


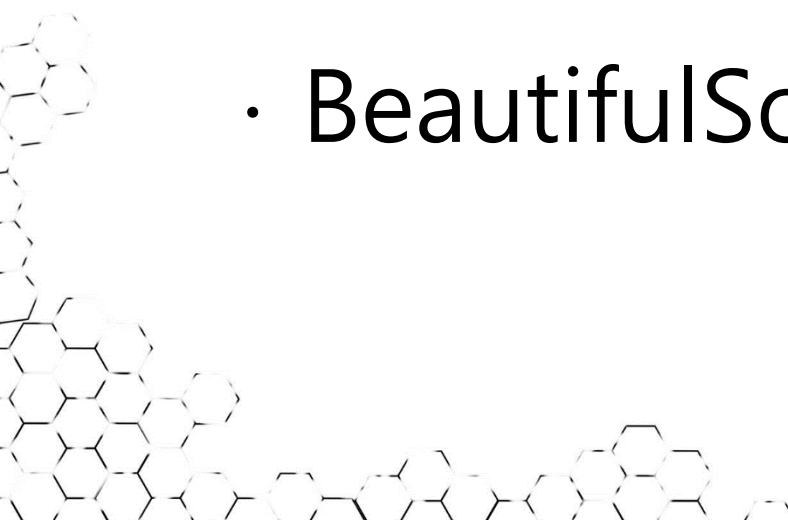


INTRODUCTION TO DATA SCIENCE

TA class V – BeautifulSoup

TA : Lee Chi-Hsuan

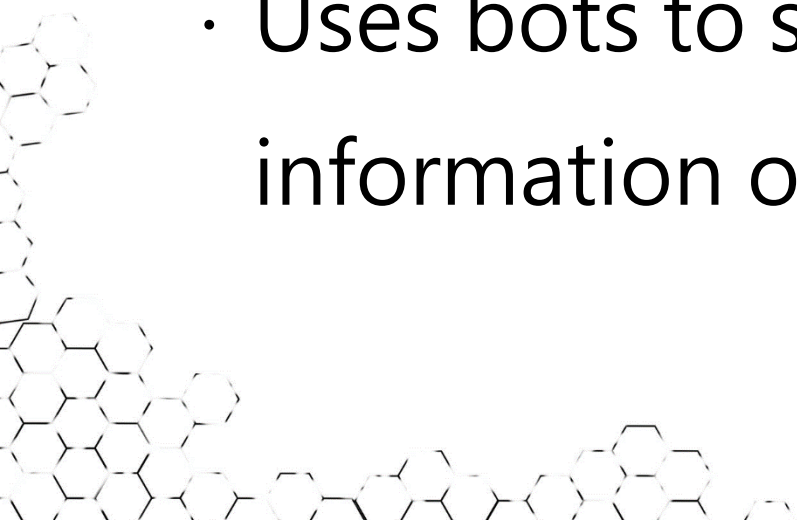


- 
- 
- Web Scraping vs Web Crawling
 - HTML
 - Regular Expression
 - Requests
 - BeautifulSoup

