

Scrapping des données

Elaboré par :

- Rekik Wiem
- Dhaouadi Mohamed Amine
- Dghaies Yasser
- Ferjani Mohamed Iheb
- Battikh Anis
- Achour Nouha

I. Introduction

Ce script Python utilise les bibliothèques Selenium et BeautifulSoup pour extraire des données à partir d'une page Web contenant un tableau de données. L'objectif principal est de récupérer les informations de la page sur les collisions de véhicules à New York depuis le portail de données ouvert de la ville.

II. Données Source

Le tableau de données présente les accidents liés aux collisions de véhicules automobiles à New York, fournissant des détails sur chaque incident spécifique. Chaque ligne de ce tableau représente un événement de collision distinct, présentant des informations telles que la date, l'heure, l'emplacement, les véhicules impliqués, et d'autres détails pertinents. Ces données sont compilées à partir des rapports officiels de la police, offrant une vue d'ensemble complète et détaillée des collisions de véhicules à moteur dans la région

III. Outils



IV. Script

```

: from selenium import webdriver
from selenium.webdriver.common.by import By
from selenium.webdriver.chrome.options import Options
from selenium.webdriver.support.ui import WebDriverWait
from selenium.webdriver.support import expected_conditions as EC
from selenium.common.exceptions import StaleElementReferenceException
from bs4 import BeautifulSoup
import pandas as pd
import time

def initialize_driver():
    options = Options()
    options.add_argument('--headless')
    return webdriver.Chrome(options=options)

def wait_for_element(driver, locator, timeout=10):
    return WebDriverWait(driver, timeout).until(EC.presence_of_element_located(locator))

def click_and_retry(driver, element, locator, retry_attempts=3):
    for _ in range(retry_attempts):
        try:
            element.click()
            return True
        except StaleElementReferenceException:
            print("Stale Element Reference Exception. Retrying...")
            element = wait_for_element(driver, locator)
        except Exception as e:
            print(f"An error occurred: {str(e)}")
            return False
    return False

def scrape_table_data(driver, table_locator):
    table = wait_for_element(driver, table_locator)
    table_html = table.get_attribute('outerHTML')
    soup = BeautifulSoup(table_html, 'html.parser')
    html_table = soup.find('table')
    return pd.read_html(str(html_table))[0]

```

```

def main():
    url = "https://data.cityofnewyork.us/Public-Safety/Motor-Vehicle-Collisions-Crashes/h9gi-nx95"
    driver = initialize_driver()

    try:
        driver.get(url)
        driver.implicitly_wait(10)

        table_locator = (By.CLASS_NAME, 'socrata-table.frozen-columns')

        if not (table := wait_for_element(driver, table_locator)):
            print("No table found on the page with class 'socrata-table frozen-columns'.")
            return

        dfs = []

        total_pages = 500

        for page in range(1, total_pages + 1):
            next_button_locator = (By.XPATH, "//button[@class='pager-button-next']")
            next_button = wait_for_element(driver, next_button_locator)

            if not click_and_retry(driver, next_button, next_button_locator):
                print("Failed to click next button. Exiting.")
                break

            time.sleep(3)

            df = scrape_table_data(driver, table_locator)

            print(df.head())

            dfs.append(df)

        result_df = pd.concat(dfs, ignore_index=True)

        result_df = result_df.dropna(axis=1, how='all')

```

ChatGPT - Google Chrome

```

        result_df = result_df.dropna(axis=1, how='all')

        result_df.to_excel("DataCollisions.xlsx", index=False)
        print("Cleaned data successfully exported DataCollisions.xlsx")

    except Exception as e:
        print(f"An error occurred: {str(e)}")

    finally:
        driver.quit()

if __name__ == "__main__":
    main()

```

2 NaN HARLEM RIVER DRIVE RAMP

NaN

Étape 1 : Initialisation du WebDriver

Le script débute en initialisant un WebDriver Chrome en mode headless à l'aide de la bibliothèque Selenium. Cette étape permet une navigation automatisée sur la page Web sans interface graphique.

