

① Initializing values directly into DataFrame - five rows of data with column headings as USN, Name, marks.

data = { "USN": ["234", "235", "236", "237", "238"],
"Name": ["Asha", "Cathy", "Rosa", "Amy", "Gina"],
"Marks": [85, 90, 78, 92, 88]}

	USN	Name	Mark
0	234	Asha	85
1	235	Carly	90
2	236	Rosa	78
3	237	Amy	92
4	238	Gina	88

```
print(df-diabetes)
```

[illegible]

③ Import datasets from a specific .csv file

```
import pandas as pd
```

```
df_csv = pd.read_csv ("/sample-sales-data.csv")
```

```
print (df_csv.head ())
```

slp		Product	Quantity	Price	Sales	Region
0		Laptop	5	1000	5000	North
1		Mouse	15	20	300	West
2		Keyboard	10	50	500	East
3		Monitor	8	200	1600	South
4		Laptop	12	950	11400	North

④ Downloading datasets from existing dataset repositories like Kaggle -

```
import pandas as pd
```

```
df_kaggle = pd.read_csv ("/Dataset of Diabetes.csv")
```

```
print (df_kaggle.head (1))
```

slp	ID	No-Patients	Gender	AGE	Urea	Cr
0	502	17975	F	50	4.7	46
1	735	34221	M	26	4.5	62
2	420	47975	F	50	4.7	

BMI CLASS

0	24.0	N
1	23.0	N
2	24.0	N
3	24.0	N
4	21.0	N

To - 20

Stock Market Data Analysis

1. ["HDFCBANK.NS", "ICICIBANK.NS", "KOTAKBANK.NS"]
2. start date and end date.
3. Plot the closing price and daily return for all three banks.

```
import yfinance as yf
import pandas as pd
import matplotlib.pyplot as plt
```

```
tickers = ["HDFCBANK.NS", "ICICIBANK.NS", "KOTAKBANK.NS"]
```

```
data = yf.download(tickers, start = "2024-01-01", end = "2024-12-30",
                    group-by = 'tickers')
```

```
print ("First 5 rows of the dataset")
```

```
print (data.head(1))
```

```
print ("\n Shape of dataset ")
```

```
print (data.shape)
```

```
print ("\n Column names")
```

```
print (data.columns)
```

```
hdfc_data = data['HDFCBANK.NS']
```

```
print ("\n Summary statistics for HDFC Industries")
```

```
print (hdfc_data.describe(1))
```

```
hdfc_data['Daily Return'] = hdfc_data['Close'].pct-change(1)
```

```
plt.figure(figsize = (12, 6))
```

```
plt.subplot(2, 1, 1)
```

```
hdfc_data['Daily Return'] = hdfc_data.plot (title = " HDFC Industries - Daily  
Returns", color = 'orange')
```

```
plt.tight_layout(1)
```

```
plt.show()
```

① To load .csv file into dataframe

```
import pandas as pd
df = pd.read_csv('/content/housing.csv')
print(df.head())
print(df.columns)
print(df.info)
```

② To display information of all columns

```
print(df.columns)
```

	population	households	median-income	median-house-value	ocean-proximity
0	322.0	126.0	322.0		
1	2401.0	1138.0	8.304		
2	496.0		8.754		
3	558.0	219.0	5.6471		
4	565.0	219.0	1.8462		

③ To display statistical information of all numerical

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

④ To display the count of unique labels for "Ocean Proximity" column

```
print(df.value_counts('ocean-proximity'))
```

⑤ To display which attributes in a dataset have missing values.

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


0 NEAR RAY
 1 NEAR RAY
 2 INLAND
 ...
 0.5 - 0.3 208

For Diabetes dataset, apply data preprocessing technique -

```

import pandas as pd
import numpy as np
import os
import sklearn.preprocessing
  
```

import LabelEncoder, MinMaxScaler, StandardScaler

```
diabetes - df = pd.read_csv("/content / Dataset - q - Diabetes.csv")
```

Check for missing value

```
print("Missing values in Diabetes dataset : \n", diabetes - df.isnull().sum())
```

Handle missing values with mean

```
numeric - cols = diabetes - df.select_dtypes(include = np.number).columns
```

```
diabetes - df[numeric - cols] = diabetes - df[numeric - cols].fillna(diabetes - df[numeric - cols].mean())
```

To fill missing categorical values with mode

```
for col in diabetes - df.select_dtypes(include = ['object']).columns:
```

```
diabetes - df[col].fillna(diabetes - df[col].mode()[0], inplace = True)
```

Identify categorical values and encode them.

```
categorical - cols = ['Gender', 'CLASS']
```

```
label - enc = LabelEncoder()
```

```
for col in categorical - cols:
```

```
diabetes - df[col] = label - enc.fit_transform(diabetes - df[col])
```

Apply min max scaling

```
min - max - scaler = MinMaxScaler()
```

```
diabetes - df - scaled = pd.DataFrame(min - max - scaler.fit_transform(diabetes - df),
```

```
columns = diabetes - df.columns)
```

Apply standardization

```
std_scaler = StandardScaler()  
diabetes_df_standardized = pd.DataFrame(std_scaler.fit_transform(diabetes_df),  
                                         column = diabetes_df.columns)  
  
i # Save the preprocessed datasets  
d diabetes_df_scaled.to_csv("/content/diabetes-preprocessed.csv", index = False)  
p os.makedirs("content", exist_ok = True)  
p diabetes_df_standardized.to_csv("/content/diabetes-standardized.csv", index = False)  
p print("Diabetes dataset preprocessing completed. Preprocessed datasets saved.")
```

Q1. Which columns in the dataset had missing values? How did you handle them?

- Numeric columns such as glucose levels, blood pressure, insulin.
- Categorical columns such as Gender and CLASS

To handle numeric data, mean of the respective column was used.
categorical data, mode

Q2. Which categorical columns did you identify in the dataset? How did you encode?
Gender and CLASS are the two categorical columns.

Q3. Encoding method - label encoding
Male \rightarrow 0
Female \rightarrow 1

Q4. What is the difference between min-max scaling and standardization?

Min Max Scaling - transforms data to fixed range $[0, 1]$
$$x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

Standardization - transforms data so that it has a mean of 0 and standard deviation of 1
$$x' = \frac{x - \mu}{\sigma}$$