



Department of Computer Engineering

Name of the Student: DIKSHA JADHAV Roll Number: C206

SAP ID: 60004240043 Class: C3 Division: C3

Batch: DV2

Subject: Data Visualization

EXPERIMENT NO: 04

AIM: Data Collection and Wrangling in Python. (CO4)

IMPLEMENTATION/SOURCE CODE WITH OUTPUT:

```
import pandas as pd
import numpy as np
import requests
from bs4 import BeautifulSoup
```

(A) Loading a CSV File

```
df = pd.read_csv("dataset_new.csv")
df.head()
```

	Country	Cases	Deaths	Recovered	Date
0	France	282587	32532	173231	2024-10-20
1	China	257357	5504	490323	2024-10-02
2	Indonesia	665376	13709	462288	2024-04-21
3	UK	176639	12831	55408	2024-10-25
4	Turkey	427238	43009	806343	2024-05-21

```
df.describe()
```

(B) Fetching Live Data from an API

```
api_key = "7a9ad05412f2af9ab39b67fa9388c0"
city = "Mumbai"
url = f"http://api.openweathermap.org/data/2.5/weather?q={city}&appid={api_key}"
response = requests.get(url)
```



(B) Fetching Live Data from an API

```
api_key = "7a9ad05412f1f2af9ab39b67fa9388c0"
city = "Mumbai"
url = f"http://api.openweathermap.org/data/2.5/weather?q={city}&appid={api_key}"

response = requests.get(url)
data = response.json()
print("data.items(), sep='\n'")
```

[6] ✓ 0.1s

Python

```
... ('coord', {'lon': 72.8479, 'lat': 19.0144})
('weather', [{'id': 800, 'main': 'Clear', 'description': 'clear sky', 'icon': '01d'}])
('base', 'stations')
('main', {'temp': 301.85, 'feels_like': 305.07, 'temp_min': 301.85, 'temp_max': 301.85, 'pressure': 1008, 'humidity': 69, 'sea_level': 1008, 'grnd_level': 100000})
('visibility', 10000)
('wind', {'speed': 4.89, 'deg': 284, 'gust': 5.35})
('clouds', {'all': 5})
('dt', 1744459814)
('sys', {'country': 'IN', 'sunrise': 1744419215, 'sunset': 1744464289})
('timezone', 19800)
('id', 1275339)
('name', 'Mumbai')
('cod', 200)
```

File Edit Selection View Go Run ... DV_Honors

index.html dv_4jpyb X

C: > Users > sham > Downloads > dv_4jpyb > url = "https://en.wikipedia.org/wiki/COVID-19_pandemic_by_country_and_territory"

Generate + Code + Markdown | Run All Restart Clear All Outputs | Jupyter Variables Outline ... Python 3.12.5

```
url = "https://en.wikipedia.org/wiki/List_of_countries_by_population_(United_Nations)"

response = requests.get(url)
soup = BeautifulSoup(response.text, "html.parser")

table = soup.find("table", {"class": "wikitable"})

df = pd.read_html(str(table))
print(df)
```

[] Python

```
... [
0          Country or territory  Population (1 July 2022) \
1          World 8021407192
2          India 1425423212
3          China[a] 1425179569
4          United States 341534046
5          Indonesia 278830529
..          ...
233          Montserrat (United Kingdom) 4453
234          Falkland Islands (United Kingdom) 3490
235          Tokelau (New Zealand) 2290
236          Niue (New Zealand) 1821
237          Vatican City[x] 505

Population (1 July 2023) Change (%) UN continental region[1] \
0          8091734930 +0.88% -
1          1438069596 +0.89% Asia
2          1425179569 +0.89% Asia
```

Launchpad 0 7 Live Share Ln 1, Col 14 CRLF Cell 18 of 18 Go Live



File Edit Selection View Go Run ... DV_Honors

index.html dv_4.ipynb X

C:\Users> soham > Downloads > dv_4.ipynb > url = "https://en.wikipedia.org/wiki/COVID-19_pandemic_by_country_and_territory"

Generate + Code + Markdown | Run All | Restart | Clear All Outputs | Jupyter Variables | Outline ... Python 3.12.5

(A) Handling Missing Values

```
df.isnull().sum() # Check missing values
df.fillna(0, inplace=True) # Replace missing values with 0
```

** Task: Identify column names and print the first five rows.**

(B) Filtering and Renaming Columns

```
print(df)
```

	Country or territory	Population (1 July 2022) \
0	World	8021407192
1	India	1425423212
2	China[a]	1425179569
3	United States	341534046
4	Indonesia	278830529
...
233	Montserrat (United Kingdom)	4453
234	Falkland Islands (United Kingdom)	3490
235	Tokelau (New Zealand)	2290

Launchpad 0 7 Live Share Ln 1, Col 14 CRLF Cell 18 of 18 Go Live

File Edit Selection View Go Run ... DV_Honors

index.html dv_4.ipynb X

C:\Users> soham > Downloads > dv_4.ipynb > url = "https://en.wikipedia.org/wiki/COVID-19_pandemic_by_country_and_territory"

Generate + Code + Markdown | Run All | Restart | Clear All Outputs | Jupyter Variables | Outline ... Python 3.12.5

```
df=df[0]
df_filtered = df[df['Population (1 July 2022)'] > 10000000]
print(df_filtered)
```

	Country or territory	Population (1 July 2022)	Population (1 July 2023) \
0	World	8021407192	8091734930
1	India	1425423212	1438069596
2	China[a]	1425179569	1422584933
3	United States	341534046	343477335
4	Indonesia	278830529	281190067
...
91	Portugal[j]	10417073	10430738
92	Tajikistan	10182222	10389799
93	Papua New Guinea	10203169	10389635
94	Azerbaijan	10295304	10318207
95	Greece	10412480	10242908

	Change (%)	UN continental region[1]	UN statistical subregion[1]
0	+0.88%	-	-
1	+0.89%	Asia	Southern Asia
2	-0.18%	Asia	Eastern Asia
3	+0.57%	Americas	Northern America
4	+0.85%	Asia	South-eastern Asia
...
91	+0.13%	Europe	Southern Europe

Launchpad 0 7 Live Share Ln 1, Col 14 CRLF Cell 18 of 18 Go Live



```
url = "https://en.wikipedia.org/wiki/COVID-19_pandemic_by_country_and_territory"

response = requests.get(url)

soup = BeautifulSoup(response.text, "html.parser")
tables = soup.find_all("table", {"class": "wikitable"})

if len(tables) >= 12:
    df = pd.read_html(str(tables[11]))[0]
    print(df.head())
else:
    print("Less than 7 tables found on the page.")
```

	Unnamed: 0	Country	Deaths / million	Deaths	Cases
0	NaN	World[a]	887	7087718	777368165
1	NaN	Peru	6601	220994	4528708
2	NaN	Bulgaria	5679	38765	1338332
3	NaN	North Macedonia	5428	9990	352060
4	NaN	Bosnia and Herzegovina	5118	16404	404142

<ipython-input-12-7ff08630b480>:20: FutureWarning: Passing literal html to 'read_html' is deprecated and will be removed in a future version. To read f
df = pd.read_html(str(tables[11]))[0] # Index 6 corresponds to the 7th table

CONCLUSION:

Data Collection offers the raw material—timely and accurate data—from various sources. Data Wrangling molds this raw material into a clean, consistent, and usable form for analysis or machine learning. These are the key steps in ensuring the quality and integrity of any data-driven project. Clean, properly prepared data frequently determines the success of analysis or model performance.