

# Collecting & Analyzing Twitter data – an Introduction

**Viktoria Spaiser**

UAF in Political Science Informatics,  
School of Politics and International Studies

# Accessing Twitter data

## 1) **Twitter Streaming API (Application Programming Interface)**

- Real-time Twitter data collection of tweets
- Spritzer sample is free (1% of all public tweets)
- Other samples or full data (e.g. Firehose) are subject to a charge

<https://dev.twitter.com/streaming/overview>

## 2) **Twitter REST APIs (in particular Twitter Search API)**

- Historic (past 7 days!) data collection of tweets (e.g. based on hashtags)
- Collection of tweets by location (place operator of the Search API)
- Collection of followers & friends data for specified Twitter user(s)
- API Rate limits apply

<https://dev.twitter.com/rest/public>

# Accessing Twitter data

Missed the date?

- No panic, there is an archive for Streamed Twitter data

<https://archive.org/details/twitterstream>

The screenshot shows the Internet Archive website interface. At the top, there's a navigation bar with links for ABOUT, CONTACT, BLOG, PROJECTS, HELP, DONATE, TERMS, JOBS, VOLUNTEER, and PEOPLE. Below this, the main heading is "Archive Team: The Twitter Stream Grab" with a sub-description: "A simple collection of JSON grabbed from the general twitter stream, for the purposes of research, history, testing and memory. This is the 'Spritzer' version, the most light and shallow of Twitter grabs. Unfortunately, we do". There are "Share" and "Favorite" buttons. Below the description, there are tabs for "About" and "Collection". The main content area shows a grid of 10 items, each with a "Download of Twitter" title and a date (e.g., Jul 1, 2016, Jun 1, 2016, May 1, 2016, Apr 1, 2016, Mar 1, 2016). To the right, there's a search bar for the collection, a "PART OF" section with links to "Archive Team" and "Web Crawls", and a "LANGUAGE" section showing "English 6".

Here you can download historic Twitter Streaming API data in JSON format

# Accessing Twitter data

## What you need to access Twitter data via Twitter APIs

1. Twitter account
2. Obtain Authentication & Authorization (OAuth):
  - this requires registration as a developer (developing an app, even if you will not) with Twitter, register here: <https://apps.twitter.com>
  - you will get: **Consumer Key**, **Consumer Secret**, **Access token**, **Access token secret**

WITHOUT THESE YOU WILL NOT BE ABLE TO ACCESS DATA VIA TWITTER APIS!!!

# Accessing Twitter data

## 1. Python

(Python 2.7 + Anaconda for Python 2.7 recommended)

useful packages: *tweepy*, *Twython*, *simplejson*, *nltk* (*Natural Language Toolkit*)

install Python 2.7: <https://www.python.org/downloads/>

install Anaconda: <https://www.continuum.io/downloads>

install packages: e.g. type “pip install tweepy” in terminal/shell

## 2. R (packages *twitteR* and *ROAuth*):

<https://www.r-bloggers.com/setting-up-the-twitter-r-package-for-text-analytics/>

## 3. Other programming languages like Java etc.

## 4. NodeXL (no coding, Windows only, for Social Network Analyses only):

<http://www.pewinternet.org/files/2014/02/How-we-analyzed-Twitter-social-media-networks.pdf>

## 5. Mecodify (new, free software for extracting & visualizing Twitter data, no coding,

soon available from: <http://www.mecodem.eu>, developed by **Walid Al-Safaq:**

[walid.al-saqaf@ims.su.se](mailto:walid.al-saqaf@ims.su.se) )

## 6. LIDA seems to have developed some software to collect tweets data, contact

**David Batty:** [d.batty@leeds.ac.uk](mailto:d.batty@leeds.ac.uk)



# Twitter data, unprocessed

```
{  
  "created_at": "Thu Apr 30 21:53:11 +0000 2015",  
  "id": 593895901623496700,  
  "id_str": "593895901623496704",  
  "text": "This is a #test tweet @LoveforTestingT with an image. http://t.co/ZvgHovKZq4",  
  "source": "<a href='\"http://twitter.com/\"' rel='\"nofollow/\">Twitter Web Client</a>",  
  "truncated": false,  
  "in_reply_to_status_id": null,  
  "in_reply_to_status_id_str": null,  
  "in_reply_to_user_id": null,  
  "in_reply_to_user_id_str": null,  
  "in_reply_to_screen_name": null,  
  "user": {  
    "id": 2993982541,  
    "id_str": "2993982541",  
    "name": "Test Demo",  
    "screen_name": "jondee_test",  
    "location": "Denver, CO",  
    "url": null,  
    "description": "this is a test account.",  
    "protected": false,  
    "verified": false,  
    "followers_count": 2,  
    "friends_count": 43,  
    "listed_count": 0,  
    "favourites_count": 0,  
    "statuses_count": 30,  
    "created_at": "Sat Jan 24 00:12:53 +0000 2015",  
  }  
}
```

...

