

# Financial Data Analysis with Python

Instructor: Luping Yu

May 10, 2022

---

## Practice 02. Web Page and Crawler

---

### Web Page

Before we start writing code, we need to understand a little bit about the structure of a web page.

When we visit a web page, our web browser makes a request to a web server. This request is called a `GET` request, since we're getting files from the server. The server then sends back files that tell our browser how to render the page for us. These files will typically include:

- HTML — the main content of the page.
- CSS — used to add styling to make the page look nicer.
- JS — Javascript files add interactivity to web pages.

After our browser receives all the files, it **renders** the page and displays it to us.

There's a lot that happens behind the scenes to render a page nicely, but we don't need to worry about most of it when we're web scraping. When we perform web scraping, we're interested in the main content of the web page, so we look primarily at the HTML.

---

### What is HTML?

- HTML stands for Hyper Text Markup Language
  - HTML is the standard **markup language** for creating Web pages
  - HTML describes the structure of a Web page
  - HTML consists of a series of **elements**
  - HTML elements tell the browser how to display the content
  - HTML elements label pieces of content such as "this is a heading", "this is a paragraph", "this is a link", etc.
- 

### How to View HTML Source?

Have you ever seen a Web page and wondered "Hey! How did they do that?"

View HTML Source Code:

- Right-click in an HTML page and select "View Page Source" (in Chrome) or "View Source" (in Edge), or similar in other browsers. This will open a window containing the HTML source code of the page.

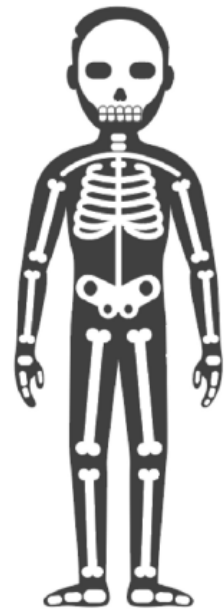
Chrome Developer Tools **F12** :

- More Tools ---> Developer Tools
- Chrome DevTools is a set of web developer tools built directly into the Google Chrome browser.

## HTML THE SKELETON

HTML provides names (tags) to describe different types of content (elements) on your website —for example: `<header>`, `<link>`, `<div>`.

This allows your browser understand what it is reading and how to render it. While your browser can render HTML by itself, you can make it dynamic and beautiful using Javascript & CSS. So think of HTML as a skeleton that comes to life with Javascript & CSS.



---

## A Simple HTML Document (a.k.a. "Page")

```
In [ ]: %%HTML

<!DOCTYPE html>
<html>

  <head>
    <title>Page Title</title>
  </head>

  <body>
    <h1>My First Heading</h1>
    <p>My first paragraph.</p>
  </body>

</html>
```

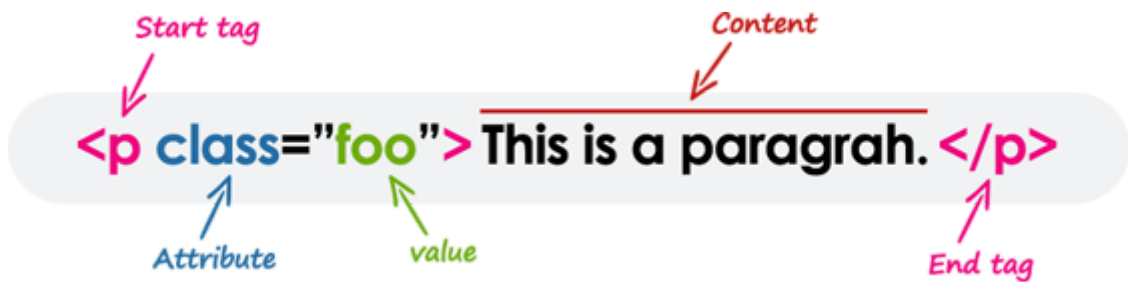
## Example Explained

- The `<!DOCTYPE html>` declaration defines that this document is an HTML document
- The `<html>` element is the root element of an HTML page
- The `<head>` element contains meta information about the HTML page