Étape 2 : Navigation vers la Page Cible

Le WebDriver accède à la page de données des collisions de véhicules à New York en utilisant l'URL spécifié. Une attente implicite garantit que la page est entièrement chargée avant de procéder à la collecte des données.

Étape 3 : Scraping de Données

Le script identifie le tableau de données sur la page (repéré par la classe CSS 'socrata-table.frozen-columns') et commence le processus d'extraction des informations. La pagination automatique est gérée pour collecter les données de plusieurs pages.

Étape 4 : Gestion des Éléments Dynamiques

Les mécanismes de gestion des éléments dynamiques, tels que l'utilisation de l'attente explicite et la gestion des exceptions Stale Element Reference, sont mis en œuvre pour assurer la fiabilité du scraping, même sur des pages dynamiques.

Étape 5 : Nettoyage des Données

Une fois les données collectées, le script les nettoie en supprimant les colonnes qui ne contiennent que des valeurs nulles, optimisant ainsi la qualité et la lisibilité du jeu de données.

Étape 6 : Exportation des Données Nettoyées

Les données nettoyées sont ensuite exportées vers un fichier Excel, facilitant leur utilisation dans des analyses ultérieures ou d'autres applications.

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	
339	04/11/202 5:28	BROOKLYN	11230	40.63529	-73.9582	(40.63529'OCEAN AV FARRAGUT ROAD					0	0	0	0	0	0	0	0	Fell Asleep Unspecified					4
340	04/11/202 5:30	QUEENS	11417	40.67917	-73.8532	(40.67917",-73.85315")	105-29 84				0	0	0	0	0	0	0	0	Driver Inat Unspecified Unspecified					4
341	04/11/202 6:05			40.61131	-74.0984	(40.61131: CLOVE RO NARROWS ROAD NOI					2	0	0	0	0	0	0	2	Reaction t Unspecified					4
342	04/11/202 6:06			40.86794	-73.8722	(40.86794' BRONX RIVER PARKWAY					3	0	0	0	0	0	0	3	Unsafe Lai Unspecified					4
343	04/11/202 6:38	MANHATT	10024	40.78423	-73.9788	(40.784225",-73.97879")	222 WEST				0	0	0	0	0	0	0	0	Driver Inat Unspecified					4
344	04/11/202 6:40	QUEENS	11691			ROCKAWA BEACH 38 STREET					0	0	0	0	0	0	0	0	Other Veh Traffic Control Disregarded					4
345	04/11/202 6:50	BROOKLYN	11236	40.63781	-73.8955	(40.63781",-73.895546")	1414 EAST				0	0	0	0	0	0	0	0	Driver Inat Unspecified					4
346	04/11/202 7:14	BROOKLYN	11212	40.66353	-73.9131	(40.66353' DUMONT AMBOY STREET					1	0	0	0	0	0	0	1	Unspecifie Unspecified Unspecified					4
347	04/11/202 7:15	QUEENS	11356	40.78322	-73.8459	(40.78322' 18 AVENU COLLEGE POINT BOU					1	0	0	0	0	0	0	1	Following Unspecified					4
