

---

# **advertools Documentation**

*Release 0.10.7*

**Elias Dabbas**

**Sep 26, 2020**



# CONTENTS

<b>1</b>	<b>Online marketing productivity and analysis tools</b>	<b>1</b>
1.1	advertools: productivity & analysis tools to scale your online marketing	1
1.2	Generate Keywords for SEM Campaigns	4
1.3	Create Ads on a Large Scale	8
1.4	Create Ads Using Long Descriptive Text (top-down approach)	9
1.5	robots.txt Tester for Large Scale Testing	12
1.6	Download, Parse, and Analyze XML Sitemaps	16
1.7	Python SEO Crawler / Spider	22
1.8	SEO Crawling & Scraping: Strategies & Recipes	31
1.9	Import Search Engine Results Pages (SERPs) for Google and YouTube	35
1.10	Import and Analyze Knowledge Graph Results on a Large Scale	42
1.11	Split, Parse, and Analyze URLs	45
1.12	Emoji: Extract, Analyze, and Get Insights	48
1.13	Extract structured entities from text lists	50
1.14	Stopwords in Several Languages	61
1.15	Text Analysis	63
1.16	Tokenize Words (N-grams)	66
1.17	Twitter Data API	67
1.18	YouTube Data API	80
<b>2</b>	<b>Indices and tables</b>	<b>99</b>
2.1	advertools	99
	<b>Python Module Index</b>	<b>107</b>
	<b>Index</b>	<b>109</b>



## ONLINE MARKETING PRODUCTIVITY AND ANALYSIS TOOLS

Crawl websites, Generate keywords for SEM campaigns, create text ads on a large scale, analyze multiple SERPs at once, gain insights from large social media posts, and get productive as an online marketer.

If these are things you are interested in, then this package might make your life a little easier.

**New:** `SEO crawler` now extracts JSON-LD canonical, alternate href, alternate hreflang, OpenGraph, and Twitter cards if available on pages

**New:** `knowledge_graph` Function for connecting to Google's Knowledge Graph Data API

**Faster:** `sitemap_to_df` Function for downloading & parsing XML sitemaps into DataFrames, is much faster now

### 1.1 `advertools`: productivity & analysis tools to scale your online marketing

A digital marketer is a data scientist.

Your job is to manage, manipulate, visualize, communicate, understand, and make decisions based on data.

You might be doing basic stuff, like copying and pasting text on spread sheets, you might be running large scale automated platforms with sophisticated algorithms, or somewhere in between. In any case your job is all about working with data.

As a data scientist you don't spend most of your time producing cool visualizations or finding great insights. The majority of your time is spent wrangling with URLs, figuring out how to stitch together two tables, hoping that the dates, won't break, without you knowing, or trying to generate the next 124,538 keywords for an upcoming campaign, by the end of the week!

`advertools` is a Python package that can hopefully make that part of your job a little easier.

### 1.1.1 Installation

```
pip install adverttools
# OR:
pip3 install adverttools
```

### SEM Campaigns

The most important thing to achieve in SEM is a proper mapping between the three main elements of a search campaign

**Keywords** (the intention) -> **Ads** (your promise) -> **Landing Pages** (your delivery of the promise) Once you have this done, you can focus on management and analysis. More importantly, once you know that you can set this up in an easy way, you know you can focus on more strategic issues. In practical terms you need two main tables to get started:

- **Keywords:** You can [generate keywords](#) (note I didn't say research) with the `kw_generate` function.
- **Ads:** There are two approaches that you can use:
  - **Bottom-up:** You can create text ads for a large number of products by simple replacement of product names, and providing a placeholder in case your text is too long. Check out the `ad_create` function for more details.
  - **Top-down:** Sometimes you have a long description text that you want to split into headlines, descriptions and whatever slots you want to split them into. `ad_from_string` helps you accomplish that.
- **Tutorials and additional resources**
  - [Setting a full SEM campaign](#) for DataCamp's website tutorial
  - [Project to practice generating SEM keywords with Python](#) on DataCamp
  - [Setting up SEM campaigns on a large scale](#) tutorial on SEMrush
  - [Visual tool to generate keywords](#) online based on the `kw_generate` function

### SEO

Probably the most comprehensive online marketing area that is both technical (crawling, indexing, rendering, redirects, etc.) and non-technical (content creation, link building, outreach, etc.). Here are some tools that can help with your SEO

- **SEO crawler:** A generic SEO crawler that can be customized, built with Scrapy, & with several features:
  - Standard SEO elements extracted by default (title, header tags, body text, status code, reponse and request headers, etc.)
  - CSS and XPath selectors: You probably have more specific needs in mind, so you can easily pass any selectors to be extracted in addition to the standard elements being extracted
  - Custom settings: full access to Scrapy's settings, allowing you to better control the crawling behavior (set custom headers, user agent, stop spider after x pages, seconds, megabytes, save crawl logs, run jobs at intervals where you can stop and resume your crawls, which is ideal for large crawls or for continuous monitoring, and many more options)
  - Following links: option to only crawl a set of specified pages or to follow and discover all pages through links
- **robots.txt downloader** A simple downloader of robots.txt files in a DataFrame format, so you can keep track of changes across crawls if any, and check the rules, sitemaps, etc.

- [XML Sitemaps downloader / parser](#) An essential part of any SEO analysis is to check XML sitemaps. This is a simple function with which you can download one or more sitemaps (by providing the URL for a robots.txt file, a sitemap file, or a sitemap index)
- [SERP importer and parser for Google & YouTube](#) Connect to Google's API and get the search data you want. Multiple search parameters supported, all in one function call, and all results returned in a DataFrame
- Tutorials and additional resources
  - A visual tool built with the `serp_goog` function to get [SERP rankings on Google](#)
  - A tutorial on [analyzing SERPs on a large scale with Python on SEMrush](#)
  - [SERP datasets on Kaggle](#) for practicing on different industries and use cases
  - [SERP notebooks on Kaggle](#) some examples on how you might tackle such data
  - [Content Analysis with XML Sitemaps and Python](#)
  - XML dataset examples: [news sites](#), [Turkish news sites](#), [Bloomberg news](#)

## Text & Content Analysis (for SEO & Social Media)

URLs, page titles, tweets, video descriptions, comments, hashtags are some examples of the types of text we deal with. `advertools` provides a few options for text analysis

- [Word frequency](#) Counting words in a text list is one of the most basic and important tasks in text mining. What is also important is counting those words by taking in consideration their relative weights in the dataset. `word_frequency` does just that.
- [URL Analysis](#) We all have to handle many thousands of URLs in reports, crawls, social media extracts, XML sitemaps and so on. `url_to_df` converts your URLs into easily readable DataFrames.
- [Emoji](#) Produced with one click, extremely expressive, highly diverse (3k+ emoji), and very popular, it's important to capture what people are trying to communicate with emoji. Extracting emoji, get their names, groups, and sub-groups is possible. The full emoji database is also available for convenience, as well as an `emoji_search` function in case you want some ideas for your next social media or any kind of communication
- [extract\\_functions](#) The text that we deal with contains many elements and entities that have their own special meaning and usage. There is a group of convenience functions to help in extracting and getting basic statistics about structured entities in text; emoji, hashtags, mentions, currency, numbers, URLs, questions and more. You can also provide a special regex for your own needs.
- [Stopwords](#) A list of stopwords in forty different languages to help in text analysis.
- [Tutorial on DataCamp for creating the word\\_frequency function and explaining the importance of the difference between absolute and weighted word frequency](#)
- [Text Analysis for Online Marketers](#) An introductory article on SEMrush

## Social Media

In addition to the text analysis techniques provided, you can also connect to the Twitter and YouTube data APIs. The main benefits of using `advertools` for this:

- Handles pagination and request limits: typically every API has a limited number of results that it returns. You have to handle pagination when you need more than the limit per request, which you typically do. This is handled by default
- DataFrame results: APIs send you back data in a format that needs to be parsed and cleaned so you can more easily start your analysis. This is also handled automatically

- Multiple requests: in YouTube's case you might want to request data for the same query across several countries, languages, channels, etc. You can specify them all in one request and get the product of all the requests in one response
- Tutorials and additional resources
- A visual tool to [check what is trending on Twitter](#) for all available locations
- A [Twitter data analysis dashboard](#) with many options
- [How to use the Twitter data API with Python](#)
- [Extracting entities from social media posts tutorial on Kaggle](#)
- [Analyzing 131k tweets by European Football clubs tutorial on Kaggle](#)
- An overview of the [YouTube data API with Python](#)

### 1.1.2 Conventions

Function names mostly start with the object you are working on, so you can use autocomplete to discover other options:

`kw_`: for keywords-related functions

`ad_`: for ad-related functions

`url_`: URL tracking and generation

`extract_`: for extracting entities from social media posts (mentions, hashtags, emoji, etc.)

`emoji_`: emoji related functions and objects

`twitter`: a module for querying the Twitter API and getting results in a DataFrame

`youtube`: a module for querying the YouTube Data API and getting results in a DataFrame

`serp_`: get search engine results pages in a DataFrame, currently available: Google and YouTube

`crawl`: a function you will probably use a lot if you do SEO

`*_to_df`: a set of convenience functions for converting to DataFrames (XML sitemaps, robots.txt files, and lists of URLs)

To install adverttools, run the following from the command line:

```
pip install adverttools
# OR:
pip3 install adverttools
```

## 1.2 Generate Keywords for SEM Campaigns

A big part of setting up SEM campaigns consists of generating keywords, and properly mapping them to landing pages and ads, as well as putting them in the right campaign and ad group structure.

Keyword research is the part of this task that takes the most time. It is very tedious, yet extremely important.

The shift here is that we are going to be *generating* keywords as opposed to researching them.

What is a keyword anyway?

It is basically a phrase that contains two things:



**Product** This is the thing that you are selling. It is simply the name of it. “barcelona”, “guitar”, “rio de janeiro”, “accounting”. The product on its own is not enough for us to understand what the user is looking for. “barcelona trips” and “barcelona football club” are completely different “keywords” for example.

**Word** To give meaning to the product, it has to come with a word. The word can be a verb like “buy” or “purchase”, and it can also be another noun, but with a clear intent expressed; “price” and “offers” for example clearly show purchase intent.

So, to *generate* keywords we need phrases that contain both, the product and the descriptive word(s). It is very easy to get the products as you know what you sell. The next thing you need to come up with are the words that work within your strategy. The most import idea here is that once you determine that you sell courses for example, there aren’t really that many words that can describe that intent; course, courses, tutorial, certification, learn, learning, education, etc. How many can you come up with? How many exist in any language? Fifteen, twenty? Once you have those are basically done.

Depending on what service you provide and what segment of the market you target it shouldn’t be difficult to come up with ideas for words (not keywords yet). You might have an e-commerce site, but want to mainly focus on cheap and discounted products. Or maybe you have luxury items, and want to exclude words that signify price sensitivity.

Let’s say you have a job site and you know that you provide jobs for engineering, graphic design, and marketing. The words are easy to come up with; “job”, “jobs”, “careers”, “vacancies”, “full time”, “part time”, “work”, and so on.

Now what we can do is use the *kw\_generate* function to come up with all possible combinations (order doesn’t matter) and/or permutations (order matters) and get a ready-to-use table to upload and start running the campaign.

```
>>> products = ['engineering', 'graphic design', 'marketing']
>>> words = ['jobs', 'careers', 'vacancies', 'full time', 'part time']
>>> adv.kw_generate(products, words)
  Campaign  Ad Group      Keyword      Criterion Type
↳ Labels
0  SEM_Campaign  Engineering      engineering jobs      Exact
↳ Jobs
1  SEM_Campaign  Engineering      engineering jobs      Phrase
↳ Jobs
2  SEM_Campaign  Engineering      +engineering +jobs      Broad
↳ Jobs
3  SEM_Campaign  Engineering      engineering careers      Exact
↳ Careers
4  SEM_Campaign  Engineering      engineering careers      Phrase
↳ Careers
..      ...      ...      ...      ...
↳ ...
625 SEM_Campaign  Marketing      part time vacancies marketing      Phrase
↳Part Time;Vacancies
626 SEM_Campaign  Marketing      +part +time +vacancies +marketing      Broad
↳Part Time;Vacancies
627 SEM_Campaign  Marketing      part time full time marketing      Exact
↳Part Time;Full Time
628 SEM_Campaign  Marketing      part time full time marketing      Phrase
↳Part Time;Full Time
629 SEM_Campaign  Marketing      +part +time +full +time +marketing      Broad
↳Part Time;Full Time
[630 rows x 5 columns]
```

And we’re done!

Check the *kw\_generate()* function for more options and details. Once you have your keywords done, you can start creating ads using either the *ad\_create* function (bottom-up approach) or the *ad\_from\_string* function (top-down

approach).

**kw\_broad** (*words*)

Return words in broad match.

**Parameters** **words** (*list*) – list of strings

**Returns formatted** words in broad match type

```
>>> keywords = ['[learn guitar]', '"guitar courses"', '+guitar +tutor']
>>> kw_broad(keywords)
['learn guitar', 'guitar courses', 'guitar tutor']
```

**kw\_exact** (*words*)

Return words in exact match.

**Parameters** **words** (*list*) – list of strings

**Returns formatted** words in exact match type

```
>>> keywords = ['learn guitar', 'guitar courses', 'guitar tutor']
>>> kw_exact(keywords)
[['learn guitar'], ['guitar courses'], ['guitar tutor']]
```

**kw\_generate** (*products, words, max\_len=3, match\_types='Exact', 'Phrase', 'Modified', capitalize\_adgroups=True, order\_matters=True, campaign\_name='SEM\_Campaign'*)  
Generate a data frame of keywords using a list of products and relevant words.

**Parameters**

- **products** (*list*) – will be used as the names of the ad groups
- **words** (*list*) – related words that make it clear that the user is interested in products
- **max\_len** (*int*) – the maximum number of words to include in each permutation of final keywords
- **match\_types** (*list*) – one or more of ('Exact', 'Phrase', 'Modified', 'Broad')
- **capitalize\_adgroups** (*bool*) – whether or not to set adgroup names in the “Ad Group” column to title case or keep them as is, default True
- **order\_matters** (*bool*) – whether or not the order of words in keywords matters, default False
- **campaign\_name** (*str*) – name of campaign

**Returns** **keywords\_df** a pandas.DataFrame ready to upload

```
>>> import advertools as adv
>>> products = ['bmw', 'toyota']
>>> words = ['buy', 'second hand']
>>> kw_df = adv.kw_generate(products, words)
>>> kw_df.head()
   Campaign Ad Group      Keyword Criterion Type      Labels
0  SEM_Campaign   Bmw    bmw buy      Exact      Buy
1  SEM_Campaign   Bmw    bmw buy      Phrase      Buy
2  SEM_Campaign   Bmw  +bmw +buy      Broad      Buy
3  SEM_Campaign   Bmw  bmw second hand      Exact  Second Hand
4  SEM_Campaign   Bmw  bmw second hand      Phrase  Second Hand
```

```
>>> kw_df.tail()
      Campaign Ad Group      Keyword Criterion Type
↳Labels
55 SEM_Campaign  Toyota  second hand toyota buy      Phrase  Second Hand;
↳Buy
56 SEM_Campaign  Toyota +second hand +toyota +buy      Broad  Second Hand;
↳Buy
57 SEM_Campaign  Toyota  second hand buy toyota      Exact  Second Hand;
↳Buy
58 SEM_Campaign  Toyota  second hand buy toyota      Phrase  Second Hand;
↳Buy
59 SEM_Campaign  Toyota +second hand +buy +toyota      Broad  Second Hand;
↳Buy
```

Sometimes you want to retain capitalization and keep it as it as is in the “Ad Group” column. This is especially important for consistency with ads DataFrames for easier integration between the two. Set *capitalize\_adgroups=False* to keep capitalization the same:

```
>>> adv.kw_generate(['SEO'], ['services', 'provider'], capitalize_adgroups=False).
↳head()
      Campaign Ad Group      Keyword Criterion Type  Labels
0 SEM_Campaign  SEO      SEO services      Exact  Services
1 SEM_Campaign  SEO      SEO services      Phrase  Services
2 SEM_Campaign  SEO +SEO +services      Broad  Services
3 SEM_Campaign  SEO      SEO provider      Exact  Provider
4 SEM_Campaign  SEO      SEO provider      Phrase  Provider
```

### **kw\_modified**(*words*)

Return words in modified broad match.

**Parameters** *words* (*list*) – list of strings

**Returns formatted** words in modified broad match type

```
>>> keywords = ['learn guitar', 'guitar courses', 'guitar tutor']
>>> kw_modified(keywords)
['+learn +guitar', '+guitar +courses', '+guitar +tutor']
```

### **kw\_neg\_broad**(*words*)

Return words in negative broad match.

**Parameters** *words* (*list*) – list of strings

**Returns formatted** words in negative broad match type

```
>>> keywords = ['learn guitar', 'guitar courses', 'guitar tutor']
>>> kw_neg_broad(keywords)
['-learn guitar', '-guitar courses', '-guitar tutor']
```

### **kw\_neg\_exact**(*words*)

Return words in negative exact match.

**Parameters** *words* (*list*) – list of strings

**Returns formatted** words in negative exact match type

```
>>> keywords = ['learn guitar', 'guitar courses', 'guitar tutor']
>>> kw_neg_exact(keywords)
['-[learn guitar]', '-[guitar courses]', '-[guitar tutor]']
```

**kw\_neg\_phrase** (*words*)

Return words in negative phrase match.

**Parameters** **words** (*list*) – list of strings

**Returns formatted** words in negative phrase match type

```
>>> keywords = ['learn guitar', 'guitar courses', 'guitar tutor']
>>> kw_neg_phrase(keywords)
['-"learn guitar"', '-"guitar courses"', '-"guitar tutor"']
```

**kw\_phrase** (*words*)

Return words in phrase match.

**Parameters** **words** (*list*) – list of strings

**Returns formatted** words in phrase match type

```
>>> keywords = ['learn guitar', 'guitar courses', 'guitar tutor']
>>> kw_phrase(keywords)
['"learn guitar"', '"guitar courses"', '"guitar tutor"']
```

## 1.3 Create Ads on a Large Scale

When creating large-scale campaigns, you also need to create ads on a large scale. For products in a similar category you typically want to use the same ads, but only replace the product name, “Get the latest <product> now”, and replace *product* as many times as you have ads.

```
>>> products = ['Dubai', 'Tokyo', 'Singapore']
>>> adv.ad_create(template='5-star Hotels in {}',
...               replacements=products,
...               max_len=30,
...               fallback='Great Cities')
['5-star Hotels In Dubai',
 '5-star Hotels In Tokyo',
 '5-star Hotels In Singapore']
```

An important thing to watch out for, is long product names. Since text ads have limits on each slot, you need to make sure you don’t exceed that limit. For this you need to provide a *fallback* text in case the product name is longer than *max\_len*.

```
>>> products = ['Dubai', 'Tokyo', 'Singapore', 'Llanfairpwllgwyngyll']
>>> adv.ad_create(template='5-star Hotels in {}',
...               replacements=products,
...               max_len=30,
...               fallback='Great Cities')
['5-star Hotels In Dubai',
 '5-star Hotels In Tokyo',
 '5-star Hotels In Singapore',
 '5-star Hotels In Great Cities']
```

**ad\_create** (*template, replacements, fallback, max\_len=30, capitalize=True*)

Insert each of the replacement strings in its place within template.

**Parameters**

- **template** (*str*) – a string format template, using braces e.g. “Get the latest {} today.”

- **replacements** (*list*) – replacement strings to be inserted in template
- **fallback** (*str*) – the string to insert in template in case replacement is longer than `max_len`
- **max\_len** (*int*) – the maximum allowed length of the full string
- **capitalize** (*bool*) – whether or not to capitalize words in the result

Returns formatted list of strings

```
>>> ad_create("Let's count {}", ['one', 'two', 'three'], 'one', 20)
["Let's Count One", "Let's Count Two", "Let's Count Three"]
```

```
>>> ad_create(template='My favorite car is {}',
...           replacements=['Toyota', 'BMW', 'Mercedes', 'Lamborghini'],
...           fallback='great',
...           max_len=28)
['My Favorite Car Is Toyota', 'My Favorite Car Is Bmw',
'My Favorite Car Is Mercedes', 'My Favorite Car Is Great']
```

```
>>> ad_create('KeEP cApITalization {}', ['As IS'],
...           fallback='fallback', max_len=50, capitalize=False)
['KeEP cApITalization As IS']
```

```
>>> ad_create('This is very long and will produce and error',
...           replacements=['something', 'long'], fallback='Very long',
...           max_len=20)
Traceback (most recent call last):
File "<input>", line 1, in <module>
File "<input>", line 26, in ad_create
ValueError: template + fallback should be <= 20 chars
```

## 1.4 Create Ads Using Long Descriptive Text (top-down approach)

Many times you have long descriptive text about your products, especially on their respective landing pages. The allowed length of text ads has become considerably long on many platforms. On Google Ads for example, you have slots of 30, 30, 30, 90, and 90 characters, for a total of 270. That's more than enough space to talk about the main features of your product.

The question is, how do you utilize that long description text that has all the details that you want, and make sure it fits correctly within the limits given by the platform you are using?

The `ad_from_string()` function does exactly that. Given a long string, it divides it into slots of any given length that you specify, and if any text remains it will be appended to the end of the returned list.

Another important benefit of this is that you can take those long descriptions (or write them) once, and then you can easily split them into different slots based on the ad format and the platform you are using.

Here is a quick overview of the available parameters and options:

**s** The string that you want to split. This would typically be available on the landing pages of each product.

**slots** The lengths that you want to split into. Note that although the default uses Google Ads' text ad template, you can change it to any other group of slots, with more or fewer slots of different lengths.

**sep** The separator by which to split the text. The default is `None` which splits the text by whitespace, but you can change it to something else if needed. Sometimes you might want the text split by hyphens (URLs for example) so you can split by that character.

**capitalize** The default is `False` which leaves the capitalization of `s` intact. If you set it to `True` then the first letter of each word would be capitalized.

### Example

Note that in any case, the returned list of characters is longer than the provided slots by one. So if you provide five slots, for example, the function will always return a list of length six.

This is to ensure that the remainder of the text is not lost if it is longer, so you know what is missing. In case you have shorter text, you will still have one element more than the provided slots to ensure consistency.

```
>>> desc_text = "Get the latest gadget online. The GX12 model comes with 13 things_
↳that do a lot of good stuff for your health. Start shopping now."
>>> len(desc_text)
130
```

Now let's see how this same description can be utilized in different scenarios

## 1.4.1 Google Text Ads

Since this is shorter than the default Google values, you will get extra empty slots (with an additional last one).

```
>>> ad_from_string(desc_text) # default values (Google text ads)
['Get the latest gadget online.',
 'The GX12 model comes with 13',
 'things that do a lot of good',
 'stuff for your health. Start shopping now.',
 '',
 '',
 '',
 '',
 '']
```

## 1.4.2 Facebook Feed Ads

In this case, it is also shorter than the default value, so you get an extra space.

```
>>> ad_from_string(desc_text, [125, 25, 30]) # Facebook feed ads
['Get the latest gadget online. The GX12 model comes with 13 things that do a lot of_
↳good stuff for your health. Start shopping',
 'now.',
 '',
 '']
```

Since it might not look good having just one word in the second slot, and an empty last one, you might want to change it as follows:

```
>>> ad_from_string(desc_text, [90, 25, 30])
['Get the latest gadget online. The GX12 model comes with 13 things that do a lot of_
↳good',
 'stuff for your health.',
 'Start shopping now.',
 '']
```

### 1.4.3 Facebook Instant Article Ad

Here is a case where our text is longer than the provided limitations, so we end up having an extra space that is not used:

```
>>> ad_from_string(desc_text, [25, 30]) # Facebook instant article ad
['Get the latest gadget',
 'online. The GX12 model comes',
 'with 13 things that do a lot of good stuff for your health. Start shopping now.']
```

**ad\_from\_string** (*s*, *slots*=30, 30, 30, 90, 90, 15, 15, *sep*=None, *capitalize*=False)

Convert string *s* to an ad by splitting it into groups of words. Each group would have a length of at most the allowed length for that slot.

If the total length of *s* exceeds the total allowed length, all remaining characters would be grouped in the last element of the returned list.

#### Parameters

- **s** (*str*) – a string of characters, with no restrictions on length
- **slots** (*list*) – an iterable of integers for the maximum lengths for each slot
- **sep** (*str*) – character(s) by which to split *s*
- **capitalize** (*bool*) – whether or not to capitalize each word after grouping. Setting it as False would not change the capitalization of the input string

**Returns text\_ad** a list of strings

```
>>> ad_from_string('this is a short ad')
['this is a short ad', '', '', '', '', '', '', '']
```

```
>>> ad_from_string('this is a longer ad and will take the first two slots')
['this as a longer ad and would', 'take the first two slots',
 '', '', '', '', '']
```

```
>>> ad_from_string("Slots can be changed the way you want", (10, 15, 10))
['Slots can', 'be changed the', 'way you', 'want']
```

```
>>> ad_from_string("The capitalization REMAInS as IS bY DefAULt",
...                (10, 15, 10))
['The', 'capitalization', 'REMAInS as', 'IS bY DefAULt']
```

```
>>> ad_from_string("set captialize=True to capitalize first letters",
...                capitalize=True)
['Set Captialize=true To', 'Capitalize First Letters',
 '', '', '', '']
```

## 1.5 robots.txt Tester for Large Scale Testing

Even though they are tiny in size, robots.txt files contain potent information that can block major sections of your site, which is what they are supposed to do. Only sometimes you might make the mistake of blocking the wrong section.

So it is very important to check if certain pages (or groups of pages) are blocked for a certain user-agent by a certain robots.txt file. Ideally, you would want to run the same check for all possible user-agents. Even more ideally, you want to be able to run the check for a large number of pages with every possible combination with user-agents!

To get the robots.txt file into an easily readable format, you can use the `robotstxt_to_df()` function to get it in a DataFrame.

```
>>> robotstxt_to_df('https://www.google.com/robots.txt')
  directive                content                robotstxt_
↳url                download_date
0    User-agent                *    https://www.google.com/robots.
↳txt 2020-06-01 14:05:16.068031+00:00
1    Disallow                /search    https://www.google.com/robots.
↳txt 2020-06-01 14:05:16.068031+00:00
2    Allow                /search/about    https://www.google.com/robots.
↳txt 2020-06-01 14:05:16.068031+00:00
3    Allow                /search/static    https://www.google.com/robots.
↳txt 2020-06-01 14:05:16.068031+00:00
4    Allow                /search/howsearchworks    https://www.google.com/robots.
↳txt 2020-06-01 14:05:16.068031+00:00
...
↳
277 User-agent                Twitterbot    https://www.google.com/robots.
↳txt 2020-06-01 14:05:16.068031+00:00
278 Allow                /imgres    https://www.google.com/robots.
↳txt 2020-06-01 14:05:16.068031+00:00
279 User-agent                facebookexternalhit    https://www.google.com/robots.
↳txt 2020-06-01 14:05:16.068031+00:00
280 Allow                /imgres    https://www.google.com/robots.
↳txt 2020-06-01 14:05:16.068031+00:00
281 Sitemap    https://www.google.com/sitemap.xml    https://www.google.com/robots.
↳txt 2020-06-01 14:05:16.068031+00:00
[282 rows x 4 columns]
```

The returned DataFrame contains columns for directives, their content, the URL of the robots.txt file, as well as the date it was downloaded. Under the *directive* column you can see the main commands; Allow, Disallow, Sitemap, Crawl-delay, User-agent, and so on. The *content* column contains the details of each of those directives (the pattern to disallow, the sitemap URL, etc.)

As for testing, the `robotstxt_test()` function runs a test for a given robots.txt file, checking which of the provided user-agents can fetch which of the provided URLs, paths, or patterns.

```
>>> robotstxt_test('https://www.example.com/robots.txt',
...                useragents=['Googlebot', 'baiduspider', 'Bingbot']
...                urls=['/', '/hello', '/some-page.html'])
```

As a result, you get a DataFrame with a row for each combination of (user-agent, URL) indicating whether or not that particular user-agent can fetch the given URL.

Some reasons why you might want to do that:

- SEO Audits: Especially for large websites with many URL patterns, and many rules for different user-agents.
- Developer or site owner about to make large changes



- Interest in strategies of certain companies

## 1.5.1 User-agents

In reality there are only two groups of user-agents that you need to worry about:

- User-agents listed in the robots.txt file: For each one of those you need to check whether or not they are blocked from fetching a certain URL (or pattern).
- \* all other user-agents: The \* includes all other user-agents, so checking the rules that apply to it should take care of the rest.

## 1.5.2 robots.txt Testing Approach

1. Get the robots.txt file that you are interested in
2. Extract the user-agents from it
3. Specify the URLs you are interested in testing
4. Run the `robotstxt_test()` function

```
>>> fb_robots = robotstxt_to_df('https://www.facebook.com/robots.txt')
>>> fb_robots
  directive                                     content
↪ robotstxt_url                download_date
0      comment Notice: Collection of data on Facebook through... https://www.
↪facebook.com/robots.txt 2020-05-31 20:12:47.576281+00:00
1      comment prohibited unless you have express written per... https://www.
↪facebook.com/robots.txt 2020-05-31 20:12:47.576281+00:00
2      comment and may only be conducted for the limited purp... https://www.
↪facebook.com/robots.txt 2020-05-31 20:12:47.576281+00:00
3      comment                                     permission. https://www.
↪facebook.com/robots.txt 2020-05-31 20:12:47.576281+00:00
4      comment See: http://www.facebook.com/apps/site_scrapin... https://www.
↪facebook.com/robots.txt 2020-05-31 20:12:47.576281+00:00
...      ...
↪      ...
461     Allow                                     /ajax/bootloader-endpoint/ https://www.
↪facebook.com/robots.txt 2020-05-31 20:12:47.576281+00:00
462     Allow /ajax/pagelet/generic.php/PagePostsSectionPagelet https://www.
↪facebook.com/robots.txt 2020-05-31 20:12:47.576281+00:00
463     Allow                                     /safetycheck/ https://www.
↪facebook.com/robots.txt 2020-05-31 20:12:47.576281+00:00
464 User-agent                                     * https://www.
↪facebook.com/robots.txt 2020-05-31 20:12:47.576281+00:00
465 Disallow                                     / https://www.
↪facebook.com/robots.txt 2020-05-31 20:12:47.576281+00:00
[466 rows x 4 columns]
```

Now that we have downloaded the file, we can easily extract the list of user-agents that it contains.

```
>>> fb_useragents = (fb_robots
...                  [fb_robots['directive']=='User-agent']
...                  ['content'].drop_duplicates()
...                  .tolist())
>>> fb_useragents
```

(continues on next page)

(continued from previous page)

```
[ 'Applebot',
  'baiduspider',
  'Bingbot',
  'Discordbot',
  'facebookexternalhit',
  'Googlebot',
  'Googlebot-Image',
  'ia_archiver',
  'LinkedInBot',
  'msnbot',
  'Naverbot',
  'Pinterestbot',
  'seznambot',
  'Slurp',
  'teoma',
  'TelegramBot',
  'Twitterbot',
  'Yandex',
  'Yeti',
  '*']
```

Quite a long list!

As a small and quick test, I'm interested in checking the home page, a random profile page (/bbc), groups and hashtag pages.

```
>>> urls_to_test = ['/', '/bbc', '/groups', '/hashtag/']
>>> fb_test = robotstxt_test('https://www.facebook.com/robots.txt',
...                          fb_useragents, urls_to_test)
>>> fb_test
```

	robotstxt_url	user_agent	url_path	can_fetch
0	https://www.facebook.com/robots.txt	*	/bbc	False
1	https://www.facebook.com/robots.txt	*	/groups	False
2	https://www.facebook.com/robots.txt	*	/	False
3	https://www.facebook.com/robots.txt	*	/hashtag/	False
4	https://www.facebook.com/robots.txt	Applebot	/	True
..	...	...	...	...
75	https://www.facebook.com/robots.txt	seznambot	/groups	True
76	https://www.facebook.com/robots.txt	teoma	/	True
77	https://www.facebook.com/robots.txt	teoma	/hashtag/	False
78	https://www.facebook.com/robots.txt	teoma	/bbc	True
79	https://www.facebook.com/robots.txt	teoma	/groups	True

```
[80 rows x 4 columns]
```

For twenty user-agents and four URLs each, we received a total of eighty test results. You can immediately see that all user-agents not listed (denoted by \* are not allowed to fetch any of the provided URLs).

