

UNIVERSIDAD NACIONAL AUTÓNOMA DE MÉXICO

FACULTAD DE CIENCIAS



Manejo de Datos



Web scraper

García Ramírez Alison

Navarro Castillo Rafael Alberto Carlos

Olivares Herrera Luis Félix

Rodríguez Hernández Fernanda

Web Scraper

Como su nombre sugiere, se le conoce como *web scraping* (*raspado web*) a una técnica de extracción de datos de páginas web mediante un software. Resulta muy útil para obtener grandes cantidades de información de manera automática.

Veremos cómo elaboramos el nuestro, bajo la temática de tiendas de ropa.

Para empezar, debemos instalar las paqueterías necesarias. En nuestra implementación, utilizamos:

- *Pandas*
- *Re*
- *Time*
- *Datetime*
- *Selenium*

Además, encontramos útil trabajar en *Jupyter Notebook*, pues en otro entorno no reconocía la última paquetería.

También, descargamos *chromedriver*. Dependiendo de la versión de *Google Chrome* que se tenga, es la versión que se debe seleccionar en el siguiente link: [ChromeDriver - WebDriver for Chrome - Downloads \(chromium.org\)](https://chromedriver.chromium.org/downloads)

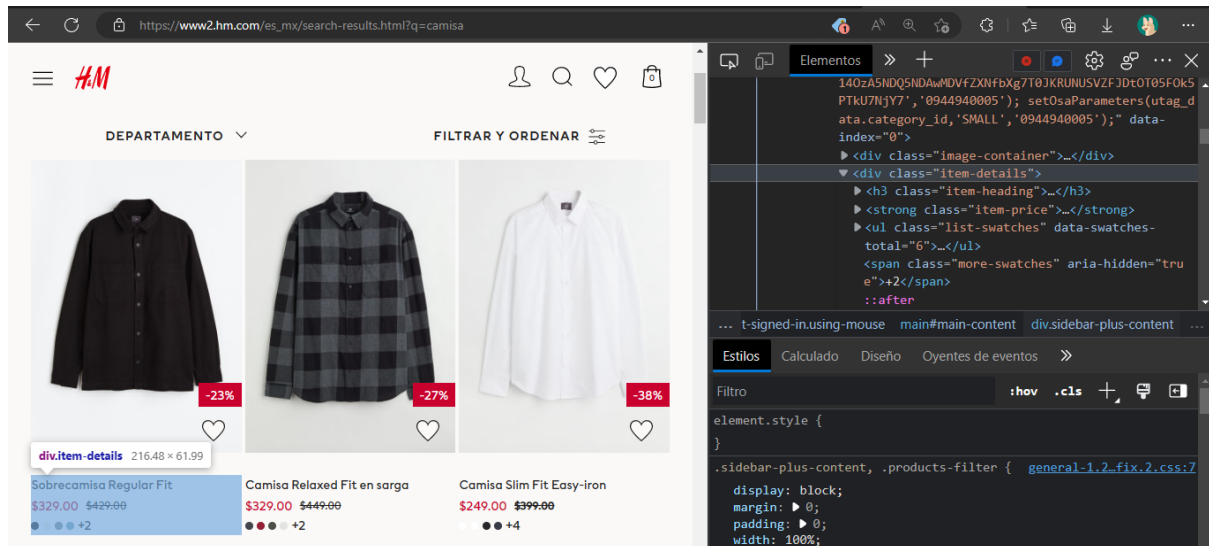
Decidimos que trabajaríamos con las marcas de *HyM*, *Bershka* y *Shein*, por lo que realizamos una función para cada sitio. Aún así, si se desea revisar diferentes páginas, el procedimiento sería muy similar.

Ya en el código, lo primero que hacemos es importar las librerías y asignar un *path* con la ruta donde se encuentra el *chromedriver* en nuestro dispositivo:

```
In [200]: 1 import pandas as pd
          2 import time
          3 from datetime import datetime
          4 import re
          5 from datetime import date
          6 from selenium import webdriver
          7 from selenium.webdriver.common.keys import Keys
          8 from selenium.webdriver.common.by import By
          9 from selenium.webdriver.chrome.options import Options
         10
         11
         12 #Ruta correspondiente|
         13 path = 'C:/Users/aragea/Documents/Rafa/MDD/PF/chromedriver.exe'
         14
         15
```

Empecemos describiendo cómo se realizó la función para *HyM*.