348	04/11/202 7:20	BROOKLYN	11226	40.65134	-73.9566	(40.65134",-73.95657")	75 MARTE				0	0	0	0	0	0	0	0	Unspecifie Unspecified					4
349	04/11/202 7:40	QUEENS	11379	40.70826	-73.8751	(40.70826",-73.87514")	69-23 76 S				0	0	0	0	0	0	0	0	Fell Asleep Unspecified Unspecified					4
350	04/11/202 7:50	BROOKLYN	11234	40.62303	-73.9254	(40.623028",-73.92538")	1490A EAS				0	0	0	0	0	0	0	0	Passing or Unspecified					4
351	04/11/202 7:59			40.74984	-73.939	(40.74984' QUEENS P 28 STREET					1	0	0	0	0	0	0	1	Following Unspecified					4
352	04/11/202 8:30			40.74901	-73.8347	(40.74901' VAN WYCK EXPWY					0	0	0	0	0	0	0	0	Driver Inexperience					4
353	04/11/202 8:30	BROOKLYN	11237	40.69426	-73.9104	(40.69426' JEFFERSON KNICKERBOCKER AVE					1	0	0	0	0	0	0	1	Fell Asleep Unspecified					4
354	04/11/202 8:35	QUEENS	11375	40.73095	-73.8486	(40.73094' 65 ROAD 108 STREET					0	0	0	0	0	0	0	0	Passing or Unspecifie Unspecifie Unspecifie					4
355	04/11/202 8:50			40.63667	-74.0257	(40.63667' 3 AVENUE 68 STREET					0	0	0	0	0	0	0	0	Alcohol Im Unspecified					4
356	04/11/202 9:00	BROOKLYN	11234	40.61826	-73.9347	(40.618256",-73.93466")	1608 COLE				0	0	0	0	0	0	0	0	Unspecified					4
357	04/11/202 9:00	QUEENS	11372	40.75048	-73.8775	(40.750477",-73.87746")	89-27 37 A				0	0	0	0	0	0	0	0	Unspecifie Unspecified					4
358	04/11/202 9:00	QUEENS	11373	40.74058	-73.8903	(40.740578",-73.890274")	45-12 74 S				0	0	0	0	0	0	0	0	Driver Inat Unspecified					4
359	04/11/202 9:15			40.70003	-73.7887	(40.70003' MERRICK (107 AVENUE					0	0	0	0	0	0	0	0	Unsafe Lai Unspecified					4
360	04/11/202 9:29	BRONX	10472	40.82737	-73.8574	(40.827374",-73.85735")	1973 CHA				0	0	0	0	0	0	0	0	Driver Inat Unspecified					4
361	04/11/202 9:30			40.7221	-73.7777	(40.72209' GRAND CENTRAL PKWY					1	0	0	0	0	0	0	1	Driver Inattention/Distractio					4
362	04/10/202 0:00	MANHATT	10018	40.75517	-73.9913	(40.75517",-73.99129")	594 8 AVE				0	0	0	0	0	0	0	0	Driver Inat Unspecifie Unspecified					4
363	04/10/202 0:00	MANHATT	10033	40.84706	-73.9382	(40.84706",-73.93818")	4162 BRO				0	0	0	0	0	0	0	0	Driver Inattention/Distractio					4
364	04/10/202 0:00	QUEENS	11414	40.66317	-73.8409	(40.663166",-73.84086")	156-71 CR				0	0	0	0	0	0	0	0	Unspecified					4
365	04/10/202 0:00			40.70644	-73.7597	(40.70644' HOLLIS AV 198 STREET					0	0	0	0	0	0	0	0	Driver Inat Unspecified					4
366	04/10/202 0:00					VICTORY BOULEVARD					0	0	0	0	0	0	0	0	Unspecifie Unspecified					4
367	04/10/202 0:00			40.81603	-73.9395	(40.81602' WEST 138 STREET					0	0	0	0	0	0	0	0	Driver Inat Unspecified					4

