

# Big Data Analytics

**Lab Practical and date** – Practical 1, Monday 19<sup>th</sup> July 2020

**Name and Roll Number-** Het Shah, 17BIT103

**Practical Objective-** Introduction to R/Python programming Language. Write program to read the data from any online website, excel file and CSV file.

## Steps Involved-

We perform data scraping using python by reading data from different file format such as excel, csv and also performing web scraping by extracting data from a website

## Background

Python is an interpreted, high-level, general-purpose programming language. Created by Guido van Rossum and first released in 1991, Python's design philosophy emphasizes code readability with its notable use of significant whitespace.

Data scraping, also known as web scraping is the process of importing information from a website into a spreadsheet or local file saved on your computer.

## File Formats Used

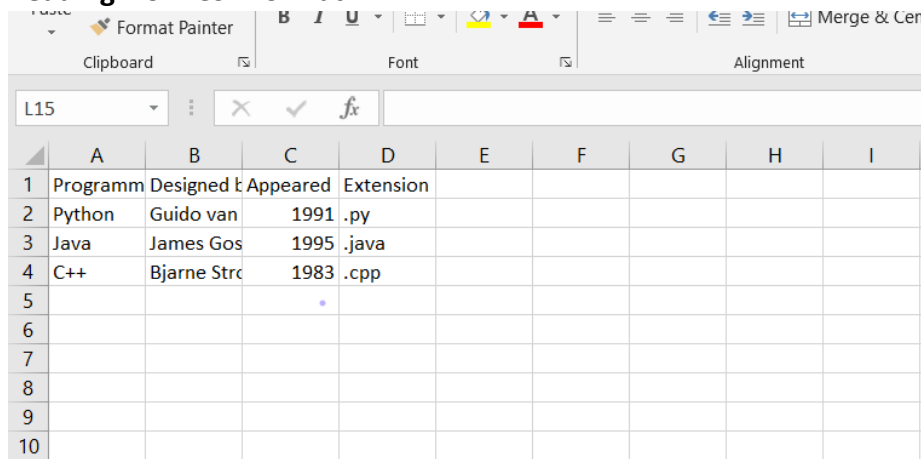
- 1) **Excel-** excel is spreadsheet application developed and maintained by MircroSoft. Excel organizes data in rows and columns format and they intersect at a space called cell.
- 2) **CSV**(comma separated values)-A minimal format compatible with many spreadsheet applications. Rows of data are represented by lines in text file, with columnar breaks delimited by a single character, usually a comma.

## Libraries used-

- 1) **Pandas**-pandas is a software library written for the Python programming language for data manipulation and analysis
- 2) **Requests**- it is a python HTTP library, allows to send HTTP requests easily
- 3) **Csv**-The csv module implements classes to read and write tabular data in CSV format.
- 4) **Lxml.html**- lxml provides a very simple and powerful API for parsing XML and HTML. It supports one-step parsing as well as step-by-step parsing using an event-driven API

## Outputs

### 1) Reading from CSV format



	A	B	C	D	E	F	G	H	I
1	Program	Designed by	Appeared	Extension					
2	Python	Guido van	1991	.py					
3	Java	James Gos	1995	.java					
4	C++	Bjarne Strou	1983	.cpp					
5									
6									
7									
8									
9									
10									

This is the CSV file from which we want to extract the data

```
In [29]: runfile('E:/Desktop/prac 1.py', wdir='E:/Desktop')
csv output
['ï»¿Programming language', 'Designed by', 'Appeared', 'Extension']
['Python', 'Guido van Rossum', '1991', '.py']
['Java', 'James Gosling', '1995', '.java']
['C++', 'Bjarne Stroustrup', '1983', '.cpp']
```

This is the python program reading the csv file and displaying it to the console

## 2) Reading from Excel Sheet

	A	B	C	D	E	F	G	H	I	J	K	L	M
1	name	roll no											
2	het	17bit103											
3													
4													
5													
6													
7													
8													
9													
10													
11													
12													

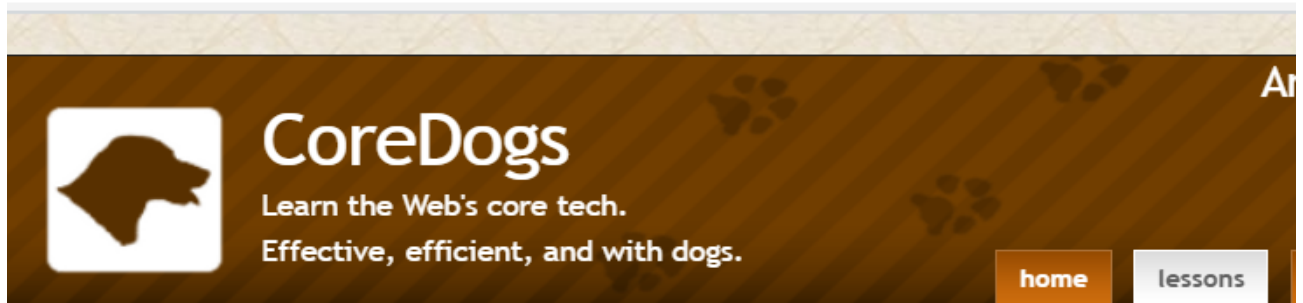
This is the Excel file from which we want to extract the data

```
excel output
name    roll no
0  het  17bit103
```

