

## 2.1. Playing around with Linux Terminal

```
lacuna@lacuna-VMware20-1: ~
Setting up gcc-aarch64-linux-gnu (4:15.2.0-4ubuntu1) ...
Setting up gcc (4:15.2.0-4ubuntu1) ...
Setting up g++-aarch64-linux-gnu (4:15.2.0-4ubuntu1) ...
Setting up g++ (4:15.2.0-4ubuntu1) ...
update-alternatives: using /usr/bin/g++ to provide /usr/bin/c++ (c++) in auto mode
Setting up build-essential (12.12ubuntu1) ...
Processing triggers for man-db (2.13.1-1) ...
Processing triggers for libc-bin (2.42-0ubuntu3) ...
lacuna@lacuna-VMware20-1:~$ python3 --version
Python 3.13.7
lacuna@lacuna-VMware20-1:~$ sudo apt install -y python3-pip
python3-pip is already the newest version (25.1.1+dfsg-1ubuntu2).
The following packages were automatically installed and are no longer required:
  linux-headers-6.17.0-5          linux-tools-6.17.0-5
  linux-headers-6.17.0-5-generic  linux-tools-6.17.0-5-generic
  linux-modules-6.17.0-5-generic
Use 'sudo apt autoremove' to remove them.

Summary:
  Upgrading: 0, Installing: 0, Removing: 0, Not Upgrading: 72
lacuna@lacuna-VMware20-1:~$ pip3 --version
pip 25.1.1 from /usr/lib/python3/dist-packages/pip (python 3.13)
lacuna@lacuna-VMware20-1:~$ 
lacuna@lacuna-VMware20-1:~$ cd ~/Desktop
mkdir -p "ZixiWang_2854187591/data" "ZixiWang_2854187591/scripts"
touch "ZixiWang_2854187591/scripts/task_1.py"
ls -l "ZixiWang_2854187591/scripts"
total 0
-rw-rw-r-- 1 lacuna lacuna 0 Jan 17 09:05 task_1.py
lacuna@lacuna-VMware20-1:~/Desktop$ ls -l "ZixiWang_2854187591"
total 8
drwxrwxr-x 2 lacuna lacuna 4096 Jan 17 09:05 data
drwxrwxr-x 2 lacuna lacuna 4096 Jan 17 09:05 scripts
lacuna@lacuna-VMware20-1:~/Desktop$ 
```

It shows I downloaded everything. I have folder data and scripts, and I created task\_1.py

## 2.2. A basic Python Script

```
(.venv) lacuna@lacuna-VMware20-1:~/Desktop/ZixiWang_2854187591$ nano task_1.py
(.venv) lacuna@lacuna-VMware20-1:~/Desktop/ZixiWang_2854187591$ nano scripts/task_1.py
(.venv) lacuna@lacuna-VMware20-1:~/Desktop/ZixiWang_2854187591$ python3 scripts/task_1.py
Enter your name: Zixi
Hello, Zixi!
```

```
lacuna@lacuna-VMware20-1: ~/Desktop/ZixiWang_2854187591 — nano scripts/task_1.py
~/Desktop/ZixiWang_2854187591
```

```
GNU nano 8.4
name = input("Enter your name: ")
print(f"Hello, {name}!")
```

I have task\_1.py and can output Hello, [name]!

