

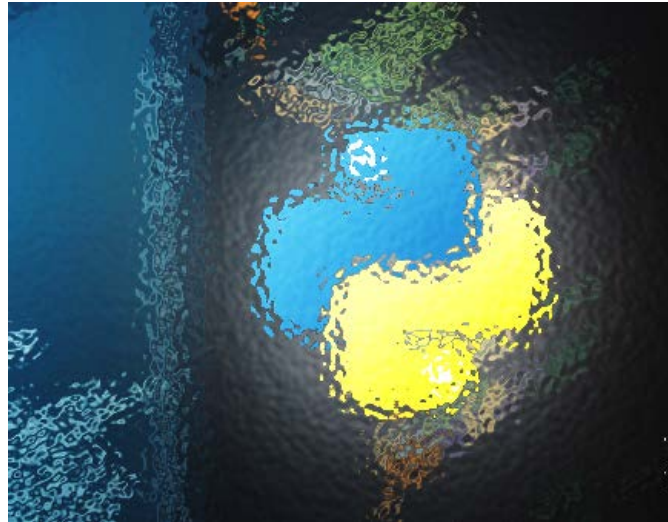


# Web Scrapping

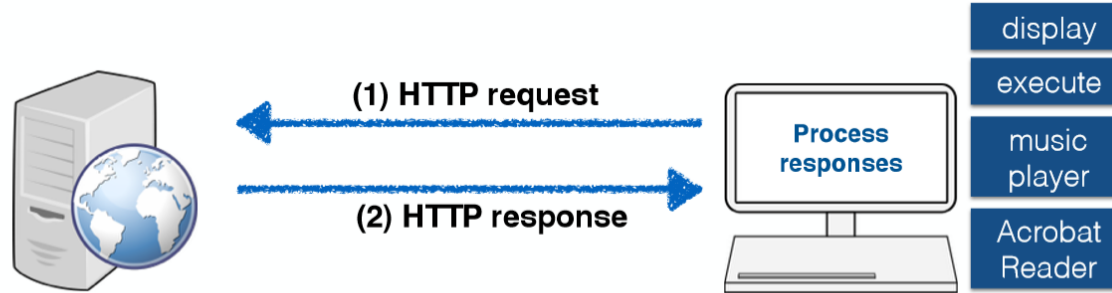
Computer-Science Projects on BigData

Elias Larbi  
Bart Gerritsen

Dec 5, 2018



# Client-server and the Web



Manages web  
resources

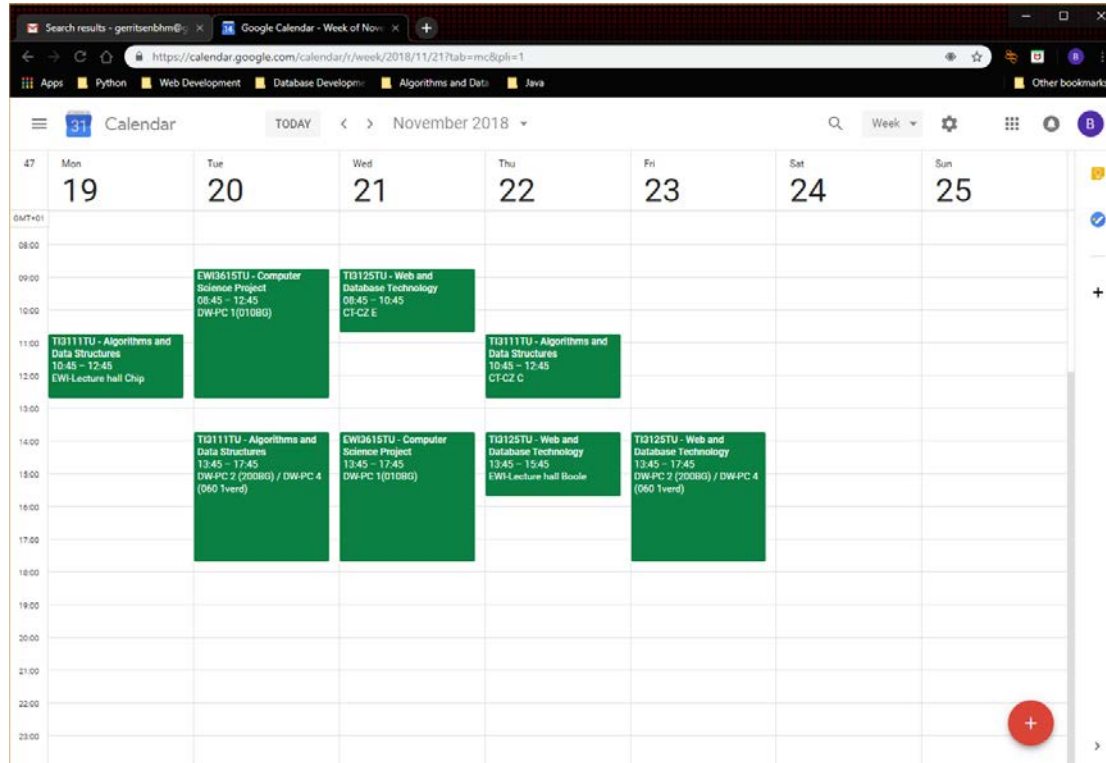
- Servers wait for data requests
- Answer thousands of clients simultaneously
- Host **web resources**