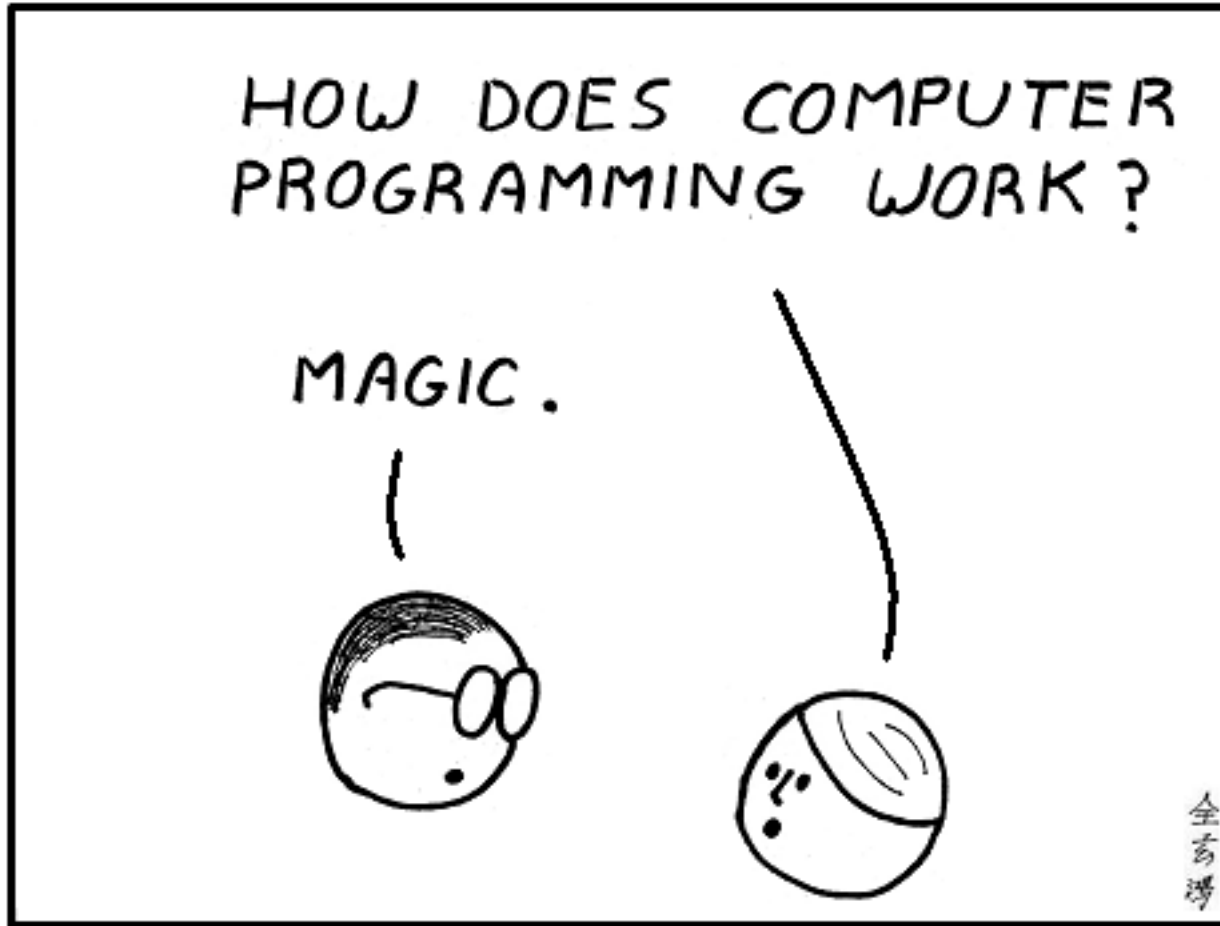
# Twitter data, key variables

Field	Description
id	Unique tweet ID number
text	Tweet text, if retweet then starts with RT @screen_name:
created_at	Timing of tweet creation, or of Twitter account creation if nested within the Twitter user field
place/coordinates	Latitude, longitude coordinates, if geo-enabled set to "true" (has to be activated by user, per default deactivated (value "false"))
user_mentions/ screen_name	Indicates whether and which Twitter user is mentioned (@) in the tweet
in_reply_to_screen_ name	Indicates whether the twitter was a reply and in that case to which Twitter user (if not a reply value "null")
user/screen_name	User name of Twitter user
user/location	Location information (e.g. name of town) as provided by Twitter user
user/name	Full name of Twitter user as provided by Twitter user
user/description	Profile description of Twitter user

