# Project Code:
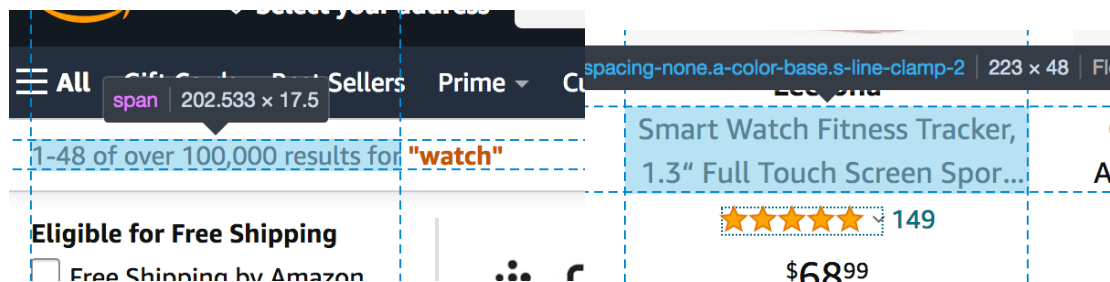
- https://github.com/munjotks/FinalProject-Munjotks.git

# Data Sources:

- User will be searching for an Amazon product through the user interface. In order to retrieve the information from amazon, I will be scraping Amazon using cache → https://www.amazon.com/
- The URL for what page to scrape will be generated from https://www.amazon.com/s?k=(SEARCHTERM)&ref=nb_sb_noss_1  and that specific page will be scraped.

```python
def create_url(searchterm):
    if ' ' in searchterm:
        searchterm = searchterm.replace(' ', '+')
    url = 'https://www.amazon.com/s?k=' + searchterm + '&ref=nb_sb_noss_1'
    return url
```

- The data I will be collecting from the specific pages will be
    - Search Term | Product Name | # of star out of 5 | Product Price | # of Reviews | # of results

```
<span class="a-offscreen">$68.99</span>
  ▼<span aria-hidden="true">
```

```
                                    ▼<i class="a-icon a-icon-star-small a-star-small-5 aok-align-bottom">
                                       <span class="a-icon-alt">5.0 out of 5 stars</span>
                                      </i>
                                      <i class="a-icon a-icon-popover"></i>
                                    </a>
                                   whitespace
                                  </span>
                                </span>
                            ▼<span aria-label="149">
                              ▼<a class="a-link-normal" href="/Fitness-Tracker-Smartwatch-Pedometer-Activ
                                 child=1&keywords=watch&qid=1607391859&sr=8-1#customerReviews">
                                 <span class="a-size-base" dir="auto">149</span>
                              </a>
```

Leefona

Smart Watch Fitness Tracker,

1. span.a-icon-alt | 80 × 18 n Spor...

⭐⭐⭐⭐⭐ ˅ 149

- Caching will be used every time the same search term is used in the user interface (Same URL)



```python
# FinalProject.py ●
# Users > munjotsingh > Documents > SI507-Python2 > FinalProject > FinalProject.py > ...
 6
 7    from bs4 import BeautifulSoup
 8    import requests
 9    import json
10
11    header = {
12        'User-Agent': 'UMSI 507 Course Final Project — Python Scraping',
13        'From': 'Munjotks@umich.edu',
14        'Course-Info': 'https://si.umich.edu/programs/courses/507'
15    }
16
17    def load_cache():
18        try:
19            cache_file = open(CACHE_FILE_NAME, 'r')
20            cache_file_contents = cache_file.read()
21            cache = json.loads(cache_file_contents)
22            cache_file.close()
23        except:
24            cache = {}
25        return cache
26
27    def save_cache(cache):
28        cache_file = open(CACHE_FILE_NAME, 'w')
29        contents_to_write = json.dumps(cache)
30        cache_file.write(contents_to_write)
31        cache_file.close()
32
33    def make_url_request_using_cache(url, cache):
34        if (url in cache.keys()):
35            print("Using Cache")
36            return cache[url]
37        else:
38            print("Fetching")
39            response = requests.get(url, headers=headers)
40            cache[url] = response.text
41            save_cache(cache)
42            return cache[url]
43
44    CACHE_DICT = load_cache()
45
```

# Database:

- I will be creating a database from the information collected from scraping the search term pages. My two tables will consist of the following fields.
    - Product Table
        - Search Term Category | Product Name | # of Stars | # of Reviews
    - Category (search term) Table
        - Search Term Category | # of Results

EXAMPLE:

| Search Term Category | Product Name | Product Price | # of Stars | # of Reviews |
|---|---|---|---|---|
| Camera | Fujifilm Instax Mini 11 Instant Camera - Lilac Purple | 69.00 | 4.8 | 2,528 |
| Camera | Digital Camera, Lecran FHD 1080P 36.0 Mega Pixels Vlogging Camera with 16X Digital Zoom, LCD Screen, Compact Portable | 48.99 | 4.3 | 83 |
| Camera | All-new Blink Outdoor – wireless, weather-resistant HD security camera with two-year battery life and motion detection | 59.99 | 4.3 | 4,619 |
| Camera | Digital Camera, Lecran FHD 1080P 36.0 Mega Pixels Vlogging Camera | 82.98 | 4.2 | 83 |

EXAMPLE:

| Search Term Category | # of Results |
|---|---|
| Camera | 10,000+ |
| Pen | 3,000 |

# Interaction and Presentation Plans

User will be asked:
- What would you like to search on Amazon?
- User Response: Camera
- How would you like the results displayed?

- User Response: top 10 by stars
- Displays top 10 results
- User can request a scatterplot;
    - Product price vs. # of stars
    - Product price vs. # of reviews