- Clients are most often Web browsers
- Telnet

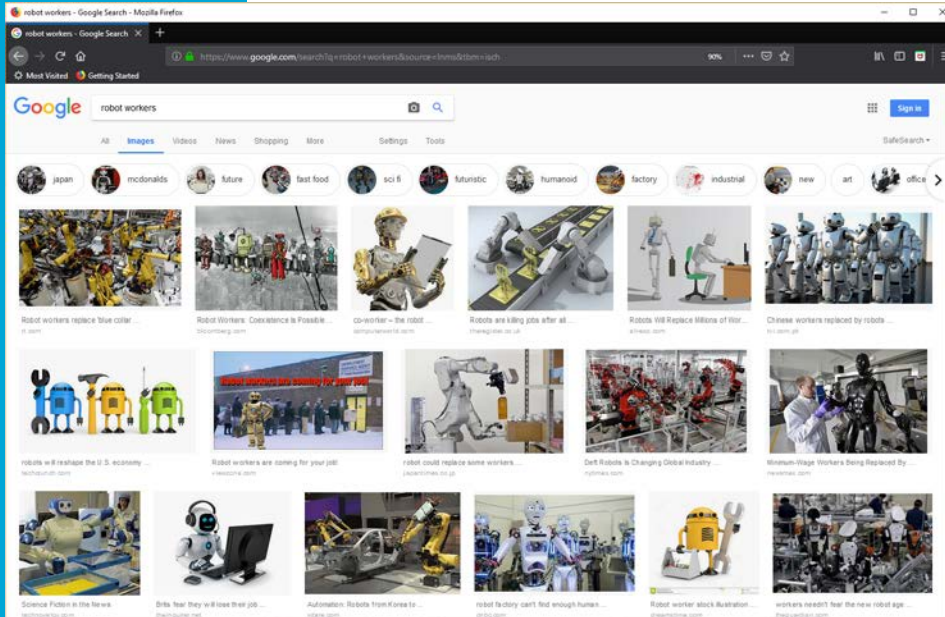
Consumes /  
manipulates web  
resources

**Web resource:** any kind of content with an identity, including static files (e.g. text, images, video), software programs, Web cam gateway, etc.

# Demo query



# Requests and web programming

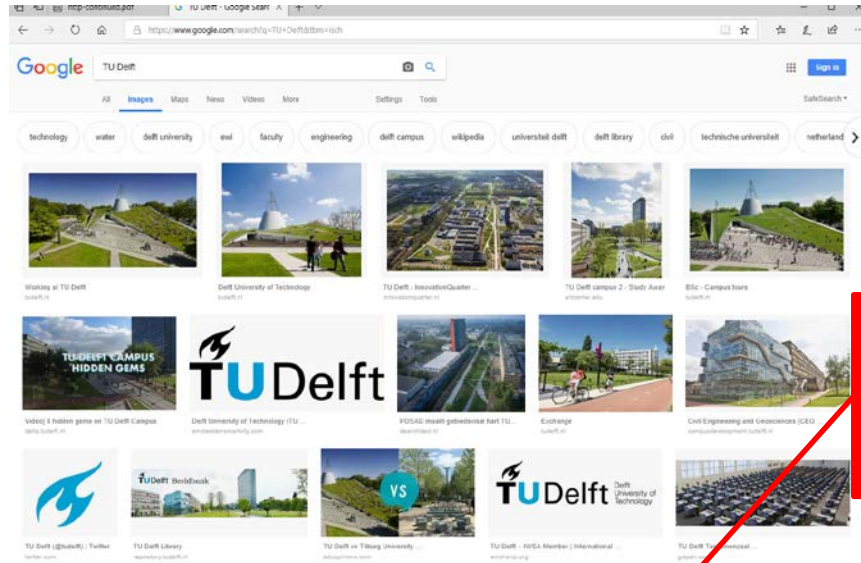


```
GENERATE CODE SNIPPETS
HTTP
1 GET /search?q=robot+workers&source=lnms&tb=isch HTTP/1.1
2 Host: www.google.com
3 cache-control: no-cache
4 Postman-Token: 30e792d6-38cf-4a08-95ee-66aa2c975a01
5
```

```
GENERATE CODE SNIPPETS
NodeJS Request
1 var request = require("request");
2
3 var options = { method: 'GET',
4   url: 'https://www.google.com/search?',
5   qs: { q: 'robot+workers', source: 'lnms', tb: 'isch' },
6   headers:
7     { 'Postman-Token': '6fda378-3a40-4330-682e-e1ea59f0e066',
8       'cache-control': 'no-cache' } };
9
10 request(options, function (error, response, body) {
11   if (error) throw new Error(error);
12
13   console.log(body);
14 });
15
```

# URL – Syntax query

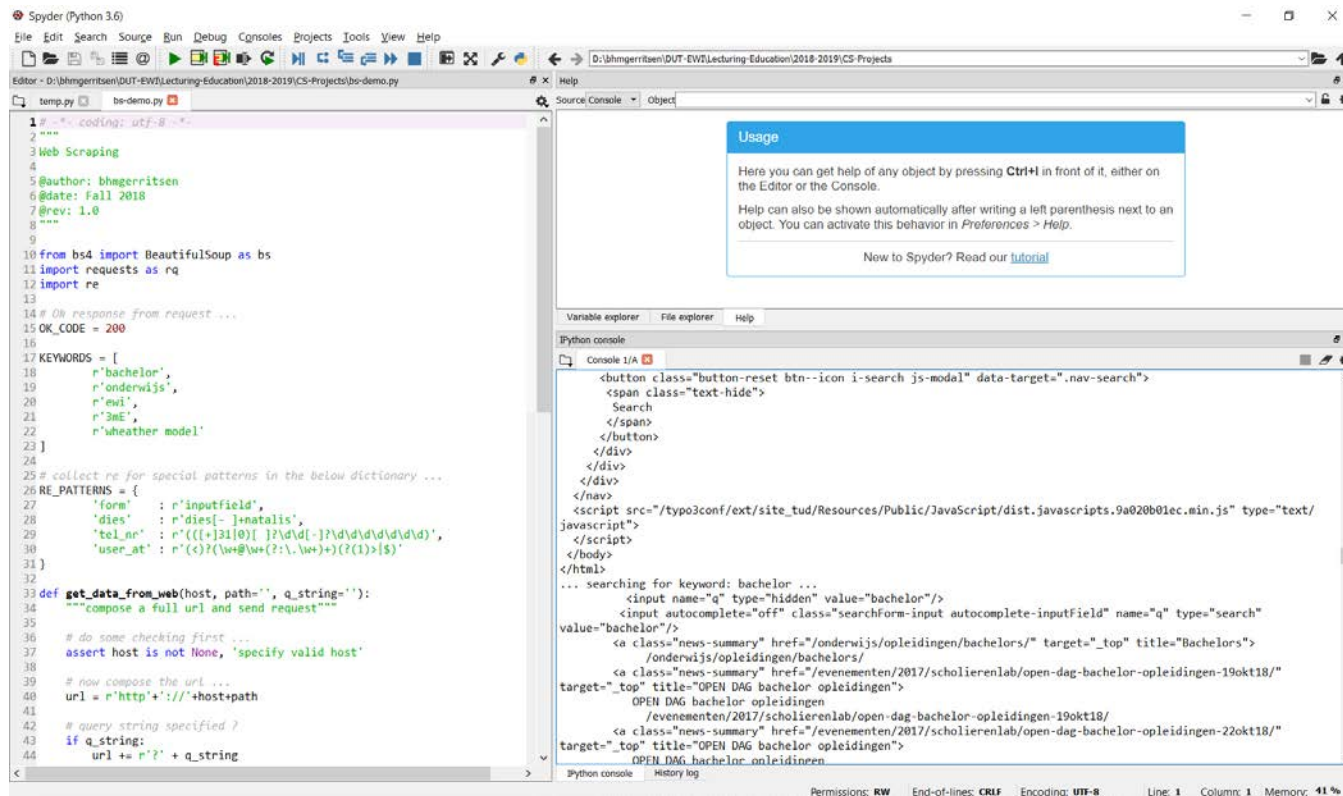
<scheme>://<user>:<password>@<host>:<port>/<path>;<params>?<query>#<frag>



