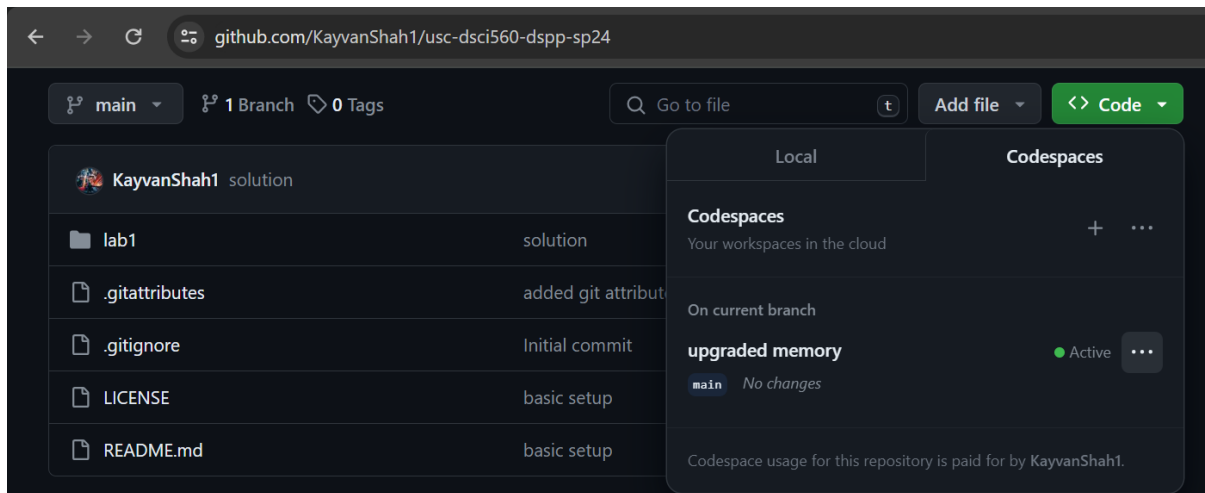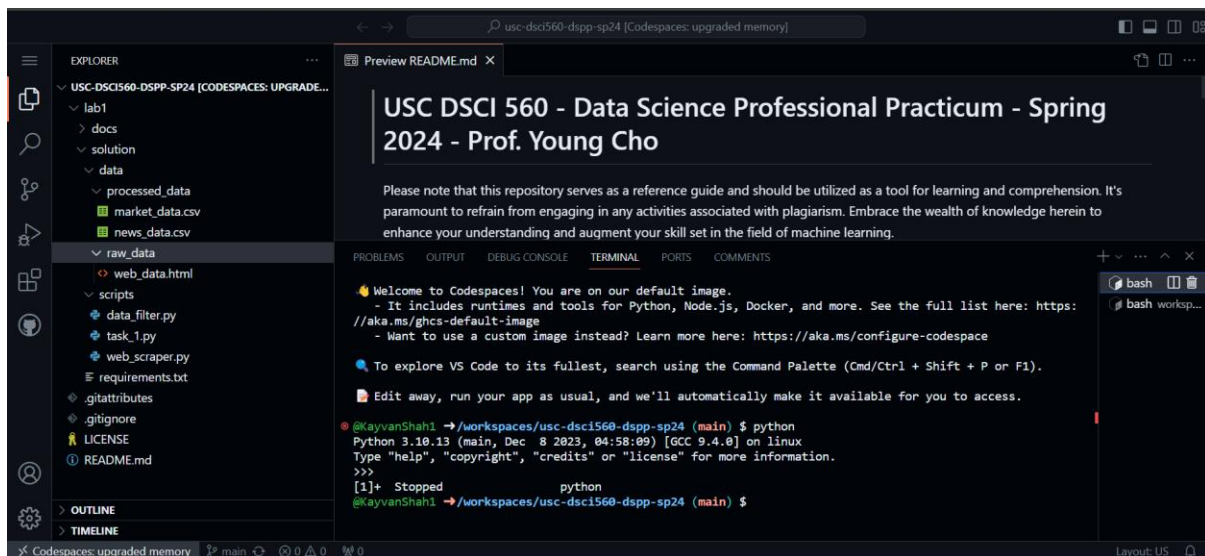# Lab 1 Solution

## Installation & Setup

- Development done locally in a local git repository on a windows machine.
- The code is synced with a remote git repository.



- A Linux Ubuntu instance is spin up using GitHub Codespace for the repository.

# Get Familiar with Linux and Python

## Playing around with Linux Terminal





## A basic Python Script



## Web Scraping Task

- Create a virtual environment and install the dependencies using the ***requirements.txt***.
    - *venv* – Created a virtual environment using command $python-m\ venv\ venv$
    - Linux has a built-in support for Chrome and using a library **chromedriver-py** we automatically manage the webdriver downloads and retrieval of the executable path while initializing the web driver.



