# Python Web Scraper

Want to pull the latest working production version (if current version is still under development)?

- ✓ Get Commit: <mark>18d45a5</mark>

Run landing_page.py to initialize the flask application and the web server. <u>(See API Information Below)</u>

**To get the part numbers:** The CSV file with the product number will be read by <mark>**read_csv()**</mark> to have an iter_ set of product numbers, *i_s_part_nums*.

**To search the website:** (in this case newegg.com) the iter_ set of part numbers, *i_s_part_nums*, is stepped through per each individual part number to create a custom url for that product search, *part_num*. This is passed into <mark>**get_custom_url()**</mark> to create a searchable page to scrape from.

**To get data off the site and store it into a Database:** <mark>**get_product_details()**</mark>- The custom url must be opened and essentially parsed using **BeautifulSoup** a Python library for pulling data out of HTML and XML files. This data is basically looped through to find the right div/section for the product list results. Then the appropriate Title, Price, Image, etc.. are pulled into the sqlite3 database.
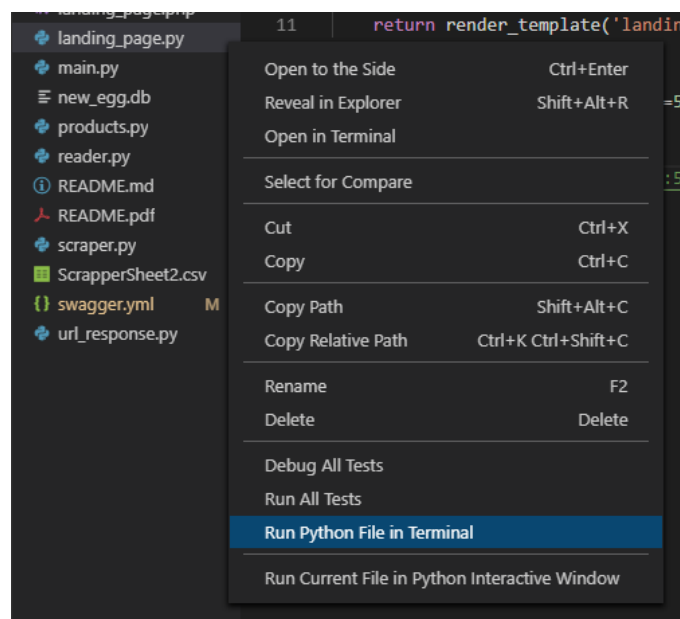
**API Information:** This web app has a built API to perform CRUD operations on it. Check out the functionality & test easily with <u>swagger UI</u>.

If the development server is not running:

*To Run The Development Server:*

→ *Right Click landing_page.py*
→ *Run Python File in Terminal*

Then navigate to the swagger user interface.

**http://localhost:5000/api/ui/#/products - For the expansion**



http://localhost:5000/api/ui/ - For Home Page

Here we can see all the operations

# Troubleshooting steps:

**VIRTUALENV**

How to run virtualenv (env/scripts/activate.ps1) on windows powershell:

1) Navigate to dir

2) RUN: PS C:\Users\Preston\web_scraping_with_python\scraper\python_web_scraper>

**set-executionpolicy remotesigned**

3) [Y] Yes ... : y [enter]

4) RUN: PS C:\Users\Preston\web_scraping_with_python\scraper\python_web_scraper>

.\env\scripts\activate.ps1

5) (env) PS C:\Users\Preston\web_scraping_with_python\scraper\python_web_scraper>

^^^ env should appear before your path.

python_web_scraper
1. Take in .csv files as searchable product data
2. Create custom url handler to search newegg.com/PRODUCT_PART_NUMBER
3. Searching newegg's results and seperating out whether products are found or not found
4. Selecting the best result-> pulling the Title, Price, and Product Image
5. Returning the found data for each .csv entry into a seperate .csv style DB- sqlite3
6. Returning data to program and printing to screen products information