

# Web Scraping and Data Extraction Report

**Author:** Rafayel

**Environment:** Ubuntu Linux

**Language:** Python 3

This report documents the complete workflow of a web scraping project, including environment setup, script execution, data extraction, and validation. Screenshots are included to demonstrate successful execution.

## 1. Project Directory Setup

This screenshot shows the creation of the project directory and subfolders using Linux terminal commands. It confirms proper navigation, folder creation, and file initialization following best practices.

```
ubuntu@ubuntu:~$ cd Desktop
ubuntu@ubuntu:~/Desktop$ ls
ubuntu-desktop-bootstrap_ubuntu-desktop-bootstrap.desktop
ubuntu@ubuntu:~/Desktop$ mkdir Rafayel_3487192016
ubuntu@ubuntu:~/Desktop$ ls
Rafayel_3487192016
ubuntu-desktop-bootstrap_ubuntu-desktop-bootstrap.desktop
ubuntu@ubuntu:~/Desktop$ cd Rafayel_3487192016/
ubuntu@ubuntu:~/Desktop/Rafayel_3487192016$ mkdir data
ubuntu@ubuntu:~/Desktop/Rafayel_3487192016$ mkdir scripts
ubuntu@ubuntu:~/Desktop/Rafayel_3487192016$ cd scripts/
ubuntu@ubuntu:~/Desktop/Rafayel_3487192016/scripts$ touch task_1.py
ubuntu@ubuntu:~/Desktop/Rafayel_3487192016/scripts$ ls
task_1.py
```

## 2. Python Script Creation and Execution

This screenshot demonstrates editing and running a Python script (task\_1.py). The successful output confirms that the Python environment is correctly configured.

```
ubuntu@ubuntu:~/Desktop/Rafayel_3487192016/scripts$ nano task_1.py
ubuntu@ubuntu:~/Desktop/Rafayel_3487192016/scripts$ python3 task_1.py
Please enter your nameRafayel
Hello, Rafayel!
ubuntu@ubuntu:~/Desktop/Rafayel_3487192016/scripts$
```

### 3. Web Scraper Execution and HTML Storage

This screenshot shows the execution of the web scraper script. A successful HTTP request is made and raw HTML data is saved locally for further processing.



## 4. HTML Content Verification

This screenshot verifies that the downloaded HTML file contains valid webpage content. The presence of metadata and page structure confirms successful data retrieval.

```
(.venv) ubuntu@ubuntu:~/Desktop/Rafayel_3487192016/scripts$ pip install requests
Collecting requests
  Downloading requests-2.32.5-py3-none-any.whl.metadata (4.9 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)
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 certifi-2026.1.4-py3-none-any.whl (152 kB)
Installing collected packages: urllib3, idna, charset_normalizer, certifi, requests
Successfully installed certifi-2026.1.4 charset_normalizer-3.4.4 idna-3.11 requests-2.32.5 urllib3-2.6.3
(.venv) ubuntu@ubuntu:~/Desktop/Rafayel_3487192016/scripts$ pip install beautifulsoup4
Collecting beautifulsoup4
  Downloading beautifulsoup4-4.14.3-py3-none-any.whl.metadata (3.8 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 beautifulsoup4-4.14.3-py3-none-any.whl (107 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: typing-extensions, soupsieve, beautifulsoup4
Successfully installed beautifulsoup4-4.14.3 soupsieve-2.8.1 typing_extensions-4.15.0
```

## 5. Data Extraction and CSV Generation

This screenshot shows the parsing process using BeautifulSoup and the creation of structured CSV files. Both market and news data are extracted successfully.

```
(.venv) ubuntu@ubuntu:~/Desktop/Rafayel_3487192016/scripts$ python3 web_scraper.py
Request successful
HTML saved to ../data/raw_data/web_data.html
(.venv) ubuntu@ubuntu:~/Desktop/Rafayel_3487192016/scripts$ head -n 10 ../data/raw_data/web_data.html
<!DOCTYPE html>
<html itemscope="" itemtype="https://schema.org/WebPage" lang="en" prefix="og:https://ogp.me/ns#">
  <head>
    <meta content="website" property="og:type"/>
    <meta content="International: Top News And Analysis" property="og:title"/>
    <meta content="CNBC International is the world leader for news on business, technology, China, trade, oil prices, the Middle East and markets." property="og:description"/>
    <meta content="https://www.cnbc.com/world/" property="og:url"/>
    <meta content="CNBC" property="og:site_name"/>
    <meta content="max-image-preview:large" name="robots"/>
    <meta content="telephone=no" name="format-detection"/>
(.venv) ubuntu@ubuntu:~/Desktop/Rafayel_3487192016/scripts$ █
```

