

*COM5507 Social Media Data
Acquisition and Processing*
Week 5. Web Technologies &
Web Scraping – Episode 1

Lectured by: Dr. Xinzhi ZHANG

Research Assistant Professor, Department of Journalism

Hong Kong Baptist University

3 Oct 2018 @ CityU M4003

Agenda

- Web scraping: an introduction
 - What & why
- The web technologies
 - The networking infrastructure & HTTP
 - HTML & CSS
- Prelude
- Web scraping: the first instances

Web scraping: an introduction

- Web scraping (also called “screen scraping,” “web harvesting,” “web data extraction,” or even “web data mining”), can be defined as “the construction of an agent to download, parse, and organize data from the web in an automated manner” (Broucke, 2017).

Why web scraping?

- The web is a rich data source of human's digital traces and social and cultural artifacts in many research domains;
- Web scraping is faster and more accurate than collecting data manually;
- It offers a solution for data collection when API (application program interface) is unavailable or is not free;
- It is a feasible empirical instance to learn a programming language.

Web scraping: an introduction

- Practitioners of web scraping
 - Academic researchers, esp. in digital humanities and social sciences
 - Business and finance analysts
 - Searching engine and web product developers
 - HR and employers
 - Digital marketing
 - Data-driven journalism professionals

The web technologies

- International Organization of Standardization (ISO) maintains the *Open Systems Interconnection (OSI) model*, in order to standardize the communication processes on the web.
- There are 7 layers

<https://medium.com/@madhavbahl10/osi-model-layers-explained-ee1d43058c1f>

An introduction [[URL](#)]

HTTP - Hypertext Transfer Protocol

- The HyperText Transfer Protocol is the set of rules to allow browsers to retrieve web documents from servers over the internet.
- It is the dominant Application Layer Protocol on the internet.
- Basic concepts: Make a Connection - Request a document - Retrieve the Document - Close the Connection
- “The web speaks HTTP.”

“Getting data” from the server

- Each time the user clicks on an anchor tag with an href= value to switch to a new page, the browser makes a connection to the web server and issues a “GET” request - to GET the content of the page at the specified URL
- The server returns the HTML document to the browser, which formats and displays the document to the user

Request

Web Server

Response

80

GET http://www.dr-chuck.com/page2.htm

```
<h1>The Second  
Page</h1><p>If you like, you  
can switch back to the <a  
href="page1.htm">First  
Page</a>.</p>
```

Browser

Click

Parse/
Render



Picture source: Charles R. Severance (2010)

HTTP in Python

- The “Requests” library
- To enable Python speaks and understands HTTP in order to “browse” the web
- URL (Uniform Resource Locator), a web address, is a reference to a web resource that specifies its location on a computer network and a mechanism for retrieving it.

<http://www.cityu.edu.hk/com/>

HTML

- HTML (Hypertext Markup Language) is the standard markup language for creating web pages and web applications.
- It is used to “construct” and “structure” the webpages.
 - [[A beginner's tutorial](#)][[Examples](#)]

HTML

The simplest HTML page looks like this:

```
<html>
  <head>
    <title>Page Title</title>
  </head>
  <body>
    <h1>Page Title</h1>
    <p>This is a really interesting paragraph.</p>
  </body>
</html>
```

source: <http://alignedleft.com/tutorials/d3/fundamentals/>

HTML

- E.g. ``
- Tags
 - Opening tag and closing tag
 - ALL HTML tags should be closed
- Attributes
 - Attributes are extra bits of information
 - Attributes appear inside the opening tag and their value is always inside quotation marks

HTML tags

- `<p>...</p>` to enclose a paragraph;
- `
` to set a line break;
- `<h1>...</h1>` to `<h6>...</h6>` for headers;
- `<div>...</div>` to indicate a “division” in an HTML document, basically used to group a set of elements;
- `<a>...` for hyperlinks;
- `...`, `...` for unordered and ordered lists respectively; inside of these, `...` is used for each list item.