Let's see who is and who is not allowed to fetch the home page.

```
>>> fb_test.query('url_path== "/"')
robotstxt_url          user_agent  url_path  can_fetch
2  https://www.facebook.com/robots.txt          *        /        False
4  https://www.facebook.com/robots.txt      Applebot  /          True
9  https://www.facebook.com/robots.txt        Bingbot   /          True
14 https://www.facebook.com/robots.txt      Discordbot /         False
18 https://www.facebook.com/robots.txt        Googlebot /          True
21 https://www.facebook.com/robots.txt      Googlebot-Image /         True
```

(continues on next page)

(continued from previous page)

26	https://www.facebook.com/robots.txt	LinkedInBot	/	False
30	https://www.facebook.com/robots.txt	Naverbot	/	True
35	https://www.facebook.com/robots.txt	Pinterestbot	/	False
39	https://www.facebook.com/robots.txt	Slurp	/	True
43	https://www.facebook.com/robots.txt	TelegramBot	/	False
47	https://www.facebook.com/robots.txt	Twitterbot	/	True
48	https://www.facebook.com/robots.txt	Yandex	/	True
55	https://www.facebook.com/robots.txt	Yeti	/	True
57	https://www.facebook.com/robots.txt	baiduspider	/	True
60	https://www.facebook.com/robots.txt	facebookexternalhit	/	False
64	https://www.facebook.com/robots.txt	ia_archiver	/	False
68	https://www.facebook.com/robots.txt	msnbot	/	True
74	https://www.facebook.com/robots.txt	seznambot	/	True
76	https://www.facebook.com/robots.txt	teoma	/	True

I'll leave it to you to figure out why LinkedIn and Pinterest are not allowed to crawl the home page but Google and Apple are, because I have no clue!

**robotstxt\_test** (*robotstxt\_url*, *user\_agents*, *urls*)

Given a `robotstxt_url` check which of the `user_agents` is allowed to fetch which of the `urls`.

All the combinations of `user_agents` and `urls` will be checked and the results returned in one `DataFrame`.

```
>>> robotstxt_test('https://facebook.com/robots.txt',
...                 user_agents=['*', 'Googlebot', 'Applebot'],
...                 urls=['/', '/bbc', '/groups', '/hashtag/'])
  robotstxt_url user_agent url_path can_fetch
0  https://facebook.com/robots.txt * / False
1  https://facebook.com/robots.txt * /bbc False
2  https://facebook.com/robots.txt * /groups False
3  https://facebook.com/robots.txt * /hashtag/ False
4  https://facebook.com/robots.txt Applebot / True
5  https://facebook.com/robots.txt Applebot /bbc True
6  https://facebook.com/robots.txt Applebot /groups True
7  https://facebook.com/robots.txt Applebot /hashtag/ False
8  https://facebook.com/robots.txt Googlebot / True
9  https://facebook.com/robots.txt Googlebot /bbc True
10 https://facebook.com/robots.txt Googlebot /groups True
11 https://facebook.com/robots.txt Googlebot /hashtag/ False
```

### Parameters

- **robotstxt\_url** (*url*) – The URL of robotx.txt file
- **user\_agents** (*str*, *list*) – One or more user agents
- **urls** (*str*, *list*) – One or more paths (relative) or URLs (absolute) to check

### Return DataFrame `robotstxt_test_df`

**robotstxt\_to\_df** (*robotstxt\_url*)

Download the contents of `robotstxt_url` into a `DataFrame`

**Parameters** **robotstxt\_url** (*url*) – The URL of the robots.txt file

**Returns DataFrame** **robotstxt\_df** A `DataFrame` containing directives, their content, the URL and time of download

## 1.6 Download, Parse, and Analyze XML Sitemaps

One of the fastest and easiest ways to get insights on a website's content is to simply download its XML sitemap(s).

Sitemaps are also important SEO tools as they reveal a lot of information about the website, and help search engines in indexing those pages. You might want to run an SEO audit and check if the URLs in the sitemap properly correspond to the actual URLs of the site, so this would be an easy way to get them.

Sitemaps basically contain a log of publishing activity, and if they have rich URLs then you can do some good analysis on their content across time as well.

The `sitemap_to_df()` function is very simple to use, and only requires the URL of a sitemap, a sitemap index, or even a robots.txt file. It goes through the sitemap(s) and returns a DataFrame containing the tags and their information.

Let's go through a quick example of what can be done with sitemaps. We can start by getting one of the BBC's sitemaps.

### 1.6.1 Regular XML Sitemaps

```
>>> bbc_sitemap = sitemap_to_df('https://www.bbc.com/sitemaps/https-sitemap-com-
↳archive-1.xml')
>>> bbc_sitemap
```

	sitemap	loc	lastmod
0	↳https://www.bbc.com/arabic/middleeast/2009/06/...	2009-06-20	14:10:48+00:00
1	↳https://www.bbc.com/arabic/middleeast/2009/06/...	2009-06-20	21:07:43+00:00
2	↳https://www.bbc.com/arabic/business/2009/06/09...	2009-06-22	12:41:48+00:00
3	↳https://www.bbc.com/arabic/multimedia/2009/06/...	2009-06-24	15:27:24+00:00
4	↳https://www.bbc.com/arabic/business/2009/06/09...	2009-06-18	15:32:54+00:00
	...		...
49994	↳https://www.bbc.com/vietnamese/world/2009/09/0...	2009-09-02	11:46:23+00:00
49995	↳https://www.bbc.com/vietnamese/world/2009/09/0...	2009-09-04	11:20:42+00:00
49996	↳https://www.bbc.com/vietnamese/world/2009/09/0...	2009-09-02	02:40:41+00:00
49997	↳https://www.bbc.com/vietnamese/football/2009/0...	2009-09-02	03:09:06+00:00
49998	↳https://www.bbc.com/vietnamese/world/2009/09/0...	2009-09-05	04:38:11+00:00

```
[49999 rows x 3 columns]
```

```
>>> bbc_sitemap.dtypes
loc                object
lastmod            datetime64[ns, UTC]
sitemap            object
dtype: object
```

Since `lastmod` is a datetime object, we can easily use it for various time-related operations. Here we look at how many articles have been published (last modified) per year.

```
>>> bbc_sitemap.set_index('lastmod').resample('A')['loc'].count()
lastmod
2008-12-31 00:00:00+00:00    2261
2009-12-31 00:00:00+00:00   47223
2010-12-31 00:00:00+00:00     0
2011-12-31 00:00:00+00:00     0
2012-12-31 00:00:00+00:00     0
2013-12-31 00:00:00+00:00     0
2014-12-31 00:00:00+00:00     0
2015-12-31 00:00:00+00:00     0
2016-12-31 00:00:00+00:00     0
2017-12-31 00:00:00+00:00     0
2018-12-31 00:00:00+00:00     0
2019-12-31 00:00:00+00:00    483
2020-12-31 00:00:00+00:00     32
Freq: A-DEC, Name: loc, dtype: int64
```

As the majority are in 2009 with a few in other years, it seems these were later updated, but we would have to check to verify (in this special case BBC's URLs contain date information, which can be compared to `lastmod` to check if there is a difference between them).

We can take a look at a sample of the URLs to get the URL template that they use.

```
>>> bbc_sitemap['loc'].sample(10).tolist()
['https://www.bbc.com/russian/rolling_news/2009/06/090628_rn_pakistani_soldiries_
↳ambush',
 'https://www.bbc.com/urdu/pakistan/2009/04/090421_mqm_speaks_rza',
 'https://www.bbc.com/arabic/middleeast/2009/07/090723_ae_silwan_tc2',
 'https://www.bbc.com/portuguese/noticias/2009/07/090729_iraquerefenbritsfn',
 'https://www.bbc.com/portuguese/noticias/2009/06/090623_egitomilitaresfn',
 'https://www.bbc.com/portuguese/noticias/2009/03/090302_gazaconferenciaml',
 'https://www.bbc.com/portuguese/noticias/2009/07/090715_hillary_iran_cq',
 'https://www.bbc.com/vietnamese/culture/2009/04/090409_machienhuu_revisiting',
 'https://www.bbc.com/portuguese/noticias/2009/05/090524_paquistaoupdateg',
 'https://www.bbc.com/arabic/worldnews/2009/06/090629_om_pakistan_report_tc2']
```

It seems the pattern is

**`https://www.bbc.com/{language}/{topic}/{YYYY}/{MM}/{YYMMDD_article_title}`**

This is quite a rich structure, full of useful information. We can easily count how many articles they have by language, by splitting by “/” and getting the elements at index three, and counting them.

```
>>> bbc_sitemap['loc'].str.split('/').str[3].value_counts()
russian    14022
persian    10968
portuguese  5403
urdu       5068
mundo      5065
vietnamese 3561
arabic     2984
hindi      1677
turkce     706
ukchina    545
Name: loc, dtype: int64
```

We can also get a subset of articles written in a certain language, and see how many articles they publish per month, week, year, etc.

```
>>> (bbc_sitemap[bbc_sitemap['loc']
...   .str.contains('/russian/')
...   .set_index('lastmod')
...   .resample('M')['loc'].count())
lastmod
2009-04-30 00:00:00+00:00    1506
2009-05-31 00:00:00+00:00    2910
2009-06-30 00:00:00+00:00    3021
2009-07-31 00:00:00+00:00    3250
2009-08-31 00:00:00+00:00    2769
...
2019-09-30 00:00:00+00:00      8
2019-10-31 00:00:00+00:00     17
2019-11-30 00:00:00+00:00     11
2019-12-31 00:00:00+00:00     24
2020-01-31 00:00:00+00:00      6
Freq: M, Name: loc, Length: 130, dtype: int64
```

The fifth element after splitting URLs is the topic or category of the article. We can do the same and count the values.

```
>>> bbc_sitemap['loc'].str.split('/').str[4].value_counts()[:30]
rolling_news      9044
world              5050
noticias          4224
iran              3682
pakistan          2103
afghanistan       1959
multimedia        1657
internacional     1555
sport             1350
international     1293
india             1285
america_latina    1274
business          1204
cultura_sociedad  913
middleeast        874
worldnews         872
russia            841
radio             769
science           755
football          674
arts              664
ciencia_tecnologia 627
entertainment     621
simp              545
vietnam           539
economia          484
haberler          424
interactivity     411
help              354
ciencia           308
Name: loc, dtype: int64
```

Finally, we can take the last element after splitting, which contains the slugs of the articles, replace underscores with spaces, split, concatenate all, put in a `pd.Series` and count the values. This way we see how many times each word occurred in an article.

```

>>> (pd.Series(
...     bbc_sitemap['loc']
...     .str.split('/')
...     .str[-1]
...     .str.replace('_', ' ')
...     .str.cat(sep=' ')
...     .split()
...     )
...     .value_counts()[:15])
rn      8808
tc2     3153
iran    1534
video   973
obama   882
us      862
china   815
ir88    727
russia  683
si      640
np      638
afghan  632
ka      565
an      556
iraq    554
dtype: int64

```

This was a quick overview and data preparation for a sample sitemap. Once you are familiar with the sitemap's structure, you can more easily start analyzing the content.

## 1.6.2 News Sitemaps

```

>>> nyt_news = sitemap_to_df('https://www.nytimes.com/sitemaps/new/news.xml.gz')
>>> nyt_news

```

	publication_name	publication_language	news_publication_date	loc	lastmod	news_title	news_keywords	image_loc	sitemap	sitemap_downloaded
0	The New York Times	en-US	2020-05-19T15:49:28Z	https://www.nytimes.com/2020/05/19/sports/hors...	2020-05-19 15:49:28+00:00	Belmont Stakes to	Run June 20 as First Leg of ... Triple Crown (Horse Racing), Horse Racing, Bel...	https://static01.nyt.com/images/2020/05/19/spo...	https://www.nytimes.com/sitemaps/new/news.xml.gz	2020-05-19 15:49:44.459267+00:00
1	The New York Times	en-US	2020-05-19T09:21:33Z	https://www.nytimes.com/2020/05/19/us/coronavi...	2020-05-19 15:49:10+00:00	Coronavirus Live News and Updates	Coronavirus (2019-nCoV)	https://static01.nyt.com/images/2020/05/19/wor...	https://www.nytimes.com/sitemaps/new/news.xml.gz	2020-05-19 15:49:44.459267+00:00
2	The New York Times	en-US	2020-04-16T22:28:14Z	https://www.nytimes.com/interactive/2020/obitu...	2020-05-19 15:48:46+00:00	Those We've Lost	Deaths (Obituaries)	NaN	https://www.nytimes.com/sitemaps/new/news.xml.gz	2020-05-19 15:49:44.459267+00:00
3	The New York Times	en-US	2020-05-19T11:24:24Z	https://www.nytimes.com/2020/05/19/nyregion/co...	2020-05-19 15:48:06+00:00	Number of N.Y.C. Students Slated for Summer Sc...	Coronavirus (2019-nCoV), New York State, New Y...	https://static01.nyt.com/images/2020/05/19/nyr...	https://www.nytimes.com/sitemaps/new/news.xml.gz	2020-05-19 15:49:44.459267+00:00

continues on next page

(continued from previous page)

```

4 https://www.nytimes.com/2020/05/19/books/coron... 2020-05-19 15:46:10+00:00
↳The New York Times en-US 2020-05-19T15:46:10Z Coronavirus
↳Shutdowns Weigh on Book Sales Books and Literature, Book Trade and Publishin...
↳https://static01.nyt.com/images/2020/05/19/boo... https://www.nytimes.com/sitemaps/
↳new/news.xml.gz 2020-05-19 15:49:44.459267+00:00
.. ...
↳ ...
↳ ...
↳ ...
↳ ...
502 https://www.nytimes.com/2020/05/14/books/revie... 2020-05-17 16:17:52+00:00
↳The New York Times en 2020-05-14T09:00:03Z The Title of Emma
↳Straub's New Novel Is Mockin... Books and Literature, Straub, Emma, All Adults...
↳https://static01.nyt.com/images/2020/04/21/boo... https://www.nytimes.com/sitemaps/
↳new/news.xml.gz 2020-05-19 15:49:44.459267+00:00
503 https://www.nytimes.com/2020/05/17/opinion/nur... 2020-05-17 16:08:29+00:00
↳The New York Times en-US 2020-05-17T15:00:07Z Coronavirus Is
↳Hitting Nursing Homes Hard. How... Nursing Homes, Coronavirus (2019-nCoV), Elderl..
↳. https://static01.nyt.com/images/2020/05/17/opi... https://www.nytimes.com/
↳sitemaps/new/news.xml.gz 2020-05-19 15:49:44.459267+00:00
504 https://www.nytimes.com/2020/05/17/business/co... 2020-05-17 16:00:08+00:00
↳The New York Times en-US 2020-05-17T16:00:08Z Autoworkers Are
↳Returning as Carmakers Gradual... Automobiles, Shutdowns (Institutional), Labor ...
↳ https://static01.nyt.com/images/2020/05/18/bus... https://www.nytimes.com/
↳sitemaps/new/news.xml.gz 2020-05-19 15:49:44.459267+00:00
505 https://www.nytimes.com/2020/05/17/opinion/let... 2020-05-17 16:00:05+00:00
↳The New York Times en-US 2020-05-17T16:00:05Z Fathers,
↳Sons, Forgiveness and Regrets Children and Childhood, Parenting
↳https://static01.nyt.com/images/2020/05/10/opi... https://www.nytimes.com/sitemaps/
↳new/news.xml.gz 2020-05-19 15:49:44.459267+00:00
506 https://www.nytimes.com/2020/05/17/opinion/let... 2020-05-17 16:00:05+00:00
↳The New York Times en-US 2020-05-17T16:00:05Z
↳ To the Class of 2020 Coronavirus (2019-nCoV), Education (K-12)
↳https://static01.nyt.com/images/2020/04/16/wor... https://www.nytimes.com/sitemaps/
↳new/news.xml.gz 2020-05-19 15:49:44.459267+00:00
[507 rows x 13 columns]

```

### 1.6.3 Video Sitemaps

```

>>> wired_video = sitemap_to_df('https://www.wired.com/video/sitemap.xml')
>>> wired_video

           loc                                video_title
↳ video_thumbnail_loc                                video_title
↳                               video_description                                video_
↳content_loc video_duration video_publication_date video_expiration_date
↳                               sitemap                                sitemap_downloaded
0 https://www.wired.com/video/watch/autocomplete... http://dwgyu36up6iuz.
↳cloudfront.net/heru80fdn/... WIRED Autocomplete Interviews - Lele Pons Answ...
↳ Lele Pons takes the WIRED Autocomplete Interv... http://dp8hsntg6do36.
↳cloudfront.net/5db75425bc... 478 2019-10-29T16:00:00+00:00
↳ NaN https://www.wired.com/video/sitemap.xml 2020-05-19 16:18:17.
↳813461+00:00
1 https://www.wired.com/video/watch/professor-ex... http://dwgyu36up6iuz.
↳cloudfront.net/heru80fdn/... Laser Expert Explains One Concept in 5 Levels ...
↳ Donna Strickland, PhD, professor at the Univer... http://dp8hsntg6do36.
↳cloudfront.net/5da6107834... 1476 2019-10-28T17:18:00+00:00
↳ NaN https://www.wired.com/video/sitemap.xml 2020-05-19 16:18:17.
↳813461+00:00

```

(continues on next page)



(continued from previous page)

```

2      https://www.wired.com/video/watch/6-levels-of-... http://dwgyu36up6iuz.
↳cloudfront.net/heru80fdn/...      6 Levels of Knife Making: Easy to Complex
↳ Knife maker Chelsea Miller explains knife maki... http://dp8hsntg6do36.
↳cloudfront.net/5db32c4a34...      963      2019-10-25T19:00:00+00:00
↳      NaN      https://www.wired.com/video/sitemap.xml 2020-05-19 16:18:17.
↳813461+00:00
3      https://www.wired.com/video/watch/mycologist-e... http://dwgyu36up6iuz.
↳cloudfront.net/heru80fdn/...      Mycologist Explains How a Slime Mold Can Solve...
↳ Physarum polycephalum is a single-celled, brai... http://dp8hsntg6do36.
↳cloudfront.net/5db31cfabc...      606      2019-10-25T16:27:00+00:00
↳      NaN      https://www.wired.com/video/sitemap.xml 2020-05-19 16:18:17.
↳813461+00:00
4      https://www.wired.com/video/watch/almost-impos... http://dwgyu36up6iuz.
↳cloudfront.net/heru80fdn/...      Why It's Almost Impossible to Do a Quintuple C...
↳ Tricking is a sport with roots in martial arts... http://dp8hsntg6do36.
↳cloudfront.net/5db2005238...      644      2019-10-24T20:34:00+00:00
↳      NaN      https://www.wired.com/video/sitemap.xml 2020-05-19 16:18:17.
↳813461+00:00
...
↳      ...
↳      ...
↳      ...
↳      ...
↳      ...
2338 https://www.wired.com/video/watch/how-to-make-... http://dwgyu36up6iuz.
↳cloudfront.net/heru80fdn/...      How To Make Wired Origami
↳ Robert Lang explains how to fold the Wired iss... http://dp8hsntg6do36.
↳cloudfront.net/5171b3cbc2...      150      2008-09-23T00:00:00+00:00
↳      NaN      https://www.wired.com/video/sitemap.xml 2020-05-19 16:18:17.
↳813461+00:00
2339 https://www.wired.com/video/watch/clover-coffe... http://dwgyu36up6iuz.
↳cloudfront.net/heru80fdn/...      Clover Coffee Machine
↳ Wired.com takes a look at the 'Clover', an $11... http://dp8hsntg6do36.
↳cloudfront.net/5171b42ec2...      147      2008-09-23T00:00:00+00:00
↳      NaN      https://www.wired.com/video/sitemap.xml 2020-05-19 16:18:17.
↳813461+00:00
2340 https://www.wired.com/video/watch/original-war... http://dwgyu36up6iuz.
↳cloudfront.net/heru80fdn/...      Original WarGames Trailer
↳      Original WarGames Trailer      http://dp8hsntg6do36.
↳cloudfront.net/5171b427c2...      140      2008-07-21T04:00:00+00:00
↳      NaN      https://www.wired.com/video/sitemap.xml 2020-05-19 16:18:17.
↳813461+00:00
2341 https://www.wired.com/video/watch/rock-band-tr... http://dwgyu36up6iuz.
↳cloudfront.net/heru80fdn/...      Rock Band Trailer
↳      Rock Band Trailer      http://dp8hsntg6do36.
↳cloudfront.net/5171b431c2...      70      2007-09-14T04:00:00+00:00
↳      NaN      https://www.wired.com/video/sitemap.xml 2020-05-19 16:18:17.
↳813461+00:00
2342 https://www.wired.com/video/watch/arrival-full... http://dwgyu36up6iuz.
↳cloudfront.net/heru80fdn/...      'Arrival' -- Full Trailer
↳ Louise Banks (Amy Adams) must learn to communi... http://dp8hsntg6do36.
↳cloudfront.net/57b344f4fd...      145      2003-10-22T04:00:00+00:00
↳      NaN      https://www.wired.com/video/sitemap.xml 2020-05-19 16:18:17.
↳813461+00:00
[2343 rows x 11 columns]

```

**sitemap\_to\_df** (*sitemap\_url*, *max\_workers*=8, *recursive*=True)

Retrieve all URLs and other available tags of a sitemap(s) and put them in a DataFrame.

You can also pass the URL of a sitemap index, or a link to a robots.txt file.

### Parameters

- **sitemap\_url** (*url*) – The URL of a sitemap, either a regular sitemap, a sitemap index, or a link to a robots.txt file. In the case of a sitemap index or robots.txt, the function will go through all the sub sitemaps and retrieve all the included URLs in one DataFrame.
- **max\_workers** (*int*) – The maximum number of workers to use for threading. The higher the faster, but with high numbers you risk being blocked and/or missing some data as you might appear like an attacker.
- **recursive** (*bool*) – Whether or not to follow and import all sub-sitemaps (in case you have a sitemap index), or to only import the given sitemap. This might be useful in case you want to explore what sitemaps are available after which you can decide which ones you are interested in.

**Return sitemap\_df** A pandas DataFrame containing all URLs, as well as other tags if available (lastmod, changefreq, priority, or others found in news, video, or image sitemaps).

## 1.7 Python SEO Crawler / Spider

A customizable crawler to analyze SEO and content of pages and websites.

This is provided by the `crawl()` function which is customized for SEO and content analysis usage, and is highly configurable. The crawler uses `Scrapy` so you get all the power that it provides in terms of performance, speed, as well as flexibility and customization.

There are two main approaches to crawl:

1. **Discovery:** You know the website to crawl, so you provide a `url_list` (one or more URLs), and you want the crawler to go through the whole website(s) by following all available links.
2. **Pre-determined a.k.a “list mode”:** You have a known set of URLs that you want to crawl and analyze, without following links or discovering new URLs.

### 1.7.1 Discovery Crawling Approach

The simplest way to use the function is to provide a list of one or more URLs and the crawler will go through all of the reachable pages.

```
>>> crawl('https://example.com', 'my_output_file.json', follow_links=True)
```

That's it! To open the file:

```
>>> import pandas as pd
>>> pd.read_json('my_output_file.json', lines=True)
```

What this does:

- Check the site's robots.txt file and get the crawl rules, which means that your crawl will be affected by these rules and the user agent you are using. Check the details below on how to change settings and user agents to control this.
- Starting with the provided URL(s) go through all links and parse pages.
- For each URL extract the most important SEO elements.
- Save them to `my_output_file.json`.

- The column headers of the output file (once you import it as a DataFrame) would be the names of the elements (title, h1, h2, etc.).

Jsonlines is the supported output format because of its flexibility in allowing different values for different scraped pages, and appending independent items to the output files.

---

**Note:** When the crawler parses pages it saves the data to the specified file by appending, and not overwriting. Otherwise it would have to store all the data in memory, which might crash your computer. A good practice is to have a separate `output_file` for every crawl with a descriptive name `sitename_crawl_YYYY_MM_DD.jl` for example. If you use the same file you will probably get duplicate data in the same file.

---

## 1.7.2 Extracted On-Page SEO Elements

The names of these elements become the headers (column names) of the `output_file`.

Element	Remarks
url	The URL requested
url_redirected_to	The actual URL that was parsed, usually but not always the same as <i>url</i>
title	The <title> tag(s)
meta_desc	Meta description
canonical	The canonical tag if available
alt_href	The <i>href</i> attribute of rel=alternate tags
alt_hreflang	The language codes of the alternate links
og:*	Open Graph data
twitter:*	Twitter card data
jsonld_*	*JSON-LD data if available. In case multiple snippets occur, the respective column names will include a number to distinguish them, <i>jsonld_1_{item_a}</i> , <i>jsonld_1_{item_b}</i> , etc. Note that the first snippet will not contain a number, so the numbering starts with “1”, starting from the second snippet. The same applies to OG and Twitter cards.
h1	<h1> tag(s)
h2	<h2> tag(s)
h3	<h3> tag(s)
body_text	The text in the <p>, <span>, and <li> tags within <body>
size	The page size in bytes
resp_metadata	Several metadata for the response download_latency, timeout etc.
status	Response status (200, 301, 302, 404, etc.)
links_url	The URLs of the links on the page
links_text	The link text (anchor text)
links_fragment	The fragment part of the link (#fragment)
links_nofollow	Boolean, whether or not the link is a nofollow link. Note that this only tells if the link itself contains a rel=“nofollow” attribute. The page might indicate “nofollow” using meta robots or X-Robots-Tag, which you have to check separately.
img_src	The <i>src</i> attribute of images
img_alt	The <i>alt</i> attribute if available or an empty string
page_depth	The depth of the crawled page
ip_address	IP address
crawl_time	Date and time the page was crawled
resp_headers_*	Available response headers (last modified, server, etc.)
request_headers_*	All available request headers (user-agent, encoding, etc.)

**Note:** All elements that may appear multiple times on a page (like header tags, or images, for example), will be joined with two “@” signs @@. For example, “**first H2 tag@@second H2 tag@@third tag**” and so on. Once you open the file, you simply have to split by @@ to get the elements as a list.

Here is a sample file of a crawl of this site (output truncated for readability):

```
>>> import pandas as pd
>>> site_crawl = pd.read_json('path/to/file.json', lines=True)
>>> site_crawl.head()

  title          url          url_redirected_to          h1
  h2          h3
body_text size download_timeout download_slot download_latency
redirect_times redirect_ttl          redirect_urls redirect_reasons
depth status          links_href          links_text
24          img_src          img_alt          ip_address
crawl_time          resp_headers_date resp_headers_content-type          resp_
headers_last-modified resp_headers_vary          resp_headers_x-ms-request-id resp_
headers_x-ms-version resp_headers_x-ms-lease-status resp_headers_x-ms-blob-type
```

(continued from previous page)

```

0 https://advertools.readthedocs https://advertools.readthedocs
↪advertools -- Python Get productive as an online ma advertools@@Indices and
↪tables Online marketing productivity NaN Generate
↪keywords for SEM camp NaN NaN advertools.readthedocs.io
↪ NaN NaN NaN https://advertools.readthedocs
↪[302] NaN NaN #@@readme.html@@advertools.kw_
↪@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@ NaN
↪ NaN 104.17.32.82 2020-05-21 10:39:35 Thu, 21 May 2020 10:39:35 GMT
↪ text/html Wed, 20 May 2020 12:26:23 GMT Accept-Encoding 720a8581-501e-
↪0043-01a2-2e77d2 2009-09-19 unlocked
↪ BlockBlob * Nginx-Proxito-Sendfile
↪ web00007c advertools master /
↪proxito/media/html/advertools advertools.readthedocs.io
↪ path subdomain max-age=31536000; includeSubDo
↪ HIT NaN Thu, 21 May 2020 11:39:35 GMT
↪public, max-age=3600 max-age=604800, report-uri="ht cloudflare
↪596daca7dbaa7e9e-BUD 02d86a3cea00007e9edb0cf2000000 text/html,application/
↪xhtml+xml en Mozilla/5.0 (Windows NT 10.0;
↪ gzip, deflate __cfduid=d76b68d148ddec1efd004
1 https://advertools.readthedocs https://advertools.readthedocs
↪advertools -- Python NaN
↪advertools Change Log - advertools 0.9.1 (2020-05-19)@@0.9.0 (202
↪Ability to specify robots.txt NaN NaN advertools.readthedocs.io
↪ NaN NaN NaN NaN NaN
↪ NaN NaN NaN index.html@@readme.html@@adver
↪@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@ NaN
↪ NaN 104.17.32.82 2020-05-21 10:39:36 Thu, 21 May 2020 10:39:35 GMT
↪ text/html Wed, 20 May 2020 12:26:23 GMT Accept-Encoding 4f7bea3b-701e-
↪0039-3f44-2f1d9f 2009-09-19 unlocked
↪ BlockBlob * Nginx-Proxito-Sendfile
↪ web00007h advertools master /
↪proxito/media/html/advertools advertools.readthedocs.io
↪ path subdomain max-age=31536000; includeSubDo
↪ HIT NaN Thu, 21 May 2020 11:39:35 GMT
↪public, max-age=3600 max-age=604800, report-uri="ht cloudflare
↪596daca9bcab7e9e-BUD 02d86a3e0e00007e9edb0d72000000 text/html,application/
↪xhtml+xml en Mozilla/5.0 (Windows NT 10.0;
↪ gzip, deflate __cfduid=d76b68d148ddec1efd004
2 https://advertools.readthedocs https://advertools.readthedocs
↪advertools -- Python Get productive as an online ma advertools@@Indices and
↪tables Online marketing productivity NaN Generate
↪keywords for SEM camp NaN NaN advertools.readthedocs.io
↪ NaN NaN NaN NaN NaN
↪NaN NaN NaN #@@readme.html@@advertools.kw_ @@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@
↪ NaN NaN 104.17.32.82 2020-
↪05-21 10:39:36 Thu, 21 May 2020 10:39:35 GMT text/html Wed, 20
↪May 2020 12:26:36 GMT Accept-Encoding 98b729fa-e01e-00bf-24c3-2e494d
↪ 2009-09-19 unlocked BlockBlob
↪ * Nginx-Proxito-Sendfile web00007c
↪ advertools latest /proxito/media/html/advertools
↪advertools.readthedocs.io path
↪subdomain max-age=31536000; includeSubDo HIT
↪ NaN Thu, 21 May 2020 11:39:35 GMT public, max-age=3600 max-
↪age=604800, report-uri="ht cloudflare 596daca9bf26d423-BUD
↪02d86a3e150000d423322742000000 text/html,application/xhtml+xml
↪ en Mozilla/5.0 (Windows NT 10.0; gzip, deflate __
↪cfduid=d76b68d148ddec1efd004

```

(continues on next page)