### 2.3. Python Web-scraping Task

```
upgrading: 0, installing: 0, removing: 0, not upgrading: 72
lacuna@lacuna-VMware20-1:~/Desktop/ZixiWang_2854187591$ cd ~/Desktop/ZixiWang_2854187591
python3 -m venv .venv
lacuna@lacuna-VMware20-1:~/Desktop/ZixiWang_2854187591$ source .venv/bin/activate
(.venv) lacuna@lacuna-VMware20-1:~/Desktop/ZixiWang_2854187591$ pip install requests beautifulsoup4
Collecting requests
  Downloading requests-2.32.5-py3-none-any.whl.metadata (4.9 kB)
Collecting beautifulsoup4
  Downloading beautifulsoup4-4.14.3-py3-none-any.whl.metadata (3.8 kB)
Collecting charset_normalizer<4,>=2 (from requests)
  Downloading charset_normalizer-3.4.4-cp313-cp313-manylinux2014_aarch64.manylinux_2_17_aarch64.manylinux_2_28_aarch64.whl.metadata (37 kB)
Collecting idna<4,>=2.5 (from requests)
  Downloading idna-3.11-py3-none-any.whl.metadata (8.4 kB)
Collecting urllib3<3,>=1.21.1 (from requests)
  Downloading urllib3-2.6.3-py3-none-any.whl.metadata (6.9 kB)
Collecting certifi=>2017.4.17 (from requests)
  Downloading certifi-2026.1.4-py3-none-any.whl.metadata (2.5 kB)
Collecting soupsieve=>1.6.1 (from beautifulsoup4)
  Downloading soupsieve-2.8.1-py3-none-any.whl.metadata (4.6 kB)
Collecting typing_extensions=>=4.0.0 (from beautifulsoup4)
  Downloading typing_extensions-4.15.0-py3-none-any.whl.metadata (3.3 kB)
Downloading requests-2.32.5-py3-none-any.whl (64 kB)
Downloading charset_normalizer-3.4.4-cp313-cp313-manylinux2014_aarch64.manylinux_2_17_aarch64.manylinux_2_28_aarch64.whl (147 kB)
Downloading idna-3.11-py3-none-any.whl (71 kB)
Downloading urllib3-2.6.3-py3-none-any.whl (131 kB)
Downloading beautifulsoup4-4.14.3-py3-none-any.whl (107 kB)
Downloading certifi-2026.1.4-py3-none-any.whl (152 kB)
Downloading soupsieve-2.8.1-py3-none-any.whl (36 kB)
Downloading typing_extensions-4.15.0-py3-none-any.whl (44 kB)
Installing collected packages: urllib3, typing_extensions, soupsieve, idna, charset_normalizer, certifi, requests, beautifulsoup4
Successfully installed beautifulsoup4-4.14.3 certifi-2026.1.4 charset_normalizer-3.4.4 idna-3.11 requests-2.32.5 soupsieve-2.8.1 typing_extensions-5.0 urllib3-2.6.3
(.venv) lacuna@lacuna-VMware20-1:~/Desktop/ZixiWang_2854187591$ python -c "import requests; import bs4; print('ok')"
ok
(.venv) lacuna@lacuna-VMware20-1:~/Desktop/ZixiWang_2854187591$
```

```
current directory: /home/lacuna/Desktop/ZixiWang_2854187591
(.venv) lacuna@lacuna-VMware20-1:~/Desktop/ZixiWang_2854187591$ ls
data  lab1_zixiwang_2854187591.zip  scripts
(.venv) lacuna@lacuna-VMware20-1:~/Desktop/ZixiWang_2854187591$ ls data
processed_data  raw_data
(.venv) lacuna@lacuna-VMware20-1:~/Desktop/ZixiWang_2854187591$
```

```
(.venv) lacuna@lacuna-VMware20-1:~/Desktop/ZixiWang_2854187591$ nano scripts/web_scraper.py
(.venv) lacuna@lacuna-VMware20-1:~/Desktop/ZixiWang_2854187591$ python3 scripts/web_scraper.py
head -n 10 data/raw_data/web_data.html
grep -n "MarketCard-container" -m 2 data/raw_data/web_data.html
grep -n "LatestNews-list" -m 1 data/raw_data/web_data.html
Fetching: https://www.cnbc.com/world/?region=world
Saved: data/raw_data/web_data.html
<div class="MarketCard-container">
<a class="MarketCard-container" href="#">
<div class="MarketCard-row">
<span class="MarketCard-symbol">
.DJI
</span>
</div>
<div class="MarketCard-row">
<span class="MarketCard-stockPosition">
49359.33
2: <a class="MarketCard-container" href="#">
19: <a class="MarketCard-container" href="#">
72:<ul class="LatestNews-list">
```

