) /$
 $(df[nc] > (Q3 + 1.5 * IQR)))$
 $any (axis=1)$

$Scaler_chose = 'minmax'$

$Scaler = MinMaxScaler()$

$df_scaled = pd.DataFrame(Scaler.fit_transform(df_clean[nc]), (axis=nc))$

$df_final = pd.concat([df_clean[cat_c], df_scaled], axis=1)$

1. No column.

2. Gender and CLASS

3. Min-Max-Scaling

if fixed range for features, bounded range, no outliers.

Standard Scaler,

\bullet outliers exist
 \bullet normal distributed data