

¿Cómo podemos capturar los datos de la web?

Alex Stewart Porras Palma

Índice

1. Proyecto de captura de datos mediante el uso de webscraping de la plataforma steam	2
2. Descripción del data set	2
3. Descripción del funcionamiento	3
4. Inspiración	4
5. Código	4
5.1. Parámetro options	4
5.2. Comando search	4
5.3. Comando sales	6
6. Robot.txt	8
7. MongoDB	9
8. Conclusiones y Líneas futuras	9
9. Anaconda	9

1. Proyecto de captura de datos mediante el uso de webscraping de la plataforma steam

En este proyecto se va a proceder a crear un bot de telegram desde el cual realizar el webscraping sobre las ofertas actuales de la plataforma steam, así como la búsqueda de un título de juego y su precio actual dado un precio objetivo.

2. Descripción del data set

El dataset será implementado en la base de datos MongoDB con los siguientes campos:

- id (string) Identificador único del elemento, es un valor autogenerado por MongoDB
- title (string) Nombre del juego
- discount (string) Descuento aplicado sobre el juego
- old_price (string) Precio original antes del descuento
- new_price (string) Precio actual después del descuento
- source (string) Origen del comando usado para extraer la información
- created_at (string) Marca de tiempo que indica cuando fue el objeto añadido a la base de datos



Figura 1: Ejemplo de dato almacenado en MongoDB

3. Descripción del funcionamiento

- Usuario de telegram envía un comando como /search Devil May Cry 5 20 o /sales
- El servidor de Telegram recibe el mensaje y lo redirige al script de python
- Selenium Navega a steam
- Recupera la información deseada de steam
 - Título del juego
 - Precios (viejo/nuevo)
 - Descuento
- Los datos recogidos se devuelven a telegram
- Se pregunta si se desea guardar los datos
 - Sí, se realiza la conexión a MongoDB y se insertan los datos
 - No. Fin



Figura 2: Esquema visual del funcionamiento del bot

4. Inspiración

La idea detrás del proyecto es la posibilidad de crear un histórico de ofertas y de precios, así como evitar la necesidad de entrar en la página web a mirar y o buscar un juego. Un ejemplo del histórico se puede encontrar en la página <https://steamdb.info/> la cual ofrece un histórico de los precios dado un intervalo de tiempo.



Figura 3: Ejemplo de steam para histórico de precios

5. Código

A modo de presentación el código se ha estructurado en un jupyter notebook

5.1. Parámetro options

```
1 options = Options()  
2 options.add_argument("start-maximized")  
3 #options.add_argument("--headless")
```

Listing 1: parámetro options function

Para ejecutar el script se pueden hacer uso de diversos parámetros en options en este caso a modo de muestra se comentará la opción de headless para que se pueda observar como se van realizando las acciones en el navegador

5.2. Comando search

```
1 def search_game_by_price(game_name, target_price):  
2     driver = webdriver.Chrome(service=Service(CHROME_DRIVER_PATH)  
3         , options=options)  
4     try:  
5  
6         driver.get('https://store.steampowered.com/')  
7  
8  
9         WebDriverWait(driver, 10).until(  
10
```

```

10         EC.presence_of_element_located((By.ID, '
11             store_nav_search_term'))
12     )
13
14     search_box = driver.find_element(By.ID, '
15         store_nav_search_term')
16     search_box.clear()
17     search_box.send_keys(game_name)
18     search_box.send_keys(Keys.RETURN)
19
20     WebDriverWait(driver, 10).until(
21         EC.presence_of_all_elements_located((By.CSS_SELECTOR,
22             'div.responsive_search_name_combined'))
23     )
24
25     first_result = driver.find_element(By.CSS_SELECTOR, 'div.
26         responsive_search_name_combined')
27
28     title = first_result.find_element(By.CSS_SELECTOR, 'span.
29         title').text.strip()
30
31     price_block = first_result.find_element(By.CSS_SELECTOR,
32         'div.search_price_discount_combined')
33
34     try:
35
36         discount_block = price_block.find_element(By.
37             CSS_SELECTOR, 'div.discount_block')
38
39         if 'no_discount' in discount_block.get_attribute('
40             class'):
41
42             final_price_element = price_block.find_element(By
43                 .CSS_SELECTOR, 'div.discount_final_price')
44             final_price = final_price_element.text.strip()
45
46             final_price = float(final_price.replace(' ', ''
47                 ).replace(',', '.').strip())
48         else:
49
50             final_price_element = discount_block.find_element
51             (By.CSS_SELECTOR, 'div.discount_final_price')
52             final_price = final_price_element.text.strip()
53
54             final_price = float(final_price.replace(' ', ''
55                 ).replace(',', '.').strip())
56
57     except Exception as e:
58         print(f"Error extracting price for {title}: {e}")