(continued from previous page)

```

3 https://adverttools.readthedocs https://adverttools.readthedocs adverttools_
↳package -- Python NaN adverttools package _
↳ Submodules@@Module contents NaN Top-level package_
↳for advertoo NaN NaN adverttools.readthedocs.io NaN _
↳ NaN NaN NaN NaN NaN NaN_
↳ NaN index.html@@readme.html@@adver @@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@
↳ NaN NaN NaN 104.17.32.82 2020-05-21_
↳10:39:36 Thu, 21 May 2020 10:39:35 GMT text/html Wed, 20 May 2020_
↳12:26:25 GMT Accept-Encoding 7a28ef3b-801e-00c2-24c3-2ed585 2009-
↳09-19 unlocked BlockBlob _
↳ * Nginx-Proxito-Sendfile web000079 _
↳ adverttools master /proxito/media/html/adverttools adverttools.
↳readthedocs.io path subdomain _
↳ max-age=31536000; includeSubDo HIT _
↳NaN Thu, 21 May 2020 11:39:35 GMT public, max-age=3600 max-age=604800,_
↳report-uri="ht cloudflare 596daca9bddb7ec2-BUD _
↳02d86a3e1300007ec2a808a2000000 text/html,application/xhtml+xml _
↳ en Mozilla/5.0 (Windows NT 10.0; gzip, deflate __
↳cfduid=d76b68d148ddeclefd004

4 https://adverttools.readthedocs https://adverttools.readthedocs Python Module_
↳Index -- Python NaN Python Module Index _
↳ NaN NaN NaN NaN NaN © Copyright_
↳2020, Eli NaN NaN NaN adverttools.readthedocs.io NaN NaN _
↳ NaN NaN NaN NaN NaN NaN NaN NaN _
↳ NaN index.html@@readme.html@@adver @@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@
↳static/minus.png - 104.17.32.82 2020-05-21 10:39:36_
↳ Thu, 21 May 2020 10:39:35 GMT text/html Wed, 20 May 2020 12:26:23_
↳GMT Accept-Encoding 75911c9e-201e-00e6-34c3-2e4ccb 2009-09-19 _
↳ unlocked BlockBlob _
↳ * Nginx-Proxito-Sendfile web00007g _
↳adverttools master /proxito/media/html/adverttools adverttools.
↳readthedocs.io path subdomain _
↳ max-age=31536000; includeSubDo HIT _
↳NaN Thu, 21 May 2020 11:39:35 GMT public, max-age=3600 max-age=604800,_
↳report-uri="ht cloudflare 596daca9b91fd437-BUD _
↳02d86a3e140000d437b81532000000 text/html,application/xhtml+xml _
↳ en Mozilla/5.0 (Windows NT 10.0; gzip, deflate __
↳cfduid=d76b68d148ddeclefd004

66 https://adverttools.readthedocs https://adverttools.readthedocs adverttools.url_
↳builders -- Pyt NaN Source code for adverttools.url _
↳ NaN NaN NaN NaN NaN © Copyright_
↳2020, Eli NaN NaN NaN adverttools.readthedocs.io NaN NaN _
↳ NaN NaN NaN NaN NaN NaN NaN NaN _
↳ NaN ../..index.html@@../..readme @@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@
↳ NaN NaN NaN 104.17.32.82 2020-05-21 10:39:39_
↳ Thu, 21 May 2020 10:39:38 GMT text/html Wed, 20 May 2020 12:26:36_
↳GMT Accept-Encoding d99f2368-c01e-006f-18c3-2ef5ef 2009-09-19 _
↳ unlocked BlockBlob _
↳ * Nginx-Proxito-Sendfile web00007a _
↳adverttools latest /proxito/media/html/adverttools adverttools.
↳readthedocs.io path subdomain _
↳ max-age=31536000; includeSubDo HIT _
↳NaN Thu, 21 May 2020 11:39:38 GMT public, max-age=3600 max-age=604800,_
↳report-uri="ht cloudflare 596dacbbb8afd437-BUD _
↳02d86a494f0000d437b828b2000000 text/html,application/xhtml+xml _
↳ en Mozilla/5.0 (Windows NT 10.0; gzip, deflate __
↳cfduid=d76b68d148ddeclefd004

```

(continues on next page)

(continued from previous page)

```

67 https://advertools.readthedocs https://advertools.readthedocs advertools.kw_
↳generate -- Pyth                                NaN Source code for advertools.kw_
↳                                             NaN
↳Copyright 2020, Eli NaN                            NaN advertools.readthedocs.io
↳ NaN NaN NaN NaN NaN NaN NaN NaN NaN NaN
↳ NaN NaN NaN .././index.html@@.././readme @@@@
↳ NaN NaN 104.17.32.82 2020-05-
↳21 10:39:39 Thu, 21 May 2020 10:39:39 GMT text/html Wed, 20 May
↳2020 12:26:36 GMT Accept-Encoding 85855c48-c01e-00ce-13c3-2e3b74
↳2009-09-19 unlocked BlockBlob
↳ * Nginx-Proxito-Sendfile web00007g
↳ advertools latest /proxito/media/html/advertools
↳advertools.readthedocs.io path
↳subdomain max-age=31536000; includeSubDo HIT
↳ NaN Thu, 21 May 2020 11:39:39 GMT public, max-age=3600 max-
↳age=604800, report-uri="ht cloudflare 596dacbd980bd423-BUD
↳02d86a4a7f0000d423323b42000000 text/html,application/xhtml+xml
↳ en Mozilla/5.0 (Windows NT 10.0; gzip, deflate
↳cfduid=d76b68d148ddeclefd004
68 https://advertools.readthedocs https://advertools.readthedocs advertools.ad_
↳from_string -- P                                NaN Source code for advertools.ad_
↳                                             NaN
↳Copyright 2020, Eli NaN                            NaN advertools.readthedocs.io
↳ NaN NaN NaN NaN NaN NaN NaN NaN NaN NaN
↳ NaN NaN NaN .././index.html@@.././readme @@@@
↳ NaN NaN 104.17.32.82 2020-05-
↳21 10:39:39 Thu, 21 May 2020 10:39:39 GMT text/html Wed, 20 May
↳2020 12:26:36 GMT Accept-Encoding b0aef497-801e-004a-1647-2f6d5c
↳2009-09-19 unlocked BlockBlob
↳ * Nginx-Proxito-Sendfile web00007k
↳ advertools latest /proxito/media/html/advertools
↳advertools.readthedocs.io path
↳subdomain max-age=31536000; includeSubDo HIT
↳ NaN Thu, 21 May 2020 11:39:39 GMT public, max-age=3600 max-
↳age=604800, report-uri="ht cloudflare 596dacbd980cd423-BUD
↳02d86a4a7f0000d423209db2000000 text/html,application/xhtml+xml
↳ en Mozilla/5.0 (Windows NT 10.0; gzip, deflate
↳cfduid=d76b68d148ddeclefd004
69 https://advertools.readthedocs https://advertools.readthedocs advertools.ad_
↳create -- Python                                NaN Source code for advertools.ad_
↳                                             NaN
↳Copyright 2020, Eli NaN                            NaN advertools.readthedocs.io
↳ NaN NaN NaN NaN NaN NaN NaN NaN NaN NaN
↳ NaN NaN NaN .././index.html@@.././readme @@@@
↳ NaN NaN 104.17.32.82 2020-05-
↳21 10:39:39 Thu, 21 May 2020 10:39:39 GMT text/html Wed, 20 May
↳2020 12:26:36 GMT Accept-Encoding 9dfdd38a-101e-00a1-7ec3-2e93a0
↳2009-09-19 unlocked BlockBlob
↳ * Nginx-Proxito-Sendfile web00007c
↳ advertools latest /proxito/media/html/advertools
↳advertools.readthedocs.io path
↳subdomain max-age=31536000; includeSubDo HIT
↳ NaN Thu, 21 May 2020 11:39:39 GMT public, max-age=3600 max-
↳age=604800, report-uri="ht cloudflare 596dacbd99847ec2-BUD
↳02d86a4a7f00007ec2a811f2000000 text/html,application/xhtml+xml
↳ en Mozilla/5.0 (Windows NT 10.0; gzip, deflate
↳cfduid=d76b68d148ddeclefd004

```

(continues on next page)

(continued from previous page)

```

70 https://advertools.readthedocs https://advertools.readthedocs      advertools.
↪emoji -- Python                                          NaN Source code for advertools.emo   ↪
↪                                          NaN                               NaN   © Copyright ↪
↪2020, Eli      NaN                               NaN  advertools.readthedocs.io       NaN   ↪
↪      NaN      NaN                               NaN                               NaN   NaN   ↪
↪ NaN  ../../index.html@@../../readme  @@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@  ↪
↪      NaN                                          NaN  104.17.32.82  2020-05-21 10:39:40 ↪
↪ Thu, 21 May 2020 10:39:39 GMT                text/html  Wed, 20 May 2020 12:26:36 ↪
↪GMT Accept-Encoding  2ad504a1-101e-000b-03c3-2e454f              2009-09-19 ↪
↪      unlocked                                           BlockBlob                               ↪
↪      * Nginx-Proxito-Sendfile                          web000079                                ↪
↪advertools                    latest /proxito/media/html/advertools  advertools.
↪readthedocs.io                path                                subdomain  ↪
↪      max-age=31536000; includeSubDo                    HIT                               ↪
↪NaN Thu, 21 May 2020 11:39:39 GMT                public, max-age=3600  max-age=604800, ↪
↪report-uri="ht                cloudflare  596dacbd9fb97e9e-BUD ↪
↪02d86a4a7f00007e9edb13a2000000 text/html,application/xhtml+xml ↪
↪      en Mozilla/5.0 (Windows NT 10.0;                gzip, deflate  __  ↪
↪cfduid=d76b68d148ddec1efd004

```

### 1.7.3 Pre-Determined Crawling Approach (List Mode)

Sometimes you might have a fixed set of URLs for which you want to scrape and analyze SEO or content performance. Some ideas:

#### SERP Data

Let's say you just ran *serp\_goog* and got a bunch of top-ranking pages that you would like to analyze, and see how that relates to their SERP ranking.

You simply provide the *url\_list* parameter and again specify the *output\_file*. This will only crawl the specified URLs, and will not follow any links.

Now you have the SERP DataFrame, as well as the crawl output file. All you have to do is to merge them by the URL columns, and end up with a richer dataset

#### News Articles

You want to follow the latest news of a certain publication, and you extract their latest news URLs from their news sitemap using *sitemap\_to\_df*. You provide those URLs and crawl them only.

#### Google Analytics / Google Search Console

Since they provide reports for URLs, you can also combine them with the ones crawled and end up with a better perspective. You might be interested in knowing more about high bounce-rate pages, pages that convert well, pages that get less traffic than you think they should and so on. You can simply export those URLs and crawl them.

Any tool that has data about a set of URLs can be used.

Again running the function is as simple as providing a list of URLs, as well as a filepath where you want the result saved.



```
>>> crawl(url_list, 'output_file.jl', follow_links=False)
```

The difference between the two approaches, is the simple parameter `follow_links`. If you keep it as `False` (the default), the crawler will only go through the provided URLs. Otherwise, it will discover pages by following links on pages that it crawls. So how do you make sure that the crawler doesn't try to crawl the whole web when `follow_links` is `True`? The `allowed_domains` parameter gives you the ability to control this, although it is an optional parameter. If you don't specify it, then it will default to only the domains in the `url_list`. It's important to note that you have to set this parameter if you have certain sub-domains that you want to crawl.

## 1.7.4 CSS and XPath Selectors

The above approaches are generic, and are useful for an exploratory SEO audit and the output is helpful for most cases.

But what if you want to extract special elements that are not included in the default output? This is extremely important, as there are key elements on pages that you need to additionally extract and analyze. Some examples might be tags, prices, social media shares, product price or availability, comments, and pretty much any element on a page that might be of interest to you.

For this you can use two special arguments for CSS and/or XPath selectors. You simply provide a dictionary `{'name_1': 'selector_1', 'name_2': 'selector_2'}` where the keys become the column names, and the values (selectors) will be used to extract the required elements.

I mostly rely on [SelectorGadget](#) which is a really great tool for getting the CSS/XPath selectors of required elements. In some pages it can get really tricky to figure that out. Other resources for learning more about selectors:

- [Scrapy's documentaion for selectors](#)
- [CSS Selector Reference on W3C](#)
- [XPath tutorial on W3C](#)

Once you have determined the elements that you want to extract and figured out what their names are going to be, you simply pass them as arguments to `css_selectors` and/or `xpath_selectors` as dictionaries, as described above.

Let's say you want to extract the links in the sidebar of this page. By default you would get all the links from the page, but you want to put those in the sidebar in a separate column. It seems that the CSS selector for them is `.toctree-l1 .internal`, and the XPath equivalent is `//*[contains(concat( " ", @class, " " ), concat( " ", "toctree-l1", " " ))]//*[contains(concat( " ", @class, " " ), concat( " ", "internal", " " ))]`. Note that this selects the *element* (the whole link object), which is not typically what you might be interested in.

So with CSS you need to append `::text` or `::attr(href)` if you want the text of the links or the `href` attribute respectively. Similarly with XPath, you will need to append `/text()` or `/@href` to the selector to get the same.

```
>>> crawl('https://advertools.readthedocs.io/en/master/advertools.spider.html',
...       'output_file.jl',
...       css_selectors={'sidebar_links': '.toctree-l1 .internal::text',
...                       'sidebar_links_url': '.toctree-l1 .internal::attr(href)'})
```

Or, instead of `css_selectors` you can add a similar dictionary for the `xpath_selectors` argument:

```
>>> crawl('https://advertools.readthedocs.io/en/master/advertools.spider.html',
...       'output_file.jl',
...       xpath_selectors={'sidebar_links': '//*[contains(concat( " ", @class, " " ),
↳ concat( " ", "toctree-l1", " " ))]//*[contains(concat( " ", @class, " " ), concat(
↳ " ", "internal", " " ))]/text()',
...                       'sidebar_links_url': '//*[contains(concat( " ", @class, "
↳ " ), concat( " ", "toctree-l1", " " ))]//*[contains(concat( " ", @class, " " ),
↳ concat( " ", "internal", " " ))]/@href'})
```

(continues on next page)

## 1.7.5 Spider Custom Settings and Additional Functionality

In addition to what you can control regarding the items you can extract, you can also customize the behaviour of the spider and set rules for crawling so you can control it even further.

This is provided by the `custom_settings` parameter. It is optional, and takes a dictionary of settings and their values. Scrapy provides a very large number of settings, and they are all available through this parameter (assuming some conditions for some of the settings).

Here are some examples that you might find interesting:

- `CONCURRENT_REQUESTS_PER_DOMAIN` Defaults to 8, and controls the number of simultaneous requests to be performed for each domain. You might want to lower this if you don't want to put too much pressure on the website's server, and you probably don't want to get blocked!
- `DEFAULT_REQUEST_HEADERS` You can change this if you need to.
- `DEPTH_LIMIT` How deep your crawl will be allowed. The default has no limit.
- `DOWNLOAD_DELAY` Similar to the first option. Controls the amount of time in seconds for the crawler to wait between consecutive pages of the same website. It can also take fractions of a second (0.4, 0.75, etc.)
- `LOG_FILE` If you want to save your crawl logs to a file, you can provide a path to it here.
- `USER_AGENT` If you want to identify yourself differently while crawling. This is affected by the robots.txt rules, so you would be potentially allowed/disallowed from certain pages based on your user-agent.
- `CLOSESPIDER_ERRORCOUNT`, `CLOSESPIDER_ITEMCOUNT`, `CLOSESPIDER_PAGECOUNT`, `CLOSESPIDER_TIMEOUT` Stop crawling after that many errors, items, pages, or seconds. These can be very useful to limit your crawling in certain cases. I particularly like to use `CLOSESPIDER_PAGECOUNT` when exploring a new website, and also to make sure that my selectors are working as expected. So for your first few crawls you might set this to one hundred for example and explore the crawled pages. Then when you are confident things are working fine, you can remove this restriction. `CLOSESPIDER_ERRORCOUNT` can also be very useful while exploring, just in case you get unexpected errors.

The next page contains a number of *strategies and recipes for crawling* with code examples and explanations.

### Usage

A very simple dictionary to be added to your function call:

```
>>> crawl('http://example.com', 'output_file.jl',
...       custom_settings={'CLOSESPIDER_PAGECOUNT': 100,
...                       'CONCURRENT_REQUESTS_PER_DOMAIN': 1,
...                       'USER_AGENT': 'custom-user-agent'})
```

Please refer to the [spider settings documentation](#) for the full details.

**crawl** (*url\_list*, *output\_file*, *follow\_links=False*, *css\_selectors=None*, *xpath\_selectors=None*, *custom\_settings=None*, *allowed\_domains=None*)  
Crawl a website's URLs based on the given *url\_list*

#### Parameters

- **url\_list** (*url*, *list*) – One or more URLs to crawl. If *follow\_links* is True, the crawler will start with these URLs and follow all links on pages recursively.
- **output\_file** (*str*) – The path to the output of the crawl. Jsonlines only is supported to allow for dynamic values. Make sure your file ends with “.jl”, e.g. *output\_file.jl*.

- **follow\_links** (*bool*) – Defaults to False. Whether or not to follow links on crawled pages.
- **css\_selectors** (*dict*) – A dictionary mapping names to CSS selectors. The names will become column headers, and the selectors will be used to extract the required data/content.
- **xpath\_selectors** (*dict*) – A dictionary mapping names to XPath selectors. The names will become column headers, and the selectors will be used to extract the required data/content.
- **custom\_settings** (*dict*) – A dictionary of optional custom settings that you might want to add to the spider’s functionality. There are over 170 settings for all kinds of options. For details please refer to the [spider settings](#) documentation.
- **allowed\_domains** (*list*) – (optional) A list of the allowed domains to crawl. This ensures that the crawler does not attempt to crawl the whole web. If not specified, it defaults to the domains of the URLs provided in `url_list`. You can use it for sub-domains if you want them to be crawled as they will not be crawled if not specified.

### Examples

Crawl a website and let the crawler discover as many pages as available

```
>>> crawl('http://example.com', 'output_file.json', follow_links=True)
>>> import pandas as pd
>>> crawl_df = pd.read_json('output_file.json', lines=True)
```

Crawl a known set of pages (on a single or multiple sites) without following links (just crawl the specified pages) or “list mode”:

```
>>> crawl(['http://example.com/product', 'http://example.com/product2',
...       'https://anotherexample.com', 'https://anotherexample.com/hello'],
...       'output_file.json', follow_links=False)
```

Crawl a website, and in addition to standard SEO elements, also get the required CSS selectors. Here we will get three additional columns *price*, *author*, and *author\_url*. Note that you need to specify if you want the text attribute or the *href* attribute if you are working with links (and all other selectors).

```
>>> crawl('http://example.com', 'output_file.json',
...       css_selectors={'price': '.a-color-price::text',
...                      'author': '.contributorNameID::text',
...                      'author_url': '.contributorNameID::attr(href)'})
```

## 1.8 SEO Crawling & Scraping: Strategies & Recipes

Once you have mastered the basics of using the `crawl` function, you probably want to achieve more with better customization and control.

These are some code strategies that might be useful to customize how you run your crawls.

Most of these options can be set using the `custom_settings` parameter that the function takes. This can be set by using a dictionary, where the keys indicate the option you want to set, and the values specify how you want to set them.

### 1.8.1 How to crawl a list of pages, and those pages only (list mode)?

Simply provide that list as the first argument, for the `url_list` parameter, and make sure that `follow_links=False`, which is the default. This simply crawls the given pages, and stops when done.

```
>>> import advertools as adv
>>> url_list = ['https://example.com/page_1',
...           'https://example.com/page_2',
...           'https://example.com/page_3',
...           'https://example.com/page_4']

>>> adv.crawl(url_list,
...           output_file='example_crawl_1.jl',
...           follow_links=False)
```

### 1.8.2 How can I crawl a website including its sub-domains?

The `crawl` function takes an optional `allowed_domains` parameter. If not provided, it defaults to the domains of the URLs in `url_list`. When the crawler goes through the pages of `example.com`, it follows links to discover pages. If it finds pages on `help.example.com` it won't crawl them (it's a different domain). The solution, therefore, is to provide a list of domains to the `allowed_domains` parameter. Make sure you also include the original domain, in this case `example.com`.

```
>>> adv.crawl('https://example.com',
...           'example_crawl_1.jl',
...           follow_links=True
...           allowed_domains=['help.example.com', 'example.com', 'community.example.
↳com'])
```

### 1.8.3 How can I save a copy of the logs of my crawl for auditing them later?

It's usually good to keep a copy of the logs of all your crawls to check for errors, exceptions, stats, etc. Pass a path of the file where you want the logs to be saved, in a dictionary to the `custom_settings` parameter. A good practice for consistency is to give the same name to the `output_file` and log file (with a different extension) for easier retrieval. For example:

```
output_file: 'website_name_crawl_1.jl'
LOG_FILE: 'website_name_crawl_1.log' (txt can also work)
output_file: 'website_name_crawl_2.jl'
LOG_FILE: 'website_name_crawl_2.log'
```

```
>>> adv.crawl('https://example.com', 'example_crawl_1.jl',
...           custom_settings={'LOG_FILE': 'example_crawl_1.log'})
```

## 1.8.4 How can I automatically stop my crawl based on a certain condition?

There are a few conditions that you can use to trigger the crawl to stop, and they mostly have descriptive names:

- `CLOSESPIDER_ERRORCOUNT`: You don't want to wait three hours for a crawl to finish, only to discover that you had errors all over the place. Set a certain number of errors to trigger the crawler to stop, so you can investigate the issue.
- `CLOSESPIDER_ITEMCOUNT`: Anything scraped from a page is an "item", h1, title, meta\_desc, etc. Set the crawler to stop after getting a certain number of items if you want that.
- `CLOSESPIDER_PAGECOUNT`: Stop the crawler after a certain number of pages have been crawled. This is useful as an exploratory technique, especially with very large websites. It might be good to crawl a few thousand pages, get an idea on its structure, and then run a full crawl with those insights in mind.
- `CLOSESPIDER_TIMEOUT`: Stop the crawler after a certain number of seconds.

```
>>> adv.crawl('https://example.com', 'example_crawl_1.jl',
...          custom_settings={'CLOSESPIDER_PAGECOUNT': 500})
```

## 1.8.5 How can I (dis)obey robots.txt rules?

The crawler obeys robots.txt rules by default. Sometimes you might want to check the results of crawls without doing that. You can set the `ROBOTSTXT_OBEY` setting under `custom_settings`:

```
>>> adv.crawl('https://example.com',
...          'example_crawl_1.jl',
...          custom_settings={'ROBOTSTXT_OBEY': False})
```

## 1.8.6 How do I set my User-agent while crawling?

Set this parameter under `custom_settings` dictionary under the key `USER_AGENT`. The default User-agent can be found by running `adv.spider.user_agent`

```
>>> adv.spider.user_agent # to get the current User-agent
>>> adv.crawl('http://example.com',
...          'example_crawl_1.jl',
...          custom_settings={'USER_AGENT': 'YOUR_USER_AGENT'})
```

## 1.8.7 How can I control the number of concurrent requests while crawling?

Some servers are set for high sensitivity to automated and/or concurrent requests, that you can quickly be blocked/banned. You also want to be polite and not kill those servers, don't you?

There are several ways to set that under the `custom_settings` parameter. The available keys are the following:

```
CONCURRENT_ITEMS: default 100
CONCURRENT_REQUESTS : default 16
CONCURRENT_REQUESTS_PER_DOMAIN: default 8
CONCURRENT_REQUESTS_PER_IP: default 0
```

```
>>> adv.crawl('https://example.com',
...           'example_crawl_1.jl',
...           custom_settings={'CONCURRENT_REQUESTS_PER_DOMAIN': 1})
```

### 1.8.8 How can I slow down the crawling so I don't hit the websites' servers too hard?

Use the `DOWNLOAD_DELAY` setting and set the interval to be waited before downloading consecutive pages from the same website (in seconds).

```
>>> adv.crawl('https://example.com', 'example_crawl_1.jl',
...           custom_settings={'DOWNLOAD_DELAY': 3}) # wait 3 seconds between pages
```

### 1.8.9 How can I set multiple settings to the same crawl job?

Simply add multiple settings to the `custom_settings` parameter.

```
>>> adv.crawl('http://example.com',
...           'example_crawl_1.jl',
...           custom_settings={'CLOSESPIDER_PAGECOUNT': 400,
...                             'CONCURRENT_ITEMS': 75,
...                             'LOG_FILE': 'output_file.log'})
```

### 1.8.10 I want to crawl a list of pages, follow links from those pages, but only to a certain specified depth

Set the `DEPTH_LIMIT` setting in the `custom_settings` parameter. A setting of 1 would follow links one level after the provided URLs in `url_list`

```
>>> adv.crawl('http://example.com',
...           'example_crawl_1.jl',
...           custom_settings={'DEPTH_LIMIT': 2}) # follow links two levels from the_
↪initial URLs, then stop
```

### 1.8.11 How do I pause/resume crawling, while making sure I don't crawl the same page twice?

There are several reasons why you might want to do this:

- You want to mainly crawl the updates to the site (you already crawled the site).
- The site is very big, and can't be crawled quickly.
- You are not in a hurry, and you also don't want to hit the servers hard, so you run your crawl across days for example.
- As an emergency measure (connection lost, battery died, etc.) you can start where you left off

Handling this is extremely simple, and all you have to do is simply provide a path to a new folder. Make sure it is new and empty, and make sure to only use it for the same crawl job reruns. That's all you have to worry about. The `JOBDIR` setting handles this.

```
>>> adv.crawl('http://example.com',
...           'example_crawl_1.jl',
...           custom_settings={'JOBDIR': '/Path/to/en/empty/folder'})
```

The first time you run the above code and then stop it. Stopping can happen by accident (lost connection, closed computer, etc.), manually (you hit ctrl+C) or you used a custom setting option to stop the crawl after a certain number of pages, seconds, etc.

The second time you want to run this, you simply run the exact same command again. If you check the folder that was created you can see a few files that manage the process. You don't need to worry about any of it. But make sure that folder doesn't get changed manually, rerun the same command as many times as you need, and the crawler should handle de-duplication for you.

## 1.9 Import Search Engine Results Pages (SERPs) for Google and YouTube

Before being able to run queries using `serp_goog()`, you will need to set up some credentials as follows (you don't need a custom search engine for `serp_youtube()`):

- **Create a custom search engine:** At first, you might be asked to enter a site to search. Enter any domain, then go to the control panel and remove it. Make sure you enable "Search the entire web" and image search. You will also need to get your search engine ID, which you can find on the control panel page.
- **Enable the custom search API:** The service will allow you to retrieve and display search results from your custom search engine programmatically. You will need to create a project for this first.
- **Create credentials for this project:** so you can get your key.
- **Enable billing for your project** if you want to run more than 100 queries per day. The first 100 queries are free; then for each additional 1,000 queries, you pay \$5.

`serp_goog(q, cx, key, c2coff=None, cr=None, dateRestrict=None, exactTerms=None, excludeTerms=None, fileType=None, filter=None, gl=None, highRange=None, hl=None, hq=None, imgColorType=None, imgDominantColor=None, imgSize=None, imgType=None, linkSite=None, lowRange=None, lr=None, num=None, orTerms=None, relatedSite=None, rights=None, safe=None, searchType=None, siteSearch=None, siteSearchFilter=None, sort=None, start=None)`

Query Google and get search results in a DataFrame.

**For each parameter, you can supply single or multiple values / arguments.** If you pass multiple arguments, all the possible combinations of arguments (the product) will be requested, and you will get one DataFrame combining all queries. See examples below.

### Parameters

- **q** – The search expression.
- **cx** – The custom search engine ID to use for this request.
- **key** – The API key of your custom search engine.
- **c2coff** – Enables or disables Simplified and Traditional Chinese Search. The default value for this parameter is 0 (zero), meaning that the feature is enabled. Supported values are: 1: Disabled 0: Enabled (default)
- **cr** – Restricts search results to documents originating in a particular country. You may use Boolean operators in the cr parameter's value. Google Search determines the country of a document by analyzing: the top-level domain (TLD) of the document's URL the geographic

location of the Web server's IP address. See the Country Parameter Values page for a list of valid values for this parameter.

- **dateRestrict** – Restricts results to URLs based on date. Supported values include: `d[number]`: requests results from the specified number of past days. `- d[number]`: requests results from the specified number of past days. `- w[number]`: requests results from the specified number of past weeks. `- m[number]`: requests results from the specified number of past months. `- y[number]`: requests results from the specified number of past years.
- **exactTerms** – Identifies a phrase that all documents in the search results must contain.
- **excludeTerms** – Identifies a word or phrase that should not appear in any documents in the search results.
- **fileType** – Restricts results to files of a specified extension. A list of file types indexable by Google can be found in Search Console Help Center.
- **filter** – Controls turning on or off the duplicate content filter. See Automatic Filtering for more information about Google's search results filters. Note that host crowding filtering applies only to multi-site searches. By default, Google applies filtering to all search results to improve the quality of those results. Acceptable values are: `"0"`: Turns off duplicate content filter. `"1"`: Turns on duplicate content filter.
- **gl** – Geolocation of end user. The `gl` parameter value is a two-letter country code. The `gl` parameter boosts search results whose country of origin matches the parameter value. See the Country Codes page for a list of valid values. Specifying a `gl` parameter value should lead to more relevant results. This is particularly true for international customers and, even more specifically, for customers in English-speaking countries other than the United States.
- **highRange** – Specifies the ending value for a search range. Use `lowRange` and `highRange` to append an inclusive search range of `lowRange...highRange` to the query.
- **hl** – Sets the user interface language. Explicitly setting this parameter improves the performance and the quality of your search results. See the Interface Languages section of Internationalizing Queries and Results Presentation for more information, and Supported Interface Languages for a list of supported languages.
- **hq** – Appends the specified query terms to the query, as if they were combined with a logical AND operator.
- **imgColorType** – Returns black and white, grayscale, or color images: `mono`, `gray`, and `color`. Acceptable values are: `"color"`: color `"gray"`: gray `"mono"`: mono
- **imgDominantColor** – Returns images of a specific dominant color. Acceptable values are: `"black"`: black `"blue"`: blue `"brown"`: brown `"gray"`: gray `"green"`: green `"orange"`: orange `"pink"`: pink `"purple"`: purple `"red"`: red `"teal"`: teal `"white"`: white `"yellow"`: yellow
- **imgSize** – Returns images of a specified size. Acceptable values are: `"huge"`: huge `"icon"`: icon `"large"`: large `"medium"`: medium `"small"`: small `"xlarge"`: xlarge `"xxlarge"`: xxlarge
- **imgType** – Returns images of a type. Acceptable values are: `"clipart"`: clipart `"face"`: face `"lineart"`: lineart `"news"`: news `"photo"`: photo
- **linkSite** – Specifies that all search results should contain a link to a particular URL
- **lowRange** – Specifies the starting value for a search range. Use `lowRange` and `highRange` to append an inclusive search range of `lowRange...highRange` to the query.
- **lr** – Restricts the search to documents written in a particular language (e.g., `lr=lang_ja`). Acceptable values are: `"lang_ar"`: Arabic `"lang_bg"`: Bulgarian `"lang_ca"`: Catalan `"lang_cs"`: Czech `"lang_da"`: Danish `"lang_de"`: German `"lang_el"`: Greek `"lang_en"`: English `"lang_es"`: Spanish `"lang_et"`: Estonian `"lang_fi"`: Finnish `"lang_fr"`: French



“lang\_hr”: Croatian “lang\_hu”: Hungarian “lang\_id”: Indonesian “lang\_is”: Icelandic “lang\_it”: Italian “lang\_iw”: Hebrew “lang\_ja”: Japanese “lang\_ko”: Korean “lang\_lt”: Lithuanian “lang\_lv”: Latvian “lang\_nl”: Dutch “lang\_no”: Norwegian “lang\_pl”: Polish “lang\_pt”: Portuguese “lang\_ro”: Romanian “lang\_ru”: Russian “lang\_sk”: Slovak “lang\_sl”: Slovenian “lang\_sr”: Serbian “lang\_sv”: Swedish “lang\_tr”: Turkish “lang\_zh-CN”: Chinese (Simplified) “lang\_zh-TW”: Chinese (Traditional)

- **num** – Number of search results to return. Valid values are integers between 1 and 10, inclusive.
- **orTerms** – Provides additional search terms to check for in a document, where each document in the search results must contain at least one of the additional search terms.
- **relatedSite** – Specifies that all search results should be pages that are related to the specified URL.
- **rights** – Filters based on licensing. Supported values include: `cc_publicdomain`, `cc_attribute`, `cc_sharealike`, `cc_noncommercial`, `cc_nonderived`, and combinations of these.
- **safe** – Search safety level. Acceptable values are: “active”: Enables SafeSearch filtering. “off”: Disables SafeSearch filtering. (default)
- **searchType** – Specifies the search type: `image`. If unspecified, results are limited to webpages. Acceptable values are: “image”: custom image search.
- **siteSearch** – Specifies all search results should be pages from a given site.
- **siteSearchFilter** – Controls whether to include or exclude results from the site named in the `siteSearch` parameter. Acceptable values are: “e”: exclude “i”: include
- **sort** – The sort expression to apply to the results.
- **start** – The index of the first result to return. Valid value are integers starting 1 (default) and the second result is 2 and so forth. For example `&start=11` gives the second page of results with the default “num” value of 10 results per page. Note: No more than 100 results will ever be returned for any query with JSON API, even if more than 100 documents match the query, so setting (start + num) to more than 100 will produce an error. Note that the maximum value for num is 10.