and many more variables...: [http://support.gnip.com/sources/twitter/data\\_format.html](http://support.gnip.com/sources/twitter/data_format.html)



Ok, let's start coding then...



# Getting data from the Streaming API

The image shows the Spyder Python IDE interface. The main editor window displays a Python script for connecting to the Twitter Streaming API. The script includes the following code:

```
1 # -*- coding: utf-8 -*-
2 """
3 Created on Thu Sep 22 08:52:33 2016
4
5 @author: viktoria
6 """
7
8 import sys
9 sys.path.insert(0, 'C:\Anaconda\Lib\site-packages')
10 import tweepy
11 from tweepy import OAuthHandler
12
13 consumer_key = 'Your Consumer Key'
14 consumer_secret = 'Your Consume Secret'
15 access_token = 'Your Access Token'
16 access_secret = 'Your Access Secret'
17
18 auth = OAuthHandler(consumer_key, consumer_secret)
19 auth.set_access_token(access_token, access_secret)
20
21 api = tweepy.API(auth)
22
23 if(api.verify_credentials):
24     print 'We sucessfully logged in'
25
26 from tweepy import Stream
27 from tweepy.streaming import StreamListener
28
29 class MyListener(StreamListener):
30
31     def on_data(self, data):
32         try:
33             with open('OutputData.json', 'a') as f:
34                 f.write(data)
35                 return True
36         except BaseException as e:
37             print("Error on_data: %s" % str(e))
38             return True
39
40     def on_error(self, status):
41         print(status)
42         return True
43
44 twitter_stream = Stream(auth, MyListener())
45 twitter_stream.filter(track=['#YourHashtag,#YourHashtag2'])
46
```

The right-hand side of the IDE shows the Object Inspector and the IPython console. The Object Inspector displays a "Usage" box with the following text:

Usage

Here you can get help of any object by pressing Ctrl+I in front of it, either on the Editor or the Console.

Help can also be shown automatically after writing a left parenthesis next to an object. You can activate this behavior in *Preferences > Object Inspector*.

The IPython console shows the following text:

```
Python 2.7.12 |Anaconda custom (x86_64)| (default, Jun 29 2016, 11:09:23)
Type "copyright", "credits" or "license" for more information.

IPython 4.1.2 -- An enhanced Interactive Python.
?          -> Introduction and overview of IPython's features.
%quickref  -> Quick reference.
help       -> Python's own help system.
object?    -> Details about 'object', use 'object??' for extra details.
%gui?     -> A brief reference about the graphical user interface.

In [1]:
```

The bottom status bar of the IDE shows: Permissions: RW, End-of-lines: LF, Encoding: UTF-8, Line: 31, Column: 29, Memory: 43 %.

# Getting data from the Search API

```
import sys
sys.path.insert(0, 'C:\Anaconda\Lib\site-packages')
import tweepy
from tweepy import OAuthHandler

consumer_key = 'Your Consumer Key'
consumer_secret = 'Your Consumer Secret'
access_token = 'Your Access Token'
access_secret = 'Your Access Secret'

auth = OAuthHandler(consumer_key, consumer_secret)
auth.set_access_token(access_token, access_secret)

api = tweepy.API(auth, wait_on_rate_limit=True, wait_on_rate_limit_notify=True)

if (not api):
    print ("Can't Authenticate")
    sys.exit(-1)

searchQuery = '#YourHashtag'
maxTweets = 10000000
tweetsPerQry = 100
fName = 'OutputData.json'

sinceId = None

max_id = -1L

tweetCount = 0
print("Downloading max {0} tweets".format(maxTweets))
with open(fName, 'w') as f:
    while tweetCount < maxTweets:
        try:
            if (max_id <= 0):
                if (not sinceId):
                    new_tweets = api.search(q=searchQuery, count=tweetsPerQry)
                else:
                    new_tweets = api.search(q=searchQuery, count=tweetsPerQry,
                                             since_id=sinceId)
            else:
                if (not sinceId):
                    new_tweets = api.search(q=searchQuery, count=tweetsPerQry,
                                             max_id=str(max_id - 1))
                else:
                    new_tweets = api.search(q=searchQuery, count=tweetsPerQry,
                                             max_id=str(max_id - 1),
                                             since_id=sinceId)

            if not new_tweets:
                print("No more tweets found")
                break
            for tweet in new_tweets:
                s = str(tweet)
                f.write(s.encode("ascii"))
            tweetCount += len(new_tweets)
            print("Downloaded {0} tweets".format(tweetCount))
            max_id = new_tweets[-1].id
        except tweepy.TweepError as e:
            print("some error : " + str(e))
            break

print ("Downloaded {0} tweets, Saved to {1}".format(tweetCount, fName))
```