Lo primero que hicimos fue realizar una búsqueda en el sitio oficial de la marca, para reconocer cómo luce el link en la pantalla de la lista de productos. Posteriormente, inspeccionamos la página para reconocer el nombre de la clase donde se almacena la información más importante, de la siguiente manera:



Como se aprecia en el link, el ejemplo fue tras buscar 'camisa'. Vemos además que la clase que nos interesa se llama *item-details*.

Revisamos a detalle cómo luce el formato de la información de cada producto:

```
In [203]: 1 #Abrimos la página
          2 driver=webdriver.Chrome(path)
          3 url1 = 'https://www2.hm.com/es_mx/search-results.html?q='+ 'camisa'
          4 driver.get(url1)
          5 time.sleep(12)
          6
          7 opciones = driver.find_elements_by_class_name("item-details")
          8
          9

In [205]: 1 opciones[0].text.split('\n')

Out[205]: ['Sobrecamisa Regular Fit',
           '$329.00 $429.00',
           'Negro',
           'Crema',
           'Azul grisáceo',
           'Azul claro',
           '+2']
```

Tras separar la información por saltos de línea al hacer *split('\n')*, vemos que los datos que nos interesan se encuentran en los primeros lugares, pues primero está el nombre, y

después los precios con y sin descuentos. Sin embargo, no podemos generalizar por completo este formato, pues en productos como el que mostraremos, hay un título adicional:



```
1 opciones[21].text.split('\n')  
  
['Premium Selection',  
 'Camisa Regular Fit de manga corta',  
 '$499.00',  
 'Beige claro',  
 'Blanco']
```

Vemos que debido al título '*Premium Selection*', la información útil se desplaza un lugar, por lo que hay que tenerlo en cuenta para no causar errores. Otra cosa importante es que si no hay descuentos, en el lugar correspondiente solo hay un dato.

Esta situación de tener etiquetas adicionales antes del nombre se presentó en los tres sitios que trabajamos, por lo que motivó a realizar un método auxiliar para reconocer si estábamos en esa situación o no. Pero primero, veamos la función principal:

```
2 def hym(producto):  
3     ...  
4     Función para obtener datos en sitio HyM  
5     -----  
6     Parameters:  
7         producto: str  
8             Tipo de producto que buscará en el sitio  
9     -----  
10    Returns:  
11        df: pandas.DataFrame()  
12            DataFrame con la información más importante de los productos  
13    ...  
14    #Abrimos la página  
15    driver=webdriver.Chrome(path)  
16    url1 = 'https://www2.hm.com/es_mx/search-results.html?q='+producto  
17    driver.get(url1)  
18    time.sleep(12)  
19  
20    #Identificamos la clase de la información de los productos  
21    opciones = driver.find_elements_by_class_name("item-details")  
22    ...
```

La primera sección es clara. Abrimos la página con el formato del link que identificamos, y guardamos las opciones de productos con la clase correspondiente.

Lo que se buscará es guardar en listas diferentes la información que necesitamos, esto es, una para los nombres de los productos, otra para el precio actual, y una última para el precio sin descuento (si es que lo hay).

Haremos un ciclo por cada producto. Llamaremos a nuestro método auxiliar *verificar()*, para saber si hay una etiqueta extra o no. Este método retorna *True* si es el orden conveniente, y un *False* si no, es decir, si hay un dato extra.

Guardamos el nombre en la lista correspondiente.

Después limpiamos el lugar de los costos, quitando el signo de peso y comas, con la intención de poder convertir los datos a flotantes. Al usar el método *split()* estaremos separando los dos precios, si es que existen. Por ello, utilizamos una excepción, para agregar un *NA* a la lista de precios sin descuento en caso de que solo existiera un dato.

```
23 #Generamos 3 listas vacías donde se añadirá la información correspondiente
24 nombres = []
25 precio_actual = []
26 precio_sin_desc = []
27
28 for i in opciones:
29     pro_info = i.text.split('\n')
30
31     #Si el orden es correcto, empezamos desde el primer índice a llenar los datos. Si no, desde el segundo
32     if verificar(pro_info):
33         nombres.append(pro_info[0])
34
35         precios = pro_info[1]
36         precios = re.sub("[$,]", "", precios).split()
37         try:
38             precio_actual.append(float(precios[0]))
39             precio_sin_desc.append(float(precios[1]))
40         except IndexError:
41             precio_sin_desc.append('NA')
```

