

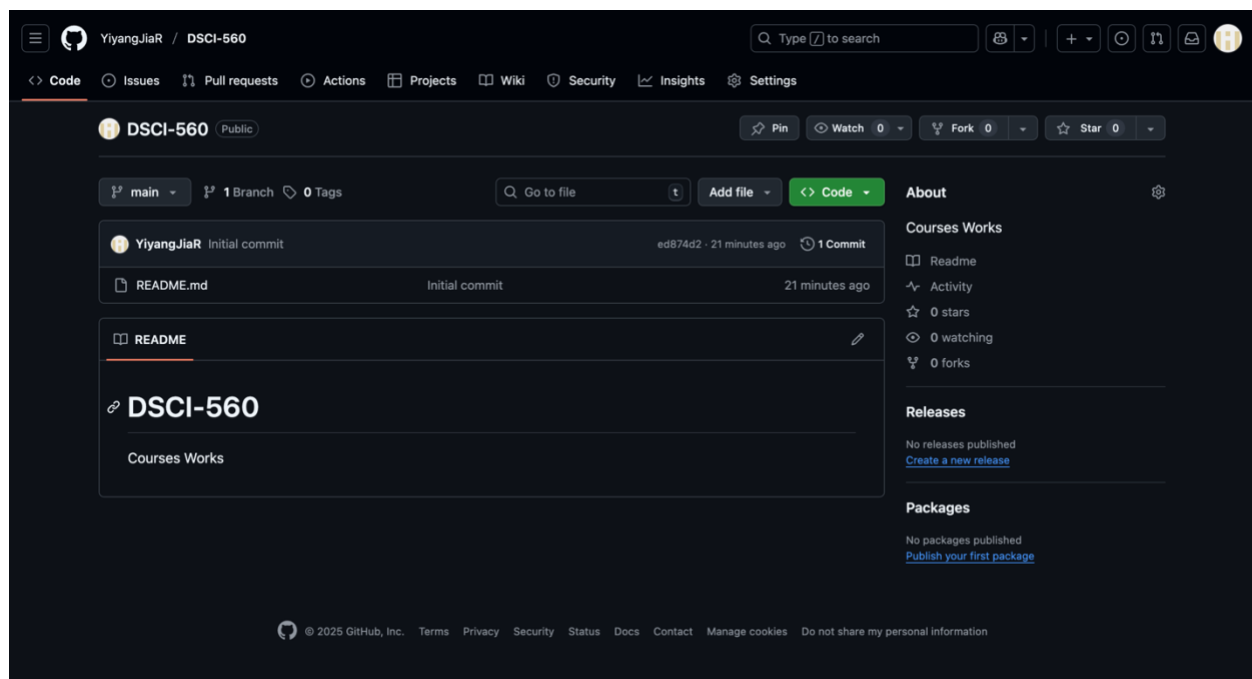
Name: Yiyang Jia, USC ID: 3882757155

Github link: <https://github.com/YiyangJiaR/DSCI-560/tree/main>

## Lab 1 Solution:

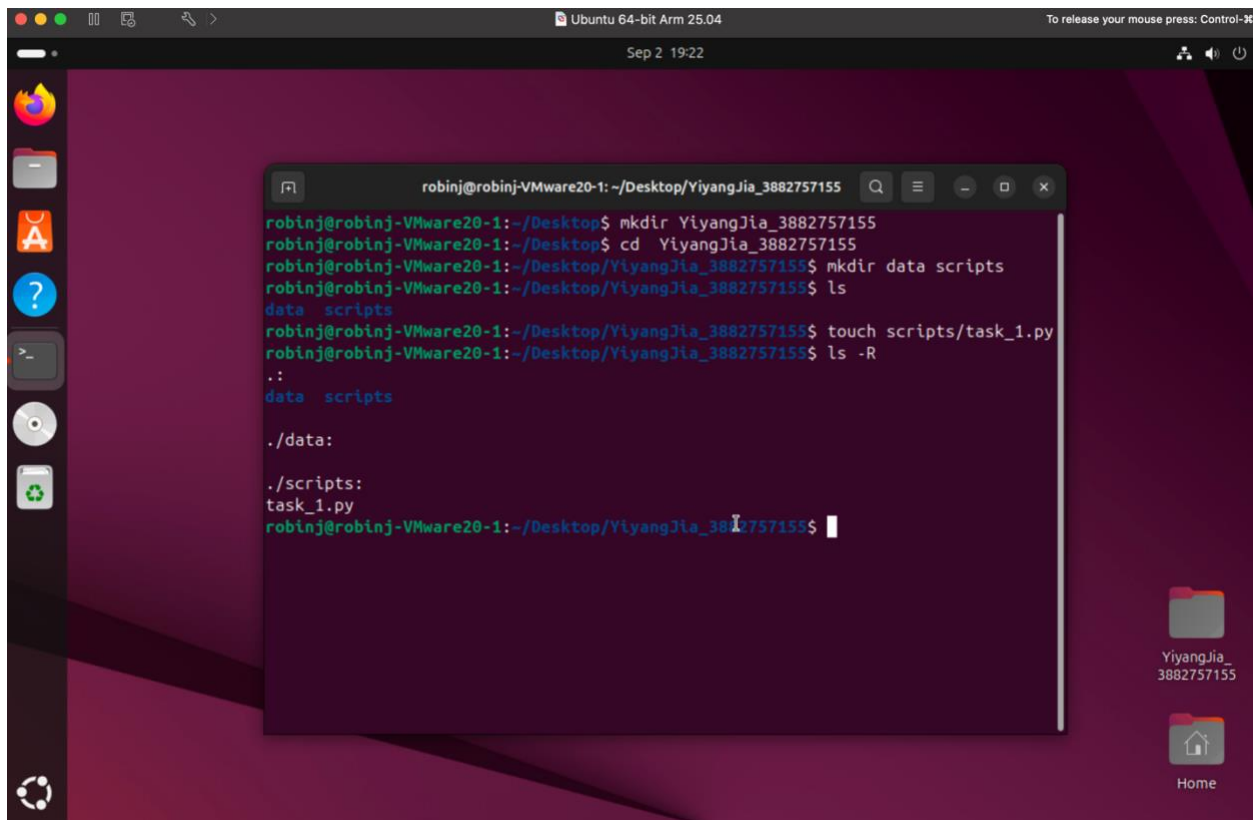
### 1. Installation & Setup

- Using Mac and M chips so download the VMware fusion and Ubuntu ARM version to create the new virtual machine.
- Create new public repo of DSCI 560



### 2. Playing Around with Linux Terminal

- Open the Linux Terminal and write the commands below to achieve the requirements. Just use the mkdir to create new directory and nano to get into the .py file to edit and cat to show what the files got after run it.



### 3. A basic python script

```
robinj@robinj-VMware20-1:~/Desktop/YiyangJia_3882757155$ cd scripts
robinj@robinj-VMware20-1:~/Desktop/YiyangJia_3882757155/scripts$ nano task_1.py
robinj@robinj-VMware20-1:~/Desktop/YiyangJia_3882757155/scripts$ cat task_1.py
name = input("Enter your name: ")
print(f"Hello, {name}!")
robinj@robinj-VMware20-1:~/Desktop/YiyangJia_3882757155/scripts$ python3 task_1.py
Enter your name: Robin
Hello, Robin!
robinj@robinj-VMware20-1:~/Desktop/YiyangJia_3882757155/scripts$
```

### 4. Python Web Scraping Task

I first installed the requests and beautifulsoup4 libraries and created the web\_scraper.py file inside the scripts and the raw\_data and processed\_data inside the web\_scraper.py file when nano inside to edit it. When first write the web\_scraper.py I just use the request.get without headers inside the function, however, it does not work so I need to add the real browser headers then it works.