The following function call will produce two queries: “hotel” in the USA, and “hotel” in France

```
>>> serp_goog(q='hotel', gl=['us', 'fr'], cx='YOUR_CX', key='YOUR_KEY')
```

The below function call will prouce four queries and make four requests:

- “flights” in UK
- “flights” in Australia
- “tickets” in UK
- “tickets” in Australia

‘cr’ here refers to ‘country restrict’, which focuses on content originating from the specified country.

```
>>> serp_goog(q=['flights', 'tickets'], cr=['countryUK', 'countryAU'],
             cx='YOUR_CX', key='YOUR_KEY')
```

**serp\_youtube** (*key*, *q=None*, *channelId=None*, *channelType=None*, *eventType=None*, *forContentOwner=None*, *forDeveloper=None*, *forMine=None*, *location=None*, *locationRadius=None*, *maxResults=None*, *onBehalfOfContentOwner=None*, *order=None*, *pageToken=None*, *publishedAfter=None*, *publishedBefore=None*, *regionCode=None*, *relatedToVideoId=None*, *relevanceLanguage=None*, *safeSearch=None*, *topicId=None*, *type=None*, *videoCaption=None*, *videoCategoryId=None*, *videoDefinition=None*, *videoDimension=None*, *videoDuration=None*, *videoEmbeddable=None*, *videoLicense=None*, *videoSyndicated=None*, *videoType=None*)

Query the YouTube API and get search results in a DataFrame. For each parameter you can supply a single or multiple value(s). Looping and merging results is handled automatically in case of multiple values.

### Parameters

- **q** – (string) The `q` parameter specifies the query term to search for. Your request can also use the Boolean NOT (-) and OR (|) operators to exclude videos or to find videos that are associated with one of several search terms. For example, to search for videos matching either “boating” or “sailing”, set the `q` parameter value to `boating|sailing`. Similarly, to search for videos matching either “boating” or “sailing” but not “fishing”, set the `q` parameter value to `boating|sailing-fishing`. Note that the pipe character must be URL-escaped when it is sent in your API request. The URL-escaped value for the pipe character is `%7C`.
- **channelId** – (string) The `channelId` parameter indicates that the API response should only contain resources created by the channel. Note: Search results are constrained to a maximum of 500 videos if your request specifies a value for the `channelId` parameter and sets the `type` parameter value to `video`, but it does not also set one of the `forContentOwner`, `forDeveloper`, or `forMine` filters.
- **channelType** – (string) The `channelType` parameter lets you restrict a search to a particular type of channel. Acceptable values are:
  - `any` – Return all channels.
  - `show` – Only retrieve shows.
- **eventType** – (string) The `eventType` parameter restricts a search to broadcast events. If you specify a value for this parameter, you must also set the `type` parameter’s value to `video`. Acceptable values are:
  - `completed` – Only include completed broadcasts.
  - `live` – Only include active broadcasts.
  - `upcoming` – Only include upcoming broadcasts.
- **forContentOwner** – (boolean) This parameter can only be used in a properly authorized request, and it is intended exclusively for YouTube content partners. The `forContentOwner` parameter restricts the search to only retrieve videos owned by the content owner identified by the `onBehalfOfContentOwner` parameter. If `forContentOwner` is set to `true`, the request must also meet these requirements: The `onBehalfOfContentOwner` parameter is required. The user authorizing the request must be using an account linked to the specified content owner. The `type` parameter value must be set to `video`. None of the following other parameters can be set: `videoDefinition`, `videoDimension`, `videoDuration`, `videoLicense`, `videoEmbeddable`, `videoSyndicated`, `videoType`.
- **forDeveloper** – (boolean) This parameter can only be used in a properly authorized request. The `forDeveloper` parameter restricts the search to only retrieve videos uploaded via the developer’s application or website. The API server uses the request’s authorization credentials to identify the developer. The `forDeveloper` parameter can be used in conjunction with optional search parameters like the `q` parameter. For this feature, each

uploaded video is automatically tagged with the project number that is associated with the developer's application in the Google Developers Console. When a search request subsequently sets the `forDeveloper` parameter to `true` the API server uses the request's authorization credentials to identify the developer. Therefore, a developer can restrict results to videos uploaded through the developer's own app or website but not to videos uploaded through other apps or sites.

- **forMine** – (boolean) This parameter can only be used in a properly authorized request. The `forMine` parameter restricts the search to only retrieve videos owned by the authenticated user. If you set this parameter to `true`, then the `type` parameter's value must also be set to `video`. In addition, none of the following other parameters can be set in the same request: `videoDefinition`, `videoDimension`, `videoDuration`, `videoLicense`, `videoEmbeddable`, `videoSyndicated`, `videoType`.
- **relatedToVideoId** – (string) The `relatedToVideoId` parameter retrieves a list of videos that are related to the video that the parameter value identifies. The parameter value must be set to a YouTube video ID and, if you are using this parameter, the `type` parameter must be set to `video`. Note that if the `relatedToVideoId` parameter is set, the only other supported parameters are `part`, `maxResults`, `pageToken`, `regionCode`, `relevanceLanguage`, `safeSearch`, `type` (which must be set to `video`), and `fields`.
- **location** – (string) The `location` parameter, in conjunction with the `locationRadius` parameter, defines a circular geographic area and also restricts a search to videos that specify, in their metadata, a geographic location that falls within that area. The parameter value is a string that specifies latitude/longitude coordinates e.g. (37.42307,-122.08427). The `location` parameter value identifies the point at the center of the area. The `locationRadius` parameter specifies the maximum distance that the location associated with a video can be from that point for the video to still be included in the search results. The API returns an error if your request specifies a value for the `location` parameter but does not also specify a value for the `locationRadius` parameter.
- **locationRadius** – (string) The `locationRadius` parameter, in conjunction with the `location` parameter, defines a circular geographic area. The parameter value must be a floating point number followed by a measurement unit. Valid measurement units are m, km, ft, and mi. For example, valid parameter values include 1500m, 5km, 10000ft, and 0.75mi. The API does not support `locationRadius` parameter values larger than 1000 kilometers. Note: See the definition of the `location` parameter for more information.
- **maxResults** – (unsigned integer) The `maxResults` parameter specifies the maximum number of items that should be returned in the result set. Acceptable values are 0 to 50, inclusive. The default value is 5.
- **onBehalfOfContentOwner** – (string) This parameter can only be used in a properly authorized request. Note: This parameter is intended exclusively for YouTube content partners. The `onBehalfOfContentOwner` parameter indicates that the request's authorization credentials identify a YouTube CMS user who is acting on behalf of the content owner specified in the parameter value. This parameter is intended for YouTube content partners that own and manage many different YouTube channels. It allows content owners to authenticate once and get access to all their video and channel data, without having to provide authentication credentials for each individual channel. The CMS account that the user authenticates with must be linked to the specified YouTube content owner.
- **order** – (string) The `order` parameter specifies the method that will be used to order resources in the API response. The default value is `relevance`. Acceptable values are:
  - `date` – Resources are sorted in reverse chronological order based on the date they were created.

rating – Resources are sorted from highest to lowest rating.

relevance – Resources are sorted based on their relevance to the search query. This is the default value for this parameter.

title – Resources are sorted alphabetically by title.

videoCount – Channels are sorted in descending order of their number of uploaded videos.

viewCount – Resources sorted from highest to lowest number of views. For live broadcasts, videos are sorted by number of concurrent viewers while the broadcasts are ongoing.

- **pageToken** – (string) The `pageToken` parameter identifies a specific page in the result set that should be returned. In an API response, the `nextPageToken` and `prevPageToken` properties identify other pages that could be retrieved.
- **publishedAfter** – (datetime) The `publishedAfter` parameter indicates that the API response should only contain resources created at or after the specified time. The value is an RFC 3339 formatted date-time value (1970-01-01T00:00:00Z).
- **publishedBefore** – (datetime) The `publishedBefore` parameter indicates that the API response should only contain resources created before or at the specified time. The value is an RFC 3339 formatted date-time value (1970-01-01T00:00:00Z).
- **regionCode** – (string) The `regionCode` parameter instructs the API to return search results for videos that can be viewed in the specified country. The parameter value is an ISO 3166-1 alpha-2 country code.
- **relevanceLanguage** – (string) The `relevanceLanguage` parameter instructs the API to return search results that are most relevant to the specified language. The parameter value is typically an ISO 639-1 two-letter language code. However, you should use the values `zh-Hans` for simplified Chinese and `zh-Hant` for traditional Chinese. Please note that results in other languages will still be returned if they are highly relevant to the search query term.
- **safeSearch** – (string) The `safeSearch` parameter indicates whether the search results should include restricted content as well as standard content. Acceptable values are:
  - moderate – YouTube will filter some content from search results and, at the least, will filter content that is restricted in your locale. Based on their content, search results could be removed from search results or demoted in search results. This is the default parameter value.
  - none – YouTube will not filter the search result set.
  - strict – YouTube will try to exclude all restricted content from the search result set.Based on their content, search results could be removed from search results or demoted in search results.
- **topicId** – (string) The `topicId` parameter indicates that the API response should only contain resources associated with the specified topic. The value identifies a Freebase topic ID.
- **type** – (string) The `type` parameter restricts a search query to only retrieve a particular type of resource. The value is a comma-separated list of resource types. The default value is `video,channel,playlist`. Acceptable values are: `channel`, `playlist`, and `video`
- **videoCaption** – (string) The `videoCaption` parameter indicates whether the API should filter video search results based on whether they have captions. If you specify a value for this parameter, you must also set the `type` parameter's value to `video`. Acceptable values are:

any – Do not filter results based on caption availability.

closedCaption – Only include videos that have captions.

none – Only include videos that do not have captions.

- **videoCategoryId** – (string) The `videoCategoryId` parameter filters video search results based on their category. If you specify a value for this parameter, you must also set the `type` parameter's value to `video`.

- **videoDefinition** – (string) The `videoDefinition` parameter lets you restrict a search to only include either high definition (HD) or standard definition (SD) videos. HD videos are available for playback in at least 720p, though higher resolutions, like 1080p, might also be available. If you specify a value for this parameter, you must also set the `type` parameter's value to `video`. Acceptable values are:

any – Return all videos, regardless of their resolution.

high – Only retrieve HD videos.

standard – Only retrieve videos in standard definition.

- **videoDimension** – (string) The `videoDimension` parameter lets you restrict a search to only retrieve 2D or 3D videos. If you specify a value for this parameter, you must also set the `type` parameter's value to `video`. Acceptable values are:

2d – Restrict search results to exclude 3D videos.

3d – Restrict search results to only include 3D videos.

any – Include both 3D and non-3D videos in returned results. This is the default value.

- **videoDuration** – (string) The `videoDuration` parameter filters video search results based on their duration. If you specify a value for this parameter, you must also set the `type` parameter's value to `video`. Acceptable values are:

any – Do not filter video search results based on their duration. This is the default value.

long – Only include videos longer than 20 minutes.

medium – Only include videos that are between four and 20 minutes long (inclusive).

short – Only include videos that are less than four minutes long.

- **videoEmbeddable** – (string) The `videoEmbeddable` parameter lets you to restrict a search to only videos that can be embedded into a webpage. If you specify a value for this parameter, you must also set the `type` parameter's value to `video`. Acceptable values are:

any – Return all videos, embeddable or not.

true – Only retrieve embeddable videos.

- **videoLicense** – (string) The `videoLicense` parameter filters search results to only include videos with a particular license. YouTube lets video uploaders choose to attach either the Creative Commons license or the standard YouTube license to each of their videos. If you specify a value for this parameter, you must also set the `type` parameter's value to `video`. Acceptable values are:

any – Return all videos, regardless of which license they have, that match the query parameters.

creativeCommon – Only return videos that have a Creative Commons license. Users can reuse videos with this license in other videos that they create.

youtube – Only return videos that have the standard YouTube license.

- **videoSyndicated** – (string) The `videoSyndicated` parameter lets you to restrict a search to only videos that can be played outside youtube.com. If you specify a value for this parameter, you must also set the `type` parameter’s value to `video`. Acceptable values are:
  - any – Return all videos, syndicated or not.
  - true – Only retrieve syndicated videos.
- **videoType** – (string) The `videoType` parameter lets you restrict a search to a particular type of videos. If you specify a value for this parameter, you must also set the `type` parameter’s value to `video`. Acceptable values are:
  - any – Return all videos.
  - episode – Only retrieve episodes of shows.
  - movie – Only retrieve movies.

**set\_logging\_level** (*level\_or\_name*)

Change the logging level during the session. Acceptable values are [0, 10, 20, 30, 40, 50, ‘NOTSET’, ‘DEBUG’, ‘INFO’, ‘WARNING’, ‘ERROR’, ‘CRITICAL’]

**youtube\_channel\_details** (*key, channel\_ids*)

Return details of channels for which the ids are given. Assumes `ids` is a comma-separated list of channel ids with no spaces.

**youtube\_video\_details** (*key, vid\_ids*)

Return details of videos for which the ids are given. Assumes `ids` is a comma-separated list of video ids with no spaces.

## 1.10 Import and Analyze Knowledge Graph Results on a Large Scale

If *analyzing SERPs* is the first step in understanding your rankings on search engines, then analyzing the knowledge graph can be thought of as step zero.

SERP positions for a certain keyword show how each page is ranked in comparison to all other eligible pages. Knowledge graph scores on the other hand, show the ranks of the different meanings that a word can take for Google (a person, a city, a brand, etc.).

**Warning:** From [Google’s documentation](#): This API is not suitable for use as a production-critical service. Your product should not form a critical dependence on this API.

It’s not clear whether this is from a technical reliability or a content correctness point of view, but it is what the docs mention. So please keep this in mind when using it.

### 1.10.1 Account Setup

In order to be able to send requests, you will need to [create a project](#), [set up billing](#), and [activate the knowledge graph API](#) for your project. You will then need to [create credentials](#) (API Key). Once you have that, you can use it as your `key` parameter when running requests, as shown below.

## 1.10.2 How to use Google’s Knowledge Graph API

What is “google”? Is it a search engine, a company, a brand, a very large number? What else is it?

And if it is all of those things, what is the relative ranking of each? What is the source of the information, its URL, images (if any)?

```
>>> key = 'YOUR_GOOGLE_DEVELOPER_KEY'
>>> google = knowledge_graph(key=key, query='google')
>>> google
  query  resultScore                                result.@type
↪ result.description    result.name
0  google      203191      ['Corporation', 'Organization', 'Thing']
↪   Technology company      Google
1  google      49462                                ['WebSite', 'Thing']
↪   Google                  Search
2  google      19142                                ['WebSite', 'Thing']
↪   nan                    Gmail
3  google      13251      ['Brand', 'WebSite', 'Thing']
↪   Website                Google Maps
4  google      7549      ['WebSite', 'SoftwareApplication', 'Thing']
↪   Website                Google Drive
5  google      6853                                ['WebSite', 'Thing']
↪   Website                Google Play
6  google      6543      ['SoftwareApplication', 'Thing']
↪   Web browser            Google Chrome
7  google      4312      ['Corporation', 'Organization', 'Thing']
↪ Multinational conglomerate company Alphabet Inc.
8  google      3395      ['SoftwareApplication', 'Thing']
↪   nan                    Google Account
9  google      1306                                ['Thing']
↪   nan                    Google

>>> google.columns
Index(['query', 'resultScore', '@type', 'result.@type', 'result.description',
      'result.image.contentUrl', 'result.image.url',
      'result.detailedDescription.articleBody',
      'result.detailedDescription.url', 'result.detailedDescription.license',
      'result.url', 'result.name', 'result.@id', 'query_time'],
      dtype='object')
```

The above table is a sample response from the `knowledge_graph()` function. Many more columns are available as you can see in the second line above. We can see that “google” is a company, with a result score of 203,191 and it is a search engine/website with a result score of 49,462. It is then understood as an email application, a mapping application, and so on, as you can see in the `result.name` column.

You can also see that we get the types under which this result falls, in the `result.@type` column. Multiple types show the type inheritance, and as you can also see, everything is a “Thing”. This is the top element in the type hierarchy under which everything belongs.

Like `serp_goog()` and `serp_youtube()`, this function works in the same manner, creating, sending, and aggregating the product of the arguments passed to it.

For example if you run

```
>>> knowledge_graph(key=key, query=['google', 'bing'], languages=['en', 'fr', 'de'])
```

The function will send 2 (queries) x 3 languages = 6 requests.

(google, en), (google, fr), (google, de), (bing, en), (bing, fr), (bing, de)

This is actually the main value of having this function, because you usually want a large sample to evaluate certain keywords across languages or types.

Let's check what "seo" and "search engine optimization" mean in different languages.

```
>>> seo = knowledge_graph(key=key, query=['seo', 'search engine optimization'],
↳languages=['en', 'es', 'de'])
>>> seo
```

	query	languages	resultScore	result.
↳name				↳result.@type
↳result.description				
0	search engine optimization	de	3587	↳Suchmaschinenoptimierung
↳				↳nan
1	search engine optimization	de	321	↳Suchmaschinenoptimierung
↳				↳nan
2	search engine optimization	de	252	↳Suchmaschinenmarketing
↳				↳nan
4	search engine optimization	en	71756	↳engine optimization
↳				↳nan
5	search engine optimization	en	5056	↳engine marketing
↳				↳nan
6	search engine optimization	en	576	↳SEOP, Inc.
↳				↳Company
13	seo	de	3313	↳Seoul
↳				↳Hauptstadt von Südkorea
14	seo	de	1509	↳Seo Ye-
↳				↳ji
↳				↳Schauspielerin
15	seo	de	584	↳Suchmaschinenoptimierung
↳				↳nan
33	seo	es	1509	↳Seo Ye-ji
↳				↳Actriz
34	seo	es	584	↳Posicionamiento en buscadores
↳				↳nan
35	seo	es	316	↳Jin
↳				↳Cantante
53	seo	en	8760	↳Search
↳				↳nan
54	seo	en	3313	↳Seoul
↳				↳Capital of South Korea
55	seo	en	1435	↳Sulli
↳				↳South
↳				↳Korean actress

```
>>> seo.columns
Index(['query', 'languages', 'resultScore', '@type', 'result.name',
      'result.@type', 'result.id', 'result.image.contentUrl',
      'result.image.url', 'result.detailedDescription.license',
      'result.detailedDescription.url',
      'result.detailedDescription.articleBody', 'result.description',
      'result.url', 'query_time'],
```

(continues on next page)



(continued from previous page)

dtype='object')

It's interesting to see how the same word can mean different things in different contexts.

**knowledge\_graph**(*key*, *query=None*, *ids=None*, *languages=None*, *types=None*, *prefix=None*, *limit=None*)

Query Google's Knowledge Graph with any combination of parameters.

Note that Google's documentation states that "This API is not suitable for use as a production-critical service." So please keep this in mind.

#### Parameters

- **key** (*string*) – Your Google developer key.
- **query** (*string*) – A literal string to search for in the Knowledge Graph.
- **ids** (*string*) – A list of entity IDs to search for in the Knowledge Graph.
- **languages** (*string*) – The list of language codes (defined in ISO 639) to run the query with, for instance *en*.
- **types** (*string*) – Restricts returned entities to those of the specified types. For example, you can specify *Person* (as defined in <http://schema.org/Person>) to restrict the results to entities representing people. If multiple types are specified, returned entities will contain one or more of these types.
- **prefix** (*boolean*) – Enables prefix (initial substring) match against names and aliases of entities. For example, a prefix *Jung* will match entities and aliases such as *Jung*, *Jungle*, and *Jung-ho Kang*.
- **limit** (*number*) – Limits the number of entities to be returned. Maximum is 500. Default is 20. Requests with high limits have a higher chance of timing out.

<https://developers.google.com/knowledge-graph/reference/rest/v1>

## 1.11 Split, Parse, and Analyze URLs

Extracting information from URLs can be a little tedious, yet very important. Using the standard for URLs we can extract a lot of information in a fairly structured manner.

There are many situations in which you have many URLs that you want to better understand:

- **Analytics reports:** Whichever analytics system you use, whether Google Analytics, search console, or any other reporting tool that reports on URLs, your reports can be enhanced by splitting URLs, and in effect becoming four or five data points as opposed to one.
- **Crawl datasets:** The result of any crawl you run typically contains the URLs, which can benefit from the same enhancement.
- **SERP datasets:** Which are basically about URLs.
- **Extracted URLs:** Extracting URLs from social media posts is one thing you might want to do to better understand those posts, and further splitting URLs can also help.
- **XML sitemaps:** Right after downloading a sitemap(s) splitting it further can help in giving a better perspective on the dataset.

The main function here is `url_to_df()`, which as the name suggests, converts URLs to DataFrames.

```

>>> urls ['https://net.location.com/path_1/path_2?price=10&color=blue#frag_1',
...       'https://net.location.com/path_1/path_2?price=15&color=red#frag_2']
>>> url_to_df(urls)

```

	netloc	path	query	fragment	url	scheme	dir_
↪1	dir_2	query_price	query_color				
0	https://net.location.com	/path_1/path_2	price=10&color=blue	frag_1	https://net.location.com/path_1/path_2?price=10&color=blue#frag_1	https	
↪	net.location.com	/path_1/path_2	price=10&color=blue	frag_1			
↪	path_1	path_2	10	blue			
1	https://net.location.com	/path_1/path_2	price=15&color=red	frag_2	https://net.location.com/path_1/path_2?price=15&color=red#frag_2	https	
↪	net.location.com	/path_1/path_2	price=15&color=red	frag_2			
↪	path_1	path_2	15	red			

- **url:** The original URLs are listed as a reference. They are decoded for easier reading, and you can set `decode=False` if you want to retain the original encoding.
- **scheme:** Self-explanatory. Note that you can also provide relative URLs `/category/sub-category?one=1&two=2` in which case the `url`, `scheme` and `netloc` columns would be empty. You can mix relative and absolute URLs as well.
- **netloc:** The network location is the sub-domain (optional) together with the domain and top-level domain and/or the country domain.
- **path:** The slug of the URL, excluding the query parameters and fragments if any. The path is also split in to directories `dir_1`, `dir_2`, `dir_3`... to make it easier to categorize and analyze the URLs.
- **query:** If query parameters are available they are given in this column, but more importantly they are parsed and included in separate columns, where each parameter has its own column (with the keys being the names). As in the example above, the query `price=10&color=blue` becomes two columns, one for price and the other for color. If any other URLs in the dataset contain the same parameters, their values will be populated in the same column, and `NA` otherwise.
- **fragment:** The final part of the URL after the hash mark #, linking to a part in the page.
- **query\_\*:** The query parameter names are prepended with `query_` to make it easy to filter them out, and to avoid any name collisions with other columns, if some URL contains a query parameter called “url” for example. In the unlikely event of having a repeated parameter in the same URL, then their values would be delimited by two “@” signs `one@@two@@three`. It’s unusual, but it happens.
- **hostname and port:** If available a column for ports will be shown, and if the hostname is different from `netloc` it would also have its own column.

### 1.11.1 Query Parameters

The great thing about parameters is that the names are descriptive (mostly!) and once given a certain column you can easily understand what data they contain. Once this is done, you can sort the products by price, filter by destination, get the red and blue items, and so on.

### 1.11.2 The URL Path (Directories):

Here things are not as straightforward, and there is no way to know what the first or second directory is supposed to indicate. In general, I can think of three main situations that you can encounter while analyzing directories.

- **Consistent URLs:** This is the simplest case, where all URLs follow the same structure. */en/product1* clearly shows that the first directory indicates the language of the page. So it can also make sense to rename those columns once you have discovered their meaning.
- **Inconsistent URLs:** This is similar to the previous situation. All URLs follow the same pattern with a few exceptions. Take the following URLs for example:
  - */topic1/title-of-article-1*
  - */es/topic1/title-of-article-2*
  - */es/topic2/title-of-article-3*
  - */topic2/title-of-artilce-4*

You can see that they follow the pattern */language/topic/article-title*, except for English, which is not explicitly mentioned, but its articles can be identified by having two instead of three directories, as we have for “/es/”. If URLs are split in this case, you will end up with *dir\_1* having “topic1”, “es”, “es”, and “topic2”, which distorts the data. Actually you want to have “en”, “es”, “es”, “en”. In such cases, after making sure you have the right rules and patterns, you might create special columns or replace/insert values to make them consistent, and get them to a state similar to the first example.

- **URLs of different types:** In many cases you will find that sites having different types of pages with completely different roles on the site.
  - */blog/post-1-title.html*
  - */community/help/topic\_1*
  - */community/help/topic\_2*

Here, once you split the directories, you will see that they don’t align properly (because of different lengths), and they can’t be compared easily. A good approach is to split your dataset into one for blog posts and another for community content for example.

The ideal case for the *path* part of the URL is to be split into directories of equal length across the dataset, having the right data in the right columns and *NA* otherwise. Or, splitting the dataset and analyzing separately.

**url\_to\_df** (*urls*, *decode=True*)

Split the given URLs into their components to a DataFrame.

Each column will have its own component, and query parameters and directories will also be parsed and given special columns each.

#### Parameters

- **urls** (*url*) – A list of URLs to split into components
- **decode** (*bool*) – Whether or not to decode the given URLs

**Return DataFrame split** A DataFrame with a column for each component

## 1.12 Emoji: Extract, Analyze, and Get Insights

An emoji is worth a thousand words! Regular expressions and helper functionality to aid in extracting and finding emoji from text.

EMOJI_ENTRIES	A dictionary of the full emoji list together with unicode code points, textual name, group, and sub-group. Based on v13.0 <a href="https://unicode.org/Public/emoji/13.0/emoji-test.txt">https://unicode.org/Public/emoji/13.0/emoji-test.txt</a>
emoji_df	The same dictionary as a pandas DataFrame
<code>extract_emoji</code>	A function for extracting and summarizing emoji in a text list, with statistics about frequencies and usage.
<code>emoji_search</code>	A function for searching across names, groups, and sub-groups to find emoji based on your keywords of choice.
EMOJI_RAW	A regular expression to extract the full list. See here on how it was developed: <a href="https://www.kaggle.com/eliasdabbas/how-to-create-a-python-regex-to-extract-emoji">https://www.kaggle.com/eliasdabbas/how-to-create-a-python-regex-to-extract-emoji</a>

### `emoji_search` (*regex*)

Return a DataFrame of all emoji entries where any description contains *regex*.

“description” can be the name of the emoji, its group, or sub-group.

**Parameters** *regex* (*str*) – regular expression (case insensitive)

```
>>> emoji_search('dog')
  codepoint      status  emoji      name      group
↳sub_group
0      1F436  fully-qualified  dog face  Animals & Nature
↳animal-mammal
1      1F415  fully-qualified      dog  Animals & Nature
↳animal-mammal
2      1F9AE  fully-qualified  guide dog  Animals & Nature
↳animal-mammal
3  1F415 200D 1F9BA  fully-qualified  service dog  Animals & Nature  animal-
↳mammal
4      1F32D  fully-qualified      hot dog      Food & Drink  food-
↳prepared
```

```
>>> blue = adv.emoji_search('blue')
  codepoint      status  emoji      name      group
↳sub_group
0      1F499  fully-qualified  blue heart  Smileys & Emotion
↳emotion
1      1FAD0  fully-qualified  blueberries  Food & Drink  food-
↳fruit
2      1F4D8  fully-qualified  blue book      Objects  book-
↳paper
3      1F535  fully-qualified  blue circle      Symbols
↳geometric
4      1F7E6  fully-qualified  blue square      Symbols
↳geometric
5      1F537  fully-qualified  large blue diamond  Symbols
↳geometric
6      1F539  fully-qualified  small blue diamond  Symbols
↳geometric
```

### `extract_emoji` (*text\_list*)

Return a summary dictionary about emoji in *text\_list*

Get a summary of the number of emoji, their frequency, the top ones, and more.

**Parameters** `text_list` (*list*) – A list of text strings.

**Returns** `summary` A dictionary with various stats about emoji

```
>>> posts = ['I am grinning ', 'A grinning cat ',
...          'hello! ', 'Just text']
```

```
>>> emoji_summary = extract_emoji(posts)
>>> emoji_summary.keys()
dict_keys(['emoji', 'emoji_text', 'emoji_flat', 'emoji_flat_text',
'emoji_counts', 'emoji_freq', 'top_emoji', 'top_emoji_text',
'top_emoji_groups', 'top_emoji_sub_groups', 'overview'])
```

```
>>> emoji_summary['emoji']
[[''], [''], [' ', ' ', ' ', ' ', ' '], []]
```

```
>>> emoji_summary['emoji_text']
[['grinning face'], ['grinning cat'], ['grinning face', 'grinning face',
'grinning face', 'yellow heart', 'yellow heart'], []]
```

A simple extract of emoji from each of the posts. An empty list if none exist

```
>>> emoji_summary['emoji_flat']
[' ', ' ', ' ', ' ', ' ', ' ', ' ']
```

```
>>> emoji_summary['emoji_flat_text']
['grinning face', 'grinning cat', 'grinning face', 'grinning face',
'grinning face', 'yellow heart', 'yellow heart']
```

All emoji in one flat list.

```
>>> emoji_summary['emoji_counts']
[1, 1, 5, 0]
```

The count of emoji per post.

```
>>> emoji_summary['emoji_freq']
[(0, 1), (1, 2), (5, 1)]
```

Shows how many posts had 0, 1, 2, 3, etc. emoji (number\_of\_emoji, count)

```
>>> emoji_summary['top_emoji']
[('', 4), (' ', 2), (' ', 1)]
```

```
>>> emoji_summary['top_emoji_text']
[('grinning face', 4), ('yellow heart', 2),
('grinning cat', 1)]
```

```
>>> emoji_summary['top_emoji_groups']
[('Smileys & Emotion', 7)]
```

```
>>> emoji_summary['top_emoji_sub_groups']
[('face-smiling', 4), ('emotion', 2), ('cat-face', 1)]
```

```
>>> emoji_summary['overview']
{'num_posts': 4,
 'num_emoji': 7,
 'emoji_per_post': 1.75,
 'unique_emoji': 3}
```

## 1.13 Extract structured entities from text lists

Structured entities are pattern matches and not inferred entities. Some example are hashtags, emoji, mentions, questions, and so on. This is in contrast to entity extraction which are inferred from the context of the sentence (people, companies, brands and so on).

All functions start with `extract_` and have a descriptive name for the type of entity that they extract.

There is also a generic `extract` function which powers all others, and it can be used for any other pattern not included. It takes a regular expression, and returns a similar dictionary to the other functions.

### 1.13.1 Extract Functions

<code>extract()</code>	A generic function that takes a regex to extract any pattern you want
<code>extract_currency()</code>	(Currency symbols together with surrounding text for context. This does not include currency abbreviations (USD, EUR, JPY, etc.), only symbols (\$, £, €, etc).
<code>adverttools.emoji.extract_emoji()</code>	All the emoji database, together with textual names, groups and sub-groups.
<code>extract_exclamation_marks()</code>	Sentences that end with an exclamation mark!
<code>extract_hashtags()</code>	Extract hashtags with descriptive statistics.
<code>extract_intense_words()</code>	Words that contain three or more repeated characters to express an intense feeling (positive or negative), “I looooooovvvvee this thing”.
<code>extract_mentions()</code>	User mentions in social media posts. Also useful for network analysis.
<code>extract_numbers()</code>	Any numbers that are included the text list. Included a modifiable list of separators to use (“,” “.”, “-”, etc.).
<code>extract_questions()</code>	Questions included in the text list.
<code>extract_urls()</code>	URLs in the text list.
<code>extract_words()</code>	Any arbitrary words that you want extracted. Works in two modes, either the word should fully match the pattern, or as part of a longer word, (“rest” can be matched from “restaurant” or not).

All functions return a dictionary with the entities extracted, along with helpful statistics. Since the entities have different meanings, most of them return additional keys depending on the context.

The recommended way of using:

```
>>> import adverttools as adv
>>> text_list = ['This is the first #text.', 'Second #sentence is here.',
... 'Hello, how are you?', 'This #sentence is the last #sentence']
>>> hashtag_summary = adv.extract_hashtags(text_list)
>>> hashtag_summary.keys()
dict_keys(['hashtags', 'hashtags_flat', 'hashtag_counts', 'hashtag_freq',
          'top_hashtags', 'overview'])
```

Now you can start exploring:

```
>>> hashtag_summary['overview']
{'num_posts': 4,
 'num_hashtags': 4,
 'hashtags_per_post': 1.0,
 'unique_hashtags': 2}
```

```
>>> hashtag_summary['hashtags']
[['#text'], ['#sentence'], [], ['#sentence', '#sentence']]
>>> hashtag_summary['hashtags_flat']
['#text', '#sentence', '#sentence', '#sentence']
>>> hashtag_summary['hashtag_counts']
[1, 1, 0, 2]
>>> hashtag_summary['hashtag_freq']
[(0, 1), (1, 2), (2, 1)]
>>> hashtag_summary['top_hashtags']
[('#sentence', 3), ('#text', 1)]
```

**extract** (*text\_list*, *regex*, *key\_name*, *extracted=None*, *\*\*kwargs*)

Return a summary dictionary about arbitrary matches in *text\_list*.

This function is used by other specialized functions to extract certain elements (hashtags, mentions, emojis, etc.). It can be used for other arbitrary elements/matches. You only need to provide your own regex.