HTML tags

- Tables

- The <table> element defines the table.
- The <tr> element defines a table row
- The <th> element defines a table head
- The <td> element defines a data cell. They must be enclosed in <tr> tags.

```
<table border="1">  
  <tr>  
    <th>Month</th>  
    <th>Savings</th>  
  </tr>  
  <tr>  
    <td>January</td>  
    <td>$100</td>  
  </tr>  
</table>
```


HTML tags

- Images
- ``
- img tag does not have a closing tag, it closes itself, ending with `"/>"`
- Types of image that are supported on HTML
 - `***.jpg`
 - `***.gif`
 - `***.png`

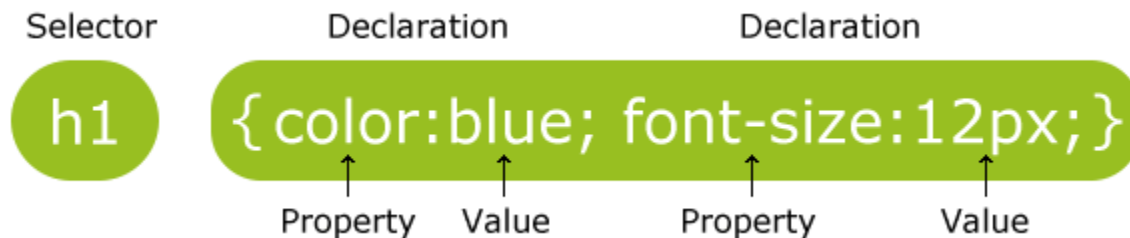
CSS

- Cascading Style Sheets are used to style the visual presentation of HTML pages.
 - Styles define how to display HTML elements
 - Styles were added to HTML 4.0 to solve a problem
- A simplest CSS reads like this (selectors and rules):

```
body {  
    background-color: white;  
    color: black;  
}
```

CSS

- Selectors and rules (also called “declarations”)
 - The selector is normally the HTML element you want to style.
 - Each declaration consists of a property and a value.
 - The property is the style attribute you want to change.
 - Each property has a value.



CSS

- Class
 - to define a group of elements
 - defined with .style
 - retrieved by: class="style":
 - can appear more than once in the same page
- id
 - to label individual, and unique element in the page
 - defined with #style
 - retrieved by: id="style"
 - can only appear once in the same page

CSS

- using id and class in the web page

class :

```
<style type="text/css">
  .footer{background:red;}
</style>
<div class="footer">footer</div>
```

id :

```
<style type="text/css">
  #footer{background:red;}
</style>
<div id="footer">footer</div>
```

- Priority in the browser:
 - style in the tag (attribute) > id > class > tag

Prelude about web scraping

- The “openness” spirit of the web
- The scientific progress
 - the on-going debate of data collection
- General purposes
 - Educational vs. academic vs. commercial

Prelude about web scraping

- Always be aware of laws and policies
 - The legal issues in the particular countries and societies
 - Terms and conditions
 - Copyright protection vs. public domain
 - “Computer Fraud and Abuse Act” (CFAA)
 - Robots.txt – the protocol set by the website for the “crawlers”

Prelude about scraping

- Revisiting the purpose
 - Governmental agencies
 - Academic institutions
 - Open data
 - Journalistic values
- Get written permission, if possible (apply for an API)
- Check the term of use
- Focusing on public information only
- Be respectful
 - “You are the guest in other people’s home”
 - Always set time intervals to avoid excessive amounts of requests
- Be open and transparent

Web scraping: the first instances

- The codes in Jupyter Notebooks
- Reading examples:
 - [Basic ideas](#)
 - [Scraping multiple pages](#)
 - [Xpath](#)
 - [\(preview\) scraping Twitter without API \(from Medium\)](#)

More preludes...