Web Scraping / Crawling

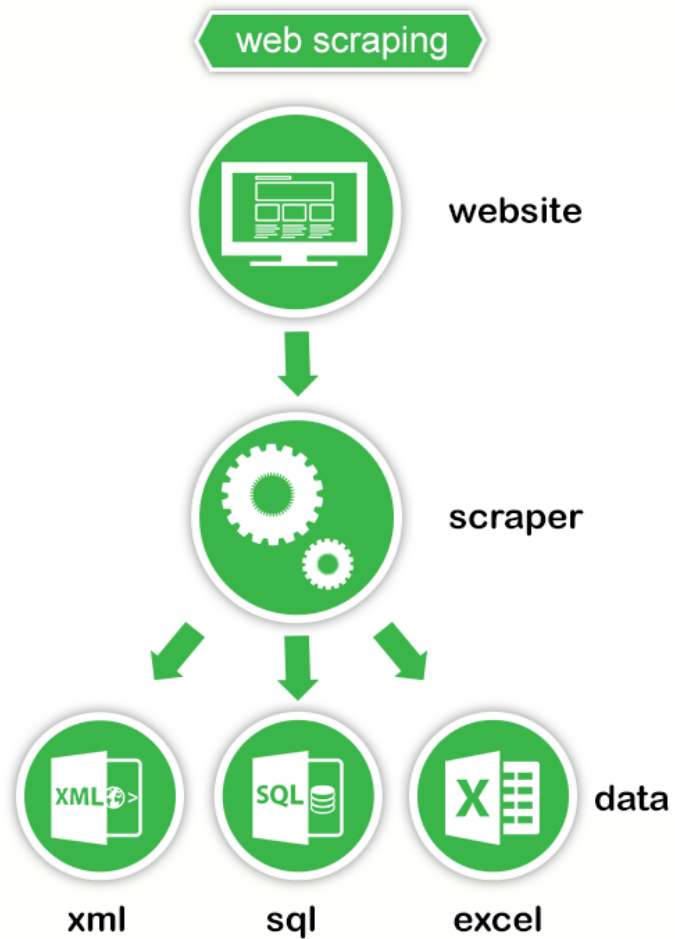


Web Scraping :

- Basically extracting data from websites in an automated manner
 - Uses bots to scrape the information or content from websites
- 

Web Crawling :

- An internet bot that systematically browses the World Wide Web, usually for the purpose of web indexing.
- It involves looking at a page in its entirety and indexing it, including its last letter and dot on the page, in the quest for information.



Steps of Web Scrapping :

1. Find the URL that you want to scrape
2. Inspecting the Page
3. Find the data you want to extract
4. Write the code
5. Run the code and extract the data
6. Store the data in the required format

HTML



- HyperText Markup Language
- The standard markup language for documents designed to be displayed in a web browser
- Can be assisted by CSS or JavaScript

Opening tag

Closing tag

`<p>My cat is very grumpy</p>`

Content

Element

Regular Expression

- A sequence of characters that specifies a search pattern.
- Such patterns are used by string-searching algorithms for "find" or "find and replace" operations on strings, or for input validation.

Metacharacter	Description
<code>^</code>	Matches the starting position within the string. In line-based tools, it matches the starting position of any line.
<code>.</code>	Matches any single character (many applications exclude newlines , and exactly which characters are considered newlines is flavor-, character-encoding-, and platform-specific, but it is safe to assume that the line feed character is included). Within POSIX bracket expressions, the dot character matches a literal dot. For example, <code>a.c</code> matches "abc", etc., but <code>[a.c]</code> matches only "a", ".", or "c".
<code>[]</code>	<p>A bracket expression. Matches a single character that is contained within the brackets. For example, <code>[abc]</code> matches "a", "b", or "c". <code>[a-z]</code> specifies a range which matches any lowercase letter from "a" to "z". These forms can be mixed: <code>[abcx-z]</code> matches "a", "b", "c", "x", "y", or "z", as does <code>[a-cx-z]</code>.</p> <p>The <code>-</code> character is treated as a literal character if it is the last or the first (after the <code>^</code>, if present) character within the brackets: <code>[abc-]</code>, <code>[-abc]</code>. Note that backslash escapes are not allowed. The <code>]</code> character can be included in a bracket expression if it is the first (after the <code>^</code>) character: <code>[]abc]</code>.</p>
<code>[^]</code>	Matches a single character that is not contained within the brackets. For example, <code>[^abc]</code> matches any character other than "a", "b", or "c". <code>[^a-z]</code> matches any single character that is not a lowercase letter from "a" to "z". Likewise, literal characters and ranges can be mixed.
<code>\$</code>	Matches the ending position of the string or the position just before a string-ending newline. In line-based tools, it matches the ending position of any line.
<code>()</code>	Defines a marked subexpression. The string matched within the parentheses can be recalled later (see the next entry, <code>\n</code>). A marked subexpression is also called a block or capturing group. BRE mode requires <code>\(\)</code> .
<code>\n</code>	Matches what the <i>n</i> th marked subexpression matched, where <i>n</i> is a digit from 1 to 9. This construct is vaguely defined in the POSIX.2 standard. Some tools allow referencing more than nine capturing groups. Also known as a backreference.
<code>*</code>	Matches the preceding element zero or more times. For example, <code>ab*c</code> matches "ac", "abc", "abbbc", etc. <code>[xyz]*</code> matches "", "x", "y", "z", "zx", "zyx", "xyzyz", and so on. <code>(ab)*</code> matches "", "ab", "abab", "ababab", and so on.
<code>{m,n}</code>	Matches the preceding element at least <i>m</i> and not more than <i>n</i> times. For example, <code>a{3,5}</code> matches only "aaa", "aaaa", and "aaaaa". This is not found in a few older instances of regexes. BRE mode requires <code>\{m,n\}</code> .

