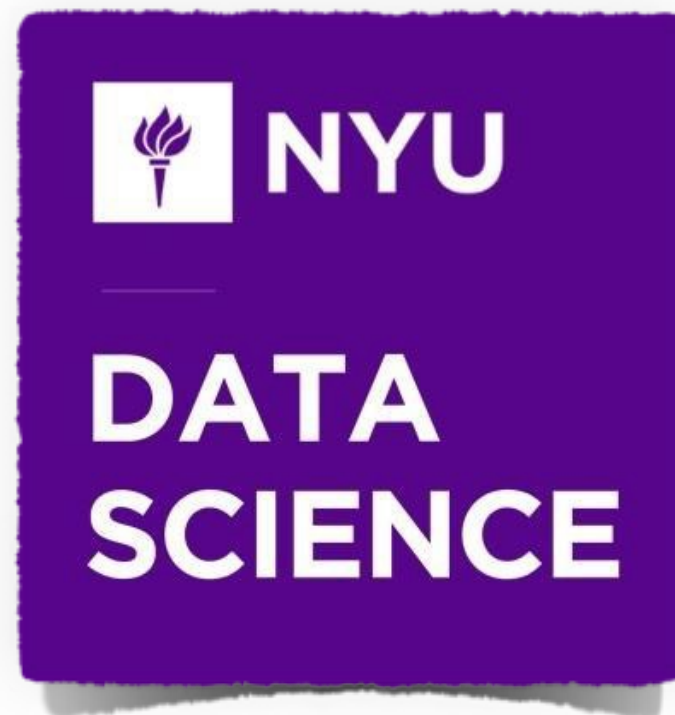


<https://github.com/bmtgoncalves/WebSci17>

Analyzing Geolocated Data with Twitter *Practice*

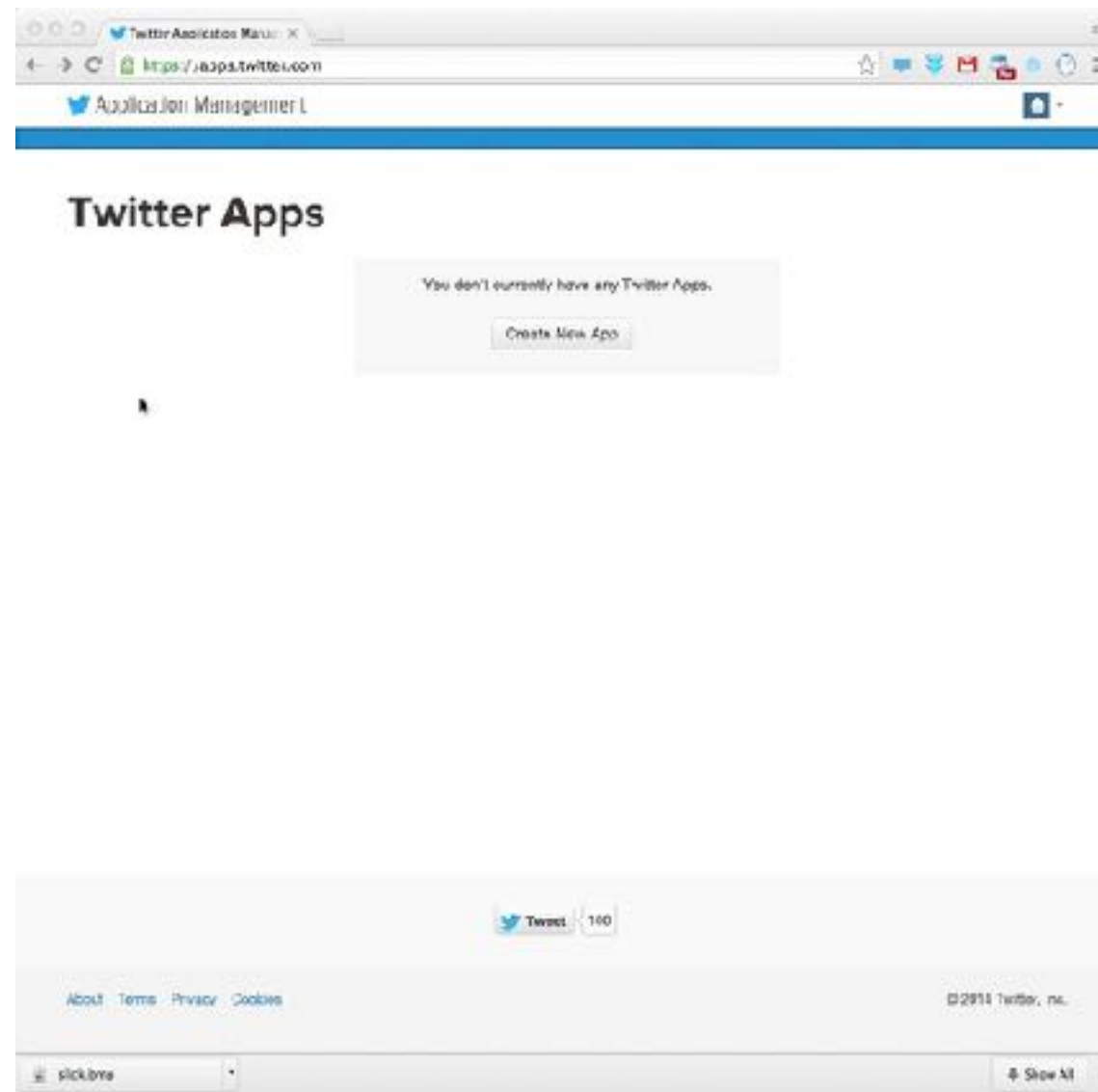
Bruno Gonçalves
www.bgoncalves.com



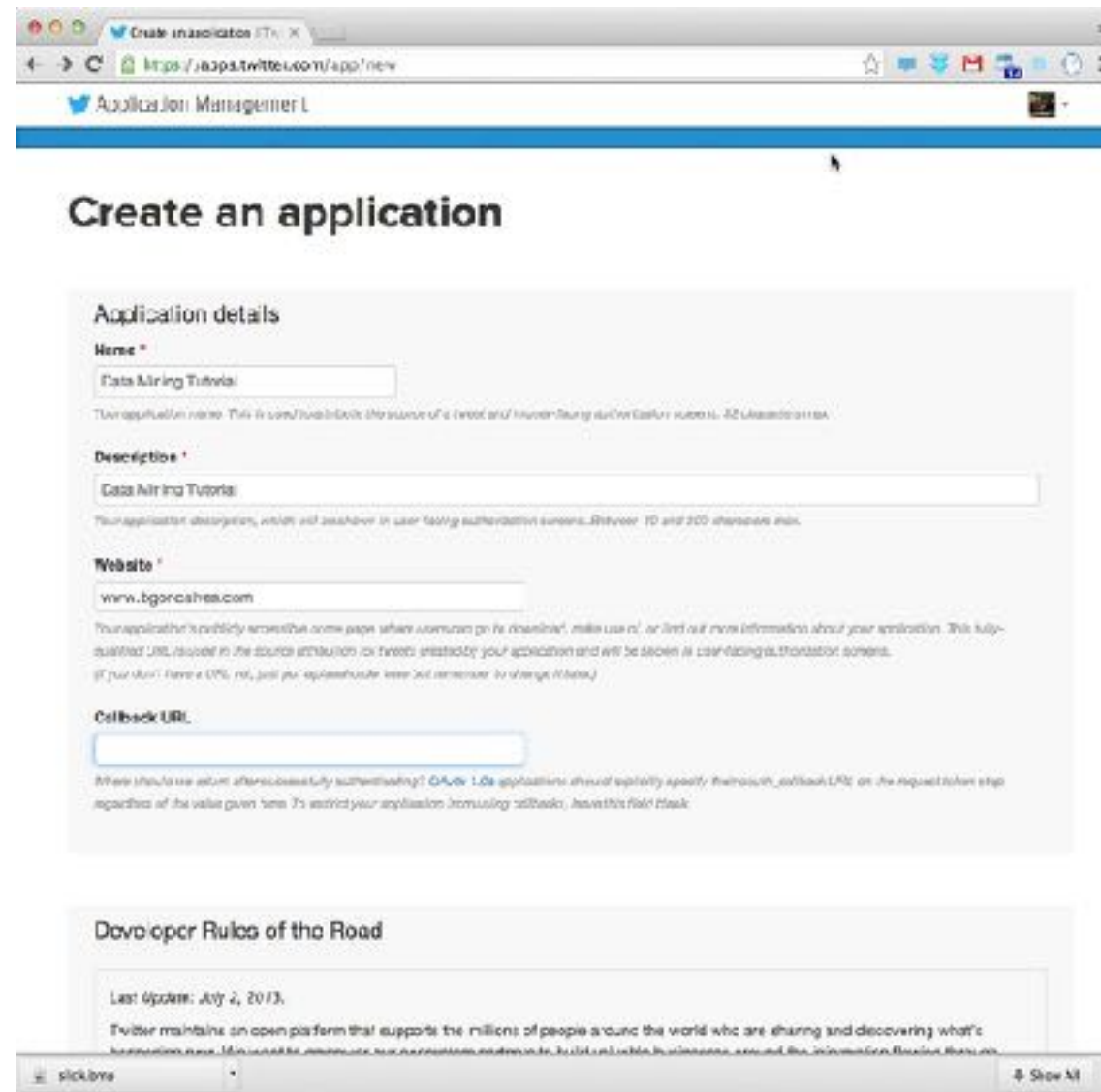
Requirements



Registering an Application



Registering an Application

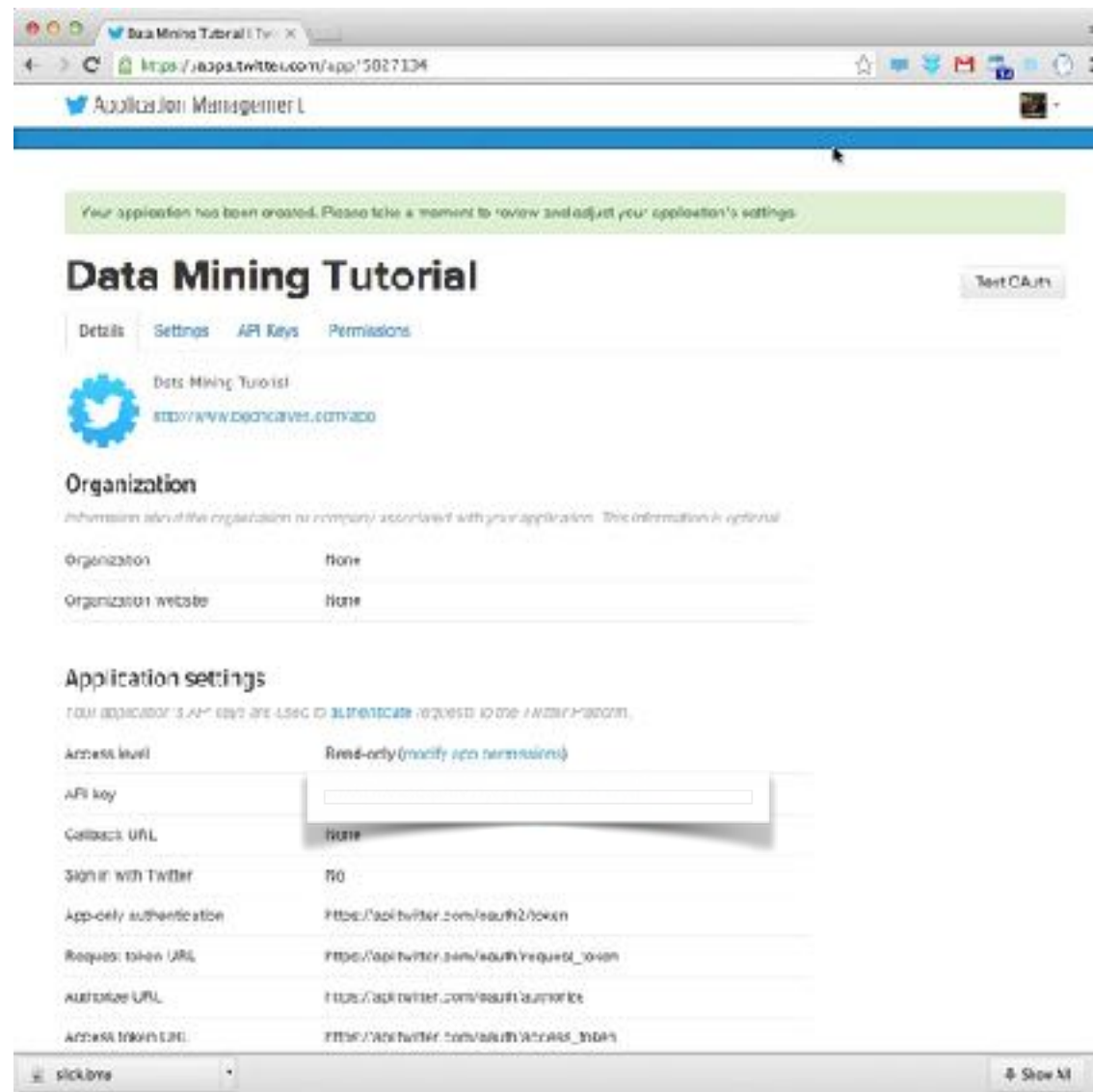


The screenshot shows a web browser window with the URL <https://apps.twitter.com/new>. The page title is "Application Manager L". The main heading is "Create an application". Below this is a form titled "Application details" with the following fields:

- Name ***: A text input field containing "Data Mining Tutorial". Below it, a small note reads: "Your application name. This is used to identify the source of a tweet and is visible during OAuth requests. 32 characters max."
- Description ***: A text input field containing "Data Mining Tutorial". Below it, a small note reads: "Your application description, which will appear in user-facing authorization screens. Between 70 and 200 characters max."
- Website ***: A text input field containing "www.bgoncalves.com". Below it, a small note reads: "Your application's publicly accessible home page where users can go to learn more about your application. This fully-qualified URL, listed in the source attribution on tweets created by your application and will be shown in user-facing authorization screens. (If you don't have a URL, just put 'placeholder text' and remember to change it later.)"
- Callback URL**: An empty text input field. Below it, a small note reads: "Where should we return after successfully authenticating? OAuth 1.0a applications should explicitly specify their oauth_callback URL on the request token step regardless of the value given here. To restrict your application (preventing callbacks), leave this field blank."

