

Proyecto 1: Plataforma de seguimiento de datos COVID-19 para Colombia

Laura Camila Blanco Gómez
Escuela de Ciencias Exactas e Ingeniería
Universidad Sergio Arboleda - Bogotá, Colombia
laura.blanco01@correo.usa.edu.co

Santiago Cáceres Linares
Escuela de Ciencias Exactas e Ingeniería
Universidad Sergio Arboleda - Bogotá, Colombia
santiago.caceres01@correo.usa.edu.co

Andrés Camilo López Ramírez
Escuela de Ciencias Exactas e Ingeniería
Universidad Sergio Arboleda - Bogotá, Colombia
andres.lopez01@correo.usa.edu.co

Resumen

En el presente proyecto se busca extraer diversos datos sobre el virus COVID-19 de una página web, guardarlos en una base de datos y representar de forma selectiva mediante gráficos de tortas, gráficos de dos dimensiones y diagramas de barras de los datos que se han extraídos de la página web.

Palabras clave:

COVID19, Colombia, Web Scraping, Python, html.

1. Marco teórico

1.1. COVID19

El COVID-19 es la enfermedad infecciosa causada por el coronavirus que se ha descubierto más recientemente y afecta a muchos países de todo el mundo. Razon por la cual fue clasificado como Pandemia por la Organización Mundial de la Salud (OMS).

La mayoría de las personas (alrededor del 80 %) se recuperan de la enfermedad sin necesidad de tratamiento hospitalario. Alrededor de 1 de cada 5 personas que contraen la COVID-19 acaba presentando un cuadro grave y experimenta dificultades para respirar.

Actualmente en Colombia se cuentan con 818.203 casos confirmados, de los cuales 68.308 son casos activos, 722.536 personas recuperadas y 25.641 fallecidos.

1.2. Web scraping

Es una técnica que permite extraer datos e información de una web. Este tutorial es una guía de inicio al web scraping con Python, utilizando para ello la librería BeautifulSoup.

Beautiful Soup es una librería Python que permite extraer información de contenido en formato HTML o XML.

Se deben seguir los siguientes pasos para un buen funcionamiento del web scraping:

- Identifica los elementos de la página de los que extraer la información
- Descarga el contenido de la página
- Crear la «sopa»
- Busca los elementos en la «sopa» y obtén la información deseada

2. Resultados

Para realizar la extracción de datos se usa la técnica **Web Scraping**, con la cual se extrajo información desde **Wikipedia**, ya que cuenta con tablas que muestran los valores de contagios, muertes, supervivencia y casos activos del COVID-19, estos se actualizan constantemente. Para esto se hizo uso de las librerías de **BeautifulSoup** y **requests**.

```

1 from bs4 import BeautifulSoup
2 import requests
3 import pandas as pd
4 import numpy as np
5 import matplotlib.pyplot as plt
6
7 url = 'https://es.wikipedia.org/wiki/Pandemia_de_enfermedad_por_coronavirus_de_2020_en_Colombia'
8 ruta = "/Users/andre/DatosCovid.CSV"
9
10 page_response = requests.get(url, timeout=5)
11 page_content = BeautifulSoup(page_response.content, "html.parser")
12 ciudadtxt = list()
13 numerosArray = list()
14
15 tablaCol = page_content.find(
16     "table", attrs={"class": "wikitable sortable collizq"})
17 ciudadRows = tablaCol.find_all("tr")
18
19 for row in ciudadRows:
20     numerostd = row.find_all("td")[1:]
21     for i in numerostd:
22         numerosArray.append(sinEspacios(i.text))
23     ciudadLink = row.find_all("a")
24     for i in ciudadLink:
25         ciudadtxt.append(i.text)

```

Listing 1: Código del web scraping

Después de hacer la extracción de datos se hace una "división" de esta información que se encuentra en **numerosArray**, esto se hizo por medio de varias listas como se ve en el siguiente código, para después guardar toda esta información de un **DataFrame**

```

1
2 for i in range(0, 330, 10):
3     total.append(numerosArray[i])
4 total.reverse()
5
6 for i in range(1, 330, 10):
7     cxMhad.append(numerosArray[i])
8 cxMhad.reverse()
9
10 for i in range(2, 330, 10):
11     muertesT.append(numerosArray[i])
12 muertesT.reverse()
13
14 for i in range(3, 330, 10):
15     porjeMT.append(numerosArray[i])
16 porjeMT.reverse()
17
18 for i in range(4, 330, 10):
19     falleMhab.append(numerosArray[i])
20 falleMhab.reverse()
21
22 for i in range(5, 330, 10):
23     recuT.append(numerosArray[i])
24 recuT.reverse()
25
26 for i in range(6, 330, 10):
27     porjeRT.append(numerosArray[i])
28 porjeRT.reverse()
29
30 for i in range(7, 330, 10):
31     CasActT.append(numerosArray[i])
32 CasActT.reverse()
33
34 for i in range(8, 330, 10):
35     porjeACT.append(numerosArray[i])
36 porjeACT.reverse()
37
38 for i in range(9, 330, 10):