This is the python program reading the csv file and displaying it to the console

### 3) Reading data from an online Website

coredogs.com/lesson/web-page-tables.html



[Home](#) » [ClientCore](#)

## A page with tables

See more about: [CSS](#) [jQuery](#) [Tables](#)

[← Exercises: A Web page with images](#)

[Basic tables →](#)

You know how to create pages, add text, interact with the user, and add images. We're on our way to making complete Web sites.

This lesson looks at the HTML tags for tables. Tables make it easy to present data in rows and columns. Here's an example:

Dog	Size	Breed
Kieran	Large	Lab
CC	Medium	Sheltie
Renata	Medium	Coonhound/lab mix

This is the website <http://coredogs.com/lesson/web-page-tables.html> from which we extracted the table

```
table parsed
  Dog      Size      Breed
0 Kieran   Large      Lab
1      CC   Medium    Sheltie
2 Renata   Medium Coonhound/lab mix
```

The python code was able to parse the table from the website and be able to display it

#### 4) Running Queries on the Parsed Data

```
table parsed
  Dog      Size      Breed
0  Kieran   Large      Lab
1    CC    Medium    Sheltie
2  Renata   Medium  Coonhound/lab mix

selecting dog breed with medium size
1    CC
2  Renata
Name: Dog, dtype: object

selecting dog breed with name size more than 5
0  Kieran
2  Renata
Name: Dog, dtype: object

selecting distinct dog breed sizes
['Large' 'Medium']

insert a row - Sunny, small, pug
  Dog      Size      Breed
0  Kieran   Large      Lab
1    CC    Medium    Sheltie
2  Renata   Medium  Coonhound/lab mix
3  Sunny    Small      pug

delete a row - with dog CC
  Dog      Size      Breed
0  Kieran   Large      Lab
2  Renata   Medium  Coonhound/lab mix
3  Sunny    Small      pug
```

In this step, we mainly performed CRUD(create, read, update and delete) and select operations using the pandas framework and we displayed the results as shown.

#### Conclusion

In this practical, we learned about data parsing from various file formats such as CSV and excel and also how to parse data from websites. Later, we used the parsed data and performed CRUD operations on them.

```
1  # -*- coding: utf-8 -*-
2  """
3  Created on Mon Jul 20 14:19:50 2020
4
5  @author: HETSHAH
6  """
7
8
9  import csv
10 import pandas as pd
11 import requests
12 import lxml.html as lh
13
14 #read from csv
15 print('csv output')
16 with open('data.csv','rt') as f:
17     data = csv.reader(f)
18     for row in data:
19         print(row)
20
21 print()
22 print()
23 print('excel output')
24 #read from excel
25 df = pd.read_excel (r'book1.xlsx')
26 print (df)
27
28
29
30 print()
31 print()
32
33 #read online
34
35 url='http://coredogs.com/lesson/web-page-tables.html'
36 page = requests.get(url)
37 doc = lh.fromstring(page.content)
38 tr_elements = doc.xpath('//tr')
39
40
41
42 tr_elements = doc.xpath('//tr')
43 col=[]
44 i=0
45
```

```

46 #getting the table headings
47 print('heading names')
48 for t in tr_elements[0]:
49     i+=1
50     name=t.text_content()
51     print('%d:"%s"'%(i,name))
52     col.append((name,[]))
53
54
55 #storing the data
56 for j in range(1,len(tr_elements)):
57     T=tr_elements[j]
58     i=0
59     for t in T.iterchildren():
60         data=t.text_content()
61         col[i][1].append(data)
62         i+=1
63
64 print()
65 print('table parsed')
66 Dict={title:column for (title,column) in col}
67 df=pd.DataFrame(Dict)
68 print(df.head())
69
70 print()
71 print('selecting dog breed with medium size')
72 print(df[df.Size == 'Medium'].Dog)
73
74 print()
75 print('selecting dog breed with name size more than 5')
76 print(df[df.Dog.str.len() >=5].Dog)
77
78
79 print()
80 print('selecting distinct dog breed sizes')
81 print(df.Size.unique())
82
83 print()
84 print('insert a row - Sunny, small, pug')
85 df = df.append({'Dog': 'Sunny', 'Size': 'Small', 'Breed': 'pug'}, ignore_index=True)
86 print(df)
87
88
89 print()

```

```
90 print( 'delete a row - with dog CC')
91 df= df[df.Dog != 'CC']
92 print(df)
93
```