Below the form is a section titled "Developer Rules of the Road" with a "Last Update: July 2, 2013" and a paragraph stating: "Twitter maintains an open platform that supports the millions of people around the world who are sharing and discovering what's interesting, new, fun and useful. We want to ensure that our platform remains open, accessible, and secure. We've created a set of rules to help you understand how to use our platform responsibly." At the bottom of the page, there are links for "pickins" and "Show All".

Registering an Application

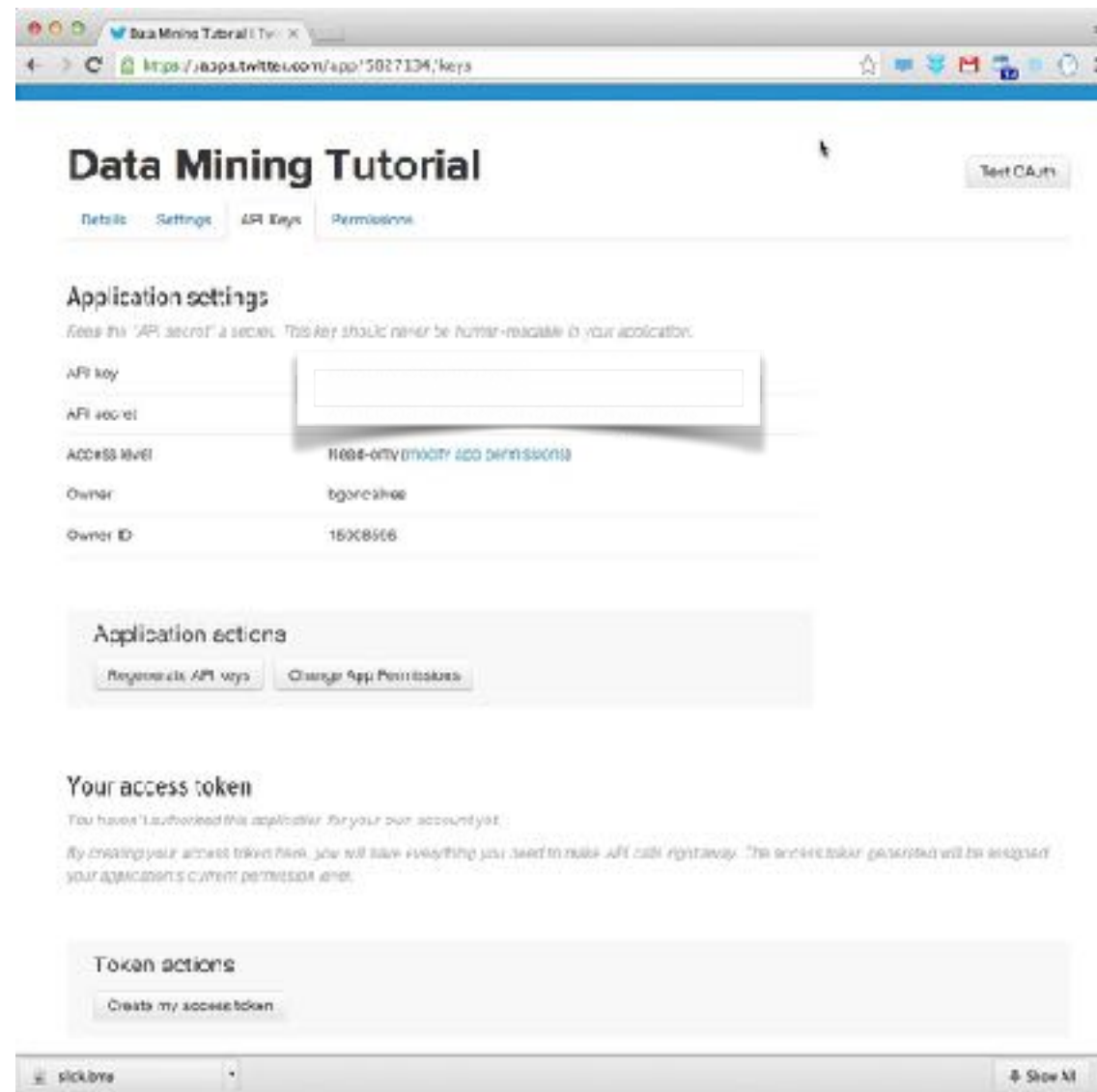


The screenshot shows the Twitter Application Management interface for an application named 'Data Mining Tutorial'. The page is titled 'Application Management' and includes a green notification bar stating 'Your application has been created. Please take a moment to review and adjust your application's settings.' The application's name is 'Data Mining Tutorial' and its website is 'http://www.bgoncalves.com/2012/01/01/data-mining-tutorial/'. The 'Organization' section is currently empty. The 'Application settings' section is expanded, showing various configuration options. A modal dialog box is open over the 'API key' field, which is currently empty. The 'API key' field is highlighted with a red border, indicating it is required. The 'API key' field is currently empty. The 'API key' field is currently empty.

Organization	Organization website
None	None

Application settings	
Access level	Read-only (modify app permissions)
API key	
Contact URL	None
Sign in with Twitter	No
App-only authentication	https://api.twitter.com/oauth2/token
Request token URL	https://api.twitter.com/oauth/request_token
Authorize URL	https://api.twitter.com/oauth/authorize
Access token URL	https://api.twitter.com/oauth/access_token

Registering an Application



Registering an Application

The screenshot shows the 'Data Mining Tutorial' application settings page on the Twitter Developer Portal. The page is divided into two main sections: 'Application settings' and 'Your access token'. Both sections have a warning box indicating that the API secret and access token secret should be kept secret. The 'Application settings' section includes fields for API key, API secret, Access level (Read-only (no write API permissions)), Owner (lgarciahe), and Owner ID (15068566). The 'Your access token' section includes fields for Access token, Access token secret, Access level (Read-only), Owner (lgarciahe), and Owner ID (15068566). There are buttons for 'Test OAuth', 'Regenerate API keys', and 'Change App Permissions'.

Data Mining Tutorial Test OAuth

[Details](#) [Settings](#) [API Keys](#) [Permissions](#)

Application settings

Keep this "API secret" a secret. This key should never be human-readable in your application.