## 6. News Data Validation

This screenshot displays the news\_data.csv file loaded into pandas. It confirms accurate extraction of timestamps, headlines, and links.

```
print(f"Storing data into {filepath}...")
with open(filepath, "w", newline="", encoding="utf-8") as f:
    writer = csv.DictWriter(f, fieldnames=data[0].keys())
    writer.writeheader()
    writer.writerows(data)

print("CSV created successfully")

def main():
    html = read_html_file(HTML_PATH)
    soup = BeautifulSoup(html, "html.parser")

    market_data = extract_market_data(soup)
    news_data = extract_latest_news(soup)

    write_csv(MARKET_CSV_PATH, market_data)
    write_csv(NEWS_CSV_PATH, news_data)

    print("Data filtering and storage completed.")

if __name__ == "__main__":
    main()

...
... Reading raw HTML data...
Filtering market banner fields...
Extracted 5 market records
Filtering latest news fields...
Extracted 30 news records
Storing data into market_data.csv...
CSV created successfully
Storing data into news_data.csv...
CSV created successfully
Data filtering and storage completed.
```

## 7. Market Data Validation

This screenshot shows market\_data.csv loaded into pandas. The table confirms correct extraction of market indices, values, and percentage changes.

LatestNews_timestamp	title	link
10 Min Ago	Who will be next to implement an Australia-style under-16s social media ban?	<a href="https://www.cnbc.com/2026/01/18/uk-australia-style-under-16s-social-media-ban">https://www.cnbc.com/2026/01/18/uk-australia-style-under-16s-social-media-ban</a>
12 Hours Ago	nan	/investingclub/
13 Hours Ago	Trump threatens to sue JPMorgan Chase for 'debanking' him	<a href="https://www.cnbc.com/2026/01/17/trump-jpmorgan-chase-debanking.html">https://www.cnbc.com/2026/01/17/trump-jpmorgan-chase-debanking.html</a>
15 Hours Ago	Trump: NATO members to face tariffs up to 25% until a Greenland deal is struck	<a href="https://www.cnbc.com/2026/01/17/trump-greenland-tariffs-nato.html">https://www.cnbc.com/2026/01/17/trump-greenland-tariffs-nato.html</a>
16 Hours Ago	Led by Texas, states race to prove they can put bitcoin on public balance sheet	<a href="https://www.cnbc.com/2026/01/17/texas-us-states-budgets-bitcoin-cryptocurrency.html">https://www.cnbc.com/2026/01/17/texas-us-states-budgets-bitcoin-cryptocurrency.html</a>
17 Hours Ago	Unshaken: Why Brazilian stocks have looked past the Venezuela attack	<a href="https://www.cnbc.com/2026/01/17/unshaken-why-brazilian-stocks-have-looked-past-the-venezuela-attack.html">https://www.cnbc.com/2026/01/17/unshaken-why-brazilian-stocks-have-looked-past-the-venezuela-attack.html</a>
17 Hours Ago	Bestselling author: How to create better habits without relying on discipline	<a href="https://www.cnbc.com/2026/01/17/james-clear-how-to-create-better-habits-without-relying-on-discipline.html">https://www.cnbc.com/2026/01/17/james-clear-how-to-create-better-habits-without-relying-on-discipline.html</a>
18 Hours Ago	Warren Buffett: To maximize your potential, ask yourself this question	<a href="https://www.cnbc.com/2026/01/17/warren-buffett-to-maximize-your-potential-ask-yourself-this-question.html">https://www.cnbc.com/2026/01/17/warren-buffett-to-maximize-your-potential-ask-yourself-this-question.html</a>
18 Hours Ago	nan	/pro/
18 Hours Ago	Buffett on parenting, giving up horse betting and why he stopped talking politics	<a href="https://www.cnbc.com/2026/01/17/warren-buffett-on-parenting-horse-betting-and-why-he-stopped-talking-politics.html">https://www.cnbc.com/2026/01/17/warren-buffett-on-parenting-horse-betting-and-why-he-stopped-talking-politics.html</a>
18 Hours Ago	Unexpected expenses take 10% of retirees' income, on average, research shows	<a href="https://www.cnbc.com/2026/01/17/unexpected-expenses-take-10-of-retirees-income-on-average-research-shows.html">https://www.cnbc.com/2026/01/17/unexpected-expenses-take-10-of-retirees-income-on-average-research-shows.html</a>
19 Hours Ago	Disney dominated the 2025 box office. Here's how it could keep the crown in 2026	<a href="https://www.cnbc.com/2026/01/17/disney-dominated-2025-box-office-how-it-could-keep-the-crown-in-2026.html">https://www.cnbc.com/2026/01/17/disney-dominated-2025-box-office-how-it-could-keep-the-crown-in-2026.html</a>