name1=value1&name  
2=value2&...

<https://www.google.com/search?q=TU+Delft&tbm=isch>

# Using Python and BeautifulSoup



The screenshot displays the Spyder Python IDE interface. The main editor window shows a Python script named `temp.py` with the following code:

```
1 # -*- coding: utf-8 -*-
2 """
3 Web Scraping
4
5 @author: bhmgerritsen
6 @date: Fall 2018
7 @rev: 1.0
8 """
9
10 from bs4 import BeautifulSoup as bs
11 import requests as rq
12 import re
13
14 # On response from request ...
15 OK_CODE = 200
16
17 KEYWORDS = [
18     r'bachelor',
19     r'onderwijs',
20     r'ewi',
21     r'2m',
22     r'weather model'
23 ]
24
25 # collect re for special patterns in the below dictionary ...
26 RE_PATTERNS = {
27     'form': r'inputfield',
28     'dies': r'dies[- ]natalis',
29     'tel_nr': r'([+3110])? \d\d(-)? \d\d\d\d\d\d\d\d',
30     'user_at': r'(<?(\w@(\w+)?(\.|\w+))?(?!(?>))$)'
31 }
32
33 def get_data_from_web(host, path='', q_string=''):
34     """compose a full url and send request"""
35
36     # do some checking first ...
37     assert host is not None, 'specify valid host'
38
39     # now compose the url ...
40     url = r'http' + '://' + host + path
41
42     # query string specified ?
43     if q_string:
44         url += r'?' + q_string
```

The console window on the right shows the output of the script, displaying a snippet of HTML from a search results page:

```
<button class="button-reset btn--icon i-search js-modal" data-target=".nav-search">
<span class="text-hide">
  Search
</span>
</button>
</div>
</div>
</div>
</nav>
<script src="/typo3conf/ext/site_tud/Resources/Public/JavaScript/dist.javascripts.9a020b01ec.min.js" type="text/
javascript">
</script>
</body>
</html>
... searching for keyword: bachelor ...
<input name="q" type="hidden" value="bachelor"/>
<input autocomplete="off" class="searchform-input autocomplete-inputfield" name="q" type="search"
value="bachelor"/>
<a class="news-summary" href="/onderwijs/opleidingen/bachelors/" target="_top" title="Bachelors">
  /onderwijs/opleidingen/bachelors/
<a class="news-summary" href="/evenementen/2017/scholierenlab/open-dag-bachelor-opleidingen-19okt18/"
target="_top" title="OPEN DAG bachelor opleidingen">
  OPEN DAG bachelor opleidingen
  /evenementen/2017/scholierenlab/open-dag-bachelor-opleidingen-19okt18/
<a class="news-summary" href="/evenementen/2017/scholierenlab/open-dag-bachelor-opleidingen-22okt18/"
target="_top" title="OPEN DAG bachelor opleidingen">
  OPEN DAG bachelor onldingen
```

The status bar at the bottom indicates the file permissions (RW), end-of-line characters (CRLF), encoding (UTF-8), and current position (Line: 1, Column: 1, Memory: 41 %).

# HTTP methods

<b>GET</b>	Get a document from the Web server.
<b>HEAD</b>	Get the header of a document from the Web server.
<b>POST</b>	Send data from the client to the server for processing.
<b>PUT</b>	Save the body of the request on the server.
<b>TRACE</b>	Trace the message through proxy servers to the server.
<b>OPTIONS</b>	Determine what methods can operate on a server.
<b>DELETE</b>	Remove a document from a Web server.

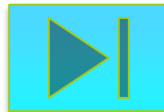
Servers may implement more or fewer methods than shown.



# Status codes

Also see:

<https://www.w3.org/Protocols/rfc2616/rfc2616-sec10.html>



## 1xx: HTTP Informational Codes

100	Continue
101	Switching Protocols
102	Processing WebDAV
103	Checkpoint draft POST PUT
122	Request-URI too long IE7

## 2xx: HTTP Successful Codes

200	OK
201	Created
202	Accepted
203	Non-Authoritative Information 1.1
204	No Content
205	Reset Content
206	Partial Content
207	Multi-Status WebDAV 4918
208	Already Reported WebDAV 5842
226	IM Used 3229 GET

## 3xx: HTTP Redirection Codes

300	Multiple Choices
301	Moved Permanently
302	Found
303	See Other 1.1
304	Not Modified
305	Use Proxy 1.1
306	Switch Proxy unused
307	Temporary Redirect 1.1
308	Permanent Redirect 7538

307 and 308 are similar to 302 and 301, but the new request method after redirect must be the same, as on initial request.