# Processing JSON data

```
1 # -*- coding: utf-8 -*-
2 """
3 Created on Thu Sep 22 15:00:59 2016
4
5 @author: viktorias
6 """
7
8 import string
9 import json
10
11 # READ IN JSON FILE
12 path = 'OutputData.json'
13 infile = open(path, 'rU')
14
15 # NAMES FOR HEADER ROW IN OUTPUT FILE
16 fields = "id screen_name name location".split()
17
18 # PREPARE YOUR OUTPUTFILE
19 outfn = "users_location.txt"
20 outfp = open(outfn, "w")
21 outfp.write(string.join(fields, "\t") + "\n") # header
22
23 # READING IN AND WRITING OUT DATA
24 #for entry in infile:
25 for line in open(path, 'r'):
26     tweet = json.loads(line)
27     # CREATE EMPTY DICTIONARY
28     r = {}
29     for f in fields:
30         r[f] = ""
31     # ASSIGN VALUE OF THE FIELDS IN JSON TO THE FIELDS IN OUR DICTIONARY
32     r['id'] = tweet['id']
33     r['screen_name'] = tweet['user']['screen_name']
34     r['name'] = tweet['user']['name']
35     r['location'] = tweet['user']['location']
36     # CREATE EMPTY LIST
37     lst = []
38     # ADD DATA FOR EACH VARIABLE
39     for f in fields:
40         lst.append(unicode(r[f]).replace("\\", "/"))
41     # WRITE ROW WITH DATA IN LIST
42     outfp.write(string.join(lst, "\t").encode("utf-8") + "\n")
43 outfp.close()
```