API key	<input type="text"/>
API secret	<input type="text"/>
Access level	Read-only (no write API permissions)
Owner	lgarciahe
Owner ID	15068566

Application actions

[Regenerate API keys](#) [Change App Permissions](#)

Your access token

This access token can be used to make API requests on your own account's behalf. Do not share your access token secret with anyone.

Access token	<input type="text"/>
Access token secret	<input type="text"/>
Access level	Read-only
Owner	lgarciahe
Owner ID	15068566

[pickins](#) [Show All](#)

API Basics

<https://dev.twitter.com/docs>

- The **twitter** module provides the oauth interface. We just need to provide the right credentials.
- Best to keep the credentials in a **dict** and parametrize our calls with the dict key. This way we can switch between different accounts easily.
- **.Twitter(auth)** takes an **OAuth** instance as argument and returns a **Twitter** object that we can use to interact with the API
- **Twitter** methods mimic API structure
- 4 basic types of objects:
 - Tweets
 - Users
 - Entities
 - Places

Authenticating with the API

```
import tweepy
from twitter_accounts import accounts

app = accounts["social"]

auth = twitter.oauth.OAuth(app["token"],
                             app["token_secret"],
                             app["api_key"],
                             app["api_secret"])

twitter_api = twitter.Twitter(auth=auth)
```

- In the remainder of this course, the **accounts** dict will live inside the `twitter_accounts.py` file
- 4 basic types of objects:
 - Tweets
 - Users
 - Entities
 - Places

Searching for Tweets

<https://dev.twitter.com/docs/api/1.1/get/search/tweets>

- `.search.tweets(query, count)`
 - `query` is the content to search for
 - `count` is the maximum number of results to return
- returns dict with a list of “`statuses`” and “`search_metadata`”

```
{u'completed_in': 0.027,  
 u'count': 15,  
 u'max_id': 438088492577345536,  
 u'max_id_str': u'438088492577345536',  
 u'next_results': u'?max_id=438088485145034752&q=soccer&include_entities=1',  
 u'query': u'soccer',  
 u'refresh_url': u'?since_id=438088492577345536&q=soccer&include_entities=1',  
 u'since_id': 0,  
 u'since_id_str': u'0'}
```
- `search_results[“search_metadata”][“next_results”]` can be used to get the next page of results

Searching for Tweets

<https://dev.twitter.com/docs/api/1.1/get/search/tweets>

```
query = "instagram"
count = 200

search_results = twitter_api.search.tweets(q=query, count=count)

statuses = search_results["statuses"]
tweet_count = 0

while True:
    try:
        next_results = search_results["search_metadata"]["next_results"]

        args = dict(parse.parse_qsl(next_results[1:]))

        search_results = twitter_api.search.tweets(**args)
        statuses = search_results["statuses"]

        print(search_results["search_metadata"]["max_id"])

        for tweet in statuses:
            tweet_count += 1

            if tweet_count % 10000 == 0:
                print(tweet_count, file=sys.stderr)

            print(tweet["text"])
    except:
        break
```

Social Connections

<https://dev.twitter.com/docs/api/1.1/get/friends/ids>
<https://dev.twitter.com/docs/api/1.1/get/followers/ids>

- `.friends.ids()` and `.followers.ids()` returns a list of up to **5000** of a users friends or followers for a given **screen_name** or **user_id**
- result is a **dict** containing multiple fields:

```
[u'next_cursor_str',  
 u'previous_cursor',  
 u'ids',  
 u'next_cursor',  
 u'previous_cursor_str']
```
- ids are contained in `results["ids"]`.
- `results["next_cursor"]` allows us to obtain the next page of results.
- `.friends.ids(screen_name=screen_name, cursor=results["next_cursor"])` will return the next page of results
- `cursor=0` means no more results

Social Connections

<https://dev.twitter.com/docs/api/1.1/get/friends/ids>
<https://dev.twitter.com/docs/api/1.1/get/followers/ids>

```
import twitter
from twitter_accounts import accounts

app = accounts["social"]

auth = twitter.oauth.OAuth(app["token"],
                             app["token_secret"],
                             app["api_key"],
                             app["api_secret"])

twitter_api = twitter.Twitter(auth=auth)

screen_name = "stephen_wolfram"

cursor = -1
followers = []

while cursor != 0:
    result = twitter_api.followers.ids(screen_name=screen_name,
                                       cursor=cursor)

    followers += result["ids"]
    cursor = result["next_cursor"]

print("Found", len(followers), "Followers")
```

User Timeline

https://dev.twitter.com/docs/api/1.1/get/statuses/user_timeline

- `.statuses.user_timeline()` returns a set of tweets posted by a single user
- Important options:
 - `include_rts='true'` to Include retweets by this user
 - `count=200` number of tweets to return in each call
 - `trim_user='true'` to not include the user information (save bandwidth and processing time)
 - `max_id=1234` to include only tweets with an id lower than `1234`
- Returns at most `200` tweets in each call. Can get all of a users tweets (up to 3200) with multiple calls using `max_id`

User Timeline

https://dev.twitter.com/docs/api/1.1/get/statuses/user_timeline

```
import twitter
from twitter_accounts import accounts

app = accounts["social"]

auth = twitter.oauth.OAuth(app["token"],
                             app["token_secret"],
                             app["api_key"],
                             app["api_secret"])

twitter_api = twitter.Twitter(auth=auth)
screen_name = "BarackObama"

args = { "count" : 200,
         "trim_user": "true",
         "include_rts": "true"
       }

tweets = twitter_api.statuses.user_timeline(screen_name = screen_name, **args)
tweets_new = tweets

while len(tweets_new) > 0:
    max_id = tweets[-1]["id"] - 1
    tweets_new = twitter_api.statuses.user_timeline(screen_name = screen_name, max_id=max_id, **args)
    tweets += tweets_new

print("Found", len(tweets), "tweets")
```

Social Interactions

```
import twitter
from twitter_accounts import accounts

app = accounts["social"]

auth = twitter.oauth.OAuth(app["token"],
                             app["token_secret"],
                             app["api_key"],
                             app["api_secret"])

twitter_api = twitter.Twitter(auth=auth)
screen_name = "BarackObama"
args = { "count" : 200, "trim_user": "true", "include_rts": "true" }

tweets = twitter_api.statuses.user_timeline(screen_name=screen_name, **args)
tweets_new = tweets

while len(tweets_new) > 0:
    max_id = tweets[-1]["id"] - 1
    tweets_new = twitter_api.statuses.user_timeline(screen_name=screen_name, max_id=max_id, **args)
    tweets += tweets_new

user = tweets[0]["user"]["id"]

for tweet in tweets:
    if "retweeted_status" in tweet:
        print(user, "->", tweet["retweeted_status"]["user"]["id"])
    elif tweet["in_reply_to_user_id"]:
        print(tweet["in_reply_to_user_id"], "->", user)
```