- The `<title>` element specifies a title for the HTML page (which is shown in the browser's title bar or in the page's tab)
- The `<body>` element defines the document's body, and is a container for all the **visible** contents, such as headings, paragraphs, images, hyperlinks, tables, lists, etc.
- The `<h1>` element defines a large heading
- The `<p>` element defines a paragraph

## HTML Element

An HTML **element** is defined by a **start tag**, some **content**, and an **end tag**:



The HTML **element** is everything from the start tag to the end tag:

```
In [ ]: <h1>My First Heading</h1>
        <p>My first paragraph.</p>
```

### 1. Empty HTML Elements

- Empty elements (also called self-closing or void elements) are not container tags.
- A typical example of an empty element, is the `<br>` element, which represents a line break. Some other common empty elements are `<img>`, `<input>`, etc.

```
In [ ]: %%HTML

        <p>This paragraph contains <br> a line break.</p>
        
        <input type="text" name="username">
```

### 2. Nesting HTML Elements

- Most HTML elements can contain any number of further elements, which are, in turn, made up of tags, attributes, and content or other elements.
- The following example shows some elements nested inside the `<p>` element.

```
In [ ]: %%HTML

        <p>Here is some <b>bold</b> text.</p>
        <p>Here is some <em>emphasized</em> text.</p>
        <p>Here is some <mark>highlighted</mark> text.</p>
```

### 3. HTML Links

HTML links are defined with the `<a>` tag:

- The link's destination is specified in the `href` attribute.
- `Attributes` are used to provide **additional information** about HTML elements.

```
In [ ]: %%HTML
<a href="https://sm.xmu.edu.cn/">This is a link</a>
```

## 4. Writing Comments in HTML

- Comments are usually added with the purpose of making the source code easier to understand.
- You can also comment out part of your HTML code for debugging purpose.
- An HTML comment begins with `<!--`, and ends with `-->`, as shown in the example below:

```
In [ ]: %%HTML
<!-- This is an HTML comment -->
<!-- This is a multi-line HTML comment
      that spans across more than one line -->
<p>This is a normal piece of text.</p>
```

## 5. HTML Elements Types

The basic elements of an HTML page are:

- A text header, denoted using the `<h1>`, `<h2>`, `<h3>`, `<h4>`, `<h5>`, `<h6>` tags.
- A paragraph, denoted using the `<p>` tag.
- A link, denoted using the `<a>` (anchor) tag.
- A list, denoted using the `<ul>` (unordered list), `<ol>` (ordered list) and `<li>` (list element) tags.
- An image, denoted using the `<img>` tag
- A divider, denoted using the `<div>` tag
- A text span, denoted using the `<span>` tag

Elements can be placed in two distinct groups: **block level** and **inline level** elements.

The former make up the document's structure, while the latter dress up the contents of a block.

- A block element occupies 100% of the available width and it is rendered with a line break before and after. Whereas, an inline element will take up only as much space as it needs.
    - The most commonly used block-level elements are `<div>`, `<p>`, `<h1>` through `<h6>`, `<form>`, `<ol>`, `<ul>`, `<li>`, and so on. Whereas, the commonly used inline-level elements are `<img>`, `<a>`, `<span>`, `<strong>`, `<b>`, `<em>`, `<i>`, `<code>`, `<input>`, `<button>`, etc.
  - The block-level elements should not be placed within inline-level elements. For example, the `<p>` element should not be placed inside the `<b>` element.
-

# HTML Attributes

**Attributes** define **additional characteristics or properties** of the element such as width and height of an image.

- Attributes are always specified in the start tag (or opening tag) and usually consists of name/value pairs like `name="value"`.
  - Some attributes are required for certain elements. For instance, an `<img>` tag must contain a `src` and `alt` attributes.

Let's take a look at some examples of the attributes usages:

```
In [ ]: %%HTML

Google</a>
<input type="text" value="Guido van Rossum">
```

In the above example `src` inside the `<img>` tag is an attribute and image path provided is its value. Similarly `href` inside the `<a>` tag is an attribute and the link provided is its value, and so on.