Para la siguiente parte, es repetir lo mismo pero modificando los índices, pues habrá que hacerlo una posición más arriba:

```
43 else:
44     nombres.append(pro_info[1]) #Se agrega el nombre
45
46     precios = pro_info[2]
47     precios = re.sub("[$,]", "", precios).split() #Se limpia los datos de los precios
48
49     #Usamos excepciones para guardar los precios
50     try:
51         precio_actual.append(float(precios[0]))
52         precio_sin_desc.append(float(precios[1]))
53     except IndexError:
54         #De no tener descuento, se llena el valor con un NA
55         precio_sin_desc.append('NA')
56
57     #Llamamos a la función auxiliar para generar dataframe
58     df = formato(nombres, precio_actual, precio_sin_desc, 'HyM', producto)
59
60
61
62     return df
```

La función *formato()* también es auxiliar. Debido a que en cada método de los tres sitios debemos generar un dataframe con el mismo orden de datos, es conveniente hacer una función por separado que dé el formato correcto, para no repetir lo mismo.

La implementación de las funciones auxiliares es la siguiente:

```
3 def verificar(pro_info):
4     """
5     Función auxiliar para verificar que los datos se encuentren en el orden deseado
6     -----
7     Parameters:
8         pro_info: list
9             Información del producto
10    """
11    #Tratamos de convertir el dato del precio a flotante. Si hay error, es porque el producto tiene un título extra,
12    #por lo que hay que cambiar los índices en la función principal
13    precio = pro_info[1]
14    precio = re.sub("[MXN$,]", "", precio).split()
15    try:
16        temp = float(precio[0])
17    except ValueError:
18        return False
19
20    return True
```

```
24 def formato(l1,l2,l3,sitio,producto):
25     """
26     Función auxiliar para ordenar los datos obtenidos y generar un DataFrame
27     -----
28     Parameters:
29         l1: list
30             Lista con nombres de los productos
31         l2: list
32             Lista con precios sin descuento
33         l3: list
34             Lista con precios sin descuento
35         sitio: str
36             Nombre de la página web donde se buscó la información
37         producto: str
38             Tipo de producto que se consultó
39     -----
40     Returns:
41         df: pandas.DataFrame()
42             Dataframe con el formato deseado
43     """
44
45     info = list(zip(l1,l2,l3))
46
47     #Creamos dataframe y añadimos columnas importantes
48     df = pd.DataFrame(info, columns=['Descripción', 'Precio Actual', 'Precio Sin Descuento'])
49     df.insert(0, "Producto", producto, allow_duplicates=False)
50     df.insert(0, "Sitio", sitio, allow_duplicates=False)
51     df.insert(5, "Fecha De Consulta", datetime.datetime.now().strftime('%Y-%m-%d'), allow_duplicates=False)
52
53
54
55     return df
56
```

En el método *formato()* lo que hacemos es ordenar las listas para crear un DataFrame que guarde los datos correctamente. Además, agregamos una columna para guardar el sitio donde se buscaron, el tipo de producto y la fecha de consulta.

Así sería la salida tras ingresar en nuestro método el producto 'Gorro':

```
64 hym('Gorro')
```

| | Sitio | Producto | Descripción | Precio Actual | Precio Sin Descuento | Fecha De Consulta |
|----|-------|----------|------------------------------|---------------|----------------------|-------------------|
| 0 | HyM | Gorro | Chamarra acolchada con gorro | 299.0 | 549.0 | 2022-12-14 |
| 1 | HyM | Gorro | Sudadera Relaxed Fit | 449.0 | NA | 2022-12-14 |
| 2 | HyM | Gorro | Gorro tejido | 179.0 | NA | 2022-12-14 |
| 3 | HyM | Gorro | Gorro con pompón | 199.0 | 249.0 | 2022-12-14 |
| 4 | HyM | Gorro | Gorro en tejido acanalado | 249.0 | NA | 2022-12-14 |
| 5 | HyM | Gorro | Sudadera Relaxed Fit | 549.0 | NA | 2022-12-14 |
| 6 | HyM | Gorro | Sudadera Regular Fit | 549.0 | NA | 2022-12-14 |
| 7 | HyM | Gorro | Sudadera con diseño | 169.0 | 249.0 | 2022-12-14 |
| 8 | HyM | Gorro | Chamarra en denim con gorro | 699.0 | NA | 2022-12-14 |
| 9 | HyM | Gorro | Sudadera con diseño | 399.0 | NA | 2022-12-14 |
| 10 | HyM | Gorro | Sudadera Relaxed Fit | 449.0 | NA | 2022-12-14 |
| 11 | HyM | Gorro | Conjunto de 2 piezas | 749.0 | NA | 2022-12-14 |
| 12 | HyM | Gorro | Sudadera Relaxed Fit | 449.0 | NA | 2022-12-14 |
| 13 | HyM | Gorro | Chaleco puffer con gorro | 749.0 | NA | 2022-12-14 |
| 14 | HyM | Gorro | Sudadera con diseño | 449.0 | NA | 2022-12-14 |
| 15 | HyM | Gorro | Sudadera Regular Fit | 499.0 | NA | 2022-12-14 |

Para los siguientes sitios, el procedimiento es análogo. Solo hay que verificar el nombre de la clase correspondiente y el formato de los datos, pero la lógica para implementar la función es la misma.

Finalmente, revisemos la función principal. Se encargará de recopilar la información de una lista de productos en los tres sitios, para después unirla en un solo DataFrame y guardarlo en un archivo de excel.

```

2 def main_excel(lis_productos):
3     ...
4     Función que genera un archivo de excel con la información de una lista de productos
5     en los sitios: HyM, Bershka y Shein
6     -----
7     Parameters:
8         lis_productos: list
9         Lista con productos que se desea analizar
10    ...
11    #Creamos 3 dataframes donde se guardarán la información de cada búsqueda
12    df1 = pd.DataFrame()
13    for i in lis_productos:
14        t = hym(i)
15        df1 = pd.concat([df1, t])
16
17    df2 = pd.DataFrame()
18    for i in lis_productos:
19        t = bka(i)
20        df2 = pd.concat([df2, t])
21
22    df3 = pd.DataFrame()
23    for i in lis_productos:
24        t = shn(i)
25        df3 = pd.concat([df3, t])
26
27    #Al finalizar los ciclos, unimos todos los dataframes y ajustamos índices
28    df_final = pd.DataFrame()
29    df_final = pd.concat([df1,df2,df3])
30    df_final.index = range(df_final.shape[0])
31
32    #Guardamos un archivo de excel con la información obtenida
33    df_final.to_excel('Proyecto.xlsx')
34
35

```

Para cada sitio, se empieza con un DataFrame vacío, de tal forma que se concatenen los que se van generando. Al terminar, guarda un archivo de excel con nombre 'Proyecto.xlsx'.

Veamos un resultado tras ingresar en esta función la lista ['Camisa', 'Falda', 'Suéter'] :

```

1 archivo = 'C:/Users/argea/Documents/Rafa/MDD/PF/Proyecto.xlsx'
2 DF = pd.read_excel(archivo)
3 DF

```

| Unnamed: 0 | Sitio | Producto | Descripción | Precio Actual | Precio Sin Descuento | Fecha De Consulta | |
|------------|-------|----------|-------------|---|----------------------|-------------------|------------|
| 0 | 0 | HyM | Camisa | Sobrecamisa recubierta | 599 | NaN | 2022-12-14 |
| 1 | 1 | HyM | Camisa | Sobrecamisa Relaxed Fit | 599 | NaN | 2022-12-14 |
| 2 | 2 | HyM | Camisa | Sobrecamisa Regular Fit | 329 | 429.0 | 2022-12-14 |
| 3 | 3 | HyM | Camisa | Camisa en sarga de algodón | 249 | 429.0 | 2022-12-14 |
| 4 | 4 | HyM | Camisa | Sobrecamisa acolchada | 599 | NaN | 2022-12-14 |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 523 | 523 | Shein | Sueter | Jersey con patrón de dibujo de hombros caídos | 470 | NaN | 2022-12-14 |
| 524 | 524 | Shein | Sueter | SHEIN Jersey de manga farol de cuello redondo | 430 | NaN | 2022-12-14 |
| 525 | 525 | Shein | Sueter | DAZY Jersey unicolor tejido de canalé | 298 | 8.0 | 2022-12-14 |
| 526 | 526 | Shein | Sueter | Jersey con patrón de Isla Fair de manga raglán... | 525 | NaN | 2022-12-14 |
| 527 | 527 | Shein | Sueter | SHEIN ICON Jersey con patrón de rayas de hombr... | 431 | NaN | 2022-12-14 |

528 rows × 7 columns

Conclusiones

Consideramos que esta herramienta es de gran utilidad, pues extraer grandes cantidades de información de manera automática resulta primordial para realizar diferentes tareas, desde consultas personales para saber dónde es más barato comprar cierto producto, hasta estudios de empresas para analizar a su competencia.

A pesar de que cada página tiene sus particularidades, en general se puede replicar el código cambiando solo unas pequeñas cosas, por lo que extender el alcance de nuestra implementación es relativamente sencillo. Esto es muy importante, pues entre más sitios sea capaz de analizar, más aplicaciones tendrá.