⌚

Sat 17 22:12

```
GNU nano 8.4
import os
import re
import json
import requests
from bs4 import BeautifulSoup

URL = "https://www.cnbc.com/world/?region=world"
OUT = "data/raw_data/web_data.html"
SYMBOLS = [".DJI", ".SPX", ".IXIC", ".VIX"]

def get_page():
    headers = {"User-Agent": "Mozilla/5.0"}
    r = requests.get(URL, headers=headers, timeout=20)
    r.raise_for_status()
    return r.text

def get_quotes(symbols):
    joined = "|".join(symbols)
    api = f"https://quote.cnbc.com/quote-html-webservice/quote.json?symbols={joined}&output=json"
    r = requests.get(api, headers={"User-Agent": "Mozilla/5.0"}, timeout=20)
    r.raise_for_status()
    txt = r.text.strip()

    try:
        return json.loads(txt)
    except Exception:
        m = re.search(r"(\{.*\})", txt, flags=re.DOTALL)
        return json.loads(m.group(1)) if m else {}

def build_market_html(quotes_json):
    soup = BeautifulSoup("", "html.parser")
    banner = soup.new_tag("div", {"class": "MarketsBanner-marketData"})

    q = quotes_json.get("QuickQuoteResult", {}).get("QuickQuote", [])
    mp = {}
    if isinstance(q, list):
        for it in q:
            syn = str(it.get("symbol", ""))
            mp[syn] = (it.get("last", ""), it.get("change_pct", ""))
    else:
        syn = str(q.get("symbol", ""))
        mp[syn] = (q.get("last", ""), q.get("change_pct", ""))

    for sym in SYMBOLS:
        last, chg = mp.get(sym, ("", ""))
        a = soup.new_tag("a", {"class": "MarketCard-container", "href": "#"})
        a.append(soup.new_tag("div", {"class": "MarketCard-row"}))
        a.append(soup.new_tag("span", {"class": "MarketCard-symbol"}))
        a.append(soup.new_tag("span", {"class": "MarketCard-stockPosition"}))
        a.append(soup.new_tag("div", {"class": "MarketCard-row"}))
        a.append(soup.new_tag("div", {"class": "MarketCard-row"}))

        a.append(soup.new_tag("div", {"class": "MarketCard-row"}))
        a.append(soup.new_tag("span", {"class": "MarketCard-symbol"}))
        a.append(soup.new_tag("span", {"class": "MarketCard-stockPosition"}))
        a.append(soup.new_tag("div", {"class": "MarketCard-row"}))
```

```

        mp[sym] = (ctt.get('last', None), ctt.get('change_pct', None))

    for sym in SYMBOLS:
        last, chg = mp.get(sym, ('', ''))
        a = soup.new_tag("a", **{"class": "MarketCard-container", "href": "#"})
        a.append(soup.new_tag("div", **[{"class": "MarketCard-row"}]))
        a.contents[-1].append(soup.new_tag("span", **{"class": "MarketCard-symbol"}))
        a.contents[-1].contents[-1].string = sym

        a.append(soup.new_tag("div", **[{"class": "MarketCard-row"})))
        a.contents[-1].append(soup.new_tag("span", **{"class": "MarketCard-stockPosition"}))
        a.contents[-1].contents[-1].string = str(last)

        a.append(soup.new_tag("div", **[{"class": "MarketCard-row"}]))
        a.contents[-1].append(soup.new_tag("span", **{"class": "MarketCard-changePct"}))
        a.contents[-1].contents[-1].string = str(chg)

        banner.append(a)

    return banner

def main():
    os.makedirs("data/raw_data", exist_ok=True)
    print("Fetching:", URL)

    html = get_page()
    soup = BeautifulSoup(html, "html.parser")
    latest_news = soup.find("ul", class_="LatestNews-list")

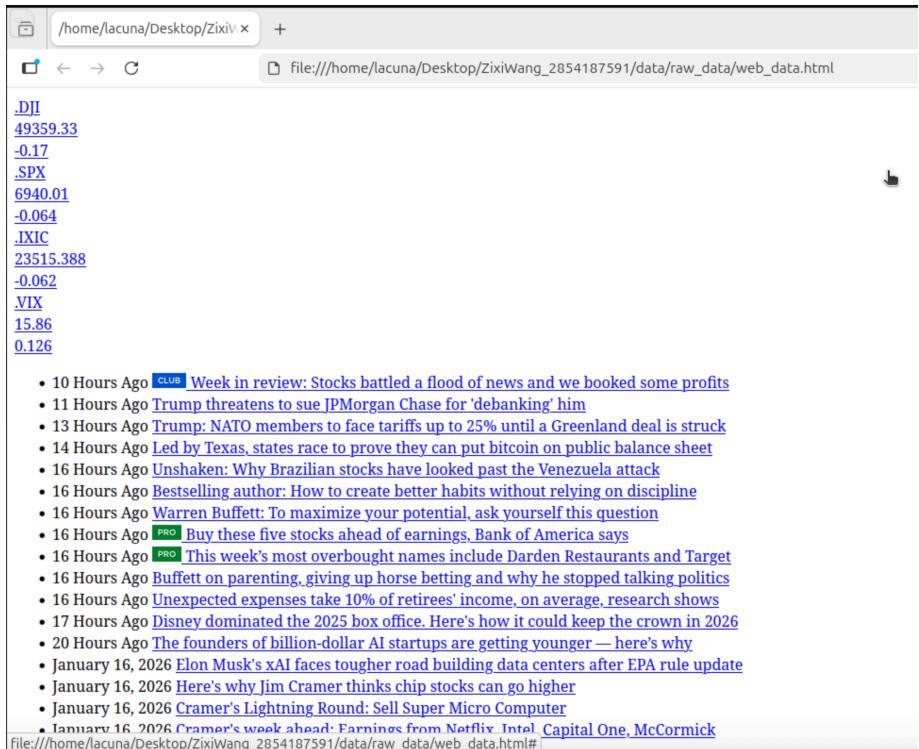
    quotes = get_quotes(SYMBOLS)
    market = build_market_html(quotes)

    with open(OUT, "w", encoding="utf-8") as f:
        f.write(market.prettify() + "\n")
        if latest_news:
            f.write(latest_news.prettify() + "\n")

    print("Saved:", OUT)

if __name__ == "__main__":
    main()