# Natural Language Processing

```
import json
import csv
import nltk
from nltk.corpus import stopwords
from nltk.collocations import BigramCollocationFinder

path = 'OutputData.json'
outpath = '/Users/macbook/Desktop/Work'
outfn = "ProcessedData.txt"

rutwout = list()
for line in open(path, 'r'):
    block = json.loads(line)
    tweet = block["text"]
    rutwout.append(tweet + "\n")

stopset = stopwords.words('english')
filter_stops = lambda w: len(w) < 2 or w in stopset

def fdist(file):
    freq = nltk.FreqDist()
    words = nltk.tokenize.regexp_tokenize(str(rutwout), "[\w]+")
    words = [word.lower() for word in words]
    words = [word for word in words if len(word) > 2]
    if words not in stopset:
        freq = nltk.FreqDist(words)
        freqwords = freq.most_common(300)
        return freqwords

def saveFreqWord(freq):
    temp_dict = dict(freq)
    writer = csv.writer(open('Wordcounts.csv', 'wb'))
    for key, value in temp_dict.items():
        writer.writerow([key, value])

saveFreqWord(fdist(path))

def bicolloc(file):
    words = nltk.tokenize.regexp_tokenize(str(rutwout), "[\w]+") #tokenising words
    words = [word.lower() for word in words]
    words = [word for word in words if len(word) > 2]
    bgm = nltk.collocations.BigramAssocMeasures() #initialising the BigramAssocMeasures
    bcf = BigramCollocationFinder.from_words(words) #initialising the BigramCollocationFinder
    bcf.apply_word_filter(filter_stops) #filtering words
    scored = bcf.score_ngrams(bgm.student_t)[:300] #computing 100 trigram-collocation with the highest scores based on T-Test
    open('/Users/macbook/Desktop/Work/Bigrams.txt', 'w').writelines(str(scored))
    temp_dict = dict(scored)
    csvData = []
    for (col1, col2), col3 in temp_dict.iteritems():
        csvData.append("%s, %s, %s" % (col1, col2, col3))
    f = open('Bigrams.csv', 'w')
    f.write("\n".join(csvData))
    f.close()

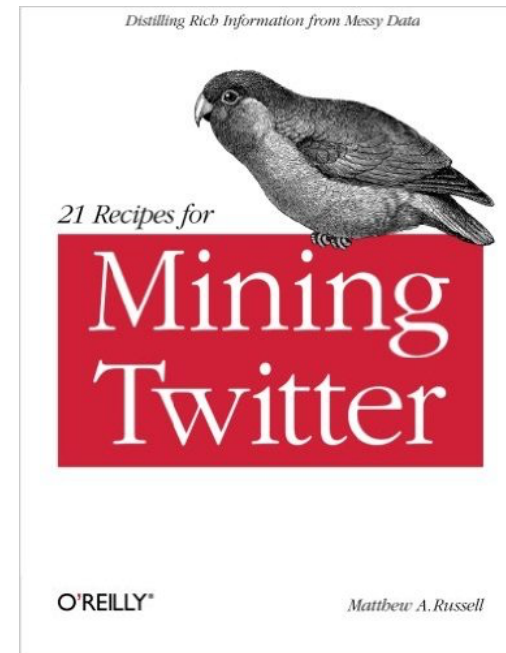
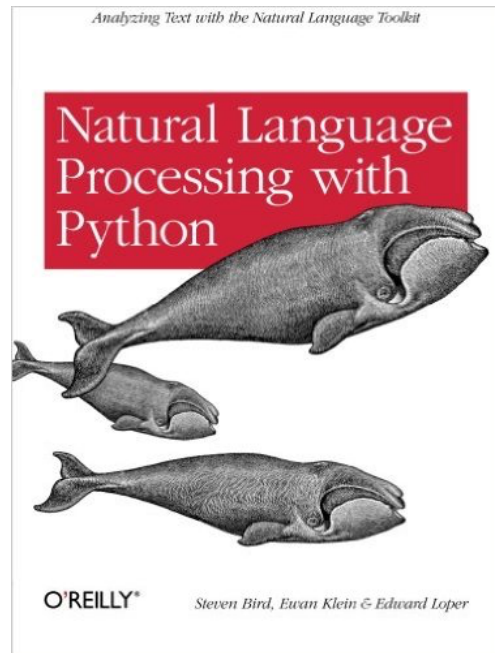
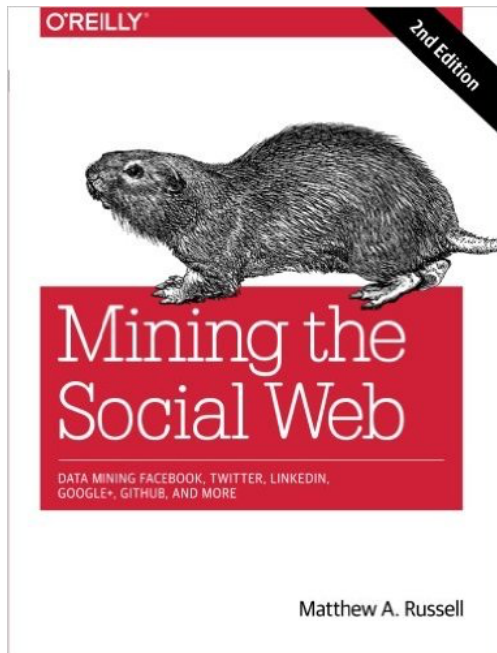
bicolloc(path)
```

# Geo-location Processing

```
1 # -*- coding: utf-8 -*-
2 """
3 Created on Thu Sep 22 19:26:59 2016
4
5 @author: viktoria
6 """
7 import json
8
9 path = 'OutputData.json'
10 with open(path, 'r') as f:
11     geo_data = {
12         "type": "FeatureCollection",
13         "features": []
14     }
15     for line in f:
16         tweet = json.loads(line)
17         if tweet['coordinates']:
18             geo_json_feature = {
19                 "type": "Feature",
20                 "geometry": tweet['coordinates'],
21                 "properties": {
22                     "text": tweet['text'],
23                     "created_at": tweet['created_at']
24                 }
25             }
26             geo_data['features'].append(geo_json_feature)
27
28 # Create GeoJSON
29 with open('geo_data.json', 'w') as fout:
30     fout.write(json.dumps(geo_data, indent=4))
31
```

You can use GeoJSON for instance in QGIS or to create interactive maps with Leaflet <http://leafletjs.com/examples/geojson.html>

# Recommended Further Reading



And many sources on the internet...