Kayvan Shah                                                                USC ID: 110665085
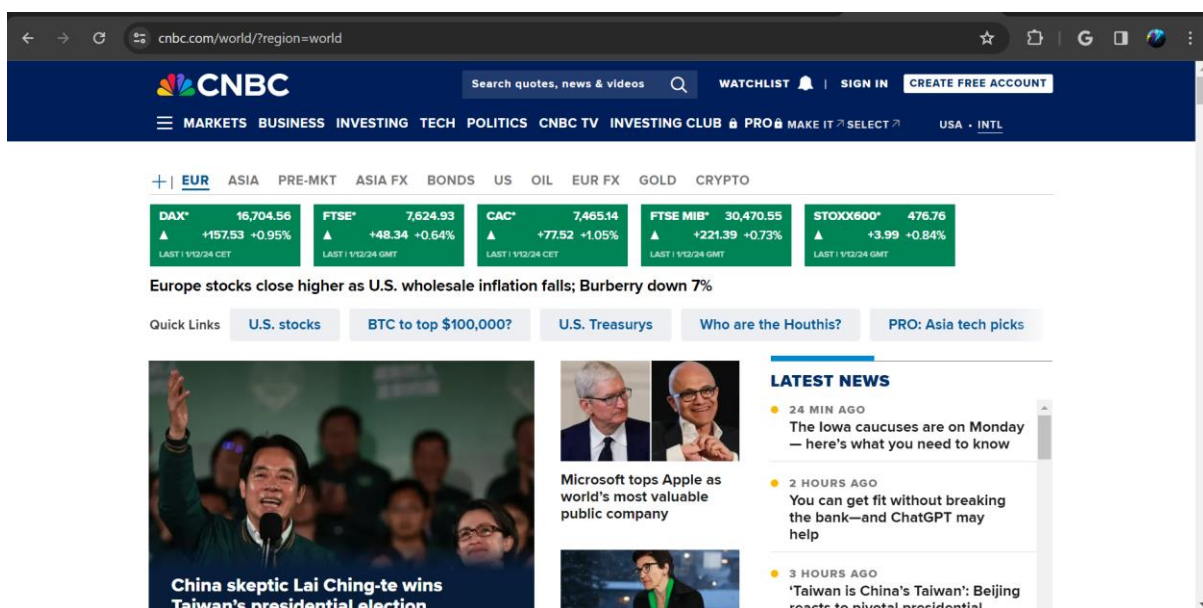
```
Installing collected packages: sortedcontainers, pytz, chromedriver-py, urllib3, tzdata, typing-extensions, soupsieve, sn
iffio, six, python-dotenv, pysocks, numpy, lxml, idna, h11, exceptiongroup, charset-normalizer, certifi, attrs, annotated
-types, wsproto, requests, python-dateutil, pydantic-core, outcome, beautifulsoup4, trio, pydantic, pandas, trio-websocke
t, pydantic-settings, selenium
Successfully installed annotated-types-0.6.0 attrs-23.2.0 beautifulsoup4-4.12.2 certifi-2023.11.17 charset-normalizer-3.3
.2 chromedriver-py-120.0.6099.109 exceptiongroup-1.2.0 h11-0.14.0 idna-3.6 lxml-5.1.0 numpy-1.26.3 outcome-1.3.0.post0 pa
ndas-2.1.4 pydantic-2.5.3 pydantic-core-2.14.6 pydantic-settings-2.1.0 pysocks-1.7.1 python-dateutil-2.8.2 python-dotenv-
1.0.0 pytz-2023.3.post1 requests-2.31.0 selenium-4.16.0 six-1.16.0 sniffio-1.3.0 sortedcontainers-2.4.0 soupsieve-2.5 tri
o-0.24.0 trio-websocket-0.11.1 typing-extensions-4.9.0 tzdata-2023.4 urllib3-2.1.0 wsproto-1.2.0

[notice] A new release of pip is available: 23.0.1 -> 23.3.2
[notice] To update, run: pip install --upgrade pip
(venv) @KayvanShah1 →/workspaces/usc-dsci560-dspp-sp24/lab1/solution (main) $ pip install -U pip
Requirement already satisfied: pip in /workspaces/usc-dsci560-dspp-sp24/venv/lib/python3.10/site-packages (23.0.1)
Collecting pip
  Downloading pip-23.3.2-py3-none-any.whl (2.1 MB)
                                    2.1/2.1 MB 6.3 MB/s eta 0:00:00
Installing collected packages: pip
  Attempting uninstall: pip
    Found existing installation: pip 23.0.1
    Uninstalling pip-23.0.1:
      Successfully uninstalled pip-23.0.1
Successfully installed pip-23.3.2
```

- A snapshot of the webpage from where the data is to be scraped.
  - The webpage loads a few fields dynamically. In this case the market banner is loaded using a JavaScript.
  - To capture the data from dynamically populated HTML elements we must use selenium with a web driver suited for the installed browser on the system.
  - While the latest news isn't loaded dynamically hence, we can get the static snapshot of the HTML tree of the scrape that its section's data.