```



I created scripts/web\_scraper.py, installed the required libraries (pip install requests beautifulsoup4), inspected the CNBC World page (<https://www.cnbc.com/world/?region=world>) to identify the Market banner and the “Latest News” section tags, and created the folders data/raw\_data and data/processed\_data. Then I wrote a script using Requests + BeautifulSoup to collect the page content and save the relevant HTML to data/raw\_data/web\_data.html. Finally, I verified the result in the terminal by printing the first 10 lines of the file (head -n 10 data/raw\_data/web\_data.html) and checking for key elements using grep (e.g., MarketCard-container and LatestNews-list).

## 2.4. Data Filtering Task

```
(.venv) lacuna@lacuna-VMware20-1:~/Desktop/ZixiWang_2854187591$ nano scripts/data_filter.py
(.venv) lacuna@lacuna-VMware20-1:~/Desktop/ZixiWang_2854187591$ python3 scripts/data_filter.py
head -n 5 data/processed_data/market_data.csv
head -n 5 data/processed_data/news_data.csv
Filtering fields...
Storing Market data...
CSV created: data/processed_data/market_data.csv (rows=4)
Storing News data...
CSV created: data/processed_data/news_data.csv (rows=30)
Done.
marketCard_symbol,marketCard_stockPosition,marketCardchangePct
.DJI,49359.33,-0.17%
.SPX,6940.01,-0.064%
.IXIC,23515.388,-0.062%
.VIX,15.86,0.126%
10 Hours Ago,Week in review: Stocks battled a flood of news and we booked some profits,https://www.cnbc.com/2026/01/17/week-in-review-stocks-battled-a-flood-of-news-and-we-booked-some-profits.html
11 Hours Ago,Trump threatens to sue JPMorgan Chase for 'debanking' him,https://www.cnbc.com/2026/01/17/trump-jpmorgan-chase-debanking.html
13 Hours Ago,Trump: NATO members to face tariffs up to 25% until a Greenland deal is struck,https://www.cnbc.com/2026/01/17/trump-greenland-tariffs-nato.html
14 Hours Ago,"Led by Texas, states race to prove they can put bitcoin on public balance sheet",https://www.cnbc.com/2026/01/17/texas-us-states-budgets-bitcoin-crypto-strategic-reserve.html
(.venv) lacuna@lacuna-VMware20-1:~/Desktop/ZixiWang_2854187591$ python3 scripts/data_filter.py
head -n 2 data/processed_data/market_data.csv
Filtering fields...
Storing Market data...
CSV created: data/processed_data/market_data.csv (rows=4)
Storing News data...
CSV created: data/processed_data/news_data.csv (rows=30)
Done.
marketCard_symbol,marketCard_stockPosition,marketCardchangePct
.DJI,49359.33,-0.17%
```

```
GNU nano 8.4
#!/usr/bin/env python3
import os
import csv
from bs4 import BeautifulSoup

INFILE = "data/raw_data/web_data.html"
MARKET_CSV = "data/processed_data/market_data.csv"
NEWS_CSV = "data/processed_data/news_data.csv"

os.makedirs("data/processed_data", exist_ok=True)

print("Filtering Fields...")

with open(INFILE, "r", encoding="utf-8", errors="ignore") as f:
    soup = BeautifulSoup(f.read(), "html.parser")
    market = []
    for card in soup.select(".MarketCard-container"):
        sym = card.select_one("span.MarketCard-symbol")
        pos = card.select_one("span.MarketCard-stockPosition")
        pct = (card.select_one("span.MarketCard-changesPct")
               or card.select_one("span.MarketCard-changePct"))

        market.append([
            sym.get_text(strip=True) if sym else '',
            pos.get_text(strip=True) if pos else '',
            pct.get_text(strip=True) if pct else '',
        ])

    market.append([
        sym.get_text(strip=True) if sym else '',
        pos.get_text(strip=True) if pos else '',