NetworkX

- High productivity software for complex networks
- Simple Python interface
- Four types of graphs supported:
 - **Graph** - UnDirected
 - **DiGraph** - Undirected
 - **MultiGraph** - Multi-edged Graph
 - **MultiDiGraph** - Directed Multigraph
- Similar interface for all types of graphs
- Nodes can be any type of Python object - Practical way to manage relationships

Growing Graphs

- `.add_node(node_id)` Add a single node with ID **node_id**
- `.add_nodes_from()` Add a list of node ids
- `.add_edge(node_i, node_j)` Adds an edge between **node_i** and **node_j**
- `.add_edges_from()` Adds a list of edges. Individual edges are represented by tuples
- `.remove_node(node_id)/.remove_nodes_from()` Removing a node removes all associated edges
- `.remove_edge(node_i, node_j)/.remove_edges_from()`

Graph Properties

- `.nodes()` Returns the list of nodes
- `.edges()` Returns the list of edges
- `.degree()` Returns a dict with each nodes degree `.in_degree()/.out_degree()` returns dicts with in/out degree for [DiGraphs](#)
- `.is_connected()` Returns true if the node is connected
- `.is_weakly_connected()/.is_strongly_connected()` for [DiGraph](#)
- `.connected_components()` A list of nodes for each connected component

NetworkX - Example

```
import networkx as NX
import numpy as np
from collections import Counter
import matplotlib.pyplot as plt

def BarabasiAlbert(N=1000000):
    G = NX.Graph()

    nodes = range(N)
    G.add_nodes_from(nodes)

    edges = [0,1,1,2,2,0]

    for node_i in range(3, N):
        pos = np.random.randint(len(edges))
        node_j = edges[pos]

        edges.append(node_i)
        edges.append(node_j)

    edges = zip(nodes, edges[1::2])

    G.add_edges_from(edges)

    return G
```

NetworkX - Example

```
import networkx as NX
import numpy as np
from collections import Counter
import matplotlib.pyplot as plt

(...)

net = BarabasiAlbert()

degrees = net.degree()
Pk = np.array(list(Counter(degrees.values()).items()))

plt.loglog(Pk.T[0], Pk.T[1], 'b*')
plt.xlabel('k')
plt.ylabel('P[k]')
plt.savefig('Pk.png')
plt.close()

print("Number of nodes:", net.number_of_nodes())
print("Number of edges:", net.number_of_edges())
```

Snowball Sampling

- Commonly used in Social Science and Computer Science
 1. Start with a single node (or small number of nodes)
 2. Get “friends” list
 3. For each friend get the “friend” list
 4. Repeat for a fixed number of layers or until enough users have been connected
- Generates a connected component from each seed
- Quickly generates a *lot* of data/API calls

Snowball Sampling

```
import networkx as NX

def snowball(net, seed, max_depth = 3, maxnodes=1000):
    seen = set()
    queue = set()

    queue.add(seed)
    queue2 = set()

    for _ in range(max_depth+1):
        while queue:
            user_id = queue.pop()
            seen.add(user_id)

            NN = net.neighbors(user_id)

            for node in NN:
                if node not in seen:
                    queue2.add(node)

            queue.update(queue2)
            queue2 = set()

    return seen

net = NX.connected_watts_strogatz_graph(10000, 4, 0.01)
neve = snowball(net, 0)

print(neve)
```

Streaming Geocoded data

<https://dev.twitter.com/streaming/overview/request-parameters#locations>

- The Streaming api provides realtime data, subject to filters
- Use `TwitterStream` instead of `Twitter` object (`.TwitterStream(auth=twitter_api.auth)`)
- `.status.filter(track=q)` will return tweets that match the query `q` in real time
- Returns generator that you can iterate over
- `.status.filter(locations=bb)` will return tweets that occur within the bounding box `bb` in real time
- `bb` is a comma separated pair of lon/lat coordinates.
 - -180,-90,180,90 - World
 - -74,40,-73,41 - NYC

Streaming Geocoded data

<https://dev.twitter.com/streaming/overview/request-parameters#locations>

```
import twitter
from twitter_accounts import accounts
import sys
import gzip

app = accounts["social"]

auth = twitter.oauth.OAuth(app["token"],
                             app["token_secret"],
                             app["api_key"],
                             app["api_secret"])

stream_api = twitter.TwitterStream(auth=auth)

query = "-74,40,-73,41" # NYC
stream_results = stream_api.statuses.filter(locations=query)
tweet_count = 0

fp = gzip.open("NYC.json.gz", "a")

for tweet in stream_results:
    try:
        tweet_count += 1
        print(tweet_count, tweet["id"])
        print(tweet, file=fp)
    except:
        pass

if tweet_count % 10000 == 0:
    print(tweet_count, file=sys.stderr)
    break
```

shapefiles

<http://www.esri.com/library/whitepapers/pdfs/shapefile.pdf>

- Open specification developed by ESRI, still the current leader in commercial GIS software
- shapefiles aren't actual (individual) files...
- but actually a set of files sharing the same name but with different extensions:

```
(py35) (master) bgoncalves@underdark:$ls -l
total 4856
-rw-r--r--@ 1 bgoncalves staff      537 Apr 17 12:40 nybb.dbf
-rw-r--r--@ 1 bgoncalves staff      562 Apr 17 12:40 nybb.prj
-rw-r--r--@ 1 bgoncalves staff 1217376 Apr 17 12:40 nybb.shp
-rw-r--r--@ 1 bgoncalves staff    12905 Apr 17 12:40 nybb.shp.xml
-rw-r--r--@ 1 bgoncalves staff      140 Apr 17 12:40 nybb.shx
-rw-r--r--  1 bgoncalves staff      536 Apr 17 12:40 nybb_wgs84.dbf
-rw-r--r--  1 bgoncalves staff      143 Apr 17 12:40 nybb_wgs84.prj
-rw-r--r--  1 bgoncalves staff      257 Apr 17 12:40 nybb_wgs84.qpj
-rw-r--r--  1 bgoncalves staff 1217376 Apr 17 12:40 nybb_wgs84.shp
-rw-r--r--  1 bgoncalves staff      140 Apr 17 12:40 nybb_wgs84.shx
(py35) (master) bgoncalves@underdark:$
```

