

Dealstream Data Scraping Project

This Scrapy project is designed to scrape business listings from the Dealstream website. It extracts detailed information about businesses for sale, including their description, price, revenue, location, and other relevant details.

Project Structure

- `dealstream_data.py`: This is the main spider file containing the scraping logic.
- `output/`: This directory is where the output JSON file will be saved.

Prerequisites

Before running the project, make sure you have the following installed:

- **Python 3.7+**: Download and install from the [Python official website](#).
- **pip**: Python's package manager, which is typically included with Python.
- **Virtual Environment (venv)**: This comes with Python 3.3+.

Setting Up the Environment

1. Navigate to the Project Directory:

Open your terminal or command prompt and navigate to the directory where the project files are located.

```
cd path/to/project/directory
```

2. Create a Virtual Environment:

Run the following command to create a virtual environment within your project directory:

```
python -m venv venv
```

3. Activate the Virtual Environment:

- On Windows:

```
venv\Scripts\activate
```

- On macOS/Linux:

```
source venv/bin/activate
```

4. Install Scrapy:

Install Scrapy using `pip`:

```
pip install scrapy
```

This command will install Scrapy and its dependencies within your virtual environment.

Running the Spider

1. Update Headers and Cookies:

To scrape Dealstream, you need to update the `headers` and `cookies` in the spider script for each run. This simulates a real user session and can help avoid being blocked.

Steps to obtain new headers and cookies:

- **Log into Dealstream:** Use a dummy account to log into Dealstream to ensure your real account is not affected if the account gets suspended.
- **Capture Network Request Headers:**
 - Open Developer Tools in your browser (usually with F12 or right-click and "Inspect").
 - Go to the **Network** tab.
 - Navigate to the Dealstream search page.
 - Right-click on a request and choose **Copy > Copy as cURL**.
- **Convert cURL to Python Headers and Cookies:** Use an online tool like curlconverter.com to convert the cURL command to Python headers and cookies format.
- **Update the Spider:** Replace the existing `headers` and `cookies` in the `dealstream_data.py` file with the newly generated values.

2. Execute the Spider:

Run the spider using the following command:

```
scrapy crawl dealstream_data
```

The scraped data will be saved to `output/dealstream.json` as specified in the custom settings of the spider.

Notes

- **Compliance:** Always comply with the website's terms of service and `robots.txt` guidelines.
- **Account Safety:** Use a separate dummy account to avoid disrupting personal accounts in case of IP bans or account suspensions.

- **Session Updates:** Regularly update headers and cookies to maintain access to the data.

Troubleshooting

- **Environment Activation:** Ensure your virtual environment is active before running any `pip` or `scrapy` commands.
- **Dependency Issues:** If you encounter dependency conflicts or issues, make sure all necessary packages are installed within your virtual environment.