

---

# **Scraper Toolkit Documentation**

***Release 1.0.0***

**Anthony Michael Pruitt**

**Jun 08, 2019**



---

## Classes:

---

<b>1</b>	<b>ScraperProject</b>	<b>3</b>
<b>2</b>	<b>PageFetcher</b>	<b>5</b>
<b>3</b>	<b>Parser</b>	<b>7</b>
<b>4</b>	<b>Exporter</b>	<b>9</b>
<b>5</b>	<b>Selector</b>	<b>11</b>
	<b>Python Module Index</b>	<b>13</b>
	<b>Index</b>	<b>15</b>



A toolkit to assist in the page-fetching, HTML-parsing, and data-exporting of a web scraping project.



**class** `scraper_toolkit.ScraperProject.ScraperProject` (*domain: str*)

Handle the page fetching, HTML parsing, and exporting of a web scraping project.

**Parameters** `domain` – Prefix to be added to scraped URLs missing the domain.

**add\_selector** (*selector: Union[str, Selector], attribute: str = None, name: str = None, post\_processing: Callable = None*)

Add the given selector to loaded CSS selectors.

**Parameters**

- **selector** – CSS selector as a string or a Selector type object.
- **attribute** – HTML attribute of the element to store
- **name** – Optional name for the parsed attribute, useful for creating the header row when exporting as a CSV file.
- **post\_processing** – Optional function called on the parsed attribute before it is stored. Useful for cleaning up and splitting data.

**add\_selectors** (*selectors: List[Selector]*)

Add multiple CSS selectors to loaded selectors.

**Parameters** `selectors` – List of Selector objects.

**export\_to\_csv** (*csv\_path: pathlib.Path, encoding: str = 'UTF-8', write\_header: bool = True*)

Export parsed data to a CSV file.

**Parameters**

- **csv\_path** – Path of the location to save the CSV file.
- **encoding** – CSV file encoding. Default is UTF-8.
- **write\_header** – If true, write a header row to the CSV file using the “name” keys in the provided data.

**fetch** (*url: str = None*) → str

Fetch HTML from the page at the given URL.

**Parameters** `url` – URL of the target page.

**Returns** HTML page as a string.

**`parse()`**

Parse HTML for elements using loaded CSS selectors and append matching elements to `self.parsed` as dictionary objects.



**class** `scraper_toolkit.components.PageFetcher.PageFetcher` (*domain: str*)

Fetch URLs and return web pages' HTML as strings.

**Parameters** `domain` – Prefix to be added to scraped URLs missing the domain.

**static** `get_full_url` (*domain: str, suffix: str*) → `str`

Return a complete URL given a domain and suffix, even if the provided suffix is the complete URL.

**Parameters**

- **domain** – The domain of the target page URL.
- **suffix** – The URL of the target page, with or without the domain prefix.

**Returns** The complete URL.

`get_html` (*url: str = None*) → `str`

Fetch the page HTML from the given URL.

**Parameters** `url` – URL of target page.

**Returns** HTML as a string.

`get_links_from_page` (*target\_url: str = None*) → `Iterable[str]`

Return a list of every href URL found from `target_url`.

**Parameters** `target_url` – URL of page to search for href links.

**Returns** List of every discovered href link on the page.

**static** `select_elements_from_html` (*html: str, selector: str*)

Return a list of HTML elements from the given html that match the provided CSS selector.

**Parameters**

- **html** – HTML of the page to parse.
- **selector** – CSS selector for target elements.

**Returns** List of HTML elements matching the CSS selector.

**select\_links\_from\_page** (*selector: str, target\_url: str = None*) → Iterable[str]

Yield the HTML of all pages linked on the *target\_url* located by the given CSS selector.

**Parameters**

- **selector** – CSS selector for elements containing href attribute
- **target\_url** – URL to search for links. If none is provided, the domain URL will be used.

**Returns** Generator for HTML as strings, fetched from the selected links.

**class** `scraper_toolkit.components.Parser.Parser` (*html: str*)

Parse HTML for specific elements or attributes

**Parameters** `html` – HTML to parse, as a string.

**add\_selector** (*selector: Union[str, scraper\_toolkit.components.Selector.Selector] = None, attribute: str = None, name: str = None, post\_processing: Callable = None*)

Add the given selector to loaded CSS selectors.