```



```

29 try:
30     cursor.execute(dataInsert)
31     db.commit()
32     print("Inserci n de datos exitosa!")
33 except:
34     print("Error en conexi n :c")
35     db.rollback()
36
37 db.close()

```

Listing 3: Código base de datos

The screenshot shows the MySQL Workbench interface. The 'Navigator' pane on the left shows the 'covid_database' schema with tables, views, stored procedures, and functions. The 'Query Editor' pane shows a query that selects all data from the 'datoscovid' table. The 'Result Grid' pane displays the query results as a table with 10 columns: Nombre, Casos_Confirmados_Totales, Casos_Confirmados_por_millon_de_habitantes, Muertes_Totales, Porcentaje_de_muertes, Fallecido_por_millon_de_habitantes, Recuperados_Totales, Porcentaje_de_recuperados, Casos, and another column. The data is sorted by 'Nombre' and shows information for various departments in Colombia.

Nombre	Casos_Confirmados_Totales	Casos_Confirmados_por_millon_de_habitantes	Muertes_Totales	Porcentaje_de_muertes	Fallecido_por_millon_de_habitantes	Recuperados_Totales	Porcentaje_de_recuperados	Casos
Vichada	527	4665.5	4	0.8%	35.4	462	87.7%	61
Vaupés	765	17109.5	10	1.3%	223.7	695	90.8%	60
Guaviare	798	9208.7	16	2%	184.6	667	83.6%	113
Guaviare	869	17161.7	12	1.4%	237	694	79.9%	163
San Andrés y Providencia	1424	22357.6	14	1%	219.8	901	63.3%	508
Arauca	1718	5839.4	47	2.7%	159.8	1518	88.4%	151
Casanare	2354	5409.1	48	2%	110.3	1804	76.6%	498
Amazonas	2742	34700.1	117	4.3%	1480.6	2615	95.4%	8
Quindío	3582	6449.4	98	2.7%	176.4	2734	76.3%	749
Putumayo	3771	10500.5	170	4.5%	473.4	3304	87.6%	291
Chocó	3996	7335.3	155	3.9%	284.5	3754	93.9%	74
Caldas	5865	5758.7	129	2.2%	126.7	4829	82.3%	883
Boyacá	7232	5819.4	147	2%	118.3	5971	82.6%	1094
La Guajira	8105	8392.7	302	3.7%	312.7	7268	89.7%	520
Coque	8476	20646.9	297	3.5%	723.5	7540	89%	618
Cauca	9404	6303.2	262	2.8%	175.6	8011	85.2%	1100
Risaralda	11139	11590.4	244	2.2%	253.9	9632	86.5%	1249
Huila	11637	10365.9	361	3.1%	321.6	9750	83.8%	1500
Tolima	12312	9188.1	333	2.7%	248.5	10879	88.4%	1058
Sucre	14072	14824.3	577	4.1%	607.8	13044	92.7%	420
Magdalena	15324	10738.4	829	5.4%	580.9	13607	88.8%	825
Norte de Santander	15690	9683.3	891	5.7%	549.9	13687	87.2%	1086
Meta	16111	15149.7	405	2.5%	380.8	14575	90.5%	1107
Nariño	18321	11256.5	682	3.7%	419	16401	89.5%	1196
Cesar	19793	15279.6	590	3%	455.5	16883	85.3%	2290
Córdoba	23946	13092.8	1539	6.4%	841.5	21270	88.8%	1018

Figura 2: Base de datos

Para representar la información extraída en graficas usamos el DataFrame para seleccionar los datos que queremos graficar y se guardan en variables *x*. y *y*. Además se hace uso de la librería **matplotlib** para definir el tipo de grafica, en este caso es *plt.pie()*.

```

1 # Grafica Barras
2 x = df['Nombre'].iloc[28:33]
3 y = df['Casos Confirmados Totales'].iloc[28:33]
4 plt.bar(x, y)
5 plt.xlabel('Departamenro')
6 plt.ylabel('Casos Confirmados Totales')
7 plt.grid()
8 plt.show()

```

Listing 4: Código gráfica de barras

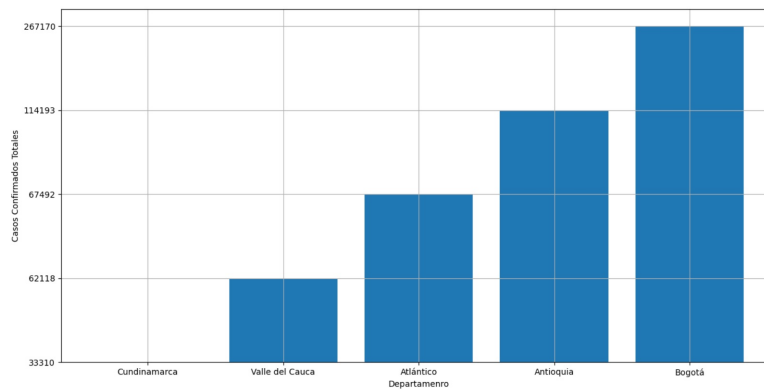


Figura 3: Grafico de barras

En el caso de la gráfica de torta se hace uso del **plt.pie()** en el cual se encontrara *autopct = "%0.1f%%"* el cual servirá para convertir los valores en porcentajes, aclarando que realizara esta conversión únicamente con los datos que se grafiquen.

```
1 #Grafico Torta
2 x = df['Muertes Totales'].iloc[28:33]
3 y = df['Nombre'].iloc[28:33]
4 plt.title('Muertes Totales')
5 plt.pie(x, labels= y, autopct = "%0.1f%%" )
6 plt.show()
```

Listing 5: Código gráfica de torta

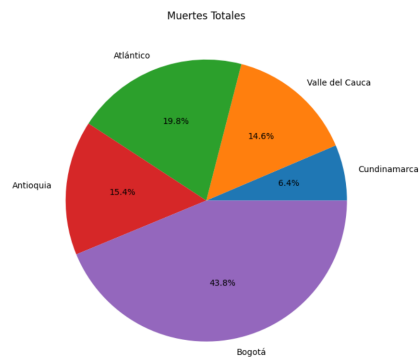


Figura 4: Grafico de torta

Y por ultimo en el grafico 2D se usa **plt.plot()** y en nuestro caso hacemos uso de **plt.grid** para que se puedan apreciar las intersecciones de los valores en la recta.

```
1 #Grafico 2D
2 x = df['Muertes Totales'].iloc[28:33]
3 y = df['Nombre'].iloc[28:33]
4 plt.plot(x,y)
5 plt.xlabel('Casos Confirmados Totales')
6 plt.ylabel('Departamento')
7 plt.grid()
8 plt.show()
```

Listing 6: Código gráfica 2D

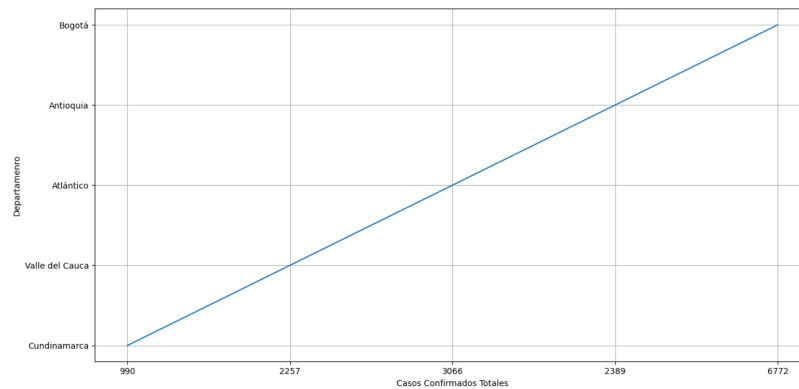


Figura 5: Grafico 2D

3. Conclusiones

- Conociendo la técnica de web scraping se facilitó para la extracción desde la página web a Python.
- Se utilizan las librerías de Python Pandas, BeautifulSoup y request, ya que nos permiten crear DataFrames y archivos CSV, buscar información y buscar información del html respectivamente.
- Se hace necesario a la hora de extraer los datos buscar páginas confiables que mantengan sus datos actualizados, ya que con esto se tendrá un panorama más realista de lo que está pasando.
- Gracias al uso de los **DataFrame** se pudo almacenar de una manera fácil la información extraída para después poder llevarla a un archivo CSV con el cual se llenó la base de datos posteriormente.

4. Link repositorio

https://github.com/ACLXRD/SA2020_G02_CovidScraper_6I3I11

Referencias

- [1] OMS, *Preguntas y respuestas sobre la enfermedad por coronavirus (COVID-19)*. Organización mundial de la salud, 2020.
- [2] J. J. L. Gómez, *Web scraping con Python. Extraer datos de una web. Guía de inicio de BeautifulSoup*. J2LOGO, 2018.
- [3] I. nacional de salud, *COVID-19 en Colombia*. MinSalud, 2020.
- [4] W. L. enciclopedia libre, *Pandemia de enfermedad por coronavirus de 2020 en Colombia*. Wikipedia, 2020.