## 4xx: HTTP Client Error Code

400	Bad Request
401	Unauthorized
402	Payment Required res
403	Forbidden
404	Not Found
405	Method Not Allowed
406	Not Acceptable
407	Proxy Authentication Required
408	Request Timeout
409	Conflict
410	Gone
411	Length Required
412	Precondition Failed
413	Request Entity Too Large
414	Request-URI Too Long
415	Unsupported Media Type
416	Requested Range Not Satisfiable
417	Expectation Failed
418	I'm a teapot 2324
422	Unprocessable Entity WebDAV 4918
423	Locked WebDAV 4918
424	Failed Dependency WebDAV 4918
425	Unordered Collection 3648
426	Upgrade Required 2817
428	Precondition Required draft
429	Too Many Requests draft
431	Request Header Fields Too Large draft
444	No Response nginx
449	Retry With MS
450	Blocked By Windows Parental Controls MS
451	Unavailable For Legal Reasons draft
499	Client Closed Request nginx

## 5xx: HTTP Server Error Codes

500	Internal Server Error
501	Not Implemented
502	Bad Gateway
503	Service Unavailable
504	Gateway Timeout
505	HTTP Version Not Supported
506	Variant Also Negotiates 2295
507	Insufficient Storage WebDAV 4918
508	Loop Detected WebDAV 5842
509	Bandwidth Limit Exceeded nostd
510	Not Extended 2774
511	Network Authentication Required draft
598	Network read timeout error nostd
599	Network connect timeout error nostd

## HTTP Code Comments

WebDAV	WebDAV extension
1.1	HTTP/1.1
GET, POST, PUT, POST	For these methods only
IE	IE extension
MS	MS extension
nginx	nginx extension
2518, 2817, 2295, 2774, 3229, 4918, 5842	RFC number
draft	Proposed draft
nostd	Non standard extension
res	Reserved for future use
unused	No more in use, deprecated

Wikipedia was used to produce all HTTP codes content:  
[http://en.wikipedia.org/wiki/HTTP\\_status](http://en.wikipedia.org/wiki/HTTP_status)



# HTTP headers

<b>Content-Type</b>	Entity type
<b>Content-Length</b>	Length/size of the message
Content-Language	Language of the entity sent (e.g. English)
<b>Content-Encoding</b>	Data transformations applied to the entity
Content-Location	Alternative location of the entity
Content-Range	For partial entities, range defines the pieces sent
<b>Content-MD5</b>	Checksum of the content
<b>Last-Modified</b>	Date on which this entity was created/modified
<b>Expires</b>	Date at which the entity will become stale
Allow	Lists the legal request methods for the entity

**Important:** Entity bodies only contain **raw** data, **header** information required to **interpret** the data.

See: [http://www.ntu.edu.sg/home/ehchua/programming/webprogramming/http\\_basics.html](http://www.ntu.edu.sg/home/ehchua/programming/webprogramming/http_basics.html)

# User-related HTTP header fields

<b>From</b>	Request	User's email address	mostly Web crawler
<b>User-Agent</b>	Request	User's browser	device customization
<b>Referer</b>	Request	Page the user came from	user interests
<b>Client-IP</b>	Request (Extension)	Client's IP address	
<b>Authorization</b>	Request	Username & password	

# Authentication

- Username + password
- HTTP headers **WWW-Authenticate** en **Authorization**
- HTTP is stateless: once logged in, the client sends the login info with each request
- More in the Security Lecture

# URL – Syntax

- Uniform resource locators offer a standardized way to point to any resource on the Internet
- Not restricted to the http scheme, syntax slightly varies from scheme to scheme
- General format (adhered to by most schemes):

```
<scheme>://<user>:<password>@<host>:<port>/<path>;<params>?<query>#<frag>
```

# URL – Syntax details

The name of a piece of a resource. Only used by the client - the fragment is not transmitted to the server.

Parameters passed to gateway resources, i.e. applications [identified by the path] such as search engines.

Additional input parameters applications may require to access a resource on the server correctly. Can be set per path segment.

the local path to the resource

the port on which the server is expecting requests for the resource

domain name (host name) or numeric IP address of the server

the username/password (may be necessary to access a resource)

determines the protocol to use when connecting to the server.