#### Parameters

- **text\_list** (*list*) – Any list of strings (social posts, page titles, etc.)
- **regex** (*str*) – The regex pattern to use for extraction.
- **key\_name** (*str*) – The name of the object extracted in singular form (hashtag, mention, etc.)
- **extracted** (*list(list)*) – List of lists, optional. If the regex is not straightforward, and matches need to be made with special code, provide the extracted words/matches as a list for each element of *text\_list*.
- **kwargs** (*mapping*) – Other kwargs that might be needed.

**Return summary** A dictionary summarizing the extracted data.

**extract\_currency** (*text\_list*, *left\_chars=20*, *right\_chars=20*)

Return a summary dictionary about currency symbols in *text\_list*

Get a summary of the number of currency symbols, their frequency, the top ones, and more.

#### Parameters

- **text\_list** (*list*) – A list of text strings.
- **left\_chars** (*int*) – The number of characters to extract, to the left of the symbol when getting *surrounding\_text*
- **right\_chars** (*int*) – The number of characters to extract, to the right of the symbol when getting *surrounding\_text*

**Returns summary** A dictionary with various stats about currencies

```
>>> posts = ['today 1 is around $4k', 'and in £ & €?', 'no idea']
>>> currency_summary = extract_currency(posts)
>>> currency_summary.keys()
dict_keys(['currency_symbols', 'currency_symbols_flat',
```

(continues on next page)

(continued from previous page)

```
'currency_symbol_counts', 'currency_symbol_freq',
'top_currency_symbols', 'overview', 'currency_symbol_names']])
```

```
>>> currency_summary['currency_symbols']
[['', '$'], ['', '£', '€'], []]
```

A simple extract of currencies from each of the posts. An empty list if none exist

```
>>> currency_summary['currency_symbols_flat']
['', '$', '', '£', '€']
```

All currency symbols in one flat list.

```
>>> currency_summary['currency_symbol_counts']
[2, 3, 0]
```

The count of currency symbols per post.

```
>>> currency_summary['currency_symbol_freq']
[(0, 1), (2, 1), (3, 1)]
```

Shows how many posts had 0, 1, 2, 3, etc. currency symbols (number\_of\_symbols, count)

```
>>> currency_summary['top_currency_symbols']
[('', 2), ('$ ', 1), ('£ ', 1), ('€ ', 1)]
```

```
>>> currency_summary['currency_symbol_names']
[['bitcoin sign', 'dollar sign'], ['bitcoin sign', 'pound sign',
'euro sign'], []]
```

```
>>> currency_summary['surrounding_text']
[['today 1 is around $4k'], ['and in £ & €?'], []]
```

```
>>> extract_currency(posts, 5, 5)['surrounding_text']
[['oday 1 is ', 'ound $4k'], ['and in £', ' & €?'], []]
```

```
>>> extract_currency(posts, 0, 3)['surrounding_text']
[['1 i', '$4k'], [' in', '£ & ', '€?'], []]
```

```
>>> currency_summary['overview']
{'num_posts': 3,
 'num_currency_symbols': 5,
 'currency_symbols_per_post': 1.6666666666666667,
 'unique_currency_symbols': 4}
```

### **extract\_exclamations** (*text\_list*)

Return a summary dictionary about exclamation (mark)s in *text\_list*

Get a summary of the number of exclamation marks, their frequency, the top ones, as well the exclamations written/said.

**Parameters** *text\_list* (*list*) – A list of text strings.

**Returns summary** A dictionary with various stats about exclamations



```
>>> posts = ['Who are you!', 'What is this!', 'No exclamation here?']
>>> exclamation_summary = extract_exclamations(posts)
>>> exclamation_summary.keys()
dict_keys(['exclamation_marks', 'exclamation_marks_flat',
'exclamation_mark_counts', 'exclamation_mark_freq',
'top_exclamation_marks', 'overview', 'exclamation_mark_names',
'exclamation_text'])
```

```
>>> exclamation_summary['exclamation_marks']
[['!'], ['!'], []]
```

A simple extract of exclamation marks from each of the posts. An empty list if none exist

```
>>> exclamation_summary['exclamation_marks_flat']
['!', '!']
```

All exclamation marks in one flat list.

```
>>> exclamation_summary['exclamation_mark_counts']
[1, 1, 0]
```

The count of exclamation marks per post.

```
>>> exclamation_summary['exclamation_mark_freq']
[(0, 1), (1, 2)]
```

Shows how many posts had 0, 1, 2, 3, etc. exclamation marks (number\_of\_symbols, count)

```
>>> exclamation_summary['top_exclamation_marks']
[(!, 2)]
```

Might be interesting if you have different types of exclamation marks

```
>>> exclamation_summary['exclamation_mark_names']
[['exclamation mark'], ['exclamation mark'], []]
```

```
>>> exclamation_summary['overview']
{'num_posts': 3,
'num_exclamation_marks': 2,
'exclamation_marks_per_post': 0.6666666666666666,
'unique_exclamation_marks': 1}
```

```
>>> posts2 = ["don't go there!", '. !', '¡Hola! ¿cómo estás?',
... 'a few different exclamation marks! make sure you see them!']
```

```
>>> exclamation_summary = extract_exclamations(posts2)
```

```
>>> exclamation_summary['exclamation_marks']
[['!'], ['!'], ['!', '!'], ['!', '!']]
# might be displayed in opposite order due to RTL exclamation mark
A simple extract of exclamation marks from each of the posts.
An empty list if none exist
```

```
>>> exclamation_summary['exclamation_marks_flat']
['!', '!', '!', '!', '!', '!']
```

All exclamation marks in one flat list.

```
>>> exclamation_summary['exclamation_mark_counts']
[1, 1, 2, 2]
```

The count of exclamation marks per post.

```
>>> exclamation_summary['exclamation_mark_freq']
[(1, 2), (2, 2)]
```

Shows how many posts had 0, 1, 2, 3, etc. exclamation marks (number\_of\_symbols, count)

```
>>> exclamation_summary['top_exclamation_marks']
[('!', 5), (';', 1)]
```

Might be interesting if you have different types of exclamation marks

```
>>> exclamation_summary['exclamation_mark_names']
[['exclamation mark'], ['exclamation mark'],
 ['inverted exclamation mark', 'exclamation mark'],
 ['exclamation mark', 'exclamation mark']]
```

```
>>> exclamation_summary['overview']
{'num_posts': 4,
 'num_exclamation_marks': 6,
 'exclamation_marks_per_post': 1.5,
 'unique_exclamation_marks': 4}
```

### **extract\_hashtags** (*text\_list*)

Return a summary dictionary about hashtags in *text\_list*

Get a summary of the number of hashtags, their frequency, the top ones, and more.

**Parameters** *text\_list* (*list*) – A list of text strings.

**Returns summary** A dictionary with various stats about hashtags

```
>>> posts = ['i like #blue', 'i like #green and #blue', 'i like all']
>>> hashtag_summary = extract_hashtags(posts)
>>> hashtag_summary.keys()
dict_keys(['hashtags', 'hashtags_flat', 'hashtag_counts', 'hashtag_freq',
 'top_hashtags', 'overview'])
```

```
>>> hashtag_summary['hashtags']
[['#blue'], ['#green', '#blue'], []]
```

A simple extract of hashtags from each of the posts. An empty list if none exist

```
>>> hashtag_summary['hashtags_flat']
['#blue', '#green', '#blue']
```

All hashtags in one flat list.

```
>>> hashtag_summary['hashtag_counts']
[1, 2, 0]
```

The count of hashtags per post.

```
>>> hashtag_summary['hashtag_freq']
[(0, 1), (1, 1), (2, 1)]
```

Shows how many posts had 0, 1, 2, 3, etc. hashtags (number\_of\_hashtags, count)

```
>>> hashtag_summary['top_hashtags']
[('#blue', 2), ('#green', 1)]
```

```
>>> hashtag_summary['overview']
{'num_posts': 3,
 'num_hashtags': 3,
 'hashtags_per_post': 1.0,
 'unique_hashtags': 2}
```

### **extract\_intense\_words** (*text\_list*, *min\_reps=3*)

Return a summary dictionary about intense words in *text\_list*

Get all instances of usage of intense words (positive or negative), using words that have *min\_reps* or more repetitions of characters. “I looooooveeee youuuuuuu”, and “I haaatttteeeee youuuuuuu” are both intense.

#### **Parameters**

- **text\_list** (*list*) – A text list from which to extract intense words
- **min\_reps** (*int*) – The number of times a character has to be repeated for the word to be considered intense.

**Returns summary** A dictionary with various stats about intense words

### **extract\_mentions** (*text\_list*)

Return a summary dictionary about mentions in *text\_list*

Get a summary of the number of mentions, their frequency, the top ones, and more.

**Parameters** **text\_list** (*list*) – A list of text strings.

**Returns summary** A dictionary with various stats about mentions

```
>>> posts = ['hello @john and @jenny', 'hi there @john', 'good morning']
>>> mention_summary = extract_mentions(posts)
>>> mention_summary.keys()
dict_keys(['mentions', 'mentions_flat', 'mention_counts', 'mention_freq',
'top_mentions', 'overview'])
```

```
>>> mention_summary['mentions']
[['@john', '@jenny'], ['@john'], []]
```

A simple extract of mentions from each of the posts. An empty list if none exist

```
>>> mention_summary['mentions_flat']
['@john', '@jenny', '@john']
```

All mentions in one flat list.

```
>>> mention_summary['mention_counts']
[2, 1, 0]
```

The count of mentions per post.

```
>>> mention_summary['mention_freq']
[(0, 1), (1, 1), (2, 1)]
```

Shows how many posts had 0, 1, 2, 3, etc. mentions (number\_of\_mentions, count)

```
>>> mention_summary['top_mentions']
[('@john', 2), ('@jenny', 1)]
```

```
>>> mention_summary['overview']
{'num_posts': 3, # number of posts
 'num_mentions': 3,
 'mentions_per_post': 1.0,
 'unique_mentions': 2}
```

**extract\_numbers** (*text\_list*, *number\_separators*=',;.-')

Return a summary dictionary about numbers in *text\_list*, separated by any of *number\_separators*

Get a summary of the number of numbers, their frequency, the top ones, and more. Typically, numbers would contain separators to make them easier to read, so these are included by default, which you can modify.

#### Parameters

- **text\_list** (*list*) – A list of text strings.
- **number\_separators** (*list(str)*) – A list of separators that you want to be included as part of the extracted numbers.

**Returns summary** A dictionary with various stats about the numbers

```
>>> posts = ['text before 123', '123,456 text after', 'phone 333-444-555',
 'multiple 123,456 and 123.456.789']
>>> number_summary = extract_numbers(posts)
>>> number_summary.keys()
dict_keys(['numbers', 'numbers_flat', 'number_counts', 'number_freq',
 'top_numbers', 'overview'])
```

```
>>> number_summary['numbers']
[['123'], ['123,456'], ['333-444-555'], ['123,456', '123.456.789']]
```

A simple extract of number from each of the posts. An empty list if none exist

```
>>> number_summary['numbers_flat']
['123', '123,456', '333-444-555', '123,456', '123.456.789']
```

All numbers in one flat list.

```
>>> number_summary['number_counts']
[1, 1, 1, 2]
```

The count of numbers per post.

```
>>> number_summary['number_freq']
[(1, 3), (2, 1)]
```

Shows how many posts had 0, 1, 2, 3, etc. numbers (number\_of\_numbers, count)

```
>>> number_summary['top_numbers']
[('123,456', 2), ('123', 1), ('333-444-555', 1), ('123.456.789', 1)]
```

```
>>> number_summary['overview']
{'num_posts': 4,
 'num_numbers': 5,
 'numbers_per_post': 1.25,
 'unique_numbers': 4}
```

**extract\_questions** (*text\_list*)

Return a summary dictionary about question(mark)s in *text\_list*

Get a summary of the number of question marks, their frequency, the top ones, as well the questions asked.

**Parameters** *text\_list* (*list*) – A list of text strings.

**Returns summary** A dictionary with various stats about questions

```
>>> posts = ['How are you?', 'What is this?', 'No question Here!']
>>> question_summary = extract_questions(posts)
>>> question_summary.keys()
dict_keys(['question_marks', 'question_marks_flat',
 'question_mark_counts', 'question_mark_freq', 'top_question_marks',
 'overview', 'question_mark_names', 'question_text'])
```

```
>>> question_summary['question_marks']
[['?'], ['?'], []]
```

A simple extract of question marks from each of the posts. An empty list if none exist

```
>>> question_summary['question_marks_flat']
['?', '?']
```

All question marks in one flat list.

```
>>> question_summary['question_mark_counts']
[1, 1, 0]
```

The count of question marks per post.

```
>>> question_summary['question_mark_freq']
[(0, 1), (1, 2)]
```

Shows how many posts had 0, 1, 2, 3, etc. question marks (*number\_of\_symbols*, *count*)

```
>>> question_summary['top_question_marks']
[( '?', 2)]
```

Might be interesting if you have different types of question marks (Arabic, Spanish, Greek, Armenian, or other)

```
>>> question_summary['question_mark_names']
[['question mark'], ['question mark'], []]
```

```
>>> question_summary['overview']
{'num_posts': 3,
 'num_question_marks': 2,
 'question_marks_per_post': 0.6666666666666666,
 'unique_question_marks': 1}
```

```
>>> posts2 = [' ', '. ', 'Hola, ¿cómo estás?',
... 'Can you see the new questions? Did you notice the different marks?']
```

```
>>> question_summary = extract_questions(posts2)
```

```
>>> question_summary['question_marks']
[[''], [''], ['¿', '?'], ['?', '?']]
# might be displayed in opposite order due to RTL question mark
A simple extract of question marks from each of the posts. An empty list if
none exist
```

```
>>> question_summary['question_marks_flat']
['', '', '¿', '?', '?', '?']
```

All question marks in one flat list.

```
>>> question_summary['question_mark_counts']
[1, 1, 2, 2]
```

The count of question marks per post.

```
>>> question_summary['question_mark_freq']
[(1, 2), (2, 2)]
```

Shows how many posts had 0, 1, 2, 3, etc. question marks (number\_of\_symbols, count)

```
>>> question_summary['top_question_marks']
[('?', 3), ('', 1), ('', 1), ('¿', 1)]
```

Might be interesting if you have different types of question marks (Arabic, Spanish, Greek, Armenian, or other)

```
>>> question_summary['question_mark_names']
[['greek question mark'], ['arabic question mark'],
['inverted question mark', 'question mark'],
['question mark', 'question mark']]
# correct order
```

```
>>> question_summary['overview']
{'num_posts': 4,
 'num_question_marks': 6,
 'question_marks_per_post': 1.5,
 'unique_question_marks': 4}
```

### **extract\_urls** (*text\_list*)

Return a summary dictionary about URLs in *text\_list*

Get a summary of the number of URLs, their frequency, the top ones, and more. This does NOT validate URLs, www.a.b would count as a URL

**Parameters** *text\_list* (*list*) – A list of text strings.

**Returns** *summary* A dictionary with various stats about URLs

```
>>> posts = ['one link http://example.com', 'two: http://a.com www.b.com',
... 'no links here',
... 'long url http://example.com/one/two/?1=one&2=two']
```

(continues on next page)

(continued from previous page)

```
>>> url_summary = extract_urls(posts)
>>> url_summary.keys()
dict_keys(['urls', 'urls_flat', 'url_counts', 'url_freq',
'top_urls', 'overview', 'top_domains', 'top_tlds'])
```

```
>>> url_summary['urls']
[['http://example.com'],
 ['http://a.com', 'http://www.b.com'],
 [],
 ['http://example.com/one/two/?1=one&2=two']]
```

A simple extract of urls from each of the posts. An empty list if none exist

```
>>> url_summary['urls_flat']
['http://example.com', 'http://a.com', 'http://www.b.com',
'http://example.com/one/two/?1=one&2=two']
```

All urls in one flat list.

```
>>> url_summary['url_counts']
[1, 2, 0, 1]
```

The count of urls per post.

```
>>> url_summary['url_freq']
[(0, 1), (1, 2), (2, 1)]
```

Shows how many posts had 0, 1, 2, 3, etc. urls (number\_of\_urls, count)

```
>>> url_summary['top_urls']
[('http://example.com', 1), ('http://a.com', 1), ('http://www.b.com', 1),
 ('http://example.com/one/two/?1=one&2=two', 1)]
```

```
>>> url_summary['top_domains']
[('example.com', 2), ('a.com', 1), ('www.b.com', 1)]
```

```
>>> url_summary['top_tlds']
[('com', 4)]
```

```
>>> url_summary['overview']
{'num_posts': 4,
 'num_urls': 4,
 'urls_per_post': 1.0,
 'unique_urls': 4}
```

**extract\_words** (*text\_list*, *words\_to\_extract*, *entire\_words\_only=False*)

Return a summary dictionary about *words\_to\_extract* in *text\_list*.

Get a summary of the number of words, their frequency, the top ones, and more.

#### Parameters

- **text\_list** (*list*) – A list of text strings.
- **words\_to\_extract** (*list*) – A list of words to extract from *text\_list*.
- **entire\_words\_only** (*bool*) – Whether or not to find only complete words (as specified by *words\_to\_find*) or find any any of the words as part of longer strings.

**Returns summary** A dictionary with various stats about the words

```
>>> posts = ['there is rain, it is raining', 'there is snow and rain',
             'there is no rain, it is snowing', 'there is nothing']
>>> word_summary = extract_words(posts, ['rain', 'snow'], True)
>>> word_summary.keys()
dict_keys(['words', 'words_flat', 'word_counts', 'word_freq',
          'top_words', 'overview'])
```

```
>>> word_summary['overview']
{'num_posts': 4,
 'num_words': 4,
 'words_per_post': 1,
 'unique_words': 2}
```

```
>>> word_summary['words']
[['rain'], ['snow', 'rain'], ['rain'], []]
```

A simple extract of mentions from each of the posts. An empty list if none exist

```
>>> word_summary['words_flat']
['rain', 'snow', 'rain', 'rain']
```

All mentions in one flat list.

```
>>> word_summary['word_counts']
[1, 2, 1, 0]
```

The count of mentions for each post.

```
>>> word_summary['word_freq']
[(0, 1) (1, 2), (2, 1)]
```

Shows how many posts had 0, 1, 2, 3, etc. words (number\_of\_words, count)

```
>>> word_summary['top_words']
[('rain', 3), ('snow', 1)]
```

Check the same posts extracting any occurrence of the specified words with `entire_words_only=False`:

```
>>> word_summary = extract_words(posts, ['rain', 'snow'], False)
```

```
>>> word_summary['overview']
{'num_posts': 4, # number of posts
 'num_words': 6,
 'words_per_post': 1.5,
 'unique_words': 4}
```

```
>>> word_summary['words']
[['rain', 'raining'], ['snow', 'rain'], ['rain', 'snowing'], []]
```

Note that the extracted words are the complete words so you can see where they occurred. In case “training” was mentioned, you would see that it is not related to rain for example.

```
>>> word_summary['words_flat']
['rain', 'raining', 'snow', 'rain', 'rain', 'snowing']
```



All mentions in one flat list.

```
>>> word_summary['word_counts']  
[2, 2, 2, 0]
```

```
>>> word_summary['word_freq']  
[(0, 1), (2, 3)]
```

Shows how many posts had 0, 1, 2, 3, etc. words (number\_of\_words, count)

```
>>> word_summary['top_words']  
[('rain', 3), ('raining', 1), ('snow', 1), ('snowing', 1)]
```

## 1.14 Stopwords in Several Languages

List of stopwords by the spaCy<sup>1</sup> package, useful in text mining, analyzing content of social media posts, tweets, web pages, keywords, etc.

Each list is accessible as part of a dictionary `stopwords` which is a normal Python dictionary.

### 1.14.1 Available Languages

- Arabic
- Azerbaijani
- Bengali
- Catalan
- Chinese
- Croatian
- Danish
- Dutch
- English
- Finnish
- French
- German
- Greek
- Hebrew
- Hindi
- Hungarian
- Indonesian
- Irish
- Italian

<sup>1</sup> Copyright (C) 2016 ExplosionAI UG (haftungsbeschränkt), 2016 spaCy GmbH, 2015 Matthew Honnibal

- Japanese
- Kazakh
- Nepali
- Norwegian
- Persian
- Polish
- Portuguese
- Romanian
- Russian
- Sinhala
- Spanish
- Swedish
- Tagalog
- Tamil
- Tatar
- Telugu
- Thai
- Turkish
- Ukrainian
- Urdu
- Vietnamese

```
>>> import advertools as adv
>>> sorted(adv.stopwords['english'])[:5]
["a", "about", "above", "across", "after"]
```

```
>>> sorted(adv.stopwords['german'])[:5]
["a", "ab", "aber", "ach", "acht"]
```

To get a list of all available languages, run

```
>>> adv.stopwords.keys()
dict_keys(['arabic', 'azerbaijani', 'bengali', 'catalan', 'chinese',
'croatian', 'danish', 'dutch', 'english', 'finnish', 'french',
'german', 'greek', 'hebrew', 'hindi', 'hungarian', 'indonesian',
'irish', 'italian', 'japanese', 'kazakh', 'nepali', 'norwegian',
'persian', 'polish', 'portuguese', 'romanian', 'russian', 'sinhala',
'spanish', 'swedish', 'tagalog', 'tamil', 'tatar', 'telugu', 'thai',
'turkish', 'ukrainian', 'urdu', 'vietnamese'])
```

## 1.15 Text Analysis

### 1.15.1 Absolute and Weighted Word Count

When analyzing a corpus of documents (I'll simply call it a text list), one of the main tasks to accomplish to start text mining is to first count the words. While there are many text mining techniques and approaches, the `word_frequency()` function works mainly by counting words in a text list. A “word” is defined as a sequence of characters split by whitespace(s), and stripped of non-word characters (commas, dots, quotation marks, etc.). A “word” is actually a phrase consisting of one word, but you have the option of getting phrases that have two words, or more. This can be done simply by providing a value for the `phrase_len` parameter.

#### Absolute vs Weighted Frequency

In social media reports, analytics, keyword reports, url and page reports, we get more information than simply the text. We get numbers describing those posts or page titles, or product names, or whatever the text list might contain. Numbers can be pageviews, shares, likes, retweets, sales, bounces, sales, etc. Since we have numbers to quantify those phrases, we can improve our counting by taking into consideration the number list that comes with the text list.

For example, if you have an e-commerce site that has two products, let's say you have bags and shoes, then your products are split 50:50 between bags and shoes. But what if you learn that shoes generate 80% of your sales? Although shoes form half your products, they generate 80% of your revenue. So the *weighted count* of your products is 80:20.

Let's say two people post two different posts on a social media platform. One of them says, “It's raining”, and the other says, “It's snowing”. As in the above example, the content is split 50:50 between “raining” and “snowing”, but we get a much more informative picture if we get the number of followers of each of those accounts (or the number of shares, likes, etc.). If one of them has a thousand followers, and other has a million (which is typical on social media, as well as in pageviews report, e-commerce and most other datasets), then you get a completely different picture about your dataset.

These two simple examples contain two posts, and a word each. The `word_frequency()` function can provide insight on hidden trends especially in large datasets, and when the sentences or phrases are also longer than a word or two each.

Let's take a look at how to use the `word_frequency()` function, and what the available parameters and options are.

**text\_list** The list of phrases or documents that you want to analyze. Here are some possible ideas that you might use this for:

- keywords, whether in a PPC or SEO report
- page titles in an analytics report
- social media posts (tweets, Facebook posts, YouTube video titles or descriptions etc.)
- e-commerce reports (where the text would be the product names)

**num\_list** Ideally, if you have more than one column describing `text_list` you should experiment with different options. Try weighting the words by pageviews, then try by bounce rate and see if you get different interesting findings. With e-commerce reports, you can see which word appears the most, and which word is associated with more revenue.

**phrase\_len** You should also experiment with different lengths of phrases. In many cases, one-word phrases might not be as meaningful as two-words or three.

**regex** The default is to simply split words by whitespace, and provide phrases of length `phrase_len`. But you may want to count the occurrences of certain patterns of text. Check out the `regex` module for the available regular

expressions that might be interesting. Some of the pre-defined ones are hashtags, mentions, questions, emoji, currencies, and more.

**rm\_words** A list of words to remove and ignore from the count. Known as stop-words these are the most frequently used words in a language, the most used, but don't add much meaning to the content (a, and, of, the, if, etc.). By default a set of English stopwords is provided (which you can check and possibly may want to modify), or run `adv.stopwords.keys()` to get a list of all the available stopwords in the available languages. In some cases (like page titles for example), you might get "words" that need to be removed as well, like the pipe "|" character for example.

**extra\_info** The returned DataFrame contains the default columns `[word, abs_freq, wtd_freq, rel_value]`. You can get extra columns for percentages and cumulative percentages that add perspective to the other columns. Set this parameter to `True` if you want that.

Below are all the columns of the returned DataFrame:

<code>word</code>	Words in the document list each on its own row. The length of these words is determined by <code>phrase_len</code> , essentially phrases if containing more than one word each.
<code>abs_freq</code>	The number of occurrences of each word in all the documents.
<code>wtd_freq</code>	Every occurrence of <code>word</code> multiplied by its respective value in <code>num_list</code> .
<code>rel_value</code>	<code>wtd_freq</code> divided by <code>abs_freq</code> , showing the value per occurrence of <code>word</code>
<code>abs_perc</code>	Absolute frequency percentage.
<code>abs_perc_cum</code>	Cumulative absolute percentage.
<code>wtd_freq_per</code>	Weighted frequency percentage.
<code>wtd_freq_per_cum</code>	Cumulative weighted frequency percentage.

**word\_frequency** (`text_list, num_list=None, phrase_len=1, regex=None, rm_words={'a', 'about', 'above', 'across', 'after', 'afterwards', 'again', 'against', 'all', 'almost', 'alone', 'along', 'already', 'also', 'although', 'always', 'am', 'among', 'amongst', 'amount', 'an', 'and', 'another', 'any', 'anyhow', 'anyone', 'anything', 'anyway', 'anywhere', 'are', 'around', 'as', 'at', 'back', 'be', 'became', 'because', 'become', 'becomes', 'becoming', 'been', 'before', 'beforehand', 'behind', 'being', 'below', 'beside', 'besides', 'between', 'beyond', 'both', 'bottom', 'but', 'by', 'ca', 'call', 'can', 'cannot', 'could', 'did', 'do', 'does', 'doing', 'done', 'down', 'due', 'during', 'each', 'eight', 'either', 'eleven', 'else', 'elsewhere', 'empty', 'enough', 'even', 'ever', 'every', 'everyone', 'everything', 'everywhere', 'except', 'few', 'fifteen', 'fifty', 'first', 'five', 'for', 'former', 'formerly', 'forty', 'four', 'from', 'front', 'full', 'further', 'get', 'give', 'go', 'had', 'has', 'have', 'he', 'hence', 'her', 'here', 'hereafter', 'hereby', 'herein', 'hereupon', 'hers', 'herself', 'him', 'himself', 'his', 'how', 'however', 'hundred', 'i', 'if', 'in', 'indeed', 'into', 'is', 'it', 'its', 'itself', 'just', 'keep', 'last', 'latter', 'latterly', 'least', 'less', 'made', 'make', 'many', 'may', 'me', 'meanwhile', 'might', 'mine', 'more', 'moreover', 'most', 'mostly', 'move', 'much', 'must', 'my', 'myself', 'name', 'namely', 'neither', 'never', 'nevertheless', 'next', 'nine', 'no', 'nobody', 'none', 'noone', 'nor', 'not', 'nothing', 'now', 'nowhere', 'of', 'off', 'often', 'on', 'once', 'one', 'only', 'onto', 'or', 'other', 'others', 'otherwise', 'our', 'ours', 'ourselves', 'out', 'over', 'own', 'part', 'per', 'perhaps', 'please', 'put', 'quite', 'rather', 're', 'really', 'regarding', 'same', 'say', 'see', 'seem', 'seemed', 'seeming', 'seems', 'serious', 'several', 'she', 'should', 'show', 'side', 'since', 'six', 'sixty', 'so', 'some', 'somehow', 'someone', 'something', 'sometime', 'sometimes', 'somewhere', 'still', 'such', 'take', 'ten', 'than', 'that', 'the', 'their', 'them', 'themselves', 'then', 'thence', 'there', 'thereafter', 'thereby', 'therefore', 'therein', 'thereupon', 'these', 'they', 'third', 'this', 'those', 'though', 'three', 'through', 'throughout', 'thru', 'thus', 'to', 'together', 'too', 'top', 'toward', 'towards', 'twelve', 'twenty', 'two', 'under', 'unless', 'until', 'up', 'upon', 'us', 'used', 'using', 'various', 'very', 'via', 'was', 'we', 'well', 'were', 'what', 'whatever', 'when', 'whence', 'whenever', 'where', 'whereafter', 'whereas', 'whereby', 'wherein', 'whereupon', 'wherever', 'whether', 'which', 'while', 'whither', 'who', 'whoever', 'whole', 'whom', 'whose', 'why', 'will', 'with', 'within', 'without', 'would', 'yet', 'you', 'your', 'yours', 'yourself', 'yourselves'}, extra_info=False)`

Count the absolute as well as the weighted frequency of words in `text_list` (based on `num_list`).

### Parameters

- **text\_list** (*list*) – Typically short phrases, but could be any list of full blown documents. Usually, you would use this to analyze tweets, book titles, URLs, etc.
- **num\_list** (*list*) – A list of numbers with the same length as `text_list`, describing a certain attribute of these ‘documents’; views, retweets, sales, etc.
- **regex** (*str*) – The regex used to split words. Doesn’t need changing in most cases.
- **phrase\_len** (*int*) – the length in words of each token the text is split into, defaults to 1.
- **rm\_words** (*set*) – Words to remove from the list a.k.a ‘stop-words’. The default uses. To get all available languages run `adv.stopwords.keys()`
- **extra\_info** (*bool*) – Whether or not to give additional metrics about the frequencies

Returns `abs_wtd_df` absolute and weighted DataFrame.

```
>>> text_list = ['apple orange', 'apple orange banana',
...             'apple kiwi', 'kiwi mango']
>>> num_list = [100, 100, 100, 400]
```

```
>>> adv.word_frequency(text_list, num_list)
   word  abs_freq  wtd_freq  rel_value
0  kiwi         2         500      250.0
1  mango         1         400      400.0
2  apple         3         300      100.0
3  orange         2         200      100.0
4  banana         1         100      100.0
```

Although “kiwi” occurred twice `abs_freq`, and “apple” occurred three times, the phrases in which “kiwi” appear have a total score of 500, so it beats “apple” on `wtd_freq` even though “apple” wins on `abs_freq`. You can sort by any of the columns of course. `rel_value` shows the value per occurrence of each word, as you can see, it is simply obtained by dividing `wtd_freq` by `abs_freq`.

```
>>> adv.word_frequency(text_list) # num_list values default to 1 each
   word  abs_freq  wtd_freq  rel_value
0  apple         3         3         1.0
1  orange         2         2         1.0
2  kiwi          2         2         1.0
3  banana         1         1         1.0
4  mango         1         1         1.0
```

```
>>> text_list2 = ['my favorite color is blue',
...              'my favorite color is green', 'the best color is green',
...              'i love the color black']
```

Setting `phrase_len` to 2, “words” become two-word phrases instead. Note that we are setting `rm_words` to the empty list so we can keep the stopwords and see if that makes sense:

```
>>> word_frequency(text_list2, phrase_len=2, rm_words=[])
   word  abs_freq  wtd_freq  rel_value
0  color is         3         3         1.0
1  my favorite         2         2         1.0
2  favorite color         2         2         1.0
3  is green          2         2         1.0
```

(continues on next page)

(continued from previous page)

4	is blue	1	1	1.0
5	the best	1	1	1.0
6	best color	1	1	1.0
7	i love	1	1	1.0
8	love the	1	1	1.0
9	the color	1	1	1.0
10	color black	1	1	1.0

The same result as above showing all possible columns by setting `extra_info` to `True`:

```
>>> adv.word_frequency(text_list, num_list, extra_info=True)
      word  abs_freq  abs_perc  abs_perc_cum  wtd_freq  wtd_freq_perc  wtd_freq_
      ←perc_cum  rel_value
0  kiwi          2  0.222222  0.222222  500      0.333333  0.
      ←333333  250.0
1  mango         1  0.111111  0.333333  400      0.266667  0.
      ←600000  400.0
2  apple         3  0.333333  0.666667  300      0.200000  0.
      ←800000  100.0
3  orange        2  0.222222  0.888889  200      0.133333  0.
      ←933333  100.0
4  banana        1  0.111111  1.000000  100      0.066667  1.
      ←000000  100.0
```

## 1.16 Tokenize Words (N-grams)

As word counting is an essential step in any text mining task, you first have to split the text into words.