        pct.get_text(strip=True) if pct else '',
    ])

print("Storing Market data...")
with open(MARKET_CSV, "w", newline="", encoding="utf-8") as f:
    w = csv.writer(f)
    w.writerow(["marketCard_symbol", "marketCard_stockPosition", "marketCardchangePct"])
    w.writerows(market)
print(f"CSV created: {MARKET_CSV} (rows={len(market)}")

news = []
for item in soup.select("ul.LatestNews-list > li.LatestNews-item"):
    t = item.select_one("time.LatestNews-timestamp")
    a = item.select_one("a.LatestNews-headline")

    ts = t.get_text(strip=True) if t else ""
    title = a.get_text(strip=True) if a else ""
    link = (a.get("href") or "").strip() if a else ""

    news.append([ts, title, link])

print("Storing News data...")
with open(NEWS_CSV, "w", newline="", encoding="utf-8") as f:
    w = csv.writer(f)
    w.writerow(["LatestNews_timestamp", "title", "link"])
    w.writerows(news)
print(f"CSV created: {NEWS_CSV} (rows={len(news)})")
print("Done.")
```

```
Open ▾ + market_data.csv
~/Desktop/ZixiWang_2854187591/data/processed_data
market_data.csv x news_data.csv chromedriver.log

marketCard_symbol,marketCard_stockPosition,marketCardchangePct
.DJI,49359.33,-0.17
.SPX,6940.01,-0.064
.IXIC,23515.388,-0.062
.VIX,15.86,0.126
```

```
Open ▾ + news_data.csv
~/Desktop/ZixiWang_2854187591/data/processed_data
market_data.csv news_data.csv x chromedriver.log

LatestNews_timestamp,title,link
10 Hours Ago,Week in review: Stocks battled a flood of news and we booked some profits,https://www.cnbc.com/2026/01/17/week-in-review-stocks-battled-a-flood-of-news-and-we-booked-some-profits.html
11 Hours Ago,Trump threatens to sue JPMorgan Chase for 'debanking' him,https://www.cnbc.com/2026/01/17/trump-jpmorgan-chase-debanking.html
13 Hours Ago,Trump: NATO members to face tariffs up to 25% until a Greenland deal is struck,https://www.cnbc.com/2026/01/17/trump-greenland-tariffs-nato.html
14 Hours Ago,"Led by Texas, states race to prove they can put bitcoin on public balance sheet",https://www.cnbc.com/2026/01/17/texas-us-states-budgets-bitcoin-crypto-strategic-reserve.html
16 Hours Ago,Unshaken: Why Brazilian stocks have looked past the Venezuela attack,https://www.cnbc.com/2026/01/17/unshaken-why-brazilian-stocks-have-looked-past-the-venezuela-attack.html
16 Hours Ago,Bestselling author: How to create better habits without relying on discipline,https://www.cnbc.com/2026/01/17/james-clear-how-to-create-better-habits-without-relying-on-discipline.html
16 Hours Ago,"Warren Buffett: To maximize your potential, ask yourself this question",https://www.cnbc.com/2026/01/17/warren-buffett-to-maximize-your-potential-ask-yourself-this-question.html
```

I created scripts/data\_filter.py to read data/raw\_data/web\_data.html into a Python list (readlines()), parse it with BeautifulSoup, and extract (marketCard\_symbol, marketCard\_stockPosition, marketCardchangePct) from the Market banner and (LatestNews-timestamp, title, link) from each Latest News entry. The script writes the extracted market data to data/processed\_data/market\_data.csv and the news data to

data/processed\_data/news\_data.csv, while printing progress messages such as “Filtering fields...”, “Storing Market data...”, and “CSV created...”. I confirmed the outputs by previewing the CSV files in the terminal using head -n 5 data/processed\_data/market\_data.csv and head -n 5 data/processed\_data/news\_data.csv.

my githublink: <https://github.com/lacunaxu/DSCI560-Lab1>