`<scheme>://<user>:<password>@<host><port>/<path>;<params>?<query>#<frag>`

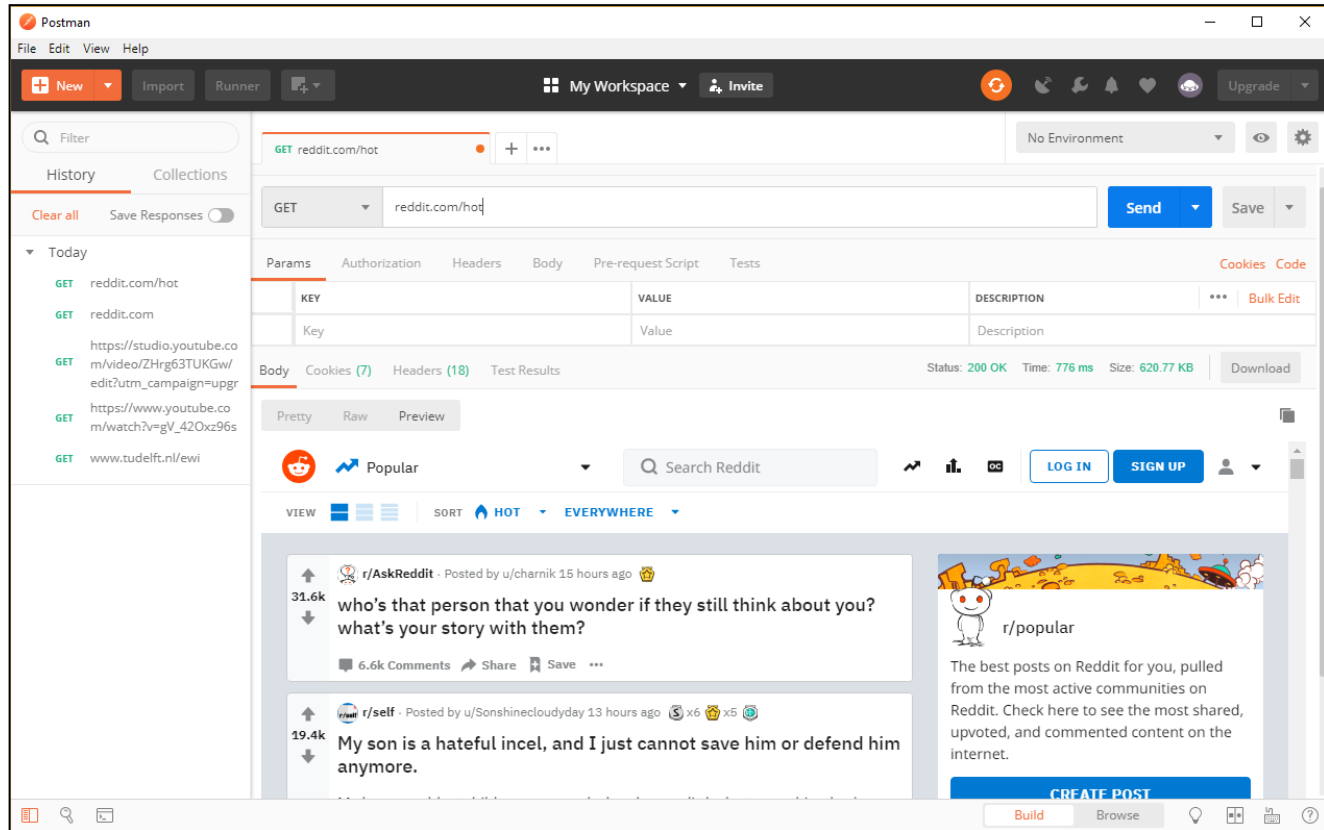
# URL– Syntax query (cont'd)

<https://www.google.com/search?q=TU+Delft&tbm=isch>

- Query component is passed to the application accessed at web server ('gateway resource')
- Enables interactive web applications
- Pattern: `name1=value1&name2=value2&...`

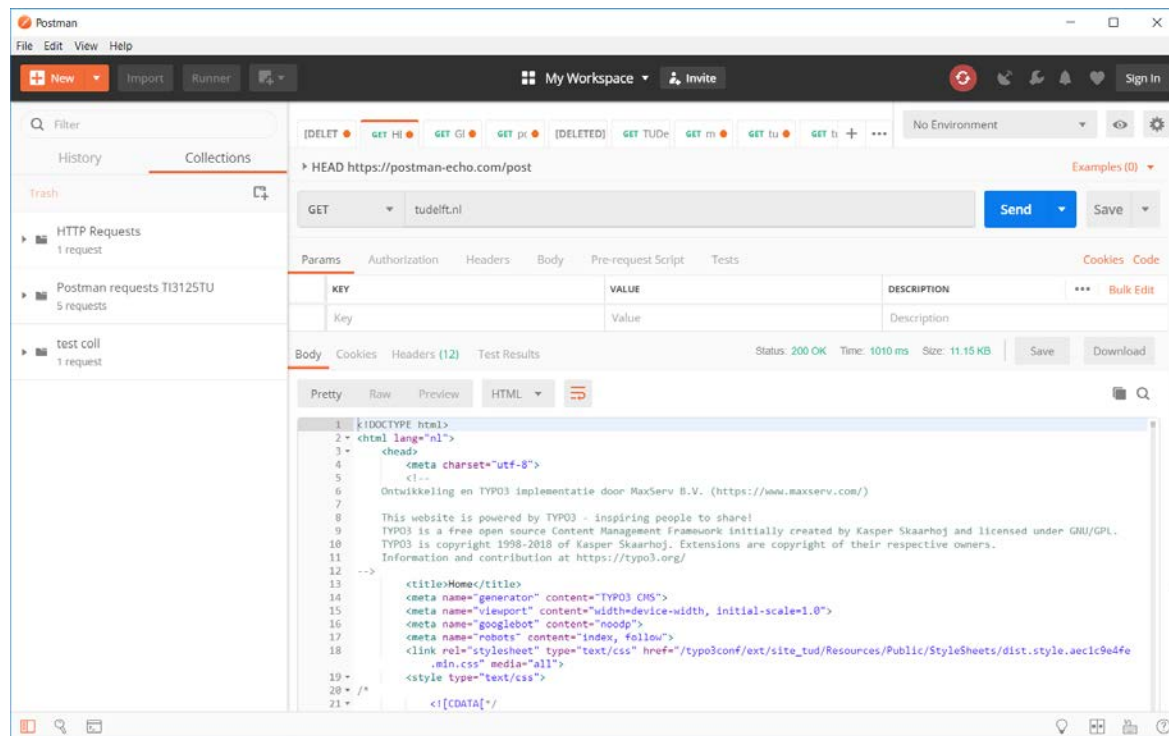


# Postman



See: <https://steelkiwi.com/blog/api-testing-useful-tools-postman-tutorial-and-hints/>

# Postman demo



# Example

## ONLINE ANALYSIS OF ENGINEERING FRAMEWORKS FOR VARIOUS LIFE CYCLE STAGES

### ONLINE ANALYSIS OF ENGINEERING FRAMEWORKS FOR VARIOUS LIFE CYCLE STAGES

Conference Paper (PDF Available) · October 2010 with 164 Reads

Conference: IRIT, VSST 2010, At Toulouse, F, Volume: 1

[Cite this publication](#)