- 授人以鱼，考人以𩚑𩚒𩚓𩚔𩚕𩚖𩚗𩚘𩚙𩚚𩚛𩚜𩚝𩚞𩚟𩚠𩚡𩚢𩚣𩚤𩚥𩚦𩚧𩚨𩚩𩚪𩚫𩚬𩚭𩚮𩚯𩚰𩚱𩚲𩚳𩚴𩚵𩚶𩚷𩚸𩚹𩚺𩚻𩚼𩚽𩚾𩚿𩛀𩛁𩛂𩛃𩛄𩛅𩛆𩛇𩛈𩛉𩛊𩛋𩛌𩛍𩛎𩛏𩛐𩛑𩛒𩛓𩛔𩛕𩛖𩛗𩛘𩛙𩛚𩛛𩛜𩛝𩛞𩛟𩛠𩛡𩛢𩛣𩛤𩛥𩛦𩛧𩛨𩛩𩛪𩛫𩛬𩛭𩛮𩛯𩛰𩛱𩛲𩛳𩛴𩛵𩛶𩛷𩛸𩛹𩛺𩛻𩛼𩛽𩛾𩛿𩜀𩜁𩜂𩜃𩜄𩜅𩜆𩜇𩜈𩜉𩜊𩜋𩜌𩜍𩜎𩜏𩜐𩜑𩜒𩜓𩜔𩜕𩜖𩜗𩜘𩜙𩜚𩜛𩜜𩜝𩜞𩜟𩜠𩜡𩜢𩜣𩜤𩜥𩜦𩜧𩜨𩜩𩜪𩜫𩜬𩜭𩜮𩜯𩜰𩜱𩜲𩜳𩜴𩜵𩜶𩜷𩜸𩜹𩜺𩜻𩜼𩜽𩜾𩜿𩝀𩝁𩝂𩝃𩝄𩝅𩝆𩝇𩝈𩝉𩝊𩝋𩝌𩝍𩝎𩝏𩝐𩝑𩝒𩝓𩝔𩝕𩝖𩝗𩝘𩝙𩝚𩝛𩝜𩝝𩝞𩝟𩝠𩝡𩝢𩝣𩝤𩝥𩝦𩝧𩝨𩝩𩝪𩝫𩝬𩝭𩝮𩝯𩝰𩝱𩝲𩝳𩝴𩝵𩝶𩝷𩝸𩝹𩝺𩝻𩝼𩝽𩝾𩝿𩞀𩞁𩞂𩞃𩞄𩞅𩞆𩞇𩞈𩞉𩞊𩞋𩞌𩞍𩞎𩞏𩞐𩞑𩞒𩞓𩞔𩞕𩞖𩞗𩞘𩞙𩞚𩞛𩞜𩞝𩞞𩞟𩞠𩞡𩞢𩞣𩞤𩞥𩞦𩞧𩞨𩞩𩞪𩞫𩞬𩞭𩞮𩞯𩞰𩞱𩞲𩞳𩞴𩞵𩞶𩞷𩞸𩞹𩞺𩞻𩞼𩞽𩞾𩞿𩟀𩟁𩟂𩟃𩟄𩟅𩟆𩟇𩟈𩟉𩟊𩟋𩟌𩟍𩟎𩟏𩟐𩟑𩟒𩟓𩟔𩟕𩟖𩟗𩟘𩟙𩟚𩟛𩟜𩟝𩟞𩟟𩟠𩟡𩟢𩟣𩟤𩟥𩟦𩟧𩟨𩟩𩟪𩟫𩟬𩟭𩟮𩟯𩟰𩟱𩟲𩟳𩟴𩟵𩟶𩟷𩟸𩟹𩟺𩟻𩟼𩟽𩟾𩟿𩠀𩠁𩠂𩠃𩠄𩠅𩠆𩠇𩠈𩠉𩠊𩠋𩠌𩠍𩠎𩠏𩠐𩠑𩠒𩠓𩠔𩠕𩠖𩠗𩠘𩠙𩠚𩠛𩠜𩠝𩠞𩠟𩠠𩠡𩠢𩠣𩠤𩠥𩠦𩠧𩠨𩠩𩠪𩠫𩠬𩠭𩠮𩠯𩠰𩠱𩠲𩠳𩠴𩠵𩠶𩠷𩠸𩠹𩠺𩠻𩠼𩠽𩠾𩠿𩡀𩡁𩡂𩡃𩡄𩡅𩡆𩡇𩡈𩡉𩡊𩡋𩡌𩡍𩡎𩡏𩡐𩡑𩡒𩡓𩡔𩡕𩡖𩡗𩡘𩡙𩡚𩡛𩡜𩡝𩡞𩡟𩡠𩡡𩡢𩡣𩡤𩡥𩡦𩡧𩡨𩡩𩡪𩡫𩡬𩡭𩡮𩡯𩡰𩡱𩡲𩡳𩡴𩡵𩡶𩡷𩡸𩡹𩡺𩡻𩡼𩡽𩡾𩡿𩢀𩢁𩢂𩢃𩢄𩢅𩢆𩢇𩢈𩢉𩢊𩢋𩢌𩢍𩢎𩢏𩢐𩢑𩢒𩢓𩢔𩢕𩢖𩢗𩢘𩢙𩢚𩢛𩢜𩢝𩢞𩢟𩢠𩢡𩢢𩢣𩢤𩢥𩢦𩢧𩢨𩢩𩢪𩢫𩢬𩢭𩢮𩢯𩢰𩢱𩢲𩢳𩢴𩢵𩢶𩢷𩢸𩢹𩢺𩢻𩢼𩢽𩢾𩢿𩣀𩣁𩣂𩣃𩣄𩣅𩣆𩣇𩣈𩣉𩣊𩣋𩣌𩣍𩣎𩣏𩣐𩣑𩣒𩣓𩣔𩣕𩣖𩣗𩣘𩣙𩣚𩣛𩣜𩣝𩣞𩣟𩣠𩣡𩣢𩣣𩣤𩣥𩣦𩣧𩣨𩣩𩣪𩣫𩣬𩣭𩣮𩣯𩣰𩣱𩣲𩣳𩣴𩣵𩣶𩣷𩣸𩣹𩣺𩣻𩣼𩣽𩣾𩣿𩤀𩤁𩤂𩤃𩤄𩤅𩤆𩤇𩤈𩤉𩤊𩤋𩤌