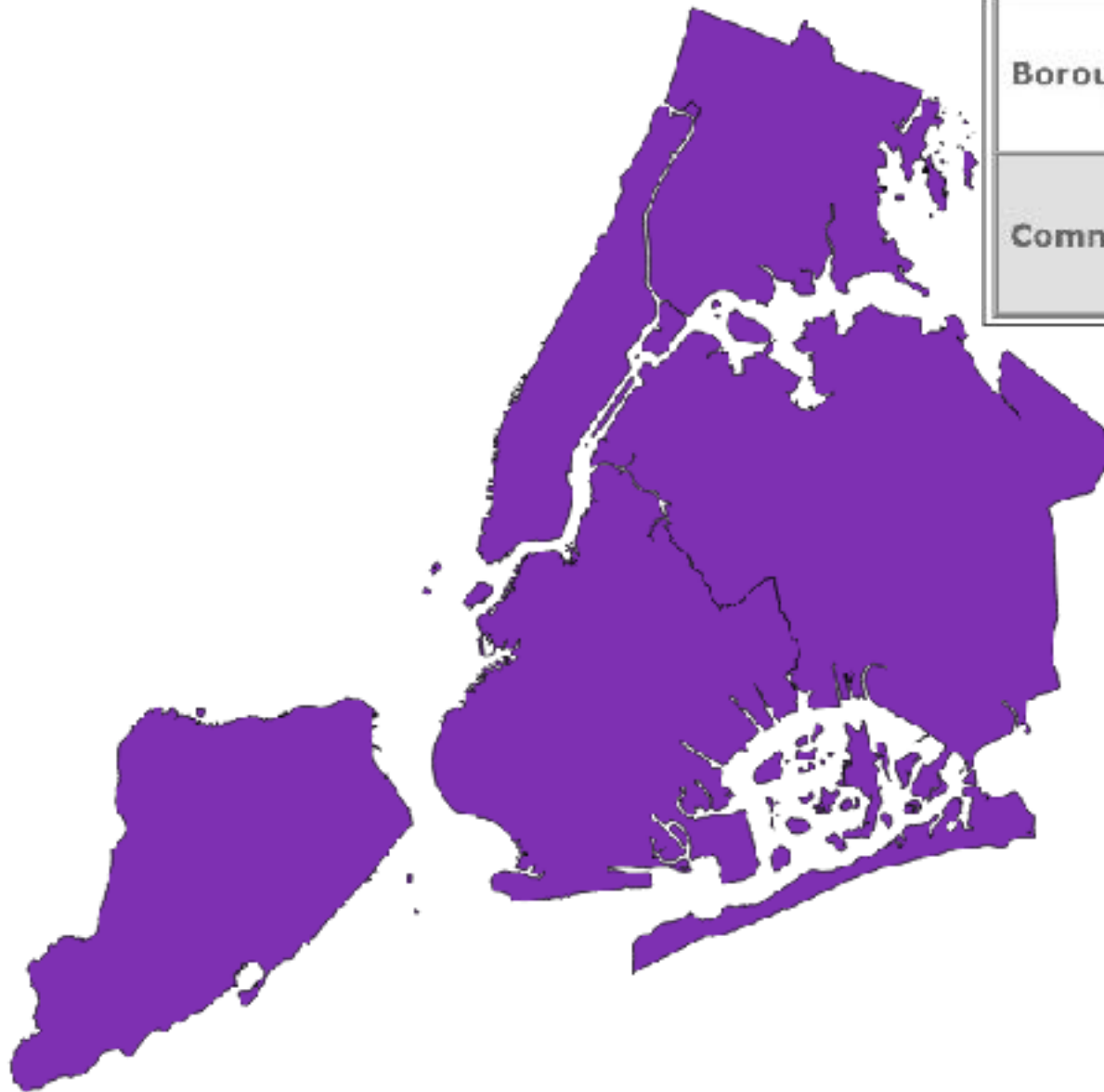
- the actual set of files changes depending on the contents, but three files are usually present:
 - **.shp** - also commonly referred to as “the” shapefile. Contains the geometric information
 - **.dbf** - a simple database containing the feature attribute table.
 - **.shx** - a spatial index, not strictly required

Shapefiles

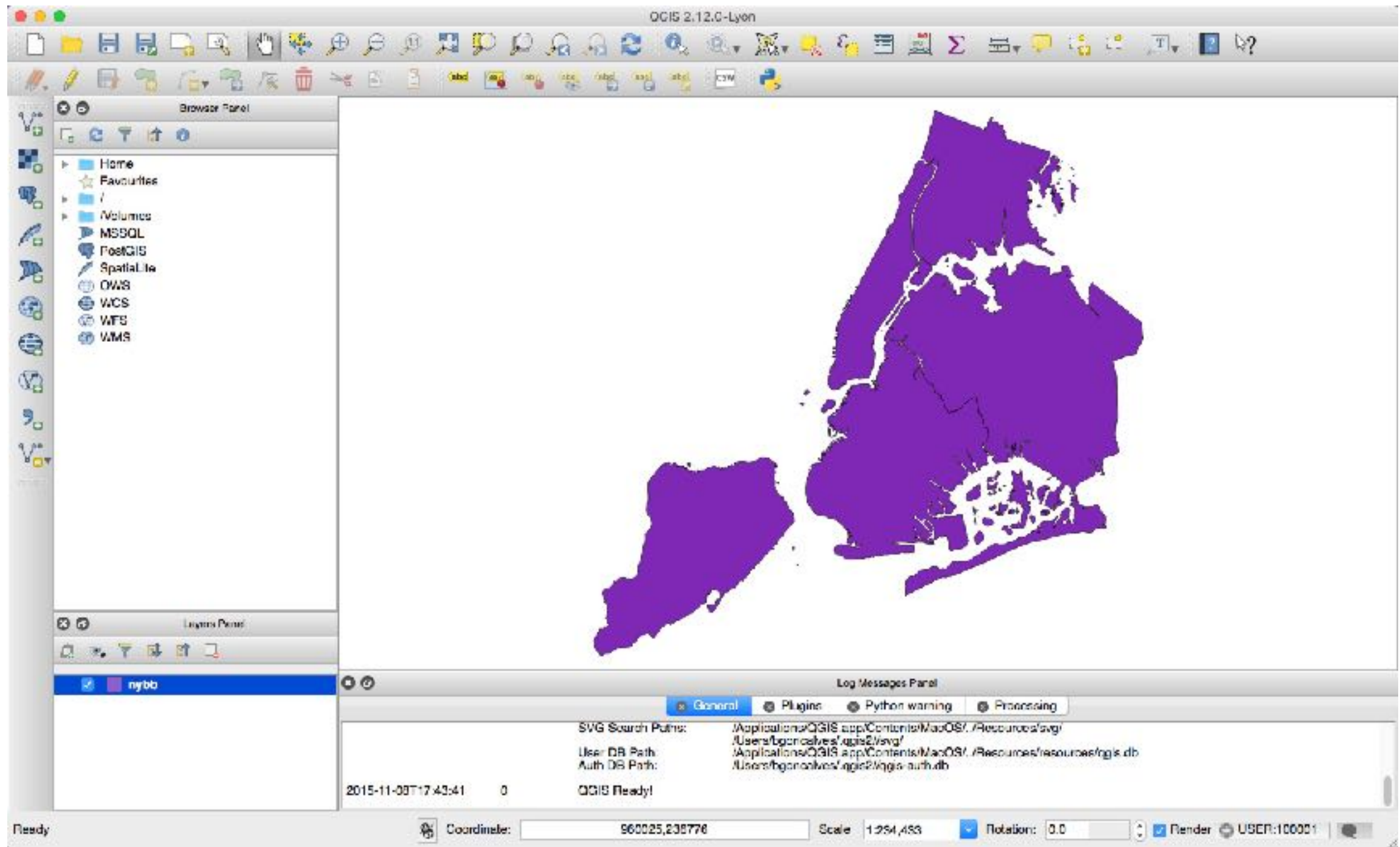
http://www.nyc.gov/html/dcp/html/bytes/districts_download_metadata.shtml#bcd