Bart H.M. Gerritsen

at 15.13 · Delft University of Technology



Bart Gerritsen@planet

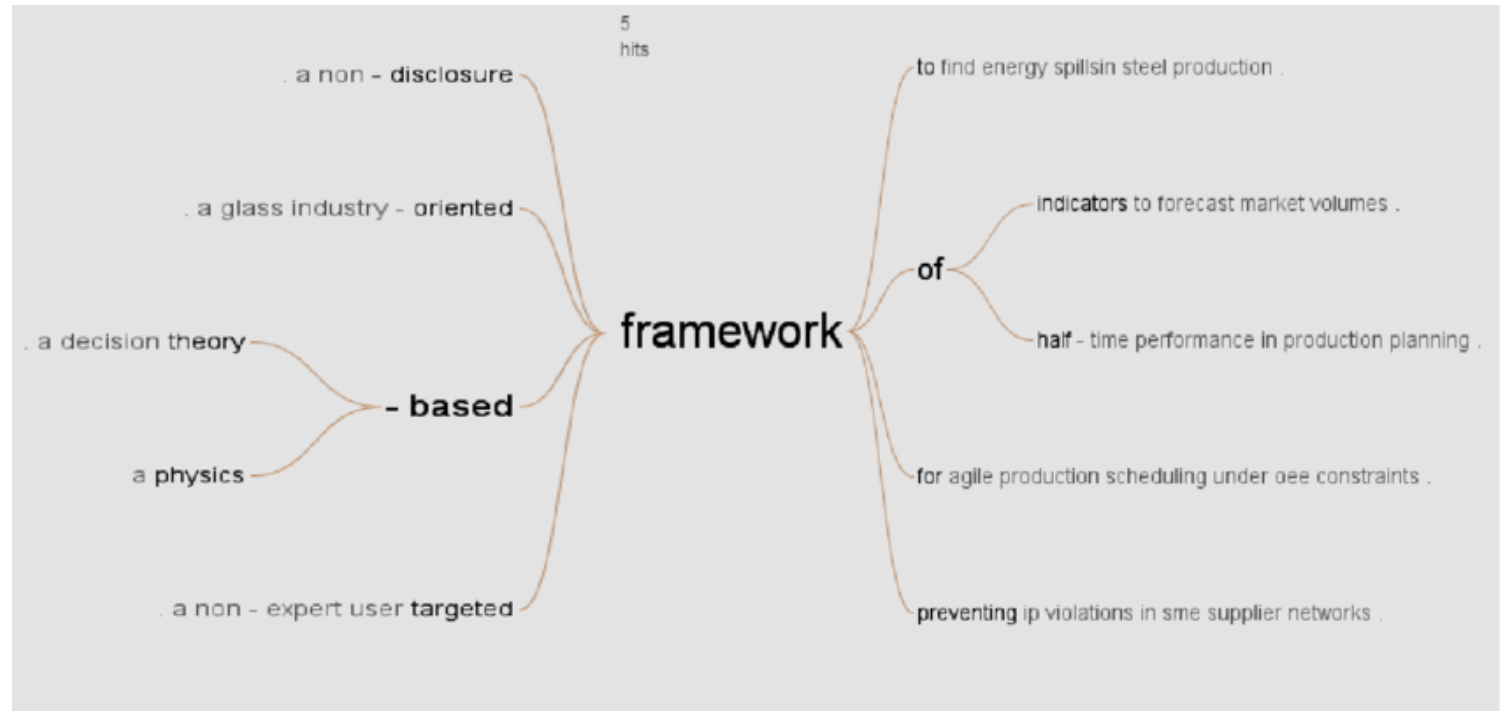


NI

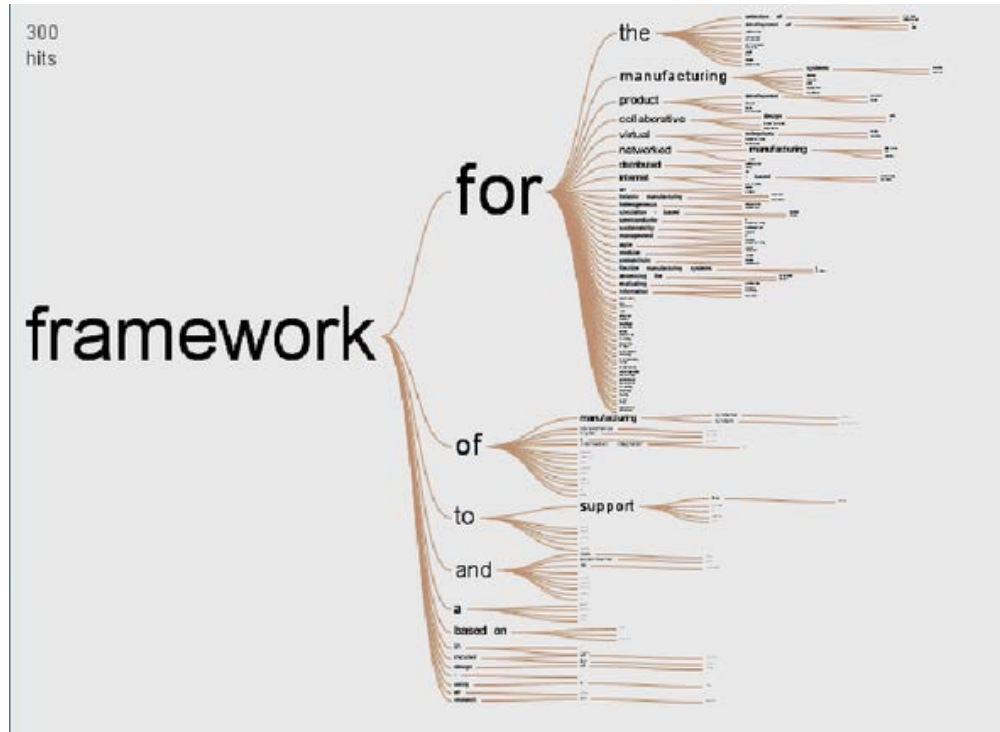
#### Abstract

Effective, adaptable and extendable frameworks can be regarded a key element for future sustainable whole-life holistic product approaches, like the development of intelligent products and the internet-of-things, agile manufacturing, smart product bundling, closed-loop life cycle management, asset management, etc. Frameworks also play a vital role in enterprise architecture and business organization. Peculiar enough and in contrast with model, technology and knowledge, framework properties form an oddly neglected engineering research field, despite the fact that frameworks are being proposed by the dozens. The question arises how to apply, reuse and extend all these frameworks in a well-formed manner so as to be able to verify, implement and build upon earlier results in a harsh industrial setting. Moreover, framework analysis, if somehow structured, can complement to standards, protocols, ontology's and other formalizations. We carried out a comprehensive bibliographic study into frameworks and applied part-of-speech-based analysis on framework properties. To that extend, we mined the available meta data of thousands of engineering frameworks disclosed on the internet in the period after 2000 and developed methods and techniques to classify them according to indicative factors for goals, resources, application context, etc. To enable searching, matching, comparison and ranking, we explored measures indicative of similarity in (part of) the framework properties. More specifically, we calculated textual energy-based association strengths to determine cross-context framework applicability. This is believed to be a necessary first step towards the formation of federations of frameworks and contract-based transformation and frameworks. Results shown in this paper are promising, but a supporting (semi-)structure like a framework ontology, is believed to further raise the online search and analysis potential of our method.

