# DSCI 560 - Lab 1 Documentation

## 1. Installation and Setup

- To keep things simple, **WSL** (Windows Subsystem for Linux) was used instead of VMware for the linux environment.

```
ariel@DESKTOP-A7HUT0A: /n    ×    +    ⌄                                    —    □    ×
Microsoft Windows [Versión 10.0.26100.2894]
(c) Microsoft Corporation. Todos los derechos reservados.

C:\Users\User\Desktop\560\DSCI-560\lab1\arielmartinez_5483611649>wsl
ariel@DESKTOP-A7HUT0A:/mnt/c/Users/User/Desktop/560/DSCI-560/lab1/arielmartinez_5483611649$ python3 --version
Python 3.10.12
ariel@DESKTOP-A7HUT0A:/mnt/c/Users/User/Desktop/560/DSCI-560/lab1/arielmartinez_5483611649$ pip3 --version
pip 22.0.2 from /usr/lib/python3/dist-packages/pip (python 3.10)
ariel@DESKTOP-A7HUT0A:/mnt/c/Users/User/Desktop/560/DSCI-560/lab1/arielmartinez_5483611649$
```

## 2.1. Playing around with Linux Terminal

- Python3 and pip were successfully installed and the basic files and repositories were created using **touch** and **mkdir** respectively.

```
ariel@DESKTOP-A7HUT0A:/mnt/c/Users/User/Desktop/560/DSCI-560/lab1/arielmartinez_5483611649$ ls
data    scripts
ariel@DESKTOP-A7HUT0A:/mnt/c/Users/User/Desktop/560/DSCI-560/lab1/arielmartinez_5483611649$ cd scripts
ariel@DESKTOP-A7HUT0A:/mnt/c/Users/User/Desktop/560/DSCI-560/lab1/arielmartinez_5483611649/scripts$ ls
task_1.py
```

## 2.2. A basic Python Script

- The **nano** command was used to write in the first python file.
- The file asks for your name (albeit a bit aggressively) and greets you in return as intended.

```
  GNU nano 6.2                              task_1.py
# Input
name = input("\nHey, you! Yeah, YOU! What's your name? Spit it out (Dumbledore asked calmly): ")

# Answer
print(f"\nHello, {name}!\n")
```

```
ariel@DESKTOP-A7HUT0A:/mnt/c/Users/User/Desktop/560/DSCI-560/lab1/arielmartinez_5483611649/scripts$ python3 task_1.py

Hey, you! Yeah, YOU! What's your name? Spit it out (Dumbledore asked calmly): Ariel

Hello, Ariel!

ariel@DESKTOP-A7HUT0A:/mnt/c/Users/User/Desktop/560/DSCI-560/lab1/arielmartinez_5483611649/scripts$
```
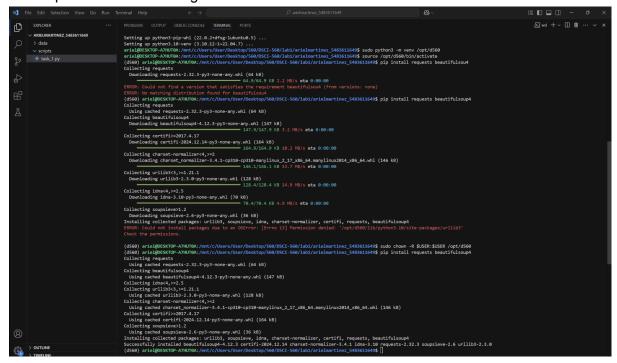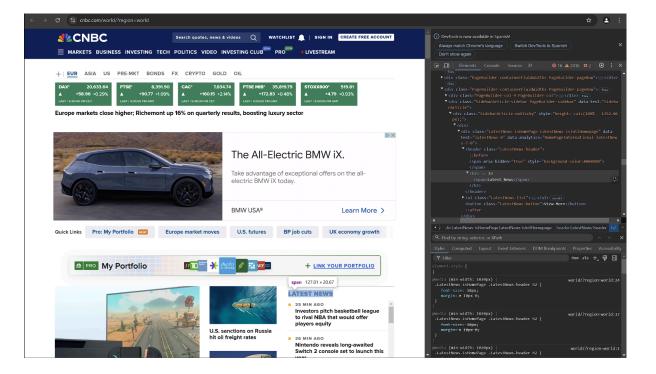
## 2.3. Python Web-scraping Task

- Command **venv** was used to create an exclusive environment (d560) for this subject which most likely will be used in future labs (inside WSL).
- Requests and BeautifulSoup were installed correctly after solving some minor problems with user rights.



- The CNBC website (https://www.cnbc.com/world/?region=world) was analyzed to find the Market Banner (`MarketsBanner-main`) and the Latest News (`LatestNews-list`) sections.

- The new folders and files were successfully created like the first time, just that this time from the VScode terminal.



- web_scraper.py successfully extracted the desired parts of the html.
- **Selenium** was used instead of **request** due to the fact that request wasn't able to extract Markets Banner data correctly due to how its data was integrated. **WebDriverWait** in particular was needed to avoid this problem by waiting to receive information on Markets Banner Row sections.
- **BeautifulSoup** was used to parse and filter the html sections and the **os** library was used to obtain the correct file path (while maintaining replicability) for the output (web_data.html).
- Prints were used to mark every section of the code.

## 2.4.  Data Filtering Task

- Script data_filter.py was successful in creating a structured csv for both sections of web_data.html.
- Similarly to the task before, **BeautifulSoup** was used to parse the html and **os** was used to correctly manage the input and output files in a replicable way.
- The Market data extraction was straightforward since it had specific classes for the 3 desired columns, but in the case of the Latest News section, the **link** was found in the href (as usual in a html) of the headline (title) class.
- Prints were used to mark every section of the code.