There are some **general purpose** attributes, such as `id`, `title`, `class`, `style`, etc. that you can use on the majority of HTML elements.

## 1.The `id` Attribute

The `id` attribute is used to give a unique identifier to an element within a document. This makes it easier to select the element using CSS or JavaScript.

- The `id` of an element must be **unique** within a single document.
  - No two elements in the same document can be named with the same `id`.
  - Each element can have only one `id`.

```
In [ ]: %%HTML

<input type="text" id="firstName">
<div id="container">Some content</div>
<p id="infoText">This is a paragraph.</p>
```

## 2. The `class` Attribute

Like `id` attribute, the `class` attribute is also used to identify elements. But unlike `id`, the class attribute **does not have to be unique** in the document.

- This means you can apply the same class to multiple elements in a document.
- Any style rules that are written to that class will be applied to all the elements having that class.

```
In [ ]: %%HTML

<input type="text" class="highlight">
<div class="box highlight">Some content</div>
<p class="highlight">This is a paragraph.</p>
```

---

### 3. The `style` Attribute

HTML is quite limited when it comes to the presentation of a web page. It was originally designed as a simple way of presenting information.

**CSS (Cascading Style Sheets)** was introduced in December 1996 by the World Wide Web Consortium (W3C) to provide a better way to style HTML elements.



The `style` attribute allows you to specify CSS styling rules such as color, font, border, etc. directly within the element. Let's check out an example to see how it works:

In [ ]: `%%HTML`

```
<p style="color: blue;">This is a paragraph.</p>

<div style="border: 1px solid red;">Some content</div>
```

### 4. HTML Attributes Types

Each element can also have attributes - each element has a different set of attributes relevant to the element.

There are a few global elements, the most common of them are:

- `id` - Denotes the unique ID of an element in a page. Used for locating elements by using links, JavaScript, and more.
- `class` - Denotes the CSS class of an element.
- `style` - Denotes the CSS styles to apply to an element.

Some websites offer data sets that are downloadable in CSV format, or accessible via an Application Programming Interface (API). But many websites with useful data don't offer these convenient options.

Web scraping is the process of **gathering information** from the Internet. Even copying and pasting the lyrics of your favorite song is a form of web scraping!

The words "web scraping" usually refer to a process that involves automation. Some websites don't like it when automatic scrapers gather their data, while others don't mind.

Common tools:

- requests
- BeautifulSoup
- Selenium
- Pandas.read\_html( )

---

## How Does Web Scraping Work?

When we scrape the web, we write code that sends a request to the server that's hosting the page we specified. The server will return the source code — HTML, mostly — for the page (or pages) we requested.

So far, we're essentially doing the same thing a web browser does — sending a server request with a specific URL and asking the server to return the code for that page.

But unlike a web browser, our web scraping code won't interpret the page's source code and display the page visually. Instead, we'll write some custom code that filters through the page's source code looking for specific elements we've specified, and extracting whatever content we've instructed it to extract.

For example, if we wanted to get all of the data from inside a table that was displayed on a web page, our code would be written to go through these steps in sequence:

- Request the content (source code) of a specific URL from the server
- Download the content that is returned
- Identify the elements of the page that are part of the table we want
- Extract and (if necessary) reformat those elements into a dataset we can analyze or use in whatever way we require.

Python and `BeautifulSoup` have built-in features designed to make this relatively straightforward.

---

## The requests library

Now that we understand the structure of a web page, it's time to get into the fun part: scraping the content we want!

The first thing we'll need to do to scrape a web page is to download the page. We can download pages using the Python `requests` library.

The requests library will make a `GET` request to a web server, which will download the HTML contents of a given web page for us.

Let's try downloading a simple [sample website](#):

```
In [ ]: import requests

page = requests.get("https://dataquestio.github.io/web-scraping-pages/simple
page
```

A `status_code` of 200 means that the page downloaded successfully.