字元	描述
<code>\</code>	將下一個字元標記為一個特殊字元（File Format Escape，清單見本表）、或一個原義字元（Identity Escape，有 <code>^\$()*+?.[\ </code> 共計12個）、或一個向後參照（backreferences）、或一個八進位跳脫符。例如，「 <code>n</code> 」符合字元「 <code>n</code> 」。「 <code>\n</code> 」符合一個換行符。序列「 <code>\\</code> 」符合「 <code>\</code> 」而「 <code>\(</code> 」則符合「 <code>(</code> 」。
<code>^</code>	符合輸入字串的開始位置。如果設定了RegExp物件的Multiline屬性， <code>^</code> 也符合「 <code>\n</code> 」或「 <code>\r</code> 」之後的位置。
<code>\$</code>	符合輸入字串的結束位置。如果設定了RegExp物件的Multiline屬性， <code>\$</code> 也符合「 <code>\n</code> 」或「 <code>\r</code> 」之前的位置。
<code>*</code>	符合前面的子表達式零次或多次。例如， <code>zo*</code> 能符合「 <code>z</code> 」、「 <code>zo</code> 」以及「 <code>zoo</code> 」。 <code>*</code> 等價於 <code>{0,}</code> 。
<code>+</code>	符合前面的子表達式一次或多次。例如，「 <code>zo+</code> 」能符合「 <code>zo</code> 」以及「 <code>zoo</code> 」，但不能符合「 <code>z</code> 」。 <code>+</code> 等價於 <code>{1,}</code> 。
<code>?</code>	符合前面的子表達式零次或一次。例如，「 <code>do(es)?</code> 」可以符合「 <code>does</code> 」中的「 <code>do</code> 」和「 <code>does</code> 」。 <code>?</code> 等價於 <code>{0,1}</code> 。
<code>{n}</code>	<code>n</code> 是一個非負整數。符合確定的 <code>n</code> 次。例如，「 <code>o{2}</code> 」不能符合「 <code>Bob</code> 」中的「 <code>o</code> 」，但是能符合「 <code>food</code> 」中的兩個 <code>o</code> 。
<code>{n,}</code>	<code>n</code> 是一個非負整數。至少符合 <code>n</code> 次。例如，「 <code>o{2,}</code> 」不能符合「 <code>Bob</code> 」中的「 <code>o</code> 」，但能符合「 <code>foooooo</code> 」中的所有 <code>o</code> 。「 <code>o{1,}</code> 」等價於「 <code>o+</code> 」。「 <code>o{0,}</code> 」則等價於「 <code>o*</code> 」。
<code>{n,m}</code>	<code>m</code> 和 <code>n</code> 均為非負整數，其中 <code>n<=m</code> 。最少符合 <code>n</code> 次且最多符合 <code>m</code> 次。例如，「 <code>o{1,3}</code> 」將符合「 <code>foooooo</code> 」中的前三個 <code>o</code> 。「 <code>o{0,1}</code> 」等價於「 <code>o?</code> 」。請注意在逗號和兩個數之間不能有空格。
<code>?</code>	非貪心量化（Non-greedy quantifiers）：當該字元緊跟在任何一個其他重複修飾詞（ <code>*,+,?,{n},{n,},{n,m}</code> ）後面時，符合模式是非貪婪的。非貪婪模式儘可能少的符合所搜尋的字串，而預設的貪婪模式則儘可能多的符合所搜尋的字串。例如，對於字串「 <code>oooo</code> 」，「 <code>o+?</code> 」將符合單個「 <code>o</code> 」，而「 <code>o+</code> 」將符合所有「 <code>o</code> 」。
<code>.</code>	符合除「 <code>\r</code> 」、「 <code>\n</code> 」之外的任何單個字元。要符合包括「 <code>\r</code> 」、「 <code>\n</code> 」在內的任何字元，請使用像「 <code>(.\ r\ n)</code> 」的模式。
<code>(pattern)</code>	符合 <code>pattern</code> 並取得這一符合的子字串。該子字串用於向後參照。所取得的符合可以從產生的Matches集合得到，在VBScript中使用SubMatches集合，在JScript中則使用\$0...\$9屬性。要符合圓括號字元，請使用「 <code>\(</code> 」或「 <code>\)</code> 」。可帶數量字尾。
<code>(?:pattern)</code>	符合 <code>pattern</code> 但不取得符合的子字串（shy groups），也就是說這是一個非取得符合，不儲存符合的子字串用於向後參照。這在使用或字元「 <code>()</code> 」來組合一個模式的各個部分是很有用。例如「 <code>industr(?:y ies)</code> 」就是一個比「 <code>industry industries</code> 」更簡略的表達式。
<code>(?=pattern)</code>	正向肯定預查（look ahead positive assert），在任何符合 <code>pattern</code> 的字串開始處符合尋找字串。這是一個非取得符合，也就是說，該符合不需要取得供以後使用。例如，「 <code>Windows(=95 98 NT 2000)</code> 」能符合「 <code>Windows2000</code> 」中的「 <code>Windows</code> 」，但不能符合「 <code>Windows3.1</code> 」中的「 <code>Windows</code> 」。預查不消耗字元，也就是說，在一個符合發生後，在最後一次符合之後立即開始下一次符合的搜尋，而不是從包含預查的字元之後開始。
<code>(?!pattern)</code>	正向否定預查（negative assert），在任何不符合 <code>pattern</code> 的字串開始處符合尋找字串。這是一個非取得符合，也就是說，該符合不需要取得供以後使用。例如「 <code>Windows(?!95 98 NT 2000)</code> 」能符合「 <code>Windows3.1</code> 」中的「 <code>Windows</code> 」，但不能符合「 <code>Windows2000</code> 」中的「 <code>Windows</code> 」。預查不消耗字元，也就是說，在一個符合發生後，在最後一次符合之後立即開始下一次符合的搜尋，而不是從包含預查的字元之後開始。