Borough Boundaries & Community Districts	Download	Metadata
Borough Boundaries (Clipped to Shoreline)	 (645k)	
Borough Boundaries (Water Areas Included)	 (31k)	
Community Districts (Clipped to Shoreline)	 (772k)	

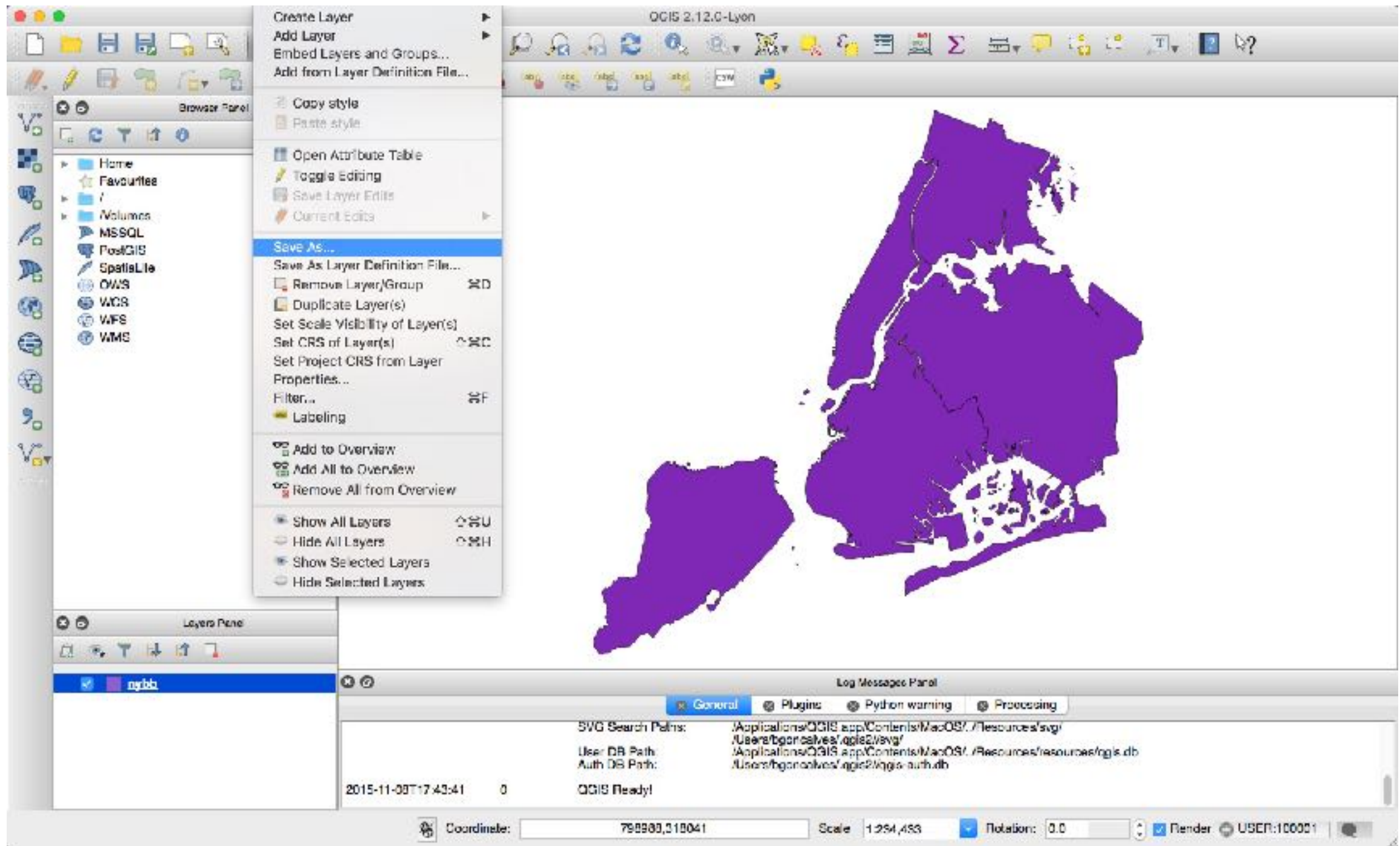
- Unfortunately it doesn't use the right reference system (**WGS84**), so we must convert it.