We won't fully dive into status codes here, but a status code starting with a 2 generally indicates success, and a code starting with a 4 or a 5 indicates an error.

We can print out the HTML content of the page using the `content` property:

```
In [ ]: page.content
```

---

## Parsing a page with BeautifulSoup

As you can see above, we now have downloaded an HTML document.

We can use the `BeautifulSoup` library to parse this document, and extract the text from the `p` tag. If we want to extract a single tag, we can use the `find_all` method, which will find all the instances of a tag on a page.

We first have to import the library, and create an instance of the `BeautifulSoup` class to parse our document:

```
In [ ]: from bs4 import BeautifulSoup

soup = BeautifulSoup(page.content, 'html.parser')

soup
```

```
In [ ]: soup.find_all('p')
```

Note that `find_all` returns a list, so we'll have to loop through, or use list indexing, it to extract text:

```
In [ ]: soup.find_all('p')[0].get_text()
```

---

## Searching for tags by class and id

We introduced classes and ids earlier, but it probably wasn't clear why they were useful.



`Classes` and `ids` are used by CSS to determine which HTML elements to apply certain styles to. But when we're scraping, we can also use them to specify the elements we want to scrape.

To illustrate this principle, we'll work with another [sample website](#):

```
In [ ]: page = requests.get("https://dataquestio.github.io/web-scraping-pages/ids_and_classes.html")
        soup = BeautifulSoup(page.content, 'html.parser')
```

Now, we can use the `find_all` method to search for items by class or by id. Let's look for any tag that has the class `outer-text`:

```
In [ ]: soup.find_all(class_='outer-text')
```

We can also search for elements by `id`:

```
In [ ]: soup.find_all(id="first")
```

---

## Using CSS Selectors

We can also search for items using `CSS` selectors. These selectors are how the CSS language allows developers to specify HTML tags to style. Here are some examples:

- `p a` — finds all `a` tags inside of a `p` tag.
- `body p a` — finds all `a` tags inside of a `p` tag inside of a `body` tag.
- `html body` — finds all `body` tags inside of an `html` tag.
- `p.outer-text` — finds all `p` tags with a class of `outer-text`.
- `p#first` — finds all `p` tags with an id of `first`.
- `body p.outer-text` — finds any `p` tags with a class of `outer-text` inside of a `body` tag.

`BeautifulSoup` objects support searching a page via CSS selectors using the `select` method. We can use CSS selectors to find all the `p` tags in our page that are inside of a `div` like this:

```
In [ ]: soup.select("div p")
```

---

## Example 1. Downloading weather data

We now know enough to proceed with extracting information about the local weather from the National Weather Service website.

The first step is to find the page we want to scrape. We'll extract weather information about downtown San Francisco from [this page](#). The page has information about the extended forecast for the next week, including time of day, temperature, and a brief description of the conditions.

Specifically, let's extract data about the extended forecast.

---

## 1. Exploring page structure with Chrome DevTools

The first thing we'll need to do is inspect the page using **Chrome Devtools**. If you're using another browser, Firefox and Safari have equivalents.

You can start the developer tools in Chrome by clicking **View -> Developer -> Developer Tools**. Make sure the Elements panel is highlighted.

The elements panel will show you all the HTML tags on the page, and let you navigate through them. It's a really handy feature!

We can then scroll up in the elements panel to find the "outermost" element that contains all of the text that corresponds to the extended forecasts. In this case, it's a `div` tag with the id `seven-day-forecast`.

If we click around on the console, and explore the div, we'll discover that each forecast item (like "Tonight", "Thursday", and "Thursday Night") is contained in a `div` with the class `tombstone-container`.