```

```

49         return None
50
51     if final_price <= target_price:
52         print(f"Game: {title}")
53         print(f"Price: {final_price}    ")
54         return title, final_price
55
56     print(f"Game: {title} is not under the target price.")
57     return None
58
59 except Exception as e:
60     print(f"Error during search: {e}")
61     return None
62
63 finally:
64     driver.quit()

```

Listing 2: search game + price function

En esta función Selenium se usa para automatizar la interacción de la página web Steam para recolectar información sobre un juego en específico. Utiliza el webdrive de Chrome, para ello es necesario descargar el chromedriver compatible con la versión de chrome instalada en el dispositivo. Esto permite la ejecución en segundo plano del navegador sin la necesidad abrirlo de manera visible. Navega a la página de steam, espera a que la página cargue busca el elemento de la barra de búsqueda utilizando el ID de HTML, borra por si hubiese texto introducido y escribe el título del juego proporcionado simulando la búsqueda como si fuera un usuario pulsando el botón enter. Tras la búsqueda espera a que los resultados carguen por completo y se recopilan todos los elementos de la búsqueda mediante un selector CSS que escoge el primer resultado de la búsqueda para extraer la información. Se evalúa si el bloque contiene información sobre descuentos, y se recupera. Finalmente se limpia la cadena de caracteres y se sustituyen las comas por puntos. (Esto último se hace por formato propio de decimal)

5.3. Comando sales

```

1 def sales_from_steam():
2     driver = webdriver.Chrome(service=Service(CHROME_DRIVER_PATH)
3     , options=options)
4     print(f"Navigating to: {URL}")
5     driver.get(URL)
6
7     try:
8         WebDriverWait(driver, 10).until(
9             EC.presence_of_element_located((By.ID, '
10                 SaleSection_13268')))
11
12         # Scroll to bottom to load all items
13         last_height = driver.execute_script("return document.body
14             .scrollHeight")
15         while True:

```

```

14     driver.execute_script("window.scrollTo(0, document.
15         body.scrollHeight);")
16     time.sleep(1)
17     new_height = driver.execute_script("return document.
18         body.scrollHeight")
19     if new_height == last_height:
20         break
21     last_height = new_height
22
23     game_containers = driver.find_elements(By.CSS_SELECTOR, '
24         div.v9uRg57bw0aPsvAnkXES0')
25
26     games = []
27
28     for game in game_containers:
29         try:
30             title = game.find_element(By.CSS_SELECTOR, 'div.
31                 StoreSaleWidgetTitle').text
32             discount = game.find_element(By.CSS_SELECTOR, '
33                 div.cnkoFkzVCby40gJ0jGGS4').text
34
35             try:
36                 old_price = game.find_element(By.CSS_SELECTOR
37                     , 'div._3fFFsvII7Y2KXNLDk_krOW').text
38                 new_price = game.find_element(By.CSS_SELECTOR
39                     , 'div._3j4dI1yA7cRfCvK8h4060B').text
40             except:
41                 old_price = new_price = "N/A"
42
43             games.append({
44                 'title': title,
45                 'discount': discount,
46                 'old_price': old_price,
47                 'new_price': new_price
48             })
49
50         except Exception as e:
51             print(f"Error parsing a game: {e}")
52             continue
53     return games
54
55 except Exception as e:
56     print(f"Error: {e}")
57 finally:
58     driver.quit()