The `word_tokenize()` function achieves that by splitting the text by whitespace. Another important thing it does after splitting is to trim the words of any non-word characters (commas, dots, exclamation marks, etc.).

You also have the option of specifying the length of the words that you want. The default is 2, which can be set through the `phrase_len` parameter.

This function is mainly a helper function for `word_frequency` to help with counting words and/or phrases.

**word\_tokenize** (*text\_list*, *phrase\_len*=2)

Split *text\_list* into phrases of length *phrase\_len* words each.

A “word” is any string between white spaces (or beginning or end of string) with delimiters stripped from both sides. Delimiters include quotes, question marks, parentheses, etc. Any delimiter contained within the string remains. See examples below.

### Parameters

- **text\_list** – List of strings.
- **phrase\_len** – Length of word tokens, defaults to 2.

**Return tokenized** List of lists, split according to `phrase_len`.

```
>>> t = ['split me into length-n-words',
... 'commas, (parentheses) get removed!',
... 'commas within text remain $1,000, but not the trailing commas.']
```

```
>>> word_tokenize(t, 1)
[['split', 'me', 'into', 'length-n-words'],
 ['commas', 'parentheses', 'get', 'removed'],
 ['commas', 'within', 'text', 'remain', '$1,000',
 'but', 'not', 'the', 'trailing', 'commas']]
```

The comma inside “\$1,000” as well as the dollar sign remain, as they are part of the “word”, but the trailing comma is stripped.

```
>>> word_tokenize(t, 2)
[['split me', 'me into', 'into length-n-words'],
 ['commas parentheses', 'parentheses get', 'get removed'],
 ['commas within', 'within text', 'text remain', 'remain $1,000',
 '$1,000 but', 'but not', 'not the', 'the trailing', 'trailing commas']]
```

```
>>> word_tokenize(t, 3)
[['split me into', 'me into length-n-words'],
 ['commas parentheses get', 'parentheses get removed'],
 ['commas within text', 'within text remain', 'text remain $1,000',
 'remain $1,000 but', '$1,000 but not', 'but not the',
 'not the trailing', 'the trailing commas']]
```

## 1.17 Twitter Data API

Easily connect to the Twitter API and start your analysis immediately.

Main Features:

**1 Get the results in a DataFrame:** With the exception of three functions that return a list of ID’s, everything else returns a pandas DataFrame, ready to use. This allows you to spend more time analyzing data, and less time figuring out the structure of the JSON response object. It’s not complicated or anything, just takes time.

**2 Manage looping and merging:** there is a limit on how many results you get per request (typically in the 100 - 200 range), several requests have to be made, and merged together. Not all responses have the same structure, so this is also handled. You only have to provide the number of responses you want through the `count` parameter where applicable (provided you are within your app’s rate limits).

**3 Unnesting nested objects:** Many response objects contain very rich embedded data, which is usually meta data about the response. For example, when you request tweets, you get a user object along with that. This is very helpful in better understanding who made the tweet, and how influential/credible they are.

**4 Documentation:** All available parameters are included in the function signatures, to make it easier to explore interactively, as well as descriptions of the parameters imported from the Twitter documentation.

### 1.17.1 Authentication

Before starting you will have to create an app through [developer.twitter.com](https://developer.twitter.com), and then you can get your authentication keys from your dashboard. Then you can authenticate as follows:

```
>>> auth_params = {
...     'app_key': 'YOUR_APP_KEY',
...     'app_secret': 'YOUR_APP_SECRET',
...     'oauth_token': 'YOUR_OAUTH_TOKEN',
...     'oauth_token_secret': 'YOUR_OAUTH_TOKEN_SECRET',
```

(continues on next page)

(continued from previous page)

```
... }
>>> import advertools as adv
>>> adv.twitter.set_auth_params(**auth_params)
```

Now every request you send will include your `auth_params` in it, and if valid you will get the respective response, for example:

```
>>> python_tweets = adv.twitter.search(q='#python', tweet_mode='extended')
```

Make sure you always specify `tweet_mode='extended'` because otherwise you will get tweets that are 140 characters long.

When you have tweets and user data in the DataFrame, the column names would be prepended with `tweet_` and `user_` to make it clear what the data belong to.

## 1.17.2 Functions

### `authenticate` (*func*)

Used internally, please use `set_auth_params` for authentication.

### `get_application_rate_limit_status` (*consumed\_only=True*)

Returns the current rate limits for methods belonging to the specified resource families.

**Parameters** `consumed_only` – Whether or not to return only items that have been consumed. Otherwise returns the full list.

<https://developer.twitter.com/en/docs/developer-utilities/rate-limit-status/api-reference/get-application-rate-limit-status>

### `get_available_trends` ()

Returns the locations that Twitter has trending topic information for.

<https://developer.twitter.com/en/docs/trends/locations-with-trending-topics/api-reference/get-trends-available>

### `get_favorites` (*user\_id=None, screen\_name=None, count=None, since\_id=None, max\_id=None, include\_entities=None, tweet\_mode=None*)

Returns the 20 most recent Tweets favorited by the authenticating or specified user.

#### Parameters

- **user\_id** – (int - optional) The ID of the user for whom to return results.
- **screen\_name** – (str - optional) The screen name of the user for whom to return results.
- **count** – (int - optional) Specifies the number of results to retrieve.
- **since\_id** – (int - optional) Returns results with an ID greater than (that is, more recent than) the specified ID. There are limits to the number of Tweets which can be accessed through the API. If the limit of Tweets has occurred since the `since_id`, the `since_id` will be forced to the oldest ID available.
- **max\_id** – (int - optional) Returns results with an ID less than (that is, older than) or equal to the specified ID.
- **include\_entities** – (bool - optional) The entities node will be omitted when set to False .



- **tweet\_mode** – (str - optional) Valid request values are compat and extended, which give compatibility mode and extended mode, respectively for Tweets that contain over 140 characters

<https://developer.twitter.com/en/docs/tweets/post-and-engage/api-reference/get-favorites-list>

**get\_followers\_ids** (*user\_id=None, screen\_name=None, cursor=None, stringify\_ids=None, count=None*)

Returns a cursored collection of user IDs for every user following the specified user.

#### Parameters

- **user\_id** – (int - optional) The ID of the user for whom to return results.
- **screen\_name** – (str - optional) The screen name of the user for whom to return results.
- **cursor** – (cursor - semi-optional) Causes the list of connections to be broken into pages of no more than 5000 IDs at a time. The number of IDs returned is not guaranteed to be 5000 as suspended users are filtered out after connections are queried. If no cursor is provided, a value of -1 will be assumed, which is the first “page.” The response from the API will include a `previous_cursor` and `next_cursor` to allow paging back and forth. See Using cursors to navigate collections for more information.
- **stringify\_ids** – (bool - optional) Some programming environments will not consume Twitter IDs due to their size. Provide this option to have IDs returned as strings instead. More about Twitter IDs.
- **count** – (int - optional) Specifies the number of results to retrieve.

<https://developer.twitter.com/en/docs/accounts-and-users/follow-search-get-users/api-reference/get-followers-ids>

**get\_followers\_list** (*user\_id=None, screen\_name=None, cursor=None, count=None, skip\_status=None, include\_user\_entities=None*)

Returns a cursored collection of user objects for users following the specified user.

#### Parameters

- **user\_id** – (int - optional) The ID of the user for whom to return results.
- **screen\_name** – (str - optional) The screen name of the user for whom to return results.
- **cursor** – (cursor - semi-optional) Causes the results to be broken into pages. If no cursor is provided, a value of -1 will be assumed, which is the first “page.” The response from the API will include a `previous_cursor` and `next_cursor` to allow paging back and forth. See Using cursors to navigate collections for more information.
- **count** – (int - optional) Specifies the number of results to retrieve.
- **skip\_status** – (bool - optional) When set to True, statuses will not be included in the returned user objects. If set to any other value, statuses will be included.
- **include\_user\_entities** – (bool - optional) The user object entities node will not be included when set to False.

<https://developer.twitter.com/en/docs/accounts-and-users/follow-search-get-users/api-reference/get-followers-list>

**get\_friends\_ids** (*user\_id=None, screen\_name=None, cursor=None, stringify\_ids=None, count=None*)

Returns a cursored collection of user IDs for every user the specified user is following (otherwise known as their “friends”).

#### Parameters

- **user\_id** – (int - optional) The ID of the user for whom to return results.
- **screen\_name** – (str - optional) The screen name of the user for whom to return results.
- **cursor** – (cursor - semi-optional) Causes the list of connections to be broken into pages of no more than 5000 IDs at a time. The number of IDs returned is not guaranteed to be 5000 as suspended users are filtered out after connections are queried. If no cursor is provided, a value of -1 will be assumed, which is the first “page.” The response from the API will include a `previous_cursor` and `next_cursor` to allow paging back and forth. See Using cursors to navigate collections for more information.
- **stringify\_ids** – (bool - optional) Some programming environments will not consume Twitter IDs due to their size. Provide this option to have IDs returned as strings instead. More about Twitter IDs.
- **count** – (int - optional) Specifies the number of results to retrieve.

<https://developer.twitter.com/en/docs/accounts-and-users/follow-search-get-users/api-reference/get-friends-ids>

**get\_friends\_list** (*user\_id=None, screen\_name=None, cursor=None, count=None, skip\_status=None, include\_user\_entities=None*)

Returns a cursored collection of user objects for every user the specified user is following (otherwise known as their “friends”).

#### Parameters

- **user\_id** – (int - optional) The ID of the user for whom to return results.
- **screen\_name** – (str - optional) The screen name of the user for whom to return results.
- **cursor** – (cursor - semi-optional) Causes the results to be broken into pages. If no cursor is provided, a value of -1 will be assumed, which is the first “page.” The response from the API will include a `previous_cursor` and `next_cursor` to allow paging back and forth. See Using cursors to navigate collections for more information.
- **count** – (int - optional) Specifies the number of results to retrieve.
- **skip\_status** – (bool - optional) When set to True statuses will not be included in the returned user objects.
- **include\_user\_entities** – (bool - optional) The user object entities node will not be included when set to False.

<https://developer.twitter.com/en/docs/accounts-and-users/follow-search-get-users/api-reference/get-friends-list>

**get\_home\_timeline** (*count=None, since\_id=None, max\_id=None, trim\_user=None, exclude\_replies=None, include\_entities=None, tweet\_mode=None*)

Returns a collection of the most recent Tweets and retweets posted by the authenticating user and the users they follow.

#### Parameters

- **count** – (int - optional) Specifies the number of results to retrieve.

- **since\_id** – (int - optional) Returns results with an ID greater than (that is, more recent than) the specified ID. There are limits to the number of Tweets which can be accessed through the API. If the limit of Tweets has occurred since the `since_id`, the `since_id` will be forced to the oldest ID available.
- **max\_id** – (int - optional) Returns results with an ID less than (that is, older than) or equal to the specified ID.
- **trim\_user** – (bool - optional) When set to True, each Tweet returned in a timeline will include a user object including only the status authors numerical ID. Omit this parameter to receive the complete user object.
- **exclude\_replies** – (bool - optional) This parameter will prevent replies from appearing in the returned timeline. Using `exclude_replies` with the `count` parameter will mean you will receive up-to count Tweets — this is because the `count` parameter retrieves that many Tweets before filtering out retweets and replies.
- **include\_entities** – (bool - optional) The entities node will not be included when set to False.
- **tweet\_mode** – (str - optional) Valid request values are `compat` and `extended`, which give compatibility mode and extended mode, respectively for Tweets that contain over 140 characters

[https://developer.twitter.com/en/docs/tweets/timelines/api-reference/get-statuses-home\\_timeline](https://developer.twitter.com/en/docs/tweets/timelines/api-reference/get-statuses-home_timeline)

**get\_list\_members** (*list\_id=None, slug=None, owner\_screen\_name=None, owner\_id=None, count=None, cursor=None, include\_entities=None, skip\_status=None*)

Returns the members of the specified list.

#### Parameters

- **list\_id** – (str - required) The numerical id of the list.
- **slug** – (str - required) You can identify a list by its slug instead of its numerical id. If you decide to do so, note that you'll also have to specify the list owner using the `owner_id` or `owner_screen_name` parameters.
- **owner\_screen\_name** – (str - optional) The screen name of the user who owns the list being requested by a slug.
- **owner\_id** – (int - optional) The user ID of the user who owns the list being requested by a slug.
- **count** – (int - optional) Specifies the number of results to retrieve.
- **cursor** – (cursor - semi-optional) Causes the collection of list members to be broken into “pages” of consistent sizes (specified by the `count` parameter). If no cursor is provided, a value of -1 will be assumed, which is the first “page.” The response from the API will include a `previous_cursor` and `next_cursor` to allow paging back and forth. See [Using cursors to navigate collections](#) for more information.
- **include\_entities** – (bool - optional) The entities node will not be included when set to False.
- **skip\_status** – (bool - optional) When set to True statuses will not be included in the returned user objects.

<https://developer.twitter.com/en/docs/accounts-and-users/create-manage-lists/api-reference/get-lists-members>

**get\_list\_memberships** (*user\_id=None, screen\_name=None, count=None, cursor=None, filter\_to\_owned\_lists=None*)

Returns the lists the specified user has been added to.

### Parameters

- **user\_id** – (int - optional) The ID of the user for whom to return results. Helpful for disambiguating when a valid user ID is also a valid screen name.
- **screen\_name** – (str - optional) The screen name of the user for whom to return results. Helpful for disambiguating when a valid screen name is also a user ID.
- **count** – (int - optional) Specifies the number of results to retrieve.
- **cursor** – (cursor - optional) Breaks the results into pages. Provide a value of -1 to begin paging. Provide values as returned in the response body's `next_cursor` and `previous_cursor` attributes to page back and forth in the list. It is recommended to always use cursors when the method supports them. See [Cursoring](#) for more information.
- **filter\_to\_owned\_lists** – (bool - optional) When True, will return just lists the authenticating user owns, and the user represented by `user_id` or `screen_name` is a member of.

<https://developer.twitter.com/en/docs/accounts-and-users/create-manage-lists/api-reference/get-lists-memberships>

**get\_list\_statuses** (*list\_id=None, slug=None, owner\_screen\_name=None, owner\_id=None, since\_id=None, max\_id=None, count=None, include\_entities=None, include\_rts=None, tweet\_mode=None*)

Returns a timeline of tweets authored by members of the specified list.

### Parameters

- **list\_id** – (str - required) The numerical id of the list.
- **slug** – (str - required) You can identify a list by its slug instead of its numerical id. If you decide to do so, note that you'll also have to specify the list owner using the `owner_id` or `owner_screen_name` parameters.
- **owner\_screen\_name** – (str - optional) The screen name of the user who owns the list being requested by a slug .
- **owner\_id** – (int - optional) The user ID of the user who owns the list being requested by a slug .
- **since\_id** – (int - optional) Returns results with an ID greater than (that is, more recent than) the specified ID. There are limits to the number of Tweets which can be accessed through the API. If the limit of Tweets has occurred since the `since_id`, the `since_id` will be forced to the oldest ID available.
- **max\_id** – (int - optional) Returns results with an ID less than (that is, older than) or equal to the specified ID.
- **count** – (int - optional) Specifies the number of results to retrieve.
- **include\_entities** – (bool - optional) Entities are ON by default in API 1.1, each tweet includes a node called "entities". This node offers a variety of metadata about the tweet in a discreet structure, including: `user_mentions`, `urls`, and `hashtags`. You can omit entities from the result by using `include_entities=False`
- **include\_rts** – (bool - optional) When set to True, the list timeline will contain native retweets (if they exist) in addition to the standard stream of tweets. The output format of retweeted tweets is identical to the representation you see in `home_timeline`.
- **tweet\_mode** – (str - optional) Valid request values are `compat` and `extended`, which give compatibility mode and extended mode, respectively for Tweets that contain over 140 characters

<https://developer.twitter.com/en/docs/accounts-and-users/create-manage-lists/api-reference/get-lists-statuses>

**get\_list\_subscribers** (*list\_id=None, slug=None, owner\_screen\_name=None, owner\_id=None, count=None, cursor=None, include\_entities=None, skip\_status=None*)

Returns the subscribers of the specified list.

#### Parameters

- **list\_id** – (str - required) The numerical id of the list.
- **slug** – (str - required) You can identify a list by its slug instead of its numerical id. If you decide to do so, note that you'll also have to specify the list owner using the `owner_id` or `owner_screen_name` parameters.
- **owner\_screen\_name** – (str - optional) The screen name of the user who owns the list being requested by a slug .
- **owner\_id** – (int - optional) The user ID of the user who owns the list being requested by a slug .
- **count** – (int - optional) Specifies the number of results to retrieve.
- **cursor** – (cursor - optional) Breaks the results into pages. A single page contains 20 lists. Provide a value of -1 to begin paging. Provide values as returned in the response body's `next_cursor` and `previous_cursor` attributes to page back and forth in the list. See Using cursors to navigate collections for more information.
- **include\_entities** – (bool - optional) When set to True, each tweet will include a node called "entities". This node offers a variety of metadata about the tweet in a discreet structure, including: `user_mentions`, `urls`, and `hashtags`. While entities are opt-in on timelines at present, they will be made a default component of output in the future. See Tweet Entities for more details.
- **skip\_status** – (bool - optional) When set to True statuses will not be included in the returned user objects.

<https://developer.twitter.com/en/docs/accounts-and-users/create-manage-lists/api-reference/get-lists-subscribers>

**get\_list\_subscriptions** (*user\_id=None, screen\_name=None, count=None, cursor=None*)

Obtain a collection of the lists the specified user is subscribed to.

#### Parameters

- **user\_id** – (int - optional) The ID of the user for whom to return results. Helpful for disambiguating when a valid user ID is also a valid screen name.
- **screen\_name** – (str - optional) The screen name of the user for whom to return results. Helpful for disambiguating when a valid screen name is also a user ID.
- **count** – (int - optional) Specifies the number of results to retrieve.
- **cursor** – (cursor - optional) Breaks the results into pages. Provide a value of -1 to begin paging. Provide values as returned in the response body's `next_cursor` and `previous_cursor` attributes to page back and forth in the list. It is recommended to always use cursors when the method supports them. See Cursoring for more information.

<https://developer.twitter.com/en/docs/accounts-and-users/create-manage-lists/api-reference/get-lists-subscriptions>

**get\_mentions\_timeline** (*count=None, since\_id=None, max\_id=None, trim\_user=None, include\_entities=None, tweet\_mode=None*)

Returns the 20 most recent mentions (tweets containing a users's @screen\_name) for the authenticating user.

#### Parameters

- **count** – (int - optional) Specifies the number of results to retrieve.
- **since\_id** – (int - optional) Returns results with an ID greater than (that is, more recent than) the specified ID. There are limits to the number of Tweets which can be accessed through the API. If the limit of Tweets has occurred since the since\_id, the since\_id will be forced to the oldest ID available.
- **max\_id** – (int - optional) Returns results with an ID less than (that is, older than) or equal to the specified ID.
- **trim\_user** – (bool - optional) When set to True, each tweet returned in a timeline will include a user object including only the status authors numerical ID. Omit this parameter to receive the complete user object.
- **include\_entities** – (bool - optional) The entities node will not be included when set to False.
- **tweet\_mode** – (str - optional) Valid request values are compat and extended, which give compatibility mode and extended mode, respectively for Tweets that contain over 140 characters

[https://developer.twitter.com/en/docs/tweets/timelines/api-reference/get-statuses-mentions\\_timeline](https://developer.twitter.com/en/docs/tweets/timelines/api-reference/get-statuses-mentions_timeline)

**get\_place\_trends** (*ids*, *exclude=None*)

Returns the top 10 trending topics for a specific WOEID, if trending information is available for it.

#### Parameters

- **id** – (int or list of ints - required) run `get_available_trends()` for the full listing. The Yahoo! Where On Earth ID of the location to return trending information for. Global information is available by using 1 as the WOEID .
- **exclude** – (str - optional) Setting this equal to hashtags will remove all hashtags from the trends list.

<https://developer.twitter.com/en/docs/trends/trends-for-location/api-reference/get-trends-place>

**get\_retweeters\_ids** (*id*, *count=None*, *cursor=None*, *stringify\_ids=None*)

Returns a collection of up to 100 user IDs belonging to users who have retweeted the tweet specified by the `id` parameter. It's better to use `get_retweets` because passing a `count > 100` will only get you duplicated data. 100 is the maximum even if there were more retweeters.

#### Parameters

- **id** – (int - required) The numerical ID of the desired status.
- **count** – (int - optional) Specifies the number of results to retrieve.
- **cursor** – (cursor - semi-optional) Causes the list of IDs to be broken into pages of no more than 100 IDs at a time. The number of IDs returned is not guaranteed to be 100 as suspended users are filtered out after connections are queried. If no cursor is provided, a value of -1 will be assumed, which is the first "page." The response from the API will include a `previous_cursor` and `next_cursor` to allow paging back and forth. See our cursor docs for more information. While this method supports the cursor parameter, the entire

result set can be returned in a single cursored collection. Using the count parameter with this method will not provide segmented cursors for use with this parameter.

- **stringify\_ids** – (bool - optional) Many programming environments will not consume Tweet ids due to their size. Provide this option to have ids returned as strings instead.

<https://developer.twitter.com/en/docs/tweets/post-and-engage/api-reference/get-statuses-retweeters-ids>

**get\_retweets** (*id*, *trim\_user=None*, *tweet\_mode=None*)

Returns up to 100 of the first retweets of a given tweet.

#### Parameters

- **id** – (int - required) The numerical ID of the desired status.
- **trim\_user** – (bool - optional) When set to True, each tweet returned in a timeline will include a user object including only the status authors numerical ID. Omit this parameter to receive the complete user object.
- **tweet\_mode** – (str - optional) Valid request values are compat and extended, which give compatibility mode and extended mode, respectively for Tweets that contain over 140 characters

<https://developer.twitter.com/en/docs/tweets/post-and-engage/api-reference/post-statuses-retweet-id>

**get\_supported\_languages** ()

Returns the list of languages supported by Twitter along with their ISO 639-1 code.

<https://developer.twitter.com/en/docs/developer-utilities/supported-languages/api-reference/get-help-languages>

**get\_user\_timeline** (*user\_id=None*, *screen\_name=None*, *since\_id=None*, *count=None*, *max\_id=None*, *trim\_user=None*, *exclude\_replies=None*, *include\_rts=None*, *tweet\_mode=None*)

Returns a collection of the most recent Tweets posted by the user indicated by the `screen_name` or `user_id` parameters.

#### Parameters

- **user\_id** – (int - optional) The ID of the user for whom to return results.
- **screen\_name** – (str - optional) The screen name of the user for whom to return results.
- **since\_id** – (int - optional) Returns results with an ID greater than (that is, more recent than) the specified ID. There are limits to the number of Tweets that can be accessed through the API. If the limit of Tweets has occurred since the `since_id`, the `since_id` will be forced to the oldest ID available.
- **count** – (int - optional) Specifies the number of results to retrieve.
- **max\_id** – (int - optional) Returns results with an ID less than (that is, older than) or equal to the specified ID.
- **trim\_user** – (bool - optional) When set to True, each Tweet returned in a timeline will include a user object including only the status authors numerical ID. Omit this parameter to receive the complete user object.
- **exclude\_replies** – (bool - optional) This parameter will prevent replies from appearing in the returned timeline. Using `exclude_replies` with the count parameter will mean you will receive up-to count tweets — this is because the count parameter retrieves that many Tweets before filtering out retweets and replies.

- **include\_rts** – (bool - optional) When set to False , the timeline will strip any native retweets (though they will still count toward both the maximal length of the timeline and the slice selected by the count parameter). Note: If you're using the trim\_user parameter in conjunction with include\_rts, the retweets will still contain a full user object.
- **tweet\_mode** – (str - optional) Valid request values are compat and extended, which give compatibility mode and extended mode, respectively for Tweets that contain over 140 characters

[https://developer.twitter.com/en/docs/tweets/timelines/api-reference/get-statuses-user\\_timeline](https://developer.twitter.com/en/docs/tweets/timelines/api-reference/get-statuses-user_timeline)

**lookup\_status** (*id, include\_entities=None, trim\_user=None, map=None, include\_ext\_alt\_text=None, include\_card\_uri=None, tweet\_mode=None*)

Returns fully-hydrated tweet objects for up to 100 tweets per request, as specified by comma-separated values passed to the `id` parameter.

#### Parameters

- **id** – (int - required) A comma separated list of Tweet IDs, up to 100 are allowed in a single request.
- **include\_entities** – (bool - optional) The entities node that may appear within embedded statuses will not be included when set to False.
- **trim\_user** – (bool - optional) When set to True, each Tweet returned in a timeline will include a user object including only the status authors numerical ID. Omit this parameter to receive the complete user object.
- **map** – (bool - optional) When using the map parameter, Tweets that do not exist or cannot be viewed by the current user will still have their key represented but with an explicitly null value paired with it
- **include\_ext\_alt\_text** – (bool - optional) If alt text has been added to any attached media entities, this parameter will return an `ext_alt_text` value in the top-level key for the media entity. If no value has been set, this will be returned as null
- **include\_card\_uri** – (bool - optional) When set to True, each Tweet returned will include a `card_uri` attribute when there is an ads card attached to the Tweet and when that card was attached using the `card_uri` value.
- **tweet\_mode** – (str - optional) Valid request values are compat and extended, which give compatibility mode and extended mode, respectively for Tweets that contain over 140 characters

<https://developer.twitter.com/en/docs/tweets/post-and-engage/api-reference/get-statuses-lookup>

**lookup\_user** (*screen\_name=None, user\_id=None, include\_entities=None, tweet\_mode=None*)

Returns fully-hydrated user objects for up to 100 users per request, as specified by comma-separated values passed to the `user_id` and/or `screen_name` parameters.

#### Parameters

- **screen\_name** – (str - optional) A comma separated list of screen names, up to 100 are allowed in a single request. You are strongly encouraged to use a POST for larger (up to 100 screen names) requests.
- **user\_id** – (int - optional) A comma separated list of user IDs, up to 100 are allowed in a single request. You are strongly encouraged to use a POST for larger requests.



- **include\_entities** – (bool - optional) The entities node that may appear within embedded statuses will not be included when set to False.
- **tweet\_mode** – (str - optional) Valid request values are compat and extended, which give compatibility mode and extended mode, respectively for Tweets that contain over 140 characters

<https://developer.twitter.com/en/docs/accounts-and-users/follow-search-get-users/api-reference/get-users-lookup>

**make\_dataframe** (*func*)

**retweeted\_of\_me** (*count=None, since\_id=None, max\_id=None, trim\_user=None, include\_entities=None, include\_user\_entities=None, tweet\_mode=None*)

Returns the most recent tweets authored by the authenticating user that have been retweeted by others.

#### Parameters

- **count** – (int - optional) Specifies the number of results to retrieve.
- **since\_id** – (int - optional) Returns results with an ID greater than (that is, more recent than) the specified ID. There are limits to the number of Tweets which can be accessed through the API. If the limit of Tweets has occurred since the `since_id`, the `since_id` will be forced to the oldest ID available.
- **max\_id** – (int - optional) Returns results with an ID less than (that is, older than) or equal to the specified ID.
- **trim\_user** – (bool - optional) When set to True, each tweet returned in a timeline will include a user object including only the status authors numerical ID. Omit this parameter to receive the complete user object.
- **include\_entities** – (bool - optional) The tweet entities node will not be included when set to False .
- **include\_user\_entities** – (bool - optional) The user entities node will not be included when set to False .
- **tweet\_mode** – (str - optional) Valid request values are compat and extended, which give compatibility mode and extended mode, respectively for Tweets that contain over 140 characters

[https://developer.twitter.com/en/docs/tweets/post-and-engage/api-reference/get-statuses-retweets\\_of\\_me](https://developer.twitter.com/en/docs/tweets/post-and-engage/api-reference/get-statuses-retweets_of_me)

**search** (*q, geocode=None, lang=None, locale=None, result\_type=None, count=None, until=None, since\_id=None, max\_id=None, include\_entities=None, tweet\_mode=None*)

Returns a collection of relevant Tweets matching a specified query.

#### Parameters

- **q** – (str - required) A UTF-8, URL-encoded search query of 500 characters maximum, including operators. Queries may additionally be limited by complexity.
- **geocode** – (lat lon dist - optional) Returns tweets by users located within a given radius of the given latitude/longitude. The location is preferentially taking from the Geotagging API, but will fall back to their Twitter profile. The parameter value is specified by " latitude,longitude,radius ", where radius units must be specified as either " mi " (miles) or " km " (kilometers). Note that you cannot use the near operator via the API to geocode arbitrary locations; however you can use this geocode parameter to search near geocodes directly. A maximum of 1,000 distinct "sub- regions" will be considered when using the radius modifier.

- **lang** – (str - optional) Restricts tweets to the given language, given by an ISO 639-1 code. Language detection is best-effort.
- **locale** – (str - optional) Specify the language of the query you are sending (only ja is currently effective). This is intended for language- specific consumers and the default should work in the majority of cases.
- **result\_type** – (str - optional) Optional. Specifies what type of search results you would prefer to receive. The current default is “mixed.” Valid values include: \* mixed : Include both popular and real time results in the response. \* recent : return only the most recent results in the response \* popular : return only the most popular results in the response.
- **count** – (int - optional) Specifies the number of results to retrieve.
- **until** – (date - optional) Returns tweets created before the given date. Date should be formatted as YYYY-MM-DD. Keep in mind that the search index has a 7-day limit. In other words, no tweets will be found for a date older than one week.
- **since\_id** – (int - optional) Returns results with an ID greater than (that is, more recent than) the specified ID. There are limits to the number of Tweets which can be accessed through the API. If the limit of Tweets has occurred since the since\_id, the since\_id will be forced to the oldest ID available.
- **max\_id** – (int - optional) Returns results with an ID less than (that is, older than) or equal to the specified ID.
- **include\_entities** – (bool - optional) The entities node will not be included when set to False.
- **tweet\_mode** – (str - optional) Valid request values are compat and extended, which give compatibility mode and extended mode, respectively for Tweets that contain over 140 characters

Operator	Finds Tweets...
watching now	containing both “watching” and “now”. This is the default operator.
“happy hour”	containing the exact phrase “happy hour”.
love OR hate	containing either “love” or “hate” (or both).
beer -root	containing “beer” but not “root”.
#haiku	containing the hashtag “haiku”.
from:interior	sent from Twitter account “interior”.
list:NASA/astronauts-in-space-now	sent from a Twitter account in the NASA list astronauts-in-space-now
to:NASA	a Tweet authored in reply to Twitter account “NASA”.
@NASA	mentioning Twitter account “NASA”.
politics filter:safe	containing “politics” with Tweets marked as potentially sensitive removed.
puppy filter:media	containing “puppy” and an image or video.
puppy -filter:retweets	containing “puppy”, filtering out retweets
puppy filter:native_video	containing “puppy” and an uploaded video, Amplify video, Periscope, or Vine.
puppy filter:periscope	containing “puppy” and a Periscope video URL.
puppy filter:vine	containing “puppy” and a Vine.
puppy filter:images	containing “puppy” and links identified as photos, including third parties such as Instagram.
puppy filter:twimg	containing “puppy” and a pic.twitter.com link representing one or more photos.
hilarious filter:links	containing “hilarious” and linking to URL.
puppy url:amazon	containing “puppy” and a URL with the word “amazon” anywhere within it.
superhero since:2015-12-21	containing “superhero” and sent since date “2015-12-21” (year-month-day).
puppy until:2015-12-21	containing “puppy” and sent before the date “2015-12-21”.
movie -scary :)	containing “movie”, but not “scary”, and with a positive attitude.
flight :(	containing “flight” and with a negative attitude.
traffic ?	containing “traffic” and asking a question.

<https://developer.twitter.com/en/docs/tweets/search/api-reference/get-search-tweets>

**search\_users** (*q*, *page=None*, *count=None*, *include\_entities=None*)

Provides a simple, relevance-based search interface to public user accounts on Twitter.

#### Parameters

- **q** – (str - required) The search query to run against people search.
- **page** – (int - optional) Specifies the page of results to retrieve.
- **count** – (int - optional) Specifies the number of results to retrieve.
- **include\_entities** – (bool - optional) The entities node will not be included in embedded Tweet objects when set to False .