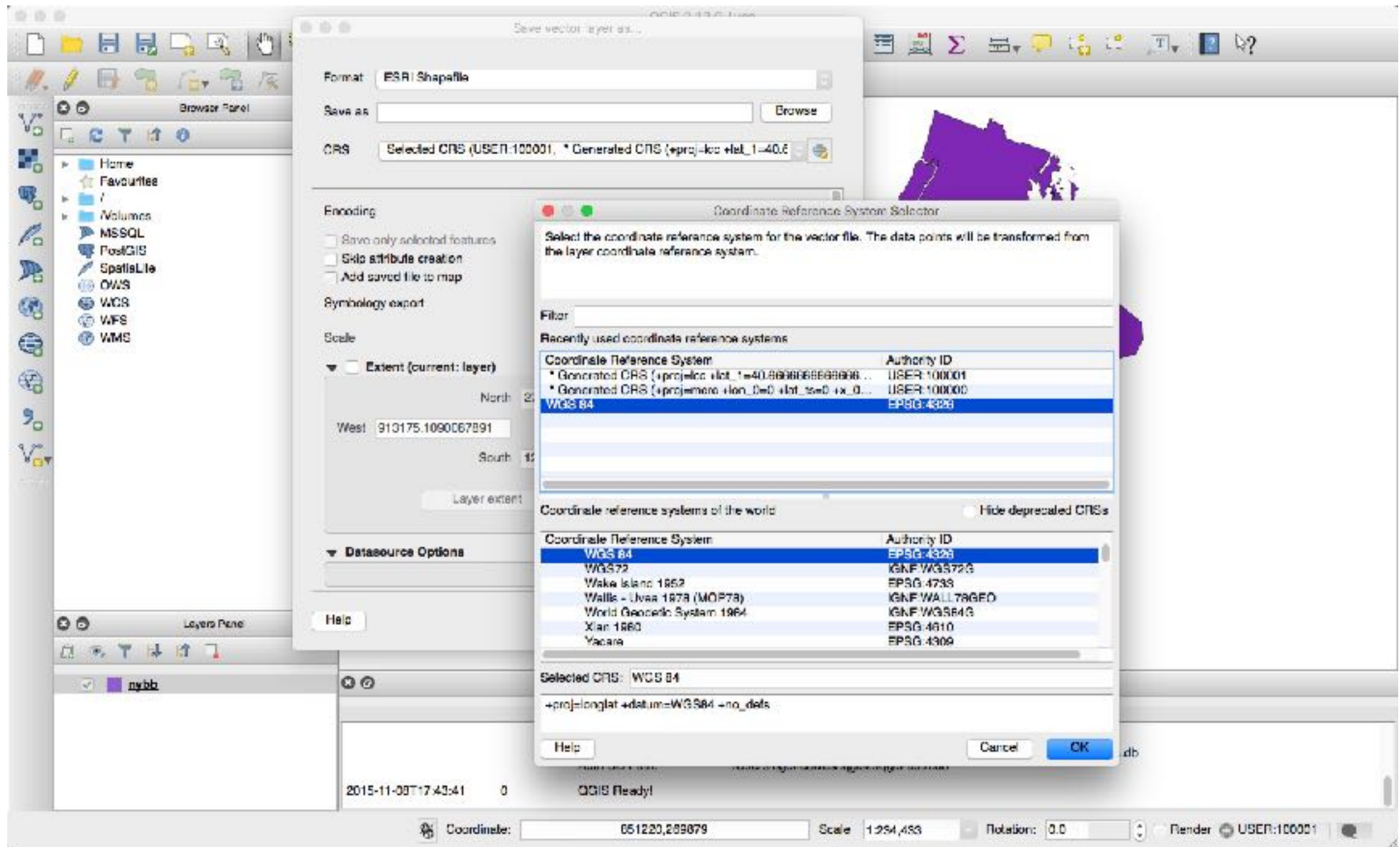
Shapefiles



Shapefiles



Shapefiles



- **pyshp** defines utility functions to load and manipulate Shapefiles programmatically.
- The **shapefile** module handles the most common operations:
 - **.Reader(filename)** - Returns a **Reader** object
- **Reader.records()/Reader.iterRecords()** returns/iterates over the different records present in the shapefile
- **Reader.shapes()/Reader.iterShapes()** - returns/iterates over the different shapes present in the shapefile
- **Reader.shapeRecords()/Reader.iterShapeRecords()** returns/iterates over both shapes and records present in the shapefile
- **Reader.record(index)/Reader.shape(index)/Reader.shapeRecord(index)** - return the record/shape/shapeRecord at index position **index**
- **Reader.numRecords** - returns the number of records in the shapefile

```
import sys
import shapefile

shp = shapefile.Reader('geofiles/nybb_15c/nybb_wgs84.shp')

print("Found", shp.numRecords, "records:")

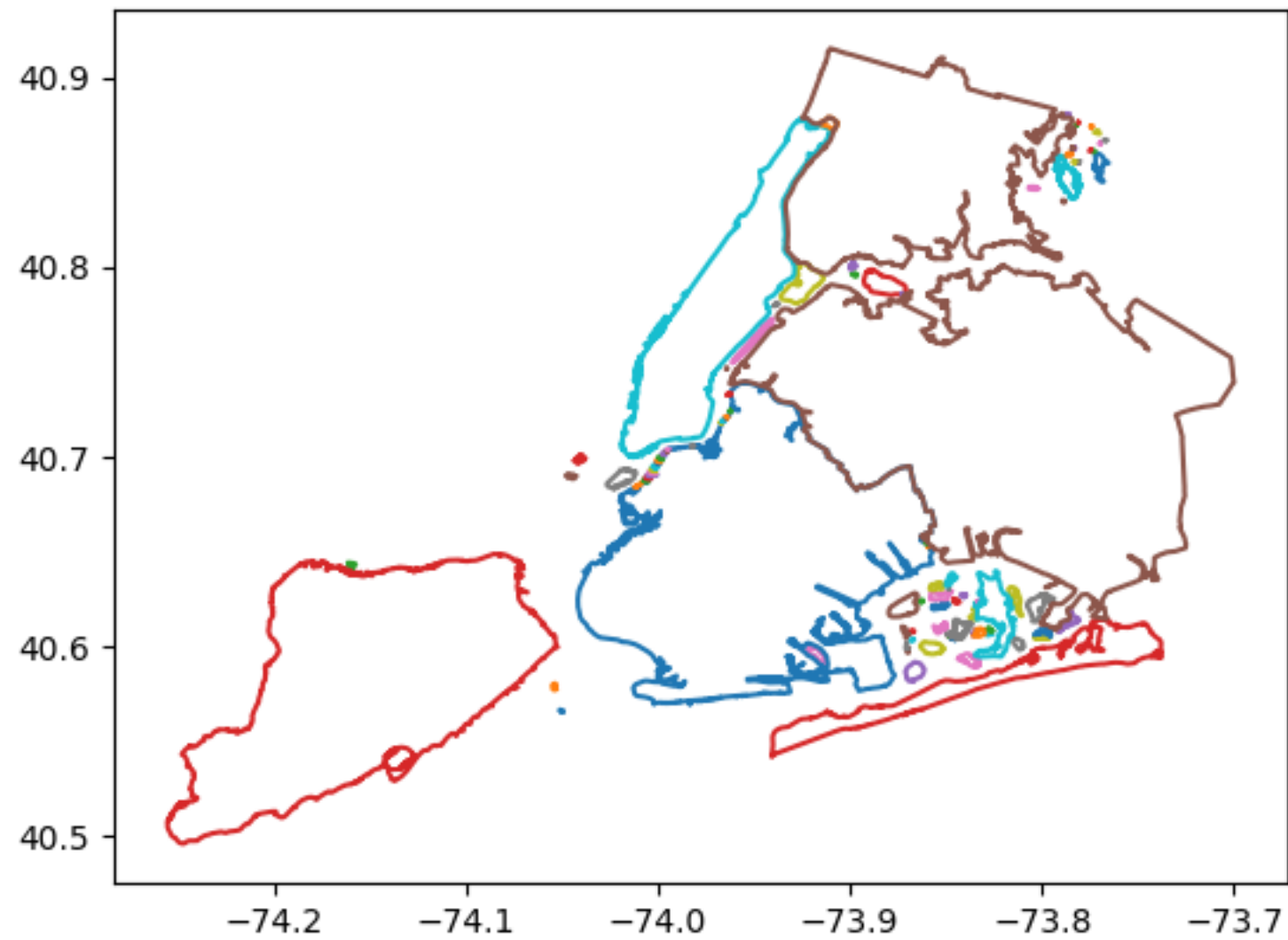
recordDict = dict(zip([record[1] for record in shp.iterRecords()], range(shp.numRecords)))

for record, id in recordDict.items():
    print(id, record)
```


- **shape** objects contain several fields:
 - **bbox** - lower left and upper right **x,y** coordinates (long/lat) - **optional**
 - **parts** - list of indexes for the first point of each of the parts making up the shape.
 - **points** - **x,y** coordinates for each point in the shape.
 - **shapeType** - integer representing the shape type - all shapes in a shapefile are required to be of the same **shapeType** or **null**.

Value	Shape Type
0	Null Shape
1	Point
3	PolyLine
5	Polygon
8	MