Tarea 3: Web Scraping

Sergio Lagos Vergara

Contents

1	Introducción														
2	Desarrollo														
		CheckName													
	2.2	ITB	5												
	2.3	ITB ajustado a la inflación	5												
	2.4	Peliculas	6												
	2.5	PeliculasMayorIngreso	9												
	2.6	PeliculasMenorIngreso	10												
3	Con	nclusión	11												
4	Bibliografía														

1 Introducción

En esta actividad se nos propuso extraer información de la página web **box office mojo**, buscando patrones con expresiones regulares en el html de la página, en donde dado un actor, debemos retornar el Ingreso Total Bruto de su carrera (ITB), el ITB ajustado a la inflación, las 5 últimas películas, las 10 películas con mayor ingreso bruto y las 10 películas con menor ingreso bruto.

A modo de ejemplo, este es el modelo de página del actor Robert Downey Jr. de donde se obtendrá la **Información**.

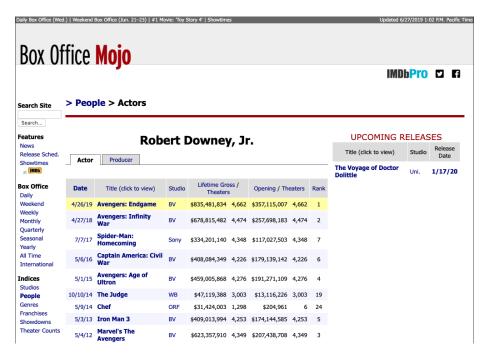


Figura 1: Robert Downey Jr. web page

A primera vista, ya podemos observar que en la tabla principal las películas están ordenadas desde la más reciente hasta la primera cinta del actor. De esta misma manera, si clickeamos en **Lifetime Gross**, nos redireccionará a una página con la misma tabla, pero las películas estan ordenadas descendentemente respecto al ingreso que tuvo y si presionamos denuevo **Lifetime Gross**, obtendremos la tabla con las películas ordenadas de manera ascendente respecto al ingreso. Esto es muy importante, ya que utilicé estos links para obtener información solicitada. Además, si seguimos bajando por la página principal del actor, nos podremos encontrar con el ITB y el ITB ajustado a la inflación.

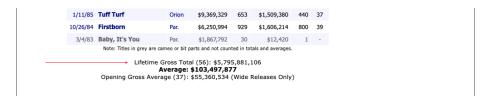


Figura 2: Lifetime Gross

54	Game 6	KMG	\$178,400	\$129,664	3/10/06
55	Friends and Lovers	Lions	\$59,500	\$33,543	4/16/99
56	Hugo Pool	N.Arts	\$26,200	\$13,330	12/12/97
Note: T	itles in grey are cameo or bit parts and not counted in t	otals and a	verages.		
	ted Total: \$6,821,775,600 age: \$121,817,400				

Figura 3: Adjusted Lifetime Gross

2 Desarrollo

Tomando en cuenta los requerimientos, decidí escribir un programa en Python 3, utilizando principalmente la librería **re**, para las expresiones regulares y **requests**, para obtener el HTML de la página.

Comenzando, dividí mi programa en 6 funciones que me ayudaron a recopilar información, pero para mejorar la explicación del código, en la figura 4 se muestra lo que se pide por consola:

```
if __name__ == '__main__':
    actor = input("Nombre del actor: ")
    datos = requests.get('https://www.boxofficemojo.com/people/chart/?id='+actor.replace(" ","")+'.htm')
    CheckName(actor,datos)
```

Figura 4: Main code

En donde 'datos' es el request.get() a la página asociada al actor.

2.1 CheckName

Esta función recibe como parámetros al actor y el requests.get() a la url asociada. Principalmente lo que hago acá es asegurarme de que el actor ingresado sea válido.

```
def CheckName(actor,datos):
    if datos.url == "https://www.boxofficemojo.com/people/":
        actor = input("Oops! Al parecer ingresaste mal el nombre, intenta denuevo: ")
        print ("\n")
        datos = requests.get('https://www.boxofficemojo.com/people/chart/?id='+actor.replace(" ","")+'.htm')
        CheckName(actor,datos)
    else:
        ITBI(datos)
        ITBInflacion(datos)
        Peliculas(actor)
        PeliculasMayorIngreso(actor)
        PeliculasMenorIngreso(actor)
        return
```

Figura 5: CheckName code

Para esto lograr esto, noté que al ingresar un actor que no se encuentra en la página, el requests.get() redirecciona a https://www.boxofficemojo.com/people. Por lo tanto, basta con hacer una comparacion entre datos.url y la página redireccionada para saber si el nombre del actor está bien o mal. Al pasar esta comparación, se llaman a las siguientes funciones.

2.2 ITB

Para obtener el ITB del actor, hice una expresión regular que haga match con la parte del HTML que se encuentra el LifetimeGross como se muestra en la Figura 2.

```
def ITB(datos):
    ITB = re.findal\(r"Lifetime\sGross\sTotal\s\\((.*?)\\):\s(\$([0-9]{1,3},([0-9]{3},)*[0-9]+)\(\.[0-9][0-9])?)", datos.text)
print ("\n")
    return
```

Figura 6: ITB code

Aquí ITB, es una lista con una lista de este estilo:

```
[('9', '$2,329,782,367', '2,329,782,367', '782,', '')]
```

Figura 7: ITB

En donde el segundo elemento de la lista de la lista, corresponde al Ingreso Total Bruto.

2.3 ITB ajustado a la inflación

Esta función es prácticamente igual a la anterior, solo que se busca el match con el Adjusted Lifetime Gross como se ve en la Figura 3.

```
def ITBInflacion(datos):
    AdjustedTotal = re.findall(r"Adjusted\sTotal\:\s(\$([0-9]{1,3},([0-9]{3},)*[0-9]{3}|[0-9]+)(\.[0-9][0-9])?)", datos.text)
    print ("Ingreso total bruto ajustado a la inflacion: ", AdjustedTotal[0][i][)
    print ("\n")
```

Figura 8: ITB ajustado code

2.4 Peliculas

En esta sección de código, utilizo la url para poder imprimir por pantalla las 5 últimas películas de algún actor.

```
def Peliculas(datos):
   Totalpeliculas = re.findall(r"<bs(.*?)</b>", datos.text)
   PeliculasGrises = re.findall(r"<font color=\"\#666666\">(.*?)</font>", datos.text)
   cont = 0

print ("Ultimas 5 peliculas: ")

for peliculas in range (1,6):
   match = re.match(r"<font color=", Totalpeliculas[peliculas])
   match2 = re.match(r"Average", Totalpeliculas[peliculas])

if match2:
   print("\n")
   return

if match:
   if PeliculasGrises[cont-1] != '(Cameo)' and PeliculasGrises[cont-1] != '(Voice)' and PeliculasGrises[cont-1] != '(Narrator)':
        print (PeliculasGrises[cont+1])
        cont = cont + 2
        else:
        print(PeliculasGrises[cont+1])
        else:
        print(Totalpeliculas[peliculas])

print ("\n")

return</pre>
```

Figura 9: Peliculas code

En este nivel se me presentaron diversos 'problemas'. Resulta que Total peliculas es una lista que contiene todos los strings encerrados entre y<\b >. Esto nos retorna una lista de estas características:

```
['Date', '<font color="#666666">The Man Who Killed Don Quixote</font', 'Absolutely Anything', 'Average: $20,169', for Ticket Price Inflation', 'font color="#666666">The Man Who Killed Don Quixote</font>', '$389,100', 'Absolute g', '$20,300', 'Average: $20,300', 'Morldwide', <font color="#666666">The Man Who Killed Don Quixote</font>', '$3
```

Figura 10: Totalpeliculas

Traducido a la página web, sería:



Figura 11: Terry Gilliam

Después de varios test cases, llegué a la conclusión de que el siempre el primer ítem de Totalpeliculas es 'Date', de que seguido de la última película del actor viene el string 'average' y que las películas marcadas en gris vienen con el tag <font.... Si solo imprimía por pantalla Totalpeliculas desde el segundo ítem, también se imprimían los tags de las películas en gris. Como solución, creé PelículasGrises, que es una lista con todas los string grises dentro del HTML de la página. Lamentablemente no solo guarda las películas grises, sino también, los estudios, el string 'voice', cuando el actor hace la voz en la película, 'narrator' cuando la narra y 'cameo' cuando hace un cameo como podemos ver a continuación:



Samuel L. Jackson

Date	Title (click to view)	Studio	Lifetime Gross / Theaters		Opening / The	Rank	
6/14/19	Shaft (2019)	WB (NL)	\$17,153,729	2,952	\$8,901,419	2,952	49
4/26/19	Avengers: Endgame	BV	\$835,481,834	4,662	\$357,115,007	4,662	-
3/8/19	Captain Marvel	BV	\$426,776,678	4,310	\$153,433,423	4,310	4
1/18/19	Glass	Uni.	\$111,035,005	3,844	\$40,328,920	3,841	17
6/15/18	Incredibles 2 (Voice)	BV	\$608,581,744	4,410	\$182,687,905	4,410	2
4/27/18	Avengers: Infinity War (Cameo)	BV	\$678,815,482	4,474	\$257,698,183	4,474	-
8/18/17	The Hitman's Bodyguard	LG/S	\$75,468,583	3,377	\$21,384,504	3,377	25
3/10/17	Kong: Skull Island	WB	\$168,052,812	3,846	\$61,025,472	3,846	10
2/3/17	I am Not Your Negro (Voice)	Magn.	\$7,123,919	320	\$686,378	43	-
1/20/17	xXx: The Return of Xander Cage	Par.	\$44,898,413	3,651	\$20,130,142	3,651	34

Figura 12: Samuel L Jackson Page



Figura 13: PeliculasGrises Samuel

Dadas las condiciones, tambien creé las variables match, que retorna true si la casilla de Totalpeliculas es un string con tag de color gris y match2, que retorna true si ya se acabaron las películas del actor en caso de que este tenga menos de 5. Match2 retorna true cuando el siguiente elemento en la lista Totalpeliculas hace match con la expresión regular "Average". Como se pudo ver anteriormente, la lista Peliculas Grises no solo almacena los string de las películas. Por lo tanto, para siempre imprimir una, se verifica que el elemento anterior en la lista no sea ni "(cameo)","(narrator)" o "(voice)" y si lo es, imprime el siguiente elemento de Peliculas Grises y no el actual.

2.5 PeliculasMayorIngreso

Para obtener esta información, como se menciona al principio, se usa la url que ordena de forma ascendente las películas respecto a sus ingresos. De esta forma, la manera de operar es casi la misma a la función Peliculas, solo que con 10 elementos. El problema aquí es que Totalpeliculas cambia un poco su estructura y hay que iterarla de diferente manera como se ve en la Figura 15.



Ben Affleck

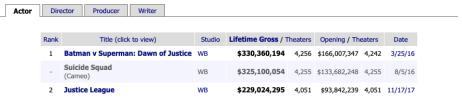


Figura 14: Ben Affleck Page



Figura 15: PeliculasGrises Ben

En síntesis, el código es casi el mismo a Peliculas, solo que hay que iterar TotalPeliculas intercaladamente desde el segundo ítem, ya que el primero siempre es 'Lifetime Gross'.

```
def PeliculasMayorIngress(actor):
    datos = requests.get('https://www.boxofficemojo.com/people/chart/?view=Actor&id='+actor.replace(" ","")+'.htm&sort=gross&order=DESC&p=.htm')
Totalpeliculas = re.findall(r"<bo(.*?)</box>, datos.text)
Peliculas&rises = re.findall(r"<font color=\"\\#666666\"\cdot\"\.*?\</font>", datos.text)
cont = 0

print ("10 peliculas con mayor ingreso: ")

for peliculas in range(0,20):
    if peliculas % 2 = 1:
    match = re.match(r"\dverage", Totalpeliculas[peliculas])
    match2 = re.match(r"\dverage", Totalpeliculas[peliculas])

if match2:
    print("\n")
    return

if match:
    if Peliculas&rises[cont-1] != '(Cameo)' and Peliculas&rises[cont-1] != '(Voice)' and Peliculas&rises[cont-1] != '(Narrator)':
        print (Peliculas&rises[cont+1])
        cont = cont + 2
    else:
        print(Peliculas&rises[cont+1])
else:
    print (Totalpeliculas[peliculas])

print ("\n")
    return
```

Figura 16: PeliculasMayorIngreso code

2.6 PeliculasMenorIngreso

En esta última función, es lo mismo que Peliculas Mayor Ingreso pero con la url de las películas ordenadas de forma ascendente respecto a su ingreso total.

Figura 17: PeliculasMenorIngreso code

3 Conclusión

En conclusión, este es un programa que dado un actor, usa expresiones regulares para identificar patrones dentro de la página web **box office mojo** y así obtener diversa información respecto al actor requerido. Quizás se puede mejorar la parte de como se imprimen las películas cuando el programa se enfrenta al problema de tener texto de otro color, pero al fin y al cabo el problema se resuelve de manera satisfactoria.

4 Bibliografía

References

- [1] BOX OFFICE MOJO
- [2] RE LIBRARY
- [3] REQUESTS LIBRARY