As requested, fetching CNBC world page HTML and saves it as web\_data.html inside the raw\_data which also inside the data directory. The results show the 10 headlines of the web\_data.html file. You will see the web\_scraper.py in the GitHub repo link.

```
robinj@robinj-VMware20-1:~/Desktop/YiyangJia_3882757155/scripts$ nano web_scraper.py
robinj@robinj-VMware20-1:~/Desktop/YiyangJia_3882757155/scripts$ python3 web_scraper.py
[OK] Saved: /home/robinj/Desktop/YiyangJia_3882757155/data/raw_data/web_data.html
robinj@robinj-VMware20-1:~/Desktop/YiyangJia_3882757155/scripts$ ls
task_1.py  web_scraper.py
robinj@robinj-VMware20-1:~/Desktop/YiyangJia_3882757155/scripts$ head -n 10 data/raw_data/web_data.html
head: cannot open 'data/raw_data/web_data.html' for reading: No such file or directory
robinj@robinj-VMware20-1:~/Desktop/YiyangJia_3882757155/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.cnn.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"/>
```

## 5. Data Filtering Task

```
robinj@robinj-VMware20-1:~/Desktop/YiyangJia_3882757155$ cat data/processed_data/market_data.csv
marketCard_symbol,marketCard_stockPosition,marketCard-changePct
robinj@robinj-VMware20-1:~/Desktop/YiyangJia_3882757155$ 
robinj@robinj-VMware20-1:~/Desktop/YiyangJia_3882757155/scripts$ nano data_filter.py
robinj@robinj-VMware20-1:~/Desktop/YiyangJia_3882757155/scripts$ python3 data_filter.py
[OK] Wrote: /home/robinj/Desktop/YiyangJia_3882757155/data/processed_data/market_data.csv and /home/robinj/Desktop/Yiya
ngJia_3882757155/data/processed_data/news_data.csv
```

I did the data filtering and when get to the market\_data.csv and the news\_data.csv it shows just the headers, and I think it's the html problem and as shown below my web\_data.html only shows the CSS class names but not the data-symbol attributes and I try to redefine the headers, but I does not work out. I do not know how to fix it, but I think the logic of the code is correct and match the lab requirements. I will

try to fix out.

```
robinj@robinj-VMware20-1:~/Desktop/YiyangJia_3882757155$ grep -iE 'data-symbol|MarketCard' data/raw_data/web_data.html | head
.MarketsBanner-container{background-color:#f2f2f2;margin-bottom:15px;padding:15px 20px;position:relative}@media (min-width:760px){.MarketsBanner-container{margin-bottom:15px;padding:15px 40px;top:16px}}@media (min-width:1020px){.MarketsBanner-container{background-color:#fff;display:flex;flex-direction:row;flex-wrap:wrap;margin:0 auto 30px;max-width:960px;padding:15px 0 0;top:16px;width:100%}}@media (min-width:1340px){.MarketsBanner-container{max-width:1290px;padding:15px 0 0}}.MarketsBanner-proMarketsBanner{top:0}@media (min-width:1340px){.MarketsBanner-proMarketsBanner{padding-bottom:8px;padding-top:8px}}@media (max-width:759px){.MarketsBanner-berkshireEvent{margin-bottom:0}}.MarketsBanner-main{border-top:1px solid #9b9b9b80;display:block;min-height:53px}.MarketsBanner-main h2{margin-bottom:0;margin-top:5px}@media (min-width:1020px){.MarketsBanner-main{flex:100%;min-height:59px}}.MarketsBanner-marketData{-webkit-overflow-scrolling:touch;align-items:flex-start;display:flex;flex-direction:row;flex-wrap:nowrap;margin:0 -20px;overflow-x:auto;overflow-y:hidden;padding:10px 0;position:relative;scrollbar-width:none}@media (min-width:760px){.MarketsBanner-marketData{margin-left:
```