<https://developer.twitter.com/en/docs/accounts-and-users/follow-search-get-users/api-reference/get-users-search>

**set\_auth\_params** (*app\_key=None*, *app\_secret=None*, *oauth\_token=None*, *oauth\_token\_secret=None*, *access\_token=None*, *token\_type='bearer'*, *oauth\_version=1*, *api\_version='1.1'*, *client\_args=None*, *auth\_endpoint='authenticate'*)

The main function for authentication. Needs to be called once in a session.

First you need to create a developer account and app: <https://developer.twitter.com/> to get your credentials.

Different ways to authenticate: [https://twython.readthedocs.io/en/latest/usage/starting\\_out.html](https://twython.readthedocs.io/en/latest/usage/starting_out.html)

**show\_lists** (*user\_id=None, screen\_name=None, reverse=None*)

Returns all lists the authenticating or specified user subscribes to, including their own.

#### Parameters

- **user\_id** – (int - optional) The ID of the user for whom to return results. Helpful for disambiguating when a valid user ID is also a valid screen name. Note: : Specifies the ID of the user to get lists from. Helpful for disambiguating when a valid user ID is also a valid screen name.
- **screen\_name** – (str - optional) The screen name of the user for whom to return results. Helpful for disambiguating when a valid screen name is also a user ID.
- **reverse** – (bool - optional) Set this to true if you would like owned lists to be returned first. See description above for information on how this parameter works.

<https://developer.twitter.com/en/docs/accounts-and-users/create-manage-lists/api-reference/get-lists-list>

**show\_owned\_lists** (*user\_id=None, screen\_name=None, count=None, cursor=None*)

Returns the lists owned by the specified Twitter user.

#### Parameters

- **user\_id** – (int - optional) The ID of the user for whom to return results. Helpful for disambiguating when a valid user ID is also a valid screen name.
- **screen\_name** – (str - optional) The screen name of the user for whom to return results. Helpful for disambiguating when a valid screen name is also a user ID.
- **count** – (int - optional) Specifies the number of results to retrieve.
- **cursor** – (cursor - optional) Breaks the results into pages. Provide a value of -1 to begin paging. Provide values as returned in the response body's `next_cursor` and `previous_cursor` attributes to page back and forth in the list. It is recommended to always use cursors when the method supports them. See [Cursoring](#) for more information.

<https://developer.twitter.com/en/docs/accounts-and-users/create-manage-lists/api-reference/get-lists-ownerships>

## 1.18 YouTube Data API

**activities\_list** (*key, part, channelId=None, home=None, mine=None, maxResults=None, pageToken=None, publishedAfter=None, publishedBefore=None, regionCode=None*)

Returns a list of channel activity events that match the request criteria. For example, you can retrieve events associated with a particular channel or with the user's own channel.

*Required parameters:*

#### Parameters

- **key** – string Your Google API key.
- **part** – string The part parameter specifies a comma-separated list of one or more activity resource properties that the API response will include. If the parameter identifies a property that contains child properties, the child properties will be included in the response. For example, in an activity resource, the `snippet` property contains other properties that identify the type of activity, a display title for the activity, and so forth. If you set `part=snippet`, the

API response will also contain all of those nested properties. The following list contains the part names that you can include in the parameter value and the quota cost for each part: contentDetails: 2 id: 0 snippet: 2

*Filters (specify exactly one of the following parameters):*

#### Parameters

- **channelId** – string The channelId parameter specifies a unique YouTube channel ID. The API will then return a list of that channel’s activities.
- **home** – boolean Note: This parameter has been deprecated. For requests that set this parameter, the API response contains items similar to those that a logged-out user would see on the YouTube home page. Note that this parameter can only be used in a properly authorized request.
- **mine** – boolean This parameter can only be used in a properly authorized request. Set this parameter’s value to true to retrieve a feed of the authenticated user’s activities.

*Optional parameters:*

#### Parameters

- **maxResults** – unsigned integer The maxResults parameter specifies the maximum number of items that should be returned in the result set.
- **pageToken** – string The pageToken parameter identifies a specific page in the result set that should be returned. In an API response, the nextPageToken and prevPageToken properties identify other pages that could be retrieved.
- **publishedAfter** – datetime The publishedAfter parameter specifies the earliest date and time that an activity could have occurred for that activity to be included in the API response. If the parameter value specifies a day, but not a time, then any activities that occurred that day will be included in the result set. The value is specified in ISO 8601 (YYYY-MM-DDThh:mm:ss.SZ) format.
- **publishedBefore** – datetime The publishedBefore parameter specifies the date and time before which an activity must have occurred for that activity to be included in the API response. If the parameter value specifies a day, but not a time, then any activities that occurred that day will be excluded from the result set. The value is specified in ISO 8601 (YYYY-MM-DDThh:mm:ss.SZ) format.
- **regionCode** – string The regionCode parameter instructs the API to return results for the specified country. The parameter value is an ISO 3166-1 alpha-2 country code. YouTube uses this value when the authorized user’s previous activity on YouTube does not provide enough information to generate the activity feed.

**captions\_list** (*key, part, videoId, id=None, onBehalfOfContentOwner=None*)

Returns a list of caption tracks that are associated with a specified video. Note that the API response does not contain the actual captions and that the captions.download method provides the ability to retrieve a caption track.

*Required parameters:*

#### Parameters

- **key** – string Your Google API key.
- **part** – string The part parameter specifies the caption resource parts that the API response will include. The list below contains the part names that you can include in the parameter value and the quota cost for each part: id: 0 snippet: 1
- **videoId** – string The videoId parameter specifies the YouTube video ID of the video for which the API should return caption tracks.

*Optional parameters:*

#### Parameters

- **id** – string The id parameter specifies a comma-separated list of IDs that identify the caption resources that should be retrieved. Each ID must identify a caption track associated with the specified video.
- **onBehalfOfContentOwner** – string This parameter can only be used in a properly authorized request. Note: This parameter is intended exclusively for YouTube content partners. The onBehalfOfContentOwner parameter indicates that the request's authorization credentials identify a YouTube CMS user who is acting on behalf of the content owner specified in the parameter value. This parameter is intended for YouTube content partners that own and manage many different YouTube channels. It allows content owners to authenticate once and get access to all their video and channel data, without having to provide authentication credentials for each individual channel. The actual CMS account that the user authenticates with must be linked to the specified YouTube content owner.

**channel\_sections\_list** (*key, part, channelId=None, id=None, mine=None, hl=None, onBehalfOfContentOwner=None*)

Returns a list of resources that match the API request criteria.

*Required parameters:*

#### Parameters

- **key** – string Your Google API key.
- **part** – string The part parameter specifies a comma-separated list of one or more channel-section resource properties that the API response will include. If the parameter identifies a property that contains child properties, the child properties will be included in the response. For example, in a channelSection resource, the snippet property contains other properties, such as a display title for the section. If you set part=snippet, the API response will also contain all of those nested properties. The following list contains the part names that you can include in the parameter value and the quota cost for each part: contentDetails: 2 id: 0 localizations: 2 snippet: 2 targeting: 2

*Filters (specify exactly one of the following parameters):*

#### Parameters

- **channelId** – string The channelId parameter specifies a YouTube channel ID. If a request specifies a value for this parameter, the API will only return the specified channel's sections.
- **id** – string The id parameter specifies a comma-separated list of IDs that uniquely identify the channelSection resources that are being retrieved. In a channelSection resource, the id property specifies the section's ID.
- **mine** – boolean This parameter can only be used in a properly authorized request. Set this parameter's value to true to retrieve a feed of the channel sections associated with the authenticated user's YouTube channel.

*Optional parameters:*

#### Parameters

- **hl** – string The hl parameter instructs the API to retrieve localized resource metadata for a specific application language that the YouTube website supports. The parameter value must be a language code included in the list returned by the `i18nLanguages.list` method. If localized resource details are available in that language, the resource's `snippet.localized` object will contain the localized values. However, if localized details are not available, the `snippet.localized` object will contain resource details in the resource's default language.

- **onBehalfOfContentOwner** – string This parameter can only be used in a properly authorized request. Note: This parameter is intended exclusively for YouTube content partners. The `onBehalfOfContentOwner` parameter indicates that the request's authorization credentials identify a YouTube CMS user who is acting on behalf of the content owner specified in the parameter value. This parameter is intended for YouTube content partners that own and manage many different YouTube channels. It allows content owners to authenticate once and get access to all their video and channel data, without having to provide authentication credentials for each individual channel. The CMS account that the user authenticates with must be linked to the specified YouTube content owner.

**channels\_list** (*key*, *part*, *categoryId=None*, *forUsername=None*, *id=None*, *managedByMe=None*, *mine=None*, *mySubscribers=None*, *hl=None*, *maxResults=None*, *onBehalfOfContentOwner=None*, *pageToken=None*)

Returns a collection of zero or more resources that match the request criteria.

*Required parameters:*

#### Parameters

- **key** – string Your Google API key.
- **part** – string The part parameter specifies a comma-separated list of one or more channel resource properties that the API response will include. If the parameter identifies a property that contains child properties, the child properties will be included in the response. For example, in a channel resource, the `contentDetails` property contains other properties, such as the uploads properties. As such, if you set `part=contentDetails`, the API response will also contain all of those nested properties. The following list contains the part names that you can include in the parameter value and the quota cost for each part: `auditDetails: 4` `brandingSettings: 2` `contentDetails: 2` `contentOwnerDetails: 2` `id: 0` `invideoPromotion: 2` `(deprecated)localizations: 2` `snippet: 2` `statistics: 2` `status: 2` `topicDetails: 2`

*Filters (specify exactly one of the following parameters):*

#### Parameters

- **categoryId** – string The `categoryId` parameter specifies a YouTube guide category, thereby requesting YouTube channels associated with that category.
- **forUsername** – string The `forUsername` parameter specifies a YouTube username, thereby requesting the channel associated with that username.
- **id** – string The `id` parameter specifies a comma-separated list of the YouTube channel ID(s) for the resource(s) that are being retrieved. In a channel resource, the `id` property specifies the channel's YouTube channel ID.
- **managedByMe** – boolean This parameter can only be used in a properly authorized request. Note: This parameter is intended exclusively for YouTube content partners. Set this parameter's value to true to instruct the API to only return channels managed by the content owner that the `onBehalfOfContentOwner` parameter specifies. The user must be authenticated as a CMS account linked to the specified content owner and `onBehalfOfContentOwner` must be provided.
- **mine** – boolean This parameter can only be used in a properly authorized request. Set this parameter's value to true to instruct the API to only return channels owned by the authenticated user.
- **mySubscribers** – boolean This parameter has been deprecated. This parameter can only be used in a properly authorized request. Use the `subscriptions.list` method and its `mySubscribers` parameter to retrieve a list of subscribers to the authenticated user's channel.

*Optional parameters:*

### Parameters

- **hl** – string The hl parameter instructs the API to retrieve localized resource metadata for a specific application language that the YouTube website supports. The parameter value must be a language code included in the list returned by the `i18nLanguages.list` method. If localized resource details are available in that language, the resource's `snippet.localized` object will contain the localized values. However, if localized details are not available, the `snippet.localized` object will contain resource details in the resource's default language.
- **maxResults** – unsigned integer The maxResults parameter specifies the maximum number of items that should be returned in the result set.
- **onBehalfOfContentOwner** – string This parameter can only be used in a properly authorized request. Note: This parameter is intended exclusively for YouTube content partners. The `onBehalfOfContentOwner` parameter indicates that the request's authorization credentials identify a YouTube CMS user who is acting on behalf of the content owner specified in the parameter value. This parameter is intended for YouTube content partners that own and manage many different YouTube channels. It allows content owners to authenticate once and get access to all their video and channel data, without having to provide authentication credentials for each individual channel. The CMS account that the user authenticates with must be linked to the specified YouTube content owner.
- **pageToken** – string The pageToken parameter identifies a specific page in the result set that should be returned. In an API response, the `nextPageToken` and `prevPageToken` properties identify other pages that could be retrieved.

**comment\_threads\_list** (*key, part, allThreadsRelatedToChannelId=None, channelId=None, id=None, videoId=None, maxResults=None, moderationStatus=None, order=None, pageToken=None, searchTerms=None, textFormat=None*)

Returns a list of comment threads that match the API request parameters.

*Required parameters:*

### Parameters

- **key** – string Your Google API key.
- **part** – string The part parameter specifies a comma-separated list of one or more comment-thread resource properties that the API response will include. The following list contains the part names that you can include in the parameter value and the quota cost for each part: `id: 0 replies: 2 snippet: 2`

*Filters (specify exactly one of the following parameters):*

### Parameters

- **allThreadsRelatedToChannelId** – string The `allThreadsRelatedToChannelId` parameter instructs the API to return all comment threads associated with the specified channel. The response can include comments about the channel or about the channel's videos.
- **channelId** – string The `channelId` parameter instructs the API to return comment threads containing comments about the specified channel. (The response will not include comments left on videos that the channel uploaded.)
- **id** – string The `id` parameter specifies a comma-separated list of comment thread IDs for the resources that should be retrieved.
- **videoId** – string The `videoId` parameter instructs the API to return comment threads associated with the specified video ID.

*Optional parameters:*

### Parameters



- **maxResults** – unsigned integer The maxResults parameter specifies the maximum number of items that should be returned in the result set.
- **moderationStatus** – string This parameter can only be used in a properly authorized request. Set this parameter to limit the returned comment threads to a particular moderation state. Note: This parameter is not supported for use in conjunction with the id parameter. The default value is published. Acceptable values are: heldForReview – Retrieve comment threads that are awaiting review by a moderator. A comment thread can be included in the response if the top-level comment or at least one of the replies to that comment are awaiting review. likelySpam – Retrieve comment threads classified as likely to be spam. A comment thread can be included in the response if the top-level comment or at least one of the replies to that comment is considered likely to be spam. published – Retrieve threads of published comments. This is the default value. A comment thread can be included in the response if its top-level comment has been published.
- **order** – string The order parameter specifies the order in which the API response should list comment threads. Valid values are: time - Comment threads are ordered by time. This is the default behavior. relevance - Comment threads are ordered by relevance. Note: This parameter is not supported for use in conjunction with the id parameter.
- **pageToken** – string The pageToken parameter identifies a specific page in the result set that should be returned. In an API response, the nextPageToken property identifies the next page of the result that can be retrieved. Note: This parameter is not supported for use in conjunction with the id parameter.
- **searchTerms** – string The searchTerms parameter instructs the API to limit the API response to only contain comments that contain the specified search terms. Note: This parameter is not supported for use in conjunction with the id parameter.
- **textFormat** – string Set this parameter's value to html or plainText to instruct the API to return the comments left by users in html formatted or in plain text. The default value is html. Acceptable values are: html – Returns the comments in HTML format. This is the default value. plainText – Returns the comments in plain text format.

**comments\_list** (*key, part, id=None, parentId=None, maxResults=None, pageToken=None, textFormat=None*)

Returns a list of comments that match the API request parameters.

*Required parameters:*

#### Parameters

- **key** – string Your Google API key.
- **part** – string The part parameter specifies a comma-separated list of one or more comment resource properties that the API response will include. The following list contains the part names that you can include in the parameter value and the quota cost for each part: id: 0 snippet: 1

*Filters (specify exactly one of the following parameters):*

#### Parameters

- **id** – string The id parameter specifies a comma-separated list of comment IDs for the resources that are being retrieved. In a comment resource, the id property specifies the comment's ID.
- **parentId** – string The parentId parameter specifies the ID of the comment for which replies should be retrieved. Note: YouTube currently supports replies only for top-level comments. However, replies to replies may be supported in the future.

*Optional parameters:*

**Parameters**

- **maxResults** – unsigned integer The maxResults parameter specifies the maximum number of items that should be returned in the result set.
- **pageToken** – string The pageToken parameter identifies a specific page in the result set that should be returned. In an API response, the nextPageToken property identifies the next page of the result that can be retrieved. Note: This parameter is not supported for use in conjunction with the id parameter.
- **textFormat** – string This parameter indicates whether the API should return comments formatted as HTML or as plain text. The default value is html. Acceptable values are: html – Returns the comments in HTML format. This is the default value. plainText – Returns the comments in plain text format.

**guide\_categories\_list** (*key, part, id=None, regionCode=None, hl=None*)

Returns a list of categories that can be associated with YouTube channels.

*Required parameters:*

**Parameters**

- **key** – string Your Google API key.
- **part** – string The part parameter specifies the guideCategory resource properties that the API response will include. Set the parameter value to snippet. The snippet part has a quota cost of 2 units.

*Filters (specify exactly one of the following parameters):*

**Parameters**

- **id** – string The id parameter specifies a comma-separated list of the YouTube channel category ID(s) for the resource(s) that are being retrieved. In a guideCategory resource, the id property specifies the YouTube channel category ID.
- **regionCode** – string The regionCode parameter instructs the API to return the list of guide categories available in the specified country. The parameter value is an ISO 3166-1 alpha-2 country code.

*Optional parameters:*

**Parameters hl** – string The hl parameter specifies the language that will be used for text values in the API response. The default value is en-US.

**i18n\_languages\_list** (*key, part, hl=None*)

Returns a list of application languages that the YouTube website supports.

*Required parameters:*

**Parameters**

- **key** – string Your Google API key.
- **part** – string The part parameter specifies the i18nLanguage resource properties that the API response will include. Set the parameter value to snippet. The snippet part has a quota cost of 1 unit.

*Optional parameters:*

**Parameters hl** – string The hl parameter specifies the language that should be used for text values in the API response. The default value is en\_US.

**i18n\_regions\_list** (*key, part, hl=None*)

Returns a list of content regions that the YouTube website supports.

*Required parameters:*

**Parameters**

- **key** – string Your Google API key.
- **part** – string The part parameter specifies the i18nRegion resource properties that the API response will include. Set the parameter value to snippet. The snippet part has a quota cost of 1 unit.

*Optional parameters:*

**Parameters hl** – string The hl parameter specifies the language that should be used for text values in the API response. The default value is en\_US.

**playlist\_items\_list** (*key, part, id=None, playlistId=None, maxResults=None, onBehalfOfContentOwner=None, pageToken=None, videoId=None*)

Returns a collection of playlist items that match the API request parameters. You can retrieve all of the playlist items in a specified playlist or retrieve one or more playlist items by their unique IDs.

*Required parameters:*

**Parameters**

- **key** – string Your Google API key.
- **part** – string The part parameter specifies a comma-separated list of one or more playlistItem resource properties that the API response will include. If the parameter identifies a property that contains child properties, the child properties will be included in the response. For example, in a playlistItem resource, the snippet property contains numerous fields, including the title, description, position, and resourceId properties. As such, if you set part=snippet, the API response will contain all of those properties. The following list contains the part names that you can include in the parameter value and the quota cost for each part: contentDetails: 2 id: 0 snippet: 2 status: 2

*Filters (specify exactly one of the following parameters):*

**Parameters**

- **id** – string The id parameter specifies a comma-separated list of one or more unique playlist item IDs.
- **playlistId** – string The playlistId parameter specifies the unique ID of the playlist for which you want to retrieve playlist items. Note that even though this is an optional parameter, every request to retrieve playlist items must specify a value for either the id parameter or the playlistId parameter.

*Optional parameters:*

**Parameters**

- **maxResults** – unsigned integer The maxResults parameter specifies the maximum number of items that should be returned in the result set.
- **onBehalfOfContentOwner** – string This parameter can only be used in a properly authorized request. Note: This parameter is intended exclusively for YouTube content partners. The onBehalfOfContentOwner parameter indicates that the request's authorization credentials identify a YouTube CMS user who is acting on behalf of the content owner specified in the parameter value. This parameter is intended for YouTube content partners that own and manage many different YouTube channels. It allows content owners to authenticate

once and get access to all their video and channel data, without having to provide authentication credentials for each individual channel. The CMS account that the user authenticates with must be linked to the specified YouTube content owner.

- **pageToken** – string The pageToken parameter identifies a specific page in the result set that should be returned. In an API response, the nextPageToken and prevPageToken properties identify other pages that could be retrieved.
- **videoId** – string The videoId parameter specifies that the request should return only the playlist items that contain the specified video.

**playlists\_list** (*key, part, channelId=None, id=None, mine=None, hl=None, maxResults=None, onBehalfOfContentOwner=None, onBehalfOfContentOwnerChannel=None, pageToken=None*)

Returns a collection of playlists that match the API request parameters. For example, you can retrieve all playlists that the authenticated user owns, or you can retrieve one or more playlists by their unique IDs.

*Required parameters:*

#### Parameters

- **key** – string Your Google API key.
- **part** – string The part parameter specifies a comma-separated list of one or more playlist resource properties that the API response will include. If the parameter identifies a property that contains child properties, the child properties will be included in the response. For example, in a playlist resource, the snippet property contains properties like author, title, description, tags, and timeCreated. As such, if you set part=snippet, the API response will contain all of those properties. The following list contains the part names that you can include in the parameter value and the quota cost for each part: contentDetails: 2 id: 0 localizations: 2 player: 0 snippet: 2 status: 2

*Filters (specify exactly one of the following parameters):*

#### Parameters

- **channelId** – string This value indicates that the API should only return the specified channel's playlists.
- **id** – string The id parameter specifies a comma-separated list of the YouTube playlist ID(s) for the resource(s) that are being retrieved. In a playlist resource, the id property specifies the playlist's YouTube playlist ID.
- **mine** – boolean This parameter can only be used in a properly authorized request. Set this parameter's value to true to instruct the API to only return playlists owned by the authenticated user.

*Optional parameters:*

#### Parameters

- **hl** – string The hl parameter instructs the API to retrieve localized resource metadata for a specific application language that the YouTube website supports. The parameter value must be a language code included in the list returned by the `i18nLanguages.list` method. If localized resource details are available in that language, the resource's snippet.localized object will contain the localized values. However, if localized details are not available, the snippet.localized object will contain resource details in the resource's default language.
- **maxResults** – unsigned integer The maxResults parameter specifies the maximum number of items that should be returned in the result set.

- **onBehalfOfContentOwner** – string This parameter can only be used in a properly authorized request. Note: This parameter is intended exclusively for YouTube content partners. The `onBehalfOfContentOwner` parameter indicates that the request's authorization credentials identify a YouTube CMS user who is acting on behalf of the content owner specified in the parameter value. This parameter is intended for YouTube content partners that own and manage many different YouTube channels. It allows content owners to authenticate once and get access to all their video and channel data, without having to provide authentication credentials for each individual channel. The CMS account that the user authenticates with must be linked to the specified YouTube content owner.
- **onBehalfOfContentOwnerChannel** – string This parameter can only be used in a properly authorized request. Note: This parameter is intended exclusively for YouTube content partners. The `onBehalfOfContentOwnerChannel` parameter specifies the YouTube channel ID of the channel to which a video is being added. This parameter is required when a request specifies a value for the `onBehalfOfContentOwner` parameter, and it can only be used in conjunction with that parameter. In addition, the request must be authorized using a CMS account that is linked to the content owner that the `onBehalfOfContentOwner` parameter specifies. Finally, the channel that the `onBehalfOfContentOwnerChannel` parameter value specifies must be linked to the content owner that the `onBehalfOfContentOwner` parameter specifies. This parameter is intended for YouTube content partners that own and manage many different YouTube channels. It allows content owners to authenticate once and perform actions on behalf of the channel specified in the parameter value, without having to provide authentication credentials for each separate channel.
- **pageToken** – string The `pageToken` parameter identifies a specific page in the result set that should be returned. In an API response, the `nextPageToken` and `prevPageToken` properties identify other pages that could be retrieved.

**search** (*key*, *part*, *forContentOwner=None*, *forDeveloper=None*, *forMine=None*, *relatedToVideoId=None*, *channelId=None*, *channelType=None*, *eventType=None*, *location=None*, *locationRadius=None*, *maxResults=None*, *onBehalfOfContentOwner=None*, *order=None*, *pageToken=None*, *publishedAfter=None*, *publishedBefore=None*, *q=None*, *regionCode=None*, *relevanceLanguage=None*, *safeSearch=None*, *topicId=None*, *type=None*, *videoCaption=None*, *videoCategoryId=None*, *videoDefinition=None*, *videoDimension=None*, *videoDuration=None*, *videoEmbeddable=None*, *videoLicense=None*, *videoSyndicated=None*, *videoType=None*)

Returns a collection of search results that match the query parameters specified in the API request. By default, a search result set identifies matching video, channel, and playlist resources, but you can also configure queries to only retrieve a specific type of resource.

*Required parameters:*

#### Parameters

- **key** – string Your Google API key.
- **part** – string The `part` parameter specifies a comma-separated list of one or more search resource properties that the API response will include. Set the parameter value to `snippet`.

*Filters (specify 0 or 1 of the following parameters):*

#### Parameters

- **forContentOwner** – boolean This parameter can only be used in a properly authorized request, and it is intended exclusively for YouTube content partners. The `forContentOwner` parameter restricts the search to only retrieve videos owned by the content owner identified by the `onBehalfOfContentOwner` parameter. If `forContentOwner` is set to `true`, the request must also meet these requirements: The `onBehalfOfContentOwner` parameter is required. The user authorizing the request must be using an account linked to the specified content owner. The `type` parameter value must be set to `video`. None of the following other param-

eters can be set: videoDefinition, videoDimension, videoDuration, videoLicense, videoEmbeddable, videoSyndicated, videoType.

- **forDeveloper** – boolean This parameter can only be used in a properly authorized request. The forDeveloper parameter restricts the search to only retrieve videos uploaded via the developer’s application or website. The API server uses the request’s authorization credentials to identify the developer. The forDeveloper parameter can be used in conjunction with optional search parameters like the q parameter. For this feature, each uploaded video is automatically tagged with the project number that is associated with the developer’s application in the Google Developers Console. When a search request subsequently sets the forDeveloper parameter to true, the API server uses the request’s authorization credentials to identify the developer. Therefore, a developer can restrict results to videos uploaded through the developer’s own app or website but not to videos uploaded through other apps or sites.
- **forMine** – boolean This parameter can only be used in a properly authorized request. The forMine parameter restricts the search to only retrieve videos owned by the authenticated user. If you set this parameter to true, then the type parameter’s value must also be set to video. In addition, none of the following other parameters can be set in the same request: videoDefinition, videoDimension, videoDuration, videoLicense, videoEmbeddable, videoSyndicated, videoType.
- **relatedToVideoId** – string The relatedToVideoId parameter retrieves a list of videos that are related to the video that the parameter value identifies. The parameter value must be set to a YouTube video ID and, if you are using this parameter, the type parameter must be set to video. Note that if the relatedToVideoId parameter is set, the only other supported parameters are part, maxResults, pageToken, regionCode, relevanceLanguage, safeSearch, type (which must be set to video), and fields.

*Optional parameters:*

#### **Parameters**

- **channelId** – string The channelId parameter indicates that the API response should only contain resources created by the channel. Note: Search results are constrained to a maximum of 500 videos if your request specifies a value for the channelId parameter and sets the type parameter value to video, but it does not also set one of the forContentOwner, forDeveloper, or forMine filters.
- **channelType** – string The channelType parameter lets you restrict a search to a particular type of channel. Acceptable values are: any – Return all channels. show – Only retrieve shows.
- **eventType** – string The eventType parameter restricts a search to broadcast events. If you specify a value for this parameter, you must also set the type parameter’s value to video. Acceptable values are: completed – Only include completed broadcasts. live – Only include active broadcasts. upcoming – Only include upcoming broadcasts.
- **location** – string The location parameter, in conjunction with the locationRadius parameter, defines a circular geographic area and also restricts a search to videos that specify, in their metadata, a geographic location that falls within that area. The parameter value is a string that specifies latitude/longitude coordinates e.g. (37.42307,-122.08427). The location parameter value identifies the point at the center of the area. The locationRadius parameter specifies the maximum distance that the location associated with a video can be from that point for the video to still be included in the search results. The API returns an error if your request specifies a value for the location parameter but does not also specify a value for the locationRadius parameter.

- **locationRadius** – string The locationRadius parameter, in conjunction with the location parameter, defines a circular geographic area. The parameter value must be a floating point number followed by a measurement unit. Valid measurement units are m, km, ft, and mi. For example, valid parameter values include 1500m, 5km, 10000ft, and 0.75mi. The API does not support locationRadius parameter values larger than 1000 kilometers. Note: See the definition of the location parameter for more information.
- **maxResults** – unsigned integer The maxResults parameter specifies the maximum number of items that should be returned in the result set.
- **onBehalfOfContentOwner** – string This parameter can only be used in a properly authorized request. Note: This parameter is intended exclusively for YouTube content partners. The onBehalfOfContentOwner parameter indicates that the request's authorization credentials identify a YouTube CMS user who is acting on behalf of the content owner specified in the parameter value. This parameter is intended for YouTube content partners that own and manage many different YouTube channels. It allows content owners to authenticate once and get access to all their video and channel data, without having to provide authentication credentials for each individual channel. The CMS account that the user authenticates with must be linked to the specified YouTube content owner.
- **order** – string The order parameter specifies the method that will be used to order resources in the API response. The default value is relevance. Acceptable values are: date – Resources are sorted in reverse chronological order based on the date they were created. rating – Resources are sorted from highest to lowest rating. relevance – Resources are sorted based on their relevance to the search query. This is the default value for this parameter. title – Resources are sorted alphabetically by title. videoCount – Channels are sorted in descending order of their number of uploaded videos. viewCount – Resources are sorted from highest to lowest number of views. For live broadcasts, videos are sorted by number of concurrent viewers while the broadcasts are ongoing.
- **pageToken** – string The pageToken parameter identifies a specific page in the result set that should be returned. In an API response, the nextPageToken and prevPageToken properties identify other pages that could be retrieved.
- **publishedAfter** – datetime The publishedAfter parameter indicates that the API response should only contain resources created at or after the specified time. The value is an RFC 3339 formatted date-time value (1970-01-01T00:00:00Z).
- **publishedBefore** – datetime The publishedBefore parameter indicates that the API response should only contain resources created before or at the specified time. The value is an RFC 3339 formatted date-time value (1970-01-01T00:00:00Z).
- **q** – string The q parameter specifies the query term to search for. Your request can also use the Boolean NOT (-) and OR (|) operators to exclude videos or to find videos that are associated with one of several search terms. For example, to search for videos matching either “boating” or “sailing”, set the q parameter value to boating|sailing. Similarly, to search for videos matching either “boating” or “sailing” but not “fishing”, set the q parameter value to boating|sailing -fishing. Note that the pipe character must be URL-escaped when it is sent in your API request. The URL-escaped value for the pipe character is %7C.
- **regionCode** – string The regionCode parameter instructs the API to return search results for videos that can be viewed in the specified country. The parameter value is an ISO 3166-1 alpha-2 country code.
- **relevanceLanguage** – string The relevanceLanguage parameter instructs the API to return search results that are most relevant to the specified language. The parameter value is typically an ISO 639-1 two-letter language code. However, you should use the values zh-Hans for simplified Chinese and zh-Hant for traditional Chinese. Please note that results

in other languages will still be returned if they are highly relevant to the search query term.