22 Hours Ago	The founders of billion-dollar AI startups are getting younger – here's why	<a href="https://www.cnbc.com/2026/01/17/billion-dollar-ai-startup-founders-are-getting-younger-here-s-why.html">https://www.cnbc.com/2026/01/17/billion-dollar-ai-startup-founders-are-getting-younger-here-s-why.html</a>
January 16, 2026	Elon Musk's xAI faces tougher road building data centers after EPA rule update	<a href="https://www.cnbc.com/2026/01/16/elon-musk-s-xai-faces-tougher-road-building-data-centers-after-epa-rule-update.html">https://www.cnbc.com/2026/01/16/elon-musk-s-xai-faces-tougher-road-building-data-centers-after-epa-rule-update.html</a>
January 16, 2026	Here's why Jim Cramer thinks chip stocks can go higher	<a href="https://www.cnbc.com/2026/01/16/heres-why-jim-cramer-thinks-chip-stocks-can-go-higher.html">https://www.cnbc.com/2026/01/16/heres-why-jim-cramer-thinks-chip-stocks-can-go-higher.html</a>
January 16, 2026	Cramer's Lightning Round: Sell Super Micro Computer	<a href="https://www.cnbc.com/2026/01/16/cramers-lightning-round-sell-super-micro-computer.html">https://www.cnbc.com/2026/01/16/cramers-lightning-round-sell-super-micro-computer.html</a>
January 16, 2026	Cramer's week ahead: Earnings from Netflix, Intel, Capital One, McCormick	<a href="https://www.cnbc.com/2026/01/16/cramers-week-ahead-earnings-from-netflix-intel-capital-one-mccormick.html">https://www.cnbc.com/2026/01/16/cramers-week-ahead-earnings-from-netflix-intel-capital-one-mccormick.html</a>
January 16, 2026	Google files to appeal search monopoly case	<a href="https://www.cnbc.com/2026/01/16/google-files-to-appeal-search-monopoly-case.html">https://www.cnbc.com/2026/01/16/google-files-to-appeal-search-monopoly-case.html</a>
January 16, 2026	More employers worry about workers' financial well-being, research shows	<a href="https://www.cnbc.com/2026/01/16/more-employers-worry-about-workers-financial-well-being-research-shows.html">https://www.cnbc.com/2026/01/16/more-employers-worry-about-workers-financial-well-being-research-shows.html</a>
January 16, 2026	Republicans want to end the 'marriage penalty' for this childcare tax credit	<a href="https://www.cnbc.com/2026/01/16/republicans-want-to-end-the-marriage-penalty-for-this-childcare-tax-credit.html">https://www.cnbc.com/2026/01/16/republicans-want-to-end-the-marriage-penalty-for-this-childcare-tax-credit.html</a>
January 16, 2026	Labor Department accused of echoing Nazi slogan in social media post	<a href="https://www.cnbc.com/2026/01/16/labor-department-accused-of-echoing-nazi-slogan-in-social-media-post.html">https://www.cnbc.com/2026/01/16/labor-department-accused-of-echoing-nazi-slogan-in-social-media-post.html</a>
January 16, 2026	Education Department to delay collections on defaulted student loans	<a href="https://www.cnbc.com/2026/01/16/education-department-to-delay-collections-on-defaulted-student-loans.html">https://www.cnbc.com/2026/01/16/education-department-to-delay-collections-on-defaulted-student-loans.html</a>
January 16, 2026	nan	/pro/
January 16, 2026	OpenAI has committed billions to recent chip deals. Some big names have been left out	<a href="https://www.cnbc.com/2026/01/16/openai-chip-deal-with-cerebras-adds-billions.html">https://www.cnbc.com/2026/01/16/openai-chip-deal-with-cerebras-adds-billions.html</a>
January 16, 2026	nan	/investingclub/
January 16, 2026	Hassett pivots to possible 'Trump cards' amid credit card battle with banks	<a href="https://www.cnbc.com/2026/01/16/white-house-hassett-trump-cards-creamy.html">https://www.cnbc.com/2026/01/16/white-house-hassett-trump-cards-creamy.html</a>
January 16, 2026	nan	/investingclub/
January 16, 2026	nan	/pro/
January 16, 2026	Coastal Virginia Offshore Wind to restart work after judge lifts Trump suspension	<a href="https://www.cnbc.com/2026/01/16/biggest-offshore-wind-project-in-us-coastal-virginia.html">https://www.cnbc.com/2026/01/16/biggest-offshore-wind-project-in-us-coastal-virginia.html</a>

## **Conclusion**

This project successfully demonstrates an end-to-end web scraping pipeline using Python. By capturing live web data, storing raw HTML, extracting structured information, and validating outputs, the project highlights practical data engineering and automation skills.