- Scraping the raw data
  - This Python script is designed to scrape data from the CNBC website related to global markets. It begins by configuring logging for informative messages and defining file and directory paths.
  - The script utilizes the `pydantic_settings` module to manage configuration settings, with a default base URL for CNBC. It then initializes a headless Chrome WebDriver using `chromedriver_py` and configurable options.
  - The main functionality involves loading the CNBC webpage, waiting for the visibility of market card rows, and extracting HTML content for the market banner and the latest news panel using BeautifulSoup.

- o The extracted HTML is saved to a file in the "raw_data" directory. As part of error handling, the script logs any exceptions that might occur during data fetching.
- o Finally, it prints the first ten lines of the saved HTML file to the standard output for quick verification.
- o Overall, this script demonstrates the use of web scraping tools like Selenium and BeautifulSoup, with logging providing insights into the different stages of the data retrieval process.

```
source /workspaces/usc-dsci560-dspp-sp24/venv/bin/activate
/workspaces/usc-dsci560-dspp-sp24/venv/bin/python /workspaces/usc-dsci560-dspp-sp24/lab1/solution/scripts/web_scraper.py
@KayvanShah1 ➜/workspaces/usc-dsci560-dspp-sp24 (main) $ source /workspaces/usc-dsci560-dspp-sp24/venv/bin/activate
(venv) @KayvanShah1 ➜/workspaces/usc-dsci560-dspp-sp24 (main) $ /workspaces/usc-dsci560-dspp-sp24/venv/bin/python /workspaces/usc-dsci560-dspp-sp
24/lab1/solution/scripts/web_scraper.py
INFO:root:Initializing the Chrome WebDriver...
INFO:root:Getting the CNBC Web Page...
INFO:root:Waiting for the Market Cards rows to be populated...
INFO:root:Extracting the market banner HTML tags...
INFO:root:Extracting the latest news panel...
INFO:root:Saving HTML page...
INFO:root:Successfully saved the HTML file to /workspaces/usc-dsci560-dspp-sp24/lab1/solution/data/raw_data/web_data.html
INFO:root:Printing the first ten lines of HTML file to stdout
<div class="MarketsBanner-marketData" id="market-data-scroll-container">
<a class="MarketCard-container MarketCard-up MarketCard-wrap" href="//www.cnbc.com/quotes/.GDAXI">
<div class="MarketCard-row">
<span class="MarketCard-symbol">
DAX*
</span>
<span class="MarketCard-stockPosition">
16,704.56
</span>
</div>
```

- Filtering the data to extract the fields of interest.

  This Python script performs the following tasks:
  1. Logging Setup: Configures logging to display INFO level messages.
  2. Path and Directory Definitions: Defines the file and directory paths using the `Path` class.
  3. HTML Parsing and Data Extraction:
     a. "read_parse_raw_data": Reads and parses an HTML file using BeautifulSoup.
     b. "NewsItem" and "MarketCard" Pydantic models are defined for structured data representation.
  4. Data Filtering Functions:
     a. "filter_latest_news": Extracts the latest news feed from the HTML, converts it into "NewsItem" Pydantic model instances, and appends them to a list.
     b. "filter_market_banner": Filters market banner data from the HTML, converts it into "MarketCard" Pydantic model instances, and appends them to a list.
  5. Main Execution:
     a. Reads and parses an HTML file ("web_data.html").
     b. Filters the latest news and market banner data using the defined functions.
     c. Logs the process at different stages (reading, filtering, and saving).
     d. Saves the filtered data into CSV files ("news_data.csv" and "market_data.csv") in the "processed_data" directory.

  Overall, the script demonstrates the use of BeautifulSoup for HTML parsing, Pydantic for structured data representation, and Pandas for data manipulation. Logging is employed to provide information about the different stages of the process, making it easier to understand and debug the execution flow.

```
(venv) @KayvanShah1 ➜/workspaces/usc-dsci560-dspp-sp24 (main) $ /workspaces/usc-dsci560-dspp-sp24/venv/bin/python /workspaces/usc-dsci560-dspp-sp
24/lab1/solution/scripts/data_filter.py
INFO:root:Reading and parsing the raw HTML file...
INFO:root:Filter the latest news feed...
INFO:root:Saving the filtered latest news feed data to CSV...
INFO:root:Filtering the market banner data...
INFO:root:Saving the filtered market data to CSV...
```

## Output Files

- Raw data output HTML file



- Processed data CSV files
  - Market data
    - Has the following fields:
      - symbol – abbreviations for a stock ticker.
      - stock_position – current/latest price of the stock.
      - change_pts – most recent change in price observed.



  - Latest news
    - Has the following fields:
      - timestamp – approximate time from the current time when the news was posted.
      - title – the headline of the news column
      - link – link to the official source of news on the website