#	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U
6476	03/20/202	14:00	BROOKLYN	11226	40.65279	-73.9469	[40.65279; NEW YORK UNDER BOULEVARD]				1	0	0	0	1	0	0				0 Unspecified
6477	03/20/202	14:00	BROOKLYN	11236	40.64141	-73.9123	[40.64140; PARK AVENUE EAST 86 STREET]				1	0	0	0	0	0	1				0 Unspecified Unspecified
6478	03/20/202	14:05	BRONX	10467			EAST FORC PELHAM PARKWAY				1	0	0	0	1	0	0				0 Failure to Unspecified
6479	03/20/202	14:08	QUEENS	11435	40.69790	-73.8115	[40.69790; 94 AVENUE; 138 PLACE]				1	0	0	0	0	0	1				0 Passing or Unspecified
6480	03/20/202	14:30	QUEENS	11434	40.68257	-73.7927	[40.68256; 136 AVENUE SUTPHIN BOULEVARD]				0	0	0	0	0	0	0				0 Unspecified Unspecified
6481	03/20/202	14:30			40.73042	-73.9518	[40.73042; MADISON AVENUE]				0	0	0	0	0	0	0				0 Lost Conn. Unspecified Unspecified
6482	03/20/202	14:34	QUEENS	11354	40.75586	-73.4362	[40.75586; -73.836155]				0	0	0	0	0	0	0				0 Turning Int. Unspecified
6483	03/20/202	14:36	QUEENS	11473	40.7418	-73.8788	[40.7418; -73.8788]				1	0	0	0	0	0	1				0 Turning Int. Unspecified
6484	03/20/202	14:38	MANHATT	10028	0	0	[40.07; 0.07]				0	0	0	0	0	0	0				0 Backing Up Unspecified
6485	03/20/202	14:43			40.68316	-73.9381	[40.68316; HALSEY STREET]				1	0	1	0	0	0	0				0 Driver Inattention/Distracted
6486	03/20/202	14:48	BROOKLYN	11229	0	0	[40.07; 0.07]				0	0	0	0	0	0	0				0 Backing Up Unspecified
6487	03/20/202	14:48	MANHATT	10037	40.75374	-73.9678	[40.75374; -73.96783]				0	0	0	0	0	0	0				0 Driver Inat. Unspecified
6488	03/20/202	14:49	BRONX	10467	0	0	[40.07; 0.07] BRONX PA WAREING AVENUE				0	0	0	0	0	0	0				0 Driver Inat. Unspecified
6489	03/20/202	14:50	BRONX	10472	40.82975	-73.8757	[40.82975; HARROD A WESTCHESTER AVENUE]				0	0	0	0	0	0	0				0 Unspecified
6490	03/20/202	14:56	QUEENS	11432	0	0	[40.07; 0.07] JAMAICA A 168 PLACE				0	0	0	0	0	0	0				0 Driver Inat. Unspecified
6491	03/20/202	15:00	BRONX	10460	40.83824	-73.8767	[40.83824; EAST 177 S BRONX PARK AVENUE]				0	0	0	0	0	0	0				0 Driver Inat. Driver Inattention/Clash
6492	03/20/202	15:00	BROOKLYN	11221	0	0	[40.07; 0.07]				0	0	0	0	0	0	0				0 Driver Inat. Unspecified
6493	03/20/202	15:00			40.69042	-73.8443	[40.69042; EAST 225 STREET]				0	0	0	0	0	0	0				0 Unsafe Spa Unspecified
6494	03/20/202	15:04			40.70245	-73.8594	[40.70245; JACOB ROBINSON PIKE]				0	0	0	0	0	0	0				0 Reaction to Uninsured Vehicle
6495	03/20/202	15:05	BRONX	10466	40.88721	-73.8608	[40.887207; -73.86082]				2	0	2	0	0	0	0				0 Driver Inexperience
6496	03/20/202	15:12			40.67269	-73.7254	[40.67269; LAURELTON PARKWAY]				0	0	0	0	0	0	0				0 Unspecified Unspecified
6497	03/20/202	15:25	BRONX	10455	40.80988	-73.9051	[40.80988; EAST 149 S BRUCKNER BOULEVARD]				1	0	0	0	0	0	1				0 Unspecified Unspecified
6498	03/20/202	15:25	QUEENS	11454	40.76582	-73.8237	[40.76582; 35 AVENUE; PARSONS BOULEVARD]				1	0	1	0	0	0	0				0 Failure to Yield Right-of-Way
6499	03/20/202	15:30	QUEENS	11372	40.7687	-73.874	[40.7687; -73.87395]				1	0	0	0	1	0	0				0 Unspecified
6500	03/20/202	15:32	BROOKLYN	11257	0	0	[40.07; 0.07]				0	0	0	0	0	0	0				0 Passenger View Obstructed/Limit
6501	03/20/202	15:40			0	0	[40.07; 0.07] CHRISTIE STREET				1	0	0	0	0	0	1				0 Physical DC Unspecified