

Creating a Scrapy project under a virtual environment is a good practice to isolate dependencies and ensure your project doesn't interfere with other Python projects on your machine.

Here's how to do it:

1. Install Python Virtual Environment Tools

If you don't already have virtual environment tools installed, you can install them using:

```
pip install virtualenv
```

Alternatively, if you're using Python 3.3+, the built-in `venv` module can be used without needing to install anything extra.

2. Create a Virtual Environment

1. Navigate to the directory where you want to create your project.
2. Create a virtual environment using either `virtualenv` or Python's `venv` module.

Using `virtualenv`:

```
virtualenv venv
```

Using `venv`:

```
python3 -m venv venv
```

This will create a folder named `venv` in your project directory, containing the Python executable and libraries specific to this virtual environment.

3. Activate the Virtual Environment

To activate the virtual environment, run the following command based on your operating system:

- **For Windows:**

```
venv\Scripts\activate
```

Once activated, you should see the virtual environment's name (`venv`) in your terminal prompt.

4. Install Scrapy Inside the Virtual Environment

Now that the virtual environment is activated, install Scrapy within this isolated environment:

```
pip install scrapy
```

This will ensure that Scrapy and its dependencies are installed only inside the virtual environment.

5. Create Your Scrapy Project

With the virtual environment activated and Scrapy installed, create your Scrapy project

```
scrapy startproject multiscraper
```

This will generate a folder structure for your project like this:

```
multiscraper/
  scrapy.cfg
  multiscraper/
    __init__.py
    items.py
    middlewares.py
    pipelines.py
    settings.py
    spiders/
      __init__.py
```

6. Define Data Structures in `items.py`

In the `items.py` file, define the data you want to scrape, such as titles, URLs, and timestamps:

```
import scrapy

class MultiscraperItem(scrapy.Item):
    # define the fields for your item here like:
    title = scrapy.Field()
    link = scrapy.Field()
    timestamp = scrapy.Field()
    source = scrapy.Field()
```

7. Create Spider for Multiple Sources

Now, create a spider in the `spiders` directory to handle multiple sources. For example, scraping headlines from **BBC** and **CNN**.

Navigate to the `spiders` folder and create a new file `multinews_spider.py`:

```
cd multiscraper/multiscraper/spiders
echo.>> multinews_spider.py
```

Edit `multinews_spider.py` with the following code:

```
import scrapy

from multiscraper.items import NewsItem

class MultiNewsSpider(scrapy.Spider):
    name = "multinews"
    allowed_domains = ['bbc.com', 'cnn.com']

    # Initial request URLs
    start_urls = [
        'https://www.bbc.com/news',
        'https://edition.cnn.com/world'
    ]

    def parse(self, response):
        if 'bbc' in response.url:
            # Scrape BBC headlines
            for article in response.css('div.gs-c-promo-body'):
                item = NewsItem()
                item['title'] = article.css('h3::text').get()
                item['link'] = response.urljoin(article.css('a::attr(href)').get())
                item['timestamp'] = article.css('time::attr(datetime)').get()
                item['source'] = 'BBC'
                yield item

            elif 'cnn' in response.url:
                # Scrape CNN headlines
                for article in response.css('article.cd__content'):
                    item = NewsItem()
                    item['title'] = article.css('span.cd__headline-text::text').get()
                    item['link'] = response.urljoin(article.css('a::attr(href)').get())
                    item['timestamp'] = article.css('span.cd__timestamp::text').get()
                    item['source'] = 'CNN'
                    yield item
```

In this spider, we scrape two different news sources: **BBC** and **CNN**. The `parse` method is used to handle responses from the URLs listed in `start_urls`.

8. Configure Settings (Optional)

You can modify settings in `settings.py` to customize how Scrapy behaves, such as controlling concurrent requests, delays, etc. For example, to slow down requests:

```
# settings.py
DOWNLOAD_DELAY = 1.0 # Add a 1 second delay between requests
CONCURRENT_REQUESTS = 4
```

9. Run the Spider

Run the spider from the terminal:

```
scrapy crawl multinews
```

You can also output data to a specific format (JSON, CSV, etc.):

```
scrapy crawl multinews -o news.json
```

This will create a `news.json` file containing the scraped headlines from both BBC and CNN in JSON format.

10. Output Example

The output will look something like this in `news.json`:

```
[
  {
    "title": "Global warming impact on weather",
    "link": "https://www.bbc.com/news/science-environment-57975048",
    "timestamp": "2024-10-12T08:12:43Z",
    "source": "BBC"
  },
  {
    "title": "New vaccines for diseases",
    "link": "https://edition.cnn.com/2024/10/12/health/new-vaccine/index.html",
    "timestamp": "2024-10-12T08:50:00",
    "source": "CNN"
  }
]
```

11. Extract from Multiple Pages

To handle pagination or scrape multiple pages, modify the spider to follow links to the next pages.

For instance, for **CNN**, you could follow the "Next page" link:

```
def parse(self, response):
    # CNN scraping logic (as before)

    # Follow pagination links
    next_page = response.css('a.pagination__next::attr(href)').get()
    if next_page is not None:
        yield response.follow(next_page, self.parse)
```

This modification will follow the "Next" page link and continue scraping.