---

## 2. Time to Start Scraping!

We now know enough to download the page and start parsing it. In the below code, we will:

- Download the web page containing the forecast.
- Create a BeautifulSoup class to parse the page.
- Find the div with id seven-day-forecast, and assign to seven\_day
- Inside seven\_day, find each individual forecast item.
- Extract and print the first forecast item.

```
In [ ]: page = requests.get("https://forecast.weather.gov/MapClick.php?lat=37.7772&lon=-122.4237")
soup = BeautifulSoup(page.content, 'html.parser')

seven_day = soup.find(id="seven-day-forecast")

forecast_items = seven_day.find_all(class_="tombstone-container")

tonight = forecast_items[0]

tonight
```

---

## 3. Extracting information from the page

As we can see, inside the forecast item **tonight** is all the information we want. There are four pieces of information we can extract:

- The name of the forecast item — in this case, **Tonight**.
- The description of the conditions — this is stored in the **title** property of **img**.
- A short description of the conditions — in this case, **Mostly Clear**.

- The temperature low — in this case, 49 degrees.

We'll extract the name of the forecast item, the short description, and the temperature first, since they're all similar:

```
In [ ]: period = tonight.find(class_="period-name").get_text()
short_desc = tonight.find(class_="short-desc").get_text()
temp = tonight.find(class_="temp").get_text()
print(period)
print(short_desc)
print(temp)
```

Now, we can extract the `title` attribute from the `img` tag. To do this, we just treat the `BeautifulSoup` object like a dictionary, and pass in the attribute we want as a key:

```
In [ ]: img = tonight.find("img")
desc = img['title']
print(desc)
```

---

## 4. Extracting all the information from the page

Now that we know how to extract each individual piece of information, we can combine our knowledge with CSS selectors and list comprehensions to **extract everything at once**.

In the below code, we will:

- Select all items with the class `period-name` inside an item with the class `tombstone-container` in `seven_day`.
- Use a list comprehension to call the `get_text` method on each `BeautifulSoup` object.

```
In [ ]: period_tags = seven_day.select(".tombstone-container .period-name")
periods = [pt.get_text() for pt in period_tags]
periods
```

As we can see above, our technique gets us each of the period names, in order. We can apply the same technique to get the other three fields:

```
In [ ]: short_descs = [sd.get_text() for sd in seven_day.select(".tombstone-container .short-desc")]
temps = [t.get_text() for t in seven_day.select(".tombstone-container .temp")]
descs = [d["title"] for d in seven_day.select(".tombstone-container img")]
print(short_descs)
print(temps)
print(descs)
```

---

## 5. Combining our data into a Pandas Dataframe

We can now combine the data into a `Pandas DataFrame` and analyze it.

```
In [ ]: import pandas as pd
```

```
weather = pd.DataFrame({"period": periods,
                        "short_desc": short_descs,
                        "temp": temps,
                        "desc": descs
})

weather
```

## Example 2. Pandas Web Scraping

Pandas makes it easy to scrape a `table` tag on a web page. You can use the function `read_html(url)` to get webpage contents.

It's only suitable for fetching `Table` type data, then let's see what kind of pages meet the conditions?

### HTML Tables

HTML tables allow web developers to arrange data into rows and columns:

```
In [ ]: %%HTML

<table>
  <tr>
    <th>Company</th>
    <th>Contact</th>
    <th>Country</th>
  </tr>
  <tr>
    <td>Apple</td>
    <td>Tim Cook</td>
    <td>United States</td>
  </tr>
  <tr>
    <td>Tencent</td>
    <td>Pony Ma</td>
    <td>China</td>
  </tr>
</table>
```

Open the web version with a browser, F12 view the HTML structure of the web page, you will find that the structure of the eligible web page has a common feature.

If you find `Table` format, you can use `pd.read_html( )`

### Sina finance - Institutional ownership

Take the aggregate shareholding data of Sina Financial institutions as an example:

```
In [ ]: url = 'https://vip.stock.finance.sina.com.cn/q/go.php/vComStockHold/kind/jjz

df = pd.read_html(url)[0]

df
```

```
In [ ]: df = pd.DataFrame()

for i in range(1, 10):
    url = 'https://vip.stock.finance.sina.com.cn/q/go.php/vComStockHold/kind
    df = pd.concat([df, pd.read_html(url)[0]])
    print("Page %s completed!" % i)

df
```