

Web Scrapping

Web Scraping

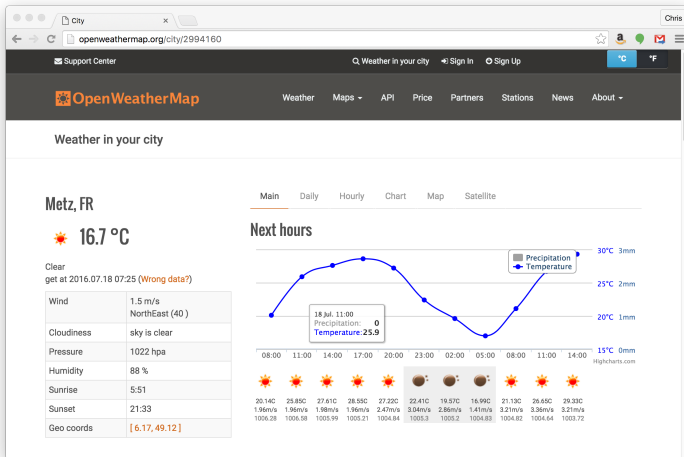
Two ways to mine data from the web

- ▶ The hard way, by web scraping
- ▶ The easy way, using web service APIs

We'll see examples of both.

Web Scrapping

Web scraping, a.k.a. screen scraping, means getting data from a web page. Suppose we want to get the current wind data for a city from [Open Weather Map](#).



Finding The Data On the Page

First you need to find the data within the HTML code for a page so you can construct a regex. Your browser's developer features can help you find the data:

The screenshot shows a web browser window displaying the OpenWeatherMap page for Metz, FR. The page shows the current temperature as 16.7°C and a table of weather data. The developer tools are open, showing the HTML structure of the page. The 'Elements' panel is selected, and the 'div' element with the class 'weather-widget' is highlighted. This element contains the weather data for Metz, FR.

OpenWeatherMap

Weather in your city

Metz, FR

☀️ 16.7°C

Clear
get at 2016.07.18 07:25 (Wrong data?)

Wind	1.5 m/s NorthEast (40)
Cloudiness	sky is clear
Pressure	1022 hpa
Humidity	88 %
Sunrise	5:51
Sunset	21:33
Geo coords	[6.17, 49.12]

amazon.fr

43% 33%

```
<div class="weather-widget">
  <h3>Metz, FR</h3>
  <h2></h2>
  <div>
    Clear
    <p></p>
    <!-- Large modal -->
    <!-- Wrong data modal -->
    <div id="wrong-data-modal" class="modal fade"
      tabindex="-1" role="dialog"></div>
    <!-- Wrong data modal -->
    <table class="table table-striped table-bordered
      table-condensed">
      <tbody>
        <tr>
          <td>Wind</td>
          <td>1.5 m/s</td>
        </tr>
      </tbody>
    </table>
  </div>
```

Getting the Web Page's HTML Code

To get the HTML code of the web page into a Python string variable that you can play with, use Python's `urllib.request` module.

```
import urllib.request
# 2994160 is the city code for Metz, FR
request = urllib.request.Request("http://www.openweathermap.com/city/2994160")
response = urllib.request.urlopen(request)
page_bytes = response.read()
page_text = page_bytes.decode()
# page_text is Python str containing the HTML code
```

Extracting the Data

Looks like the wind data is in the second `<td>` element after the `<div class="weather-widget">` tag, following a `<td>Wind</td>` element. We can play around with the HTML text in the Python REPL. We eventually end up with:

```
wind = re.findall(r'<td>Wind</td><td>(.*?)</td>', page_text.replace("\n",""))[0]
```

Notice that we used a capture group to get the element data.

Aside: Parsing HTML

HTML is context free language, which roughly means that it supports arbitrary nesting of elements. For example, you could have arbitrarily nested `div` elements with "leaf" elements containing text data, e.g.:

```
<div>
  <div>
    <div>some text</div>
  </div>
</div>
```

By the rules of HTML, you could nest `div` tags as deeply as you want. Regular expressions match regular languages, which don't support arbitrary nesting. So how can we use regexes to "parse" HTML?

Regex Matching in HTML Code

Parsing means scanning the linear sequence of symbols in a string to determine its structure (usually by putting the symbols in a tree). We don't need to parse HTML to find data on a web page. While the HTML **language** supports arbitrary nesting, a particular web page will be nested to a particular depth, resulting a simple linear sequence of symbols that we can match with a regular expression.