Requests



- An elegant and simple HTTP library for Python, built for human beings.
- Get
- Post
- Header

BeautifulSoup

- A Python library for pulling data out of HTML and XML files.






Development Design DIY Giveaway Support Contributors

Python Tutorial: Zip Files – Creating and Extracting Zip Archives

November 19, 2019 by Corey Schafer – 0 Comments

In this video, we will be learning how to create and extract zip archives. We will start by using the zipfile module, and then we will see how to do this using the shutil module. We will learn how to do this with single files and directories, as well as learning how to use glob as well. Let's get started...

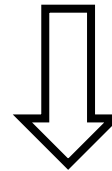


Watch on YouTube

Filed Under: Development, Python
Tagged With: glob, shutil, zip, zipfile

Main Contributor
Andre Nevores

Top Contributors (17)
Andre Nevores
Cyril Hadger
Abhishek Rajan
Hien Nguyen
Alan Haise - totacart.com
Justin Gou
Alex Canady
Sergey Trubin
Justin Prasley
Sriks
phila
Jerome Massey
Robert Butler
Jonathan Lovett



	headline	summary	video_link
0	Python Tutorial: Zip Files – Creating and Extr...	In this video, we will be learning how to crea...	https://youtube.com/watch?v=z0gguhEmWiY
1	Python Data Science Tutorial: Analyzing the 20...	In this Python Programming video, we will be l...	https://youtube.com/watch?v=_P7X8tMplsw
2	Python Multiprocessing Tutorial: Run Code in P...	In this Python Programming video, we will be l...	https://youtube.com/watch?v=fKl2JW_qrso
3	Python Threading Tutorial: Run Code Concurrent...	In this Python Programming video, we will be l...	https://youtube.com/watch?v=IEEhzQoKtQU
4	Update (2019-09-03)	Hey everyone. I wanted to give you an update o...	NaN
5	Python Quick Tip: The Difference Between "==" ...	In this Python Programming Tutorial, we will b...	https://youtube.com/watch?v=mO_dS3rXDIs
6	Python Tutorial: Calling External Commands Usi...	In this Python Programming Tutorial, we will b...	https://youtube.com/watch?v=2Fp1N6dof0Y
7	Visual Studio Code (Windows) – Setting up a Py...	In this Python Programming Tutorial, we will b...	https://youtube.com/watch?v=nh9rCzPJ20
8	Visual Studio Code (Mac) – Setting up a Python...	In this Python Programming Tutorial, we will b...	https://youtube.com/watch?v=06l63_p-2A4
9	Clarifying the Issues with Mutable Default Arg...	In this Python Programming Tutorial, we will b...	https://youtube.com/watch?v=_JGmemulNww

Let's have a quick walkthrough

Imports

```
1 import pandas as pd
2 import numpy as np
3 import requests
4 import re
5 import csv
6 from bs4 import BeautifulSoup
```

Web Scraping vs Web Crawling

Web scraping :

The process of processing a web document and extracting information out of it.
You can do web scraping without doing web crawling.

Web crawling :

The process of iteratively finding and fetching web links starting from a list of seed URL's.
Strictly speaking, to do web crawling, you have to do some degree of web scraping (to extract the URL's.)