- **safeSearch** – string The safeSearch parameter indicates whether the search results should include restricted content as well as standard content. Acceptable values are: moderate – YouTube will filter some content from search results and, at the least, will filter content that is restricted in your locale. Based on their content, search results could be removed from search results or demoted in search results. This is the default parameter value. none – YouTube will not filter the search result set. strict – YouTube will try to exclude all restricted content from the search result set. Based on their content, search results could be removed from search results or demoted in search results.
- **topicId** – string The topicId parameter indicates that the API response should only contain resources associated with the specified topic. The value identifies a Freebase topic ID. Important: Due to the deprecation of Freebase and the Freebase API, the topicId parameter started working differently as of February 27, 2017. At that time, YouTube started supporting a small set of curated topic IDs, and you can only use that smaller set of IDs as values for this parameter. See topic IDs supported as of February 15, 2017 Topics Music topics /m/04rlf Music (parent topic) /m/02mscn Christian music /m/0ggq0m Classical music /m/01lyv Country /m/02lkt Electronic music /m/0glt670 Hip hop music /m/05rwpb Independent music /m/03\_d0 Jazz /m/028sqc Music of Asia /m/0g293 Music of Latin America /m/064t9 Pop music /m/06cqb Reggae /m/06j6l Rhythm and blues /m/06by7 Rock music /m/0gywn Soul music Gaming topics /m/0bzvm2 Gaming (parent topic) /m/025zzc Action game /m/02ntfj Action-adventure game /m/0b1vjn Casual game /m/02hygl Music video game /m/04q1x3q Puzzle video game /m/01sjng Racing video game /m/040313g Role-playing video game /m/021bp2 Simulation video game /m/022dc6 Sports game /m/03hf\_rm Strategy video game Sports topics /m/06ntj Sports (parent topic) /m/0jm\_ American football /m/018jz Baseball /m/018w8 Basketball /m/01cgz Boxing /m/09xp\_ Cricket /m/02vx4 Football /m/037hz Golf /m/03tmr Ice hockey /m/01h7lh Mixed martial arts /m/0410tth Motorsport /m/07bs0 Tennis /m/07\_53 Volleyball Entertainment topics /m/02jtt Entertainment (parent topic) /m/09kqc Humor /m/02vxn Movies /m/05qjc Performing arts /m/066wd Professional wrestling /m/0f2f9 TV shows Lifestyle topics /m/019\_rr Lifestyle (parent topic) /m/032tl Fashion /m/027x7n Fitness /m/02wbm Food /m/03glg Hobby /m/068hy Pets /m/041xxh Physical attractiveness [Beauty] /m/07c1v Technology /m/07bxq Tourism /m/07yv9 Vehicles Society topics /m/098wr Society (parent topic) /m/09s1f Business /m/0kt51 Health /m/01h6rj Military /m/05qt0 Politics /m/06bvp Religion Other topics /m/01k8wb Knowledge
- **type** – string The type parameter restricts a search query to only retrieve a particular type of resource. The value is a comma-separated list of resource types. The default value is video,channel,playlist. Acceptable values are: channel,playlist,video
- **videoCaption** – string The videoCaption parameter indicates whether the API should filter video search results based on whether they have captions. If you specify a value for this parameter, you must also set the type parameter's value to video. Acceptable values are: any – Do not filter results based on caption availability. closedCaption – Only include videos that have captions. none – Only include videos that do not have captions.
- **videoCategoryId** – string The videoCategoryId parameter filters video search results based on their category. If you specify a value for this parameter, you must also set the type parameter's value to video.
- **videoDefinition** – string The videoDefinition parameter lets you restrict a search to only include either high definition (HD) or standard definition (SD) videos. HD videos are available for playback in at least 720p, though higher resolutions, like 1080p, might also be available. If you specify a value for this parameter, you must also set the type parameter's value to video. Acceptable values are: any – Return all videos, regardless of their resolution. high – Only retrieve HD videos. standard – Only retrieve videos in standard definition.



- **videoDimension** – string The videoDimension parameter lets you restrict a search to only retrieve 2D or 3D videos. If you specify a value for this parameter, you must also set the type parameter's value to video. Acceptable values are: 2d – Restrict search results to exclude 3D videos. 3d – Restrict search results to only include 3D videos. any – Include both 3D and non-3D videos in returned results. This is the default value.
- **videoDuration** – string The videoDuration parameter filters video search results based on their duration. If you specify a value for this parameter, you must also set the type parameter's value to video. Acceptable values are: any – Do not filter video search results based on their duration. This is the default value. long – Only include videos longer than 20 minutes. medium – Only include videos that are between four and 20 minutes long (inclusive). short – Only include videos that are less than four minutes long.
- **videoEmbeddable** – string The videoEmbeddable parameter lets you to restrict a search to only videos that can be embedded into a webpage. If you specify a value for this parameter, you must also set the type parameter's value to video. Acceptable values are: any – Return all videos, embeddable or not. true – Only retrieve embeddable videos.
- **videoLicense** – string The videoLicense parameter filters search results to only include videos with a particular license. YouTube lets video uploaders choose to attach either the Creative Commons license or the standard YouTube license to each of their videos. If you specify a value for this parameter, you must also set the type parameter's value to video. Acceptable values are: any – Return all videos, regardless of which license they have, that match the query parameters. creativeCommon – Only return videos that have a Creative Commons license. Users can reuse videos with this license in other videos that they create. Learn more. youtube – Only return videos that have the standard YouTube license.
- **videoSyndicated** – string The videoSyndicated parameter lets you to restrict a search to only videos that can be played outside youtube.com. If you specify a value for this parameter, you must also set the type parameter's value to video. Acceptable values are: any – Return all videos, syndicated or not. true – Only retrieve syndicated videos.
- **videoType** – string The videoType parameter lets you restrict a search to a particular type of videos. If you specify a value for this parameter, you must also set the type parameter's value to video. Acceptable values are: any – Return all videos. episode – Only retrieve episodes of shows. movie – Only retrieve movies.

**subscriptions\_list** (*key, part, channelId=None, id=None, mine=None, myRecentSubscribers=None, mySubscribers=None, forChannelId=None, maxResults=None, onBehalfOfContentOwner=None, onBehalfOfContentOwnerChannel=None, order=None, pageToken=None*)

Returns subscription resources that match the API request criteria.

*Required parameters:*

#### Parameters

- **key** – string Your Google API key.
- **part** – string The part parameter specifies a comma-separated list of one or more subscription resource properties that the API response will include. If the parameter identifies a property that contains child properties, the child properties will be included in the response. For example, in a subscription resource, the snippet property contains other properties, such as a display title for the subscription. If you set part=snippet, the API response will also contain all of those nested properties. The following list contains the part names that you can include in the parameter value and the quota cost for each part: contentDetails: 2 id: 0 snippet: 2 subscriberSnippet: 2

*Filters (specify exactly one of the following parameters):*

### Parameters

- **channelId** – string The channelId parameter specifies a YouTube channel ID. The API will only return that channel's subscriptions.
- **id** – string The id parameter specifies a comma-separated list of the YouTube subscription ID(s) for the resource(s) that are being retrieved. In a subscription resource, the id property specifies the YouTube subscription ID.
- **mine** – boolean This parameter can only be used in a properly authorized request. Set this parameter's value to true to retrieve a feed of the authenticated user's subscriptions.
- **myRecentSubscribers** – boolean This parameter can only be used in a properly authorized request. Set this parameter's value to true to retrieve a feed of the subscribers of the authenticated user in reverse chronological order (newest first). Note that this parameter only supports retrieval of the most recent 1000 subscribers to the authenticated user's channel. To retrieve a complete list of subscribers, use the mySubscribers parameter. That parameter, which does not return subscribers in a particular order, does not limit the number of subscribers that can be retrieved.
- **mySubscribers** – boolean This parameter can only be used in a properly authorized request. Set this parameter's value to true to retrieve a feed of the subscribers of the authenticated user in no particular order.

*Optional parameters:*

### Parameters

- **forChannelId** – string The forChannelId parameter specifies a comma-separated list of channel IDs. The API response will then only contain subscriptions matching those channels.
- **maxResults** – unsigned integer The maxResults parameter specifies the maximum number of items that should be returned in the result set.
- **onBehalfOfContentOwner** – string Note: This parameter is intended exclusively for YouTube content partners. The onBehalfOfContentOwner parameter indicates that the request's authorization credentials identify a YouTube CMS user who is acting on behalf of the content owner specified in the parameter value. This parameter is intended for YouTube content partners that own and manage many different YouTube channels. It allows content owners to authenticate once and get access to all their video and channel data, without having to provide authentication credentials for each individual channel. The CMS account that the user authenticates with must be linked to the specified YouTube content owner.
- **onBehalfOfContentOwnerChannel** – string This parameter can only be used in a properly authorized request. Note: This parameter is intended exclusively for YouTube content partners. The onBehalfOfContentOwnerChannel parameter specifies the YouTube channel ID of the channel to which a video is being added. This parameter is required when a request specifies a value for the onBehalfOfContentOwner parameter, and it can only be used in conjunction with that parameter. In addition, the request must be authorized using a CMS account that is linked to the content owner that the onBehalfOfContentOwner parameter specifies. Finally, the channel that the onBehalfOfContentOwnerChannel parameter value specifies must be linked to the content owner that the onBehalfOfContentOwner parameter specifies. This parameter is intended for YouTube content partners that own and manage many different YouTube channels. It allows content owners to authenticate once and perform actions on behalf of the channel specified in the parameter value, without having to provide authentication credentials for each separate channel.
- **order** – string The order parameter specifies the method that will be used to sort resources in the API response. The default value is SUBSCRIP-

TION\_ORDER\_RELEVANCE. Acceptable values are: alphabetical – Sort alphabetically. relevance – Sort by relevance. unread – Sort by order of activity.

- **pageToken** – string The pageToken parameter identifies a specific page in the result set that should be returned. In an API response, the nextPageToken and prevPageToken properties identify other pages that could be retrieved.

**video\_categories\_list** (*key, part, id=None, regionCode=None, hl=None*)

Returns a list of categories that can be associated with YouTube videos.

*Required parameters:*

#### Parameters

- **key** – string Your Google API key.
- **part** – string The part parameter specifies the videoCategory resource properties that the API response will include. Set the parameter value to snippet. The snippet part has a quota cost of 2 units.

*Filters (specify exactly one of the following parameters):*

#### Parameters

- **id** – string The id parameter specifies a comma-separated list of video category IDs for the resources that you are retrieving.
- **regionCode** – string The regionCode parameter instructs the API to return the list of video categories available in the specified country. The parameter value is an ISO 3166-1 alpha-2 country code.

*Optional parameters:*

**Parameters hl** – string The hl parameter specifies the language that should be used for text values in the API response. The default value is en\_US.

**videos\_list** (*key, part, chart=None, id=None, myRating=None, hl=None, maxHeight=None, maxResults=None, maxWidth=None, onBehalfOfContentOwner=None, pageToken=None, regionCode=None, videoCategoryId=None*)

Returns a list of videos that match the API request parameters.

*Required parameters:*

#### Parameters

- **key** – string Your Google API key.
- **part** – string The part parameter specifies a comma-separated list of one or more video resource properties that the API response will include. If the parameter identifies a property that contains child properties, the child properties will be included in the response. For example, in a video resource, the snippet property contains the channelId, title, description, tags, and categoryId properties. As such, if you set part=snippet, the API response will contain all of those properties. The following list contains the part names that you can include in the parameter value and the quota cost for each part: contentDetails: 2 fileDetails: 1 id: 0 liveStreamingDetails: 2 localizations: 2 player: 0 processingDetails: 1 recordingDetails: 2 snippet: 2 statistics: 2 status: 2 suggestions: 1 topicDetails: 2

*Filters (specify exactly one of the following parameters):*

#### Parameters

- **chart** – string The chart parameter identifies the chart that you want to retrieve. Acceptable values are: mostPopular – Return the most popular videos for the specified content region and video category.

- **id** – string The id parameter specifies a comma-separated list of the YouTube video ID(s) for the resource(s) that are being retrieved. In a video resource, the id property specifies the video's ID.
- **myRating** – string This parameter can only be used in a properly authorized request. Set this parameter's value to like or dislike to instruct the API to only return videos liked or disliked by the authenticated user. Acceptable values are: dislike – Returns only videos disliked by the authenticated user. like – Returns only video liked by the authenticated user.

*Optional parameters:*

### Parameters

- **hl** – string The hl parameter instructs the API to retrieve localized resource metadata for a specific application language that the YouTube website supports. The parameter value must be a language code included in the list returned by the `i18nLanguages.list` method. If localized resource details are available in that language, the resource's `snippet.localized` object will contain the localized values. However, if localized details are not available, the `snippet.localized` object will contain resource details in the resource's default language.
- **maxHeight** – unsigned integer The maxHeight parameter specifies the maximum height of the embedded player returned in the `player.embedHtml` property. You can use this parameter to specify that instead of the default dimensions, the embed code should use a height appropriate for your application layout. If the `maxWidth` parameter is also provided, the player may be shorter than the `maxHeight` in order to not violate the maximum width. Acceptable values are 72 to 8192, inclusive.
- **maxResults** – unsigned integer The maxResults parameter specifies the maximum number of items that should be returned in the result set.
- **maxWidth** – unsigned integer The maxWidth parameter specifies the maximum width of the embedded player returned in the `player.embedHtml` property. You can use this parameter to specify that instead of the default dimensions, the embed code should use a width appropriate for your application layout. If the `maxHeight` parameter is also provided, the player may be narrower than `maxWidth` in order to not violate the maximum height. Acceptable values are 72 to 8192, inclusive.
- **onBehalfOfContentOwner** – string This parameter can only be used in a properly authorized request. Note: This parameter is intended exclusively for YouTube content partners. The `onBehalfOfContentOwner` parameter indicates that the request's authorization credentials identify a YouTube CMS user who is acting on behalf of the content owner specified in the parameter value. This parameter is intended for YouTube content partners that own and manage many different YouTube channels. It allows content owners to authenticate once and get access to all their video and channel data, without having to provide authentication credentials for each individual channel. The CMS account that the user authenticates with must be linked to the specified YouTube content owner.
- **pageToken** – string The pageToken parameter identifies a specific page in the result set that should be returned. In an API response, the `nextPageToken` and `prevPageToken` properties identify other pages that could be retrieved. Note: This parameter is supported for use in conjunction with the `myRating` parameter, but it is not supported for use in conjunction with the `id` parameter.
- **regionCode** – string The regionCode parameter instructs the API to select a video chart available in the specified region. This parameter can only be used in conjunction with the `chart` parameter. The parameter value is an ISO 3166-1 alpha-2 country code.
- **videoCategoryId** – string The videoCategoryId parameter identifies the video category for which the chart should be retrieved. This parameter can only be used in conjunction with

the chart parameter. By default, charts are not restricted to a particular category. The default value is 0.



## INDICES AND TABLES

- `genindex`
- `modindex`
- `search`

### 2.1 advertools

#### 2.1.1 advertools package

##### Subpackages

`advertools.code_recipes` package

##### Submodules

##### Module contents

##### Submodules

##### Regular Expressions for Extracting Structured Entities

A collection of regular expressions for use in different contexts. Each one is available in two formats:

- **REGEX\_RAW: (HASHTAG\_RAW, MENTION\_RAW, etc.) raw string only, for sharing** and combining with other regexes
- **REGEX: (HASHTAG, MENTION, etc.) compiled regex, readable, and annotated**

Based on Unicode database v11.0.0

URL regex from Regular Expressions Cookbook 2nd Ed. O'Reilly

## URL Builders

`url_utm_ga(url, utm_source, utm_medium=None, utm_campaign=None, utm_content=None, utm_term=None)`

Generate a URL with UTM codes for your campaigns.

### Parameters

- `url` (*str*) – a valid URL, required
- `utm_source` (*str*) – the referrer of the traffic (e.g. facebook, twitter)
- `utm_medium` (*str*) – marketing medium (e.g. banner, email)
- `utm_campaign` (*str*) – the name of the campaign (e.g. summer\_promo, 20pct\_off)
- `utm_content` (*str*) – ad name / differentiator (e.g. 728x90, mpu, square\_banner)
- `utm_term` (*str*) – search terms bid on (only relevant for search campaigns)

**Returns** URL-encoded string for the campaign

```
>>> url_utm_ga('mysite.com', utm_source='the source')
'mysite.com?utm_source=the+source'
```

```
>>> url_utm_ga('mysite.com', utm_source='the source',
...           utm_medium='THE MEDIUM!!',
...           utm_campaign='campaign*name&^%', utm_content='728x90')
'mysite.com?utm_content=728x90&utm_campaign=campaign%2Aname%26%5E%25&utm_
↵medium=THE+MEDIUM%21%21&utm_source=the+source'
```

## Module contents

Top-level package for advertools.

### 2.1.2 Change Log - advertools

#### Unreleased

- **Added**
  - New parameter *recursive* for `sitemap_to_df` to control whether or not to get all sub sitemaps (default), or to only get the current one
  - Strategies and recipes for crawling and scraping

#### 0.10.7 (2020-09-18)

- **Added**
  - New function `knowledge_graph` for querying Google's API
  - Faster `sitemap_to_df` with threads
  - New parameter *max\_workers* for `sitemap_to_df` to determine how fast it could go
  - New parameter *capitalize\_adgroups* for `kw_generate` to determine whether or not to keep ad groups as is, or set them to title case (the default)



- **Fixed**
  - Remove restrictions on the number of URLs provided to `crawl`, assuming `follow_links` is set to `False` (list mode)
  - JSON-LD issue breaking crawls when it's invalid (now skipped)
- **Removed**
  - Deprecate the `youtube.guide_categories_list` (no longer supported by the API)

### 0.10.6 (2020-06-30)

- **Added**
  - JSON-LD support in crawling. If available on a page, JSON-LD items will have special columns, and multiple JSON-LD snippets will be numbered for easy filtering
- **Changed**
  - Stricter parsing for rel attributes, making sure they are in link elements as well
  - Date column names for `robotstxt_to_df` and `sitemap_to_df` unified as “download\_date”
  - Numbering OG, Twitter, and JSON-LD where multiple elements are present in the same page, follows a unified approach: no numbering for the first element, and numbers start with “1” from the second element on. “element”, “element\_1”, “element\_2” etc.

### 0.10.5 (2020-06-14)

- **Added**
  - **New features for the `crawl` function:**
    - \* Extract canonical tags if available
    - \* Extract alternate `href` and `hreflang` tags if available
    - \* Open Graph data “og:title”, “og:type”, “og:image”, etc.
    - \* Twitter cards data “twitter:site”, “twitter:title”, etc.
- **Fixed**
  - **Minor fixes to `robotstxt_to_df`:**
    - \* Allow whitespace in fields
    - \* Allow case-insensitive fields
- **Changed**
  - `crawl` now only supports `output_file` with the extension “.jl”
  - `word_frequency` drops `wtd_freq` and `rel_value` columns if `num_list` is not provided

#### 0.10.4 (2020-06-07)

- **Added**
  - New function `url_to_df`, splitting URLs into their components and to a DataFrame
  - Slight speed up for `robotstxt_test`

#### 0.10.3 (2020-06-03)

- **Added**
  - New function `robotstxt_test`, testing URLs and whether they can be fetched by certain user-agents
- **Changed**
  - Documentation main page relayout, grouping of topics, & sidebar captions
  - Various documentation clarifications and new tests

#### 0.10.2 (2020-05-25)

- **Added**
  - User-Agent info to requests getting sitemaps and robotstxt files
  - CSS/XPath selectors support for the `crawl` function
  - Support for custom spider settings with a new parameter `custom_settings`
- **Fixed**
  - Update changed supported search operators and values for CSE

#### 0.10.1 (2020-05-23)

- **Changed**
  - Links are better handled, and new output columns are available: `links_url`, `links_text`, `links_fragment`, `links_nofollow`
  - `body_text` extraction is improved by containing `<p>`, `<li>`, and `<span>` elements

#### 0.10.0 (2020-05-21)

- **Added**
  - New function `crawl` for crawling and parsing websites
  - New function `robotstxt_to_df` downloading robots.txt files into DataFrames

### 0.9.1 (2020-05-19)

- **Added**
  - Ability to specify robots.txt file for `sitemap_to_df`
  - Ability to retrieve any kind of sitemap (news, video, or images)
  - Errors column to the returned DataFrame if any errors occur
  - A new `sitemap_downloaded` column showing datetime of getting the sitemap
- **Fixed**
  - Logging issue causing `sitemap_to_df` to log the same action twice
  - Issue preventing URLs not ending with xml or gz from being retrieved
  - Correct sitemap URL showing in the `sitemap` column

### 0.9.0 (2020-04-03)

- **Added**
  - New function `sitemap_to_df` imports an XML sitemap into a DataFrame

### 0.8.1 (2020-02-08)

- **Changed**
  - Column `query_time` is now named `queryTime` in the `youtube` functions
  - Handle `json_normalize` import from pandas based on pandas version

### 0.8.0 (2020-02-02)

- **Added**
  - New module `youtube` connecting to all GET requests in API
  - `extract_numbers` new function
  - `emoji_search` new function
  - `emoji_df` new variable containing all emoji as a DataFrame
- **Changed**
  - Emoji database updated to v13.0
  - `serp_goog` with expanded `pagemap` and metadata
- **Fixed**
  - `serp_goog` errors, some parameters not appearing in result df
  - `extract_numbers` issue when providing dash as a separator in the middle

### 0.7.3 (2019-04-17)

- **Added**
  - New function *extract\_exclamations* very similar to *extract\_questions*
  - New function *extract\_urls*, also counts top domains and top TLDs
  - New keys to *extract\_emoji*; *top\_emoji\_categories* & *top\_emoji\_sub\_categories*
  - Groups and sub-groups to *emoji db*

### 0.7.2 (2019-03-29)

- **Changed**
  - *Emoji regex* updated
  - Simpler extraction of Spanish *questions*

### 0.7.1 (2019-03-26)

- **Fixed**
  - Missing `__init__` imports.

### 0.7.0 (2019-03-26)

- **Added**
  - New *extract\_* functions:
    - \* Generic *extract* used by all others, and takes arbitrary regex to extract text.
    - \* *extract\_questions* to get question mark statistics, as well as the text of questions asked.
    - \* *extract\_currency* shows text that has currency symbols in it, as well as surrounding text.
    - \* *extract\_intense\_words* gets statistics about, and extract words with any character repeated three or more times, indicating an intense feeling (+ve or -ve).
  - New function *word\_tokenize*:
    - \* Used by *word\_frequency* to get tokens of 1,2,3-word phrases (or more).
    - \* Split a list of text into tokens of a specified number of words each.
  - New stop-words from the `spaCy` package:

**current:** Arabic, Azerbaijani, Danish, Dutch, English, Finnish, French, German, Greek, Hungarian, Italian, Kazakh, Nepali, Norwegian, Portuguese, Romanian, Russian, Spanish, Swedish, Turkish.

**new:** Bengali, Catalan, Chinese, Croatian, Hebrew, Hindi, Indonesian, Irish, Japanese, Persian, Polish, Sinhala, Tagalog, Tamil, Tatar, Telugu, Thai, Ukrainian, Urdu, Vietnamese
- **Changed**
  - *word\_frequency* takes new parameters:
    - \* *regex* defaults to words, but can be changed to anything ‘S+’ to split words and keep punctuation for example.
    - \* *sep* not longer used as an option, the above *regex* can be used instead

- \* *num\_list* now optional, and defaults to counts of 1 each if not provided. Useful for counting *abs\_freq* only if data not available.
- \* *phrase\_len* the number of words in each split token. Defaults to 1 and can be set to 2 or higher. This helps in analyzing phrases as opposed to words.
- Parameters supplied to *serp\_goog* appear at the beginning of the result df
- *serp\_youtube* now contains *nextPageToken* to make paginating requests easier

### 0.6.0 (2019-02-11)

- **New function**
  - *extract\_words* to extract an arbitrary set of words
- **Minor updates**
  - *ad\_from\_string* slots argument reflects new text ad lengths
  - *hashtag* regex improved

### 0.5.3 (2019-01-31)

- **Fix minor bugs**
  - Handle Twitter search queries with 0 results in final request

### 0.5.2 (2018-12-01)

- **Fix minor bugs**
  - Properly handle requests for >50 items (*serp\_youtube*)
  - Rewrite test for *\_dict\_product*
  - Fix issue with string printing error msg

### 0.5.1 (2018-11-06)

- **Fix minor bugs**
  - *\_dict\_product* implemented with lists
  - Missing keys in some YouTube responses

### 0.5.0 (2018-11-04)

- **New function *serp\_youtube***
  - Query YouTube API for videos, channels, or playlists
  - Multiple queries (product of parameters) in one function call
  - Reponse looping and merging handled, one DataFrame
- *serp\_goog* return Google's original error messages
- twitter responses with entities, get the entities extracted, each in a separate column

#### 0.4.1 (2018-10-13)

- **New function *serp\_goog* (based on Google CSE)**
  - Query Google search and get the result in a DataFrame
  - Make multiple queries / requests in one function call
  - All responses merged in one DataFrame
- `twitter.get_place_trends` results are ranked by town and country

#### 0.4.0 (2018-10-08)

- **New Twitter module based on `twython`**
  - Wraps 20+ functions for getting Twitter API data
  - Gets data in a pandas DataFrame
  - Handles looping over requests higher than the defaults
- Tested on Python 3.7

#### 0.3.0 (2018-08-14)

- Search engine marketing cheat sheet.
- **New set of `extract_` functions with summary stats for each:**
  - `extract_hashtags`
  - `extract_mentions`
  - `extract_emoji`
- Tests and bug fixes

#### 0.2.0 (2018-07-06)

- New set of `kw_<match-type>` functions.
- Full testing and coverage.

#### 0.1.0 (2018-07-02)

- First release on PyPI.
- **Functions available:**
  - `ad_create`: create a text ad place words in placeholders
  - **`ad_from_string`: split a long string to shorter string that fit into** given slots
  - `kw_generate`: generate keywords from lists of products and words
  - `url_utm_ga`: generate a UTM-tagged URL for Google Analytics tracking
  - **`word_frequency`: measure the absolute and weighted frequency of words in** collection of documents

## PYTHON MODULE INDEX

### a

- advertools, 100
- advertools.ad\_create, 8
- advertools.ad\_from\_string, 9
- advertools.code\_recipes, 99
- advertools.code\_recipes.spider\_strategies,  
31
- advertools.emoji, 47
- advertools.extract, 50
- advertools.knowledge\_graph, 42
- advertools.kw\_generate, 4
- advertools.regex, 99
- advertools.robotstxt, 11
- advertools.serp, 35
- advertools.sitemaps, 15
- advertools.spider, 22
- advertools.stopwords, 61
- advertools.twitter, 67
- advertools.url\_builders, 99
- advertools.urlytics, 45
- advertools.word\_frequency, 62
- advertools.word\_tokenize, 66
- advertools.youtube, 80





## A

activities\_list() (in module *advertools.youtube*), 80  
 ad\_create() (in module *advertools.ad\_create*), 8  
 ad\_from\_string() (in module *advertools.ad\_from\_string*), 11  
 advertools  
     module, 100  
 advertools.ad\_create  
     module, 8  
 advertools.ad\_from\_string  
     module, 9  
 advertools.code\_recipes  
     module, 99  
 advertools.code\_recipes.spider\_strategies  
     module, 31  
 advertools.emoji  
     module, 47  
 advertools.extract  
     module, 50  
 advertools.knowledge\_graph  
     module, 42  
 advertools.kw\_generate  
     module, 4  
 advertools.regex  
     module, 99  
 advertools.robotstxt  
     module, 11  
 advertools.serp  
     module, 35  
 advertools.sitemaps  
     module, 15  
 advertools.spider  
     module, 22  
 advertools.stopwords  
     module, 61  
 advertools.twitter  
     module, 67  
 advertools.url\_builders  
     module, 99  
 advertools.urlytics  
     module, 45

advertools.word\_frequency  
     module, 62  
 advertools.word\_tokenize  
     module, 66  
 advertools.youtube  
     module, 80  
 authenticate() (in module *advertools.twitter*), 68

## C

capitalize, 10  
 captions\_list() (in module *advertools.youtube*), 81  
 channel\_sections\_list() (in module *advertools.youtube*), 82  
 channels\_list() (in module *advertools.youtube*), 83  
 comment\_threads\_list() (in module *advertools.youtube*), 84  
 comments\_list() (in module *advertools.youtube*), 85  
 crawl() (in module *advertools.spider*), 30

## E

emoji\_search() (in module *advertools.emoji*), 48  
 extra\_info, 64  
 extract() (in module *advertools.extract*), 51  
 extract\_currency() (in module *advertools.extract*), 51  
 extract\_emoji() (in module *advertools.emoji*), 48  
 extract\_exclamations() (in module *advertools.extract*), 52  
 extract\_hashtags() (in module *advertools.extract*), 54  
 extract\_intense\_words() (in module *advertools.extract*), 55  
 extract\_mentions() (in module *advertools.extract*), 55  
 extract\_numbers() (in module *advertools.extract*), 56  
 extract\_questions() (in module *advertools.extract*), 57  
 extract\_urls() (in module *advertools.extract*), 58

`extract_words()` (in module `advertools.extract`), 59

## G

`get_application_rate_limit_status()` (in module `advertools.twitter`), 68

`get_available_trends()` (in module `advertools.twitter`), 68

`get_favorites()` (in module `advertools.twitter`), 68

`get_followers_ids()` (in module `advertools.twitter`), 69

`get_followers_list()` (in module `advertools.twitter`), 69

`get_friends_ids()` (in module `advertools.twitter`), 69

`get_friends_list()` (in module `advertools.twitter`), 70

`get_home_timeline()` (in module `advertools.twitter`), 70

`get_list_members()` (in module `advertools.twitter`), 71

`get_list_memberships()` (in module `advertools.twitter`), 71

`get_list_statuses()` (in module `advertools.twitter`), 72

`get_list_subscribers()` (in module `advertools.twitter`), 73

`get_list_subscriptions()` (in module `advertools.twitter`), 73

`get_mentions_timeline()` (in module `advertools.twitter`), 73

`get_place_trends()` (in module `advertools.twitter`), 74

`get_retweeters_ids()` (in module `advertools.twitter`), 74

`get_retweets()` (in module `advertools.twitter`), 75

`get_supported_languages()` (in module `advertools.twitter`), 75

`get_user_timeline()` (in module `advertools.twitter`), 75

`guide_categories_list()` (in module `advertools.youtube`), 86

## I

`i18n_languages_list()` (in module `advertools.youtube`), 86

`i18n_regions_list()` (in module `advertools.youtube`), 86

## K

`knowledge_graph()` (in module `advertools.knowledge_graph`), 45

`kw_broad()` (in module `advertools.kw_generate`), 6

`kw_exact()` (in module `advertools.kw_generate`), 6

`kw_generate()` (in module `advertools.kw_generate`), 6

`kw_modified()` (in module `advertools.kw_generate`), 7

`kw_neg_broad()` (in module `advertools.kw_generate`), 7

`kw_neg_exact()` (in module `advertools.kw_generate`), 7

`kw_neg_phrase()` (in module `advertools.kw_generate`), 7

`kw_phrase()` (in module `advertools.kw_generate`), 8

## L

`lookup_status()` (in module `advertools.twitter`), 76

`lookup_user()` (in module `advertools.twitter`), 76

## M

`make_dataframe()` (in module `advertools.twitter`), 77

module

`advertools`, 100

`advertools.ad_create`, 8

`advertools.ad_from_string`, 9

`advertools.code_recipes`, 99

`advertools.code_recipes.spider_strategies`, 31

`advertools.emoji`, 47

`advertools.extract`, 50

`advertools.knowledge_graph`, 42

`advertools.kw_generate`, 4

`advertools.regex`, 99

`advertools.robotstxt`, 11

`advertools.serp`, 35

`advertools.sitemaps`, 15

`advertools.spider`, 22

`advertools.stopwords`, 61

`advertools.twitter`, 67

`advertools.url_builders`, 99

`advertools.urlytics`, 45

`advertools.word_frequency`, 62

`advertools.word_tokenize`, 66

`advertools.youtube`, 80

## N

`num_list`, 63

## P

`phrase_len`, 63

`playlist_items_list()` (in module `advertools.youtube`), 87

`playlists_list()` (in module `advertools.youtube`), 88

## R

regex, [63](#)

retweeted\_of\_me() (in module *advertools.twitter*),  
[77](#)

rm\_words, [64](#)

robotstxt\_test() (in module *advertools.robotstxt*),  
[15](#)

robotstxt\_to\_df() (in module *advertools.robotstxt*), [15](#)

## S

s, [9](#)

search() (in module *advertools.twitter*), [77](#)

search() (in module *advertools.youtube*), [89](#)

search\_users() (in module *advertools.twitter*), [79](#)

sep, [10](#)

serp\_goog() (in module *advertools.serp*), [35](#)

serp\_youtube() (in module *advertools.serp*), [37](#)

set\_auth\_params() (in module *advertools.twitter*),  
[79](#)

set\_logging\_level() (in module *advertools.serp*),  
[42](#)

show\_lists() (in module *advertools.twitter*), [80](#)

show\_owned\_lists() (in module *advertools.twitter*), [80](#)

sitemap\_to\_df() (in module *advertools.sitemaps*),  
[21](#)

slots, [9](#)

subscriptions\_list() (in module *advertools.youtube*), [93](#)

## T

text\_list, [63](#)

## U

url\_to\_df() (in module *advertools.urlytics*), [47](#)

url\_utm\_ga() (in module *advertools.url\_builders*),  
[100](#)

## V

video\_categories\_list() (in module *advertools.youtube*), [95](#)

videos\_list() (in module *advertools.youtube*), [95](#)

## W

word\_frequency() (in module *advertools.word\_frequency*), [64](#)

word\_tokenize() (in module *advertools.word\_tokenize*), [66](#)

## Y

youtube\_channel\_details() (in module *advertools.serp*), [42](#)

youtube\_video\_details() (in module *advertools.serp*), [42](#)