Demonstrating Various data Preprocessing techniques for

1) Data cleaning: Handling missing values, Handling
Categorical data, handling outliners

2) Data Transformations: min-max Scalar /Normalisation,
Standard Scalar.

1) Housing dataset

(1) Column information    (ii) Statistical Summary of Columns

(3)    Unique values in Ocean Proximity column

(4)    Columns with missing values.

Code    import Pandas as pd

```
housing_dt = pd. read_csv ('/path/to/housing.csv')
Print ("Column Information:")
Print (housing_dt.info())

Print ("\n Statistical Summary:")
Print (housing_dt. describe())

Print ("\n unique values in 'Ocean Proximity':")
Print (housing_dt ['ocean Proximity']. value_counts())
Print ("\n Columns With missing values:")
Print (housing_dt. isnull().sum())
```

output

Index ([ 'longitude', 'latitude', 'housing-median-age',
'total-rooms', 'total-bedrooms', 'population', 'households',
'median-income', 'median-house-value', 'ocean-proximity' ],

Column Information:

Data Columns (total 10 Columns):

| Column | non-null Count | Dtype |
|---|---|---|
| 0 longitude | 20640 non-null | float 64 |
| 1 latitude | 20640 non-null | float 64 |
| 2 housing-media-age | 2063 non-null | float 64. |
| : | | |
| 9 Ocean-Proximity | 20640 non-null | Object. |

Statistical Summary

| | longitude | latitude | housingmedian-age | totalrooms |
|---|---|---|---|---|
| Count | 20640.0000 | 20640.0000 | 20640.0000 | 20640.0000 |
| mean | -119.5697 | 35.6318 | 28.6394 | 2635.54 |
| std | 2.0035 | 2.1359 | 12.5843 | 2181.53 |
| min | -124.35000 | 3.25400 | 18.0000 | 1447.543 |
| max | -121.8000 | 33.9300 | 28.100 | 1447.430 |
| 25% | -118.542 | 34.2600 | 29.010 | 2127.250 |
| 50% | -118.0100 | 37.7100 | 37.200 | 3148.100 |
| 75% | | | | |
| max | -114.3100 | 41.9500 | 52.153 | 39320.010 |

unique values in 'Ocean Proximity'

ocean-Proximity

1H ocean      9136

INland        6551

Near ocean    2658

Near Bay      2290

Island        5

↳ dtype : int

columns with missing Values

longitude         0

Latitude          0

housing           0
median-age

Ocean Proximity   0

dtype : int

2) Python code to implement the following data Preprocessing techniques for Diabetes & adult income dataset.

```python
import Pandas as Pd
import numpg as np
import matplottlib pyplot as Plt

df = pd-read_csv ('/content/dataset of diabetg.csv')
  df.head(100)

df = pd.Ocod-csv("/content/adult data1.csv")
  df.head(10)

import numpy as np
# introduce some missing values for demonstration.

df.loc[5, 'AGE'] = np.nan
df.loc[9, 'BMI'] = np.nan
df.head(10)
```

output

| | ID | No_Pation | Gender | Age | Urea | Cr | BMI | Cloy |
|---|---|---|---|---|---|---|---|---|
| | | | | 50.0 | 4-7 | 46 | NaN | N |
| 0 | 50 2 | 17975 | F | | | | | |
| 1 | 412 | 47975 | M | 43.0 | 4.5 | 62 | 0.17913 | N |
| 2 | 327 | 87656 | F | 32.0 | 4.2 | 72 | 0.17916 | N |
| 3 | 634 | 67643 | M | NaN | 4.7 | 2.3 | 0.69 | N |

i) Which columns in dataset had missing values? How did you handle them?

- Adult Income dataset : (adult-csv)

  No columns had missing values

- Diabetes dataset (Dataset of diabetes. csv)

  No Columns has missing values

  - Handling method Since no missing value were found, no computation has performed

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

adult Income Dataset (adult.csv)

⟶ Categorical columns

  ⟶ workclass

  → Education

  → occupation

  → race

  → gender

  → income

  → relationship

Encoding method; one-hot encoding has applied using pd.get_dummies (db, drop_first = True)

→ diabetes Dataset

    Categorical column

        → Gender
        → class

encoding: one-hot encoding applied
pd.get_dummies (db, drop_first = True)

③ difference b/w min-max Scaling & Standardisation

| feature | min max scaling | Standardisation |
|---|---|---|
| definition | Sales values b/w in fixed range (usually 0 & 1) | center data or mean (0) with SD 1 |
| formula | $X_{scaled} = \dfrac{X_{max} - X_{min} \times}{X_{in}}$ | $X_{stand}$ $\quad d = \sigma x - \mu$ |
| Effect | Preserving original data distribution but scales it within a lin range | Changes data distribution to have zero mean & unit variance |

④ when would you use one over other

minmax — he need data to be within a fixed range
The dataset doesn't contain extreme outliers

Standardisation is used when

- The dataset has anything varying units &
  a Gaussian distribution (
  normal
- There are significant outliers,
  as Standardisation is less
  sensitive to them.

Sen
10-03-2024