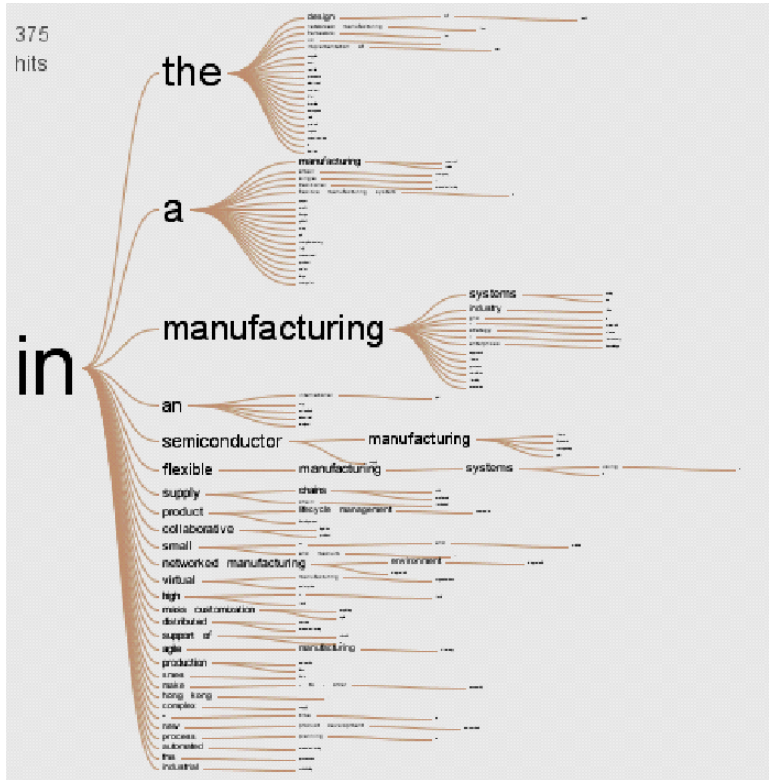
# Looking for patterns ...