**Parameters**

- **selector** – CSS selector as a string or a Selector type object.
- **attribute** – HTML attribute of the element to store
- **name** – Optional name for the parsed attribute, useful for creating the header row when exporting as a CSV file.
- **post\_processing** – Optional function called on the parsed attribute before it is stored. Useful for cleaning up and splitting data.

**parse** ()

Parse HTML for elements using loaded CSS selectors and append matching elements to `self.parsed` as dictionary objects.



```
class scraper_toolkit.components.Exporter.Exporter (data: Union[Parser, dict,  
List[dict]])
```

Export data from parsers.

**Param** data: The data to export as a Parser object, a dictionary, or a list of dictionaries.

```
export_to_csv (csv_path: Union[pathlib.Path, pathlib.PurePath, str], encoding: str = 'UTF-8',  
write_header: bool = True)
```

Export parsed data to a CSV file.

#### Parameters

- **csv\_path** – Path of the location to save the CSV file.
- **encoding** – CSV file encoding. Default is UTF-8.
- **write\_header** – If true, write a header row to the CSV file using the “name” keys in the provided data.



```
class scraper_toolkit.components.Selector.Selector(selector_str: str, name: str =  
                                                None, attribute: str = None,  
                                                post_processing: Callable =  
                                                None)
```

Represent a CSS selector with an optional name and target attribute, with an optional post-processing function.

#### Parameters

- **selector\_str** – CSS selector as a string.
- **name** – Optional name for the parsed attribute, useful for creating the header row when exporting as a CSV file.
- **attribute** – HTML attribute of the element to store
- **post\_processing** – Optional function called on the parsed attribute before it is stored. Useful for cleaning up and splitting data.





### S

`scraper_toolkit.components.Exporter`, [9](#)  
`scraper_toolkit.components.PageFetcher`,  
[5](#)  
`scraper_toolkit.components.Parser`, [7](#)  
`scraper_toolkit.components.Selector`, [11](#)  
`scraper_toolkit.ScraperProject`, [3](#)



## A

add\_selector() (*scraper\_toolkit.components.Parser.Parser*  
*method*), 7

add\_selector() (*scraper\_toolkit.ScraperProject.ScraperProject*  
*method*), 3

add\_selectors() (*scraper\_toolkit.ScraperProject.ScraperProject*  
*method*), 3

## S

scraper\_toolkit.components.Exporter  
*(module)*, 9

scraper\_toolkit.components.PageFetcher  
*(module)*, 5

scraper\_toolkit.components.Parser (*module*), 7

scraper\_toolkit.components.Selector  
*(module)*, 11

export\_to\_csv() (*scraper\_toolkit.components.Exporter.Exporter*  
*method*), 9

export\_to\_csv() (*scraper\_toolkit.ScraperProject.ScraperProject*  
*method*), 3

Exporter (*class in scraper\_toolkit.components.Exporter*),  
9

scraper\_toolkit.ScraperProject (*module*), 3

ScraperProject (*class in scraper\_toolkit.ScraperProject*), 3

select\_elements\_from\_html() (*scraper\_toolkit.components.PageFetcher.PageFetcher*  
*static method*), 5

select\_links\_from\_page() (*scraper\_toolkit.components.PageFetcher.PageFetcher*  
*method*), 5

Selector (*class in scraper\_toolkit.components.Selector*),  
11

fetch() (*scraper\_toolkit.ScraperProject.ScraperProject*  
*method*), 3

## G

get\_full\_url() (*scraper\_toolkit.components.PageFetcher.PageFetcher*  
*static method*), 5

get\_html() (*scraper\_toolkit.components.PageFetcher.PageFetcher*  
*method*), 5

get\_links\_from\_page() (*scraper\_toolkit.components.PageFetcher.PageFetcher*  
*method*), 5

## P

PageFetcher (*class in scraper\_toolkit.components.PageFetcher*),  
5

parse() (*scraper\_toolkit.components.Parser.Parser*  
*method*), 7

parse() (*scraper\_toolkit.ScraperProject.ScraperProject*  
*method*), 4

Parser (*class in scraper\_toolkit.components.Parser*), 7