𩤍𩤎𩤏𩤐𩤑𩤒𩤓𩤔𩤕𩤖𩤗𩤘𩤙𩤚𩤛𩤜𩤝𩤞𩤟𩤠𩤡𩤢𩤣𩤤𩤥𩤦𩤧𩤨𩤩𩤪𩤫𩤬𩤭𩤮𩤯𩤰𩤱𩤲𩤳𩤴𩤵𩤶𩤷𩤸𩤹𩤺𩤻𩤼𩤽𩤾𩤿𩥀𩥁𩥂𩥃𩥄𩥅𩥆𩥇𩥈𩥉𩥊𩥋𩥌𩥍𩥎𩥏𩥐𩥑𩥒𩥓𩥔𩥕𩥖𩥗𩥘𩥙𩥚𩥛𩥜𩥝𩥞𩥟𩥠𩥡𩥢𩥣𩥤𩥥𩥦𩥧𩥨𩥩𩥪𩥫𩥬𩥭𩥮𩥯𩥰𩥱𩥲𩥳𩥴𩥵𩥶𩥷𩥸𩥹𩥺𩥻𩥼𩥽𩥾𩥿𩦀𩦁𩦂𩦃𩦄𩦅𩦆𩦇𩦈𩦉𩦊𩦋𩦌𩦍𩦎𩦏𩦐𩦑𩦒𩦓𩦔𩦕𩦖𩦗𩦘𩦙𩦚𩦛𩦜𩦝𩦞𩦟𩦠𩦡𩦢𩦣𩦤𩦥𩦦𩦧𩦨𩦩𩦪𩦫𩦬𩦭𩦮𩦯𩦰𩦱𩦲𩦳𩦴𩦵𩦶𩦷𩦸𩦹𩦺𩦻𩦼𩦽𩦾𩦿𩧀𩧁𩧂𩧃𩧄𩧅𩧆𩧇𩧈𩧉𩧊𩧋𩧌𩧍𩧎𩧏𩧐𩧑𩧒𩧓𩧔𩧕𩧖𩧗𩧘𩧙𩧚𩧛𩧜𩧝𩧞𩧟𩧠𩧡𩧢𩧣𩧤𩧥𩧦𩧧𩧨𩧩𩧪𩧫𩧬𩧭𩧮𩧯𩧰𩧱𩧲𩧳𩧴𩧵𩧶𩧷𩧸𩧹𩧺𩧻𩧼𩧽𩧾𩧿𩨀𩨁𩨂𩨃𩨄𩨅𩨆𩨇𩨈𩨉𩨊𩨋𩨌𩨍𩨎𩨏𩨐𩨑𩨒𩨓𩨔𩨕𩨖𩨗𩨘𩨙𩨚𩨛𩨜𩨝𩨞𩨟𩨠𩨡𩨢𩨣𩨤𩨥𩨦𩨧𩨨𩨩𩨪𩨫𩨬𩨭𩨮𩨯𩨰𩨱𩨲𩨳𩨴𩨵𩨶𩨷𩨸𩨹𩨺𩨻𩨼𩨽𩨾𩨿𩩀𩩁𩩂𩩃𩩄𩩅𩩆𩩇𩩈𩩉𩩊𩩋𩩌𩩍𩩎𩩏𩩐𩩑𩩒𩩓𩩔𩩕𩩖𩩗𩩘𩩙𩩚𩩛𩩜𩩝𩩞𩩟𩩠𩩡𩩢𩩣𩩤𩩥𩩦𩩧𩩨𩩩𩩪𩩫𩩬𩩭𩩮𩩯𩩰𩩱𩩲𩩳𩩴𩩵𩩶𩩷𩩸𩩹𩩺𩩻𩩼𩩽𩩾𩩿𩪀𩪁𩪂𩪃𩪄𩪅𩪆𩪇𩪈𩪉𩪊𩪋𩪌𩪍

More preludes...

- “The web is messy.”
- There will never be an “one-size-fits-all” solution to all the web scraping.
- Solutions
 - Test & rework
 - Search online tutorials and try to replicate their codes
 - Compare the different outputs
 - Consult a friend