# Looking for patterns (2)

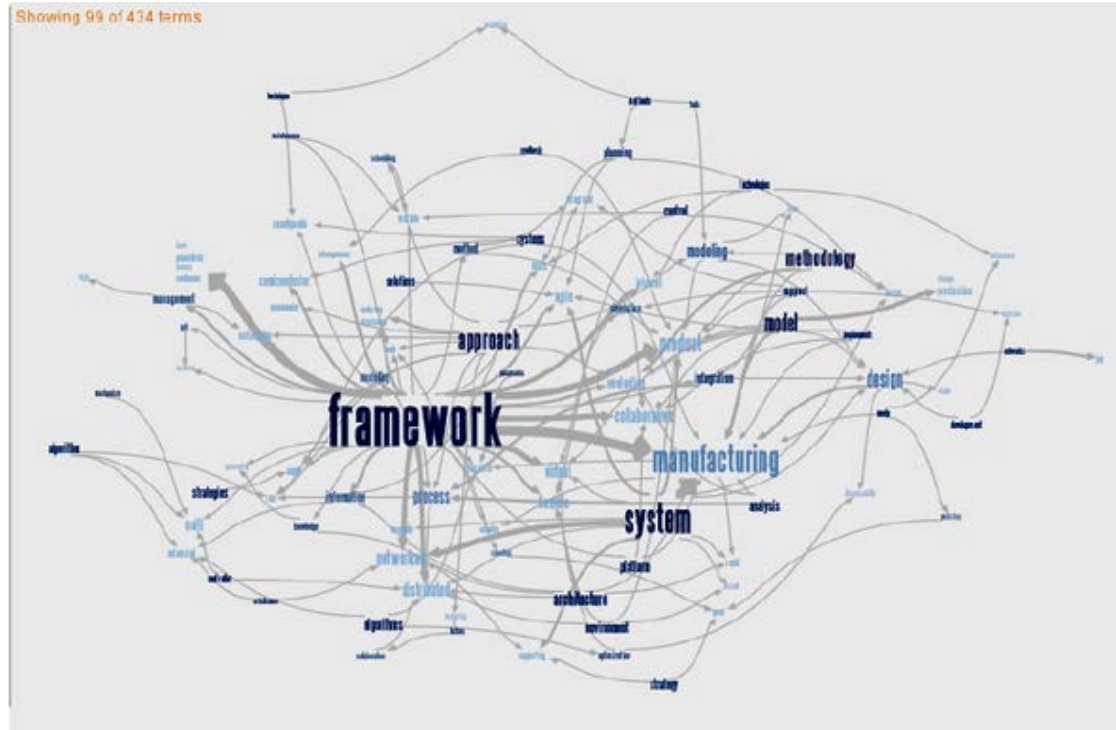


# Pattern element analysis





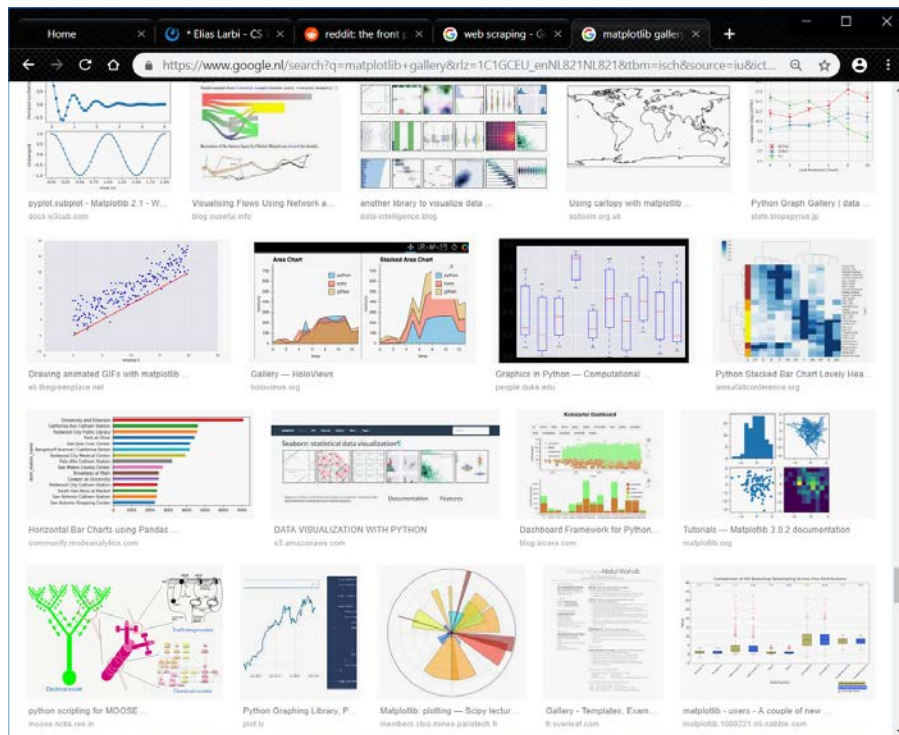
# Looking for patterns (3)



# Math analysis (Markov model)

$d(t_0) = 3$			$d(t_0) = 2$			$d(t_0) = 1$			$t_0 = CT$	$d(t_1) = 1$			$d(t_1) = 2$			$d(t_1) = 3$			$d(t_1)$
Term	freq	%	Term	freq	%	Term	freq	%		Term	freq	%	Term	freq	%	Term	freq	%	Term
	13	6.2%	A	27	12.6%	A	34	16.1%	F	for	100	47.4%	the	16	7.6%	design	19	8.0%	of
	7	3.3%	An	10	4.7%	modeling	8	3.3%	F	to	19	9.0%	product	12	5.7%	product	8	3.8%	and
	6	2.4%	and	7	3.3%	integration	8	3.3%	F	of	17	8.1%	support	11	5.2%	and	6	2.8%	design
	5	2.4%	product	7	3.3%	a	7	3.3%	F	and	13	6.2%	on	5	2.4%	the	5	2.4%	information
ign	5	2.4%	an	5	2.4%	design	7	3.3%	F	based	7	3.3%	collaborative	5	2.4%	of	5	2.4%	development
	4	1.9%	design	5	2.4%	modeling	5	2.4%	F	with	4	1.9%	concurrent	3	1.4%	integration	5	2.4%	information
	3	1.4%	the	4	1.9%	system	5	2.4%	F	in	3	1.4%	Web-based	3	1.4%	development	4	1.9%	for
em	3	1.4%	information	4	1.9%	conceptual	5	2.4%	F	model	2	0.9%	evaluation	3	1.4%	a	4	1.9%	product
	3	1.4%	a	4	1.9%	integrated	5	2.4%	F	and	2	0.9%	for	3	1.4%	management	3	1.4%	platform
verio	2	0.9%	-	3	1.4%	theoretical	4	1.9%	F	supporting	2	0.9%	development	2	0.9%	to	2	0.9%	engineering
I	2	0.9%	of	3	1.4%	engineering	4	1.9%	F				to	2	0.9%	function	2	0.9%	concurrent
ly	2	0.9%	on	3	1.4%	software	4	1.9%	F				integrating	2	0.9%	support	2	0.9%	knowledge
duct	2	0.9%	data	3	1.4%	and	3	1.4%	F				its	2	0.9%	communication	2	0.9%	sharing
ributed	2	0.9%	function	2	0.9%	the	3	1.4%	F				virtual	2	0.9%	tool	2	0.9%	integration
manufacturing	2	0.9%	via	2	0.9%	recovery	3	1.4%	F				optimal	2	0.9%	distributed	2	0.9%	data
	2	0.9%	processes	2	0.9%	unified	3	1.4%	F				knowledge	2	0.9%	planning	2	0.9%	environment
			knowledge	2	0.9%	optimization	3	1.4%	F				new	2	0.9%	knowledge	2	0.9%	management
			Part	2	0.9%	management	3	1.4%	F				integrate	2	0.9%	engineering	2	0.9%	
			Research	2	0.9%	development	2	0.9%	F				design	2	0.9%	conceptual	2	0.9%	
			cost	2	0.9%	generic	2	0.9%	F							lifecycle	2	0.9%	
			models	2	0.9%	decision-making	2	0.9%	F							rapid	2	0.9%	
			design-decision	2	0.9%	making	2	0.9%	F							collaborative	2	0.9%	
			decision	2	0.9%	support	2	0.9%	F							products	2	0.9%	
						information	2	0.9%	F							enterprises	2	0.9%	
						on	2	0.9%	F										
						Internet-based	2	0.9%	F										
						component	2	0.9%	F										
						reference	2	0.9%	F										
						exchange	2	0.9%	F										
						collaborative	2	0.9%	F										
						general	2	0.9%	F										
									211										
SING	79	37.4%	MISSING	10	7.6%	MISSING	9	0.0%	-	MISSING	33	15.6%	MISSING	33	15.6%	MISSING	38	18.0%	MISSING

# Matplotlib



# Some pointers to start ...

- <http://duspviz.mit.edu/tutorials/python-scraping>
- <https://www.dataquest.io/blog/web-scraping-tutorial-python/>
- <https://www.dataquest.io/blog/web-scraping-beautifulsoup/>

