

## LAB-1 :

write python code, consider filename as "housing.csv"

i) To load .csv file into the dataframe

```
import pandas as pd
filename = "housing.csv"
df = pd.read_csv(filename)
```

ii) To display information of all columns

```
print(df.info())
```

iii) To display statistical information of all numerical

```
print(df.describe())
```

iv) To display the count of unique labels for "Ocean Proximity" column

```
if "Ocean Proximity" in df.columns:
    print(df["Ocean Proximity"].value_counts())
```

else: print("Not found")

v) To display which attributes in dataset have missing values count greater than zero

```
m-values = df.isnull().sum()
m-columns = m-values[m-values > 0]
```

```
if not m-columns.empty:
    print(m-columns)
```

else: print("No missing values found")

## Output :

ii) 0	longitude	3000	non-null	float64
1	latitude	3000	non-null	float64
2	housing-median-age	3000	non-null	float64
3	total-rooms	3000	non-null	float64
4	total-bedrooms	3000	non-null	float64
5	population	3000	non-null	float64
6	house-holds	3000	non-null	float64
7	median-income	3000	non-null	float64
8	median-house-value	3000	non-null	float64

iii)	2000
count	-119.5
mean	1.99
std	-124.1
min	-121.
25%	-118
50%	-118
75%	-118
max	-118

iv) Ocean  
v) No n

## Question

Q1) which  
you  
son: nu  
Adult  
Diabe  
Hond  
Adult

Diab

Q2)

son

Car

D

Co



	Longitude	Latitude
count	2000	3000
mean	-119.5	35.6
std	1.99	2.12
min	-124.18	30.5
max	-121.81	33.9
25%	-121.81	34.2
50%	-118.98	37.6
75%	-118.02	41.9
max	-114.49	41.9

- iv) 'Ocean Proximity' column not found in the dataset  
v) No missing values found in the dataset

### Questions:

Q1) Which columns in the datasets had missing values? How did you handle them?

Ans: Missing value columns:

Adult Income dataset → Age, Salary

Diabetes dataset → Glucose, BMI

Handling approach:

Adult Income dataset → for age → used median since it's less sensitive to outliers  
→ for salary → used mean as salaries typically follow normal distribution

Diabetes dataset → Glucose → used median since glucose levels may have outliers  
→ BMI → used mean assuming normal distribution.

Q2) Which categorical columns did you identify in the dataset? How did you encode them?

Ans: Adult Income Dataset:

Categorical columns: Gender → original encoding  
City → One-Hot Encoding

Diabetes Dataset:

Categorical columns: Gender → original encoding  
Outcome → One-Hot Encoding.



Q3) What is the difference between Min-Max scaling and standardization? When would you use one over the other?

Soln:

Min-Max Scaling:

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

- scales values between 0 & 1.
- Sensitive to outliers

Standardization

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

- Transform data to have mean = 0 and variance = 1
- less affected by outliers.

When data is not normally distributed and has known bounds, min-max scaling is used.  
 When data follows a normal distribution, standardization is used.