```

Listing 3: sales function

El comando sales utiliza selenium par acceder y extraer los datos de las ofertas actuales de steam desde su apartado ofertas especiales. De la misma manera que el anterior, comienza iniciando el navegador con el webdriver de chrome. Se dirige directamente a la página de las ofertas especiales y espera a que cargue la sección de las ofertas, éste es localizado mediante HTML ID SalesSection_13268, debido a que no todo el contenido carga a primera vista es necesario desplazarse hacia abajo, el script simula el desplazamiento manual ejecutando JS que se desplaza hasta el final de la página. Una vez hecho esto, se localizan todos los elementos del contenedor (12 en este caso) mediante HTML CSS selector. Se itera sobre cada elemento y se extrae la información.

6. Robot.txt

```
1 Host: store.steampowered.com
2 User-Agent: *
3 Disallow: /share/
4 Disallow: /news/externalpost/
5 Disallow: /account/emailoptout/?*token=
6 Disallow: /login/?*guestpasskey=
7 Disallow: /join/?*redir=
8 Disallow: /account/ackgift/
9 Disallow: /email/
10 Disallow: /widget/
```

Listing 4: Robots.txt steam

Antes de realizar el scraping es importante revisar el archivo robots.txt de steam, desde aquí se informa a los usuarios a qué partes se pueden acceder mediante el uso de bots. User-agent: * indica que las reglas especificadas a continuación se aplican a todos los agentes de usuario automatizados, sin excepción. Las entradas "Disallow."^{en}umeran rutas específicas en el sitio web que los bots no deben visitar. Por ejemplo, los directorios relacionados con cuentas de usuario, operaciones de correo electrónico etc. Cabe destacar que en el caso del uso de selenium la idea es la de simular un comportamiento humano por lo que con un modelo robusto y eficaz sería posible saltarse las alertas de bots, (los captcha). El bot realizado para este proyecto no infringe la política de robots así que se podría continuar con su uso.

7. MongoDB

Como base de datos se ha escogido MongoDB, el motivo es el de aprovechar la tecnologías proporcionadas por la asignatura, así como sacar partido a la estructura de documentos, para almacenar los datos. En este caso al ser un ejemplo sencillo, se tiene en cuenta de que los juegos pueden tener descuento o no, por lo que el esquema flexible de MongoDB puede gestionar estas inconsistencias sin ningún tipo de problema. Además admite operaciones de lectura y escritura de alta velocidad, lo que lo hace ideal para aplicaciones en tiempo real, como un bot de Telegram que recopila y almacena información del juego dinámicamente según las solicitudes de los usuarios o el análisis de ventas. Aunque con la versión gratuita se obtiene una base de datos limitada es más que suficiente para gestionar este proyecto. Para realizar la conexión se hace uso de su sistema en la nube MongoDB atlas lo que es especialmente práctico a la hora de implementar el bot en servidores en la nube.

8. Conclusiones y Líneas futuras

El proyecto demuestra como la automatización del web scraping se puede utilizar para la creación de una herramienta para conseguir y almacenar datos en tiempo real se ha creado un sistema escalable y flexible además los usuarios pueden buscar juegos específicos y recibir información en tiempo real según sus precios deseados, además de ver las ventas en curso directamente desde la tienda de Steam. Además, MongoDB permite almacenar y recuperar posteriormente los juegos consultados, lo que facilita el análisis o el seguimiento histórico.

9. Anaconda

Para hacer un mejor uso tanto de las librerías como del entorno y su espacio, se ha hecho uso de Anaconda, para crear un entorno virtual, el cual se usará como kernel para compilar tanto el jupyter notebook, como el script del bot, de esta manera, una vez finalizado el proyecto, se puede eliminar la carpeta y de esta manera conservar el espacio, manteniendo una práctica limpia y eficaz.

Como mejoras al proyecto sería interesante implementar un comando en segundo plano que revise periódicamente el precio de un juego o de varios juegos y que envíe una notificación a telegram cuando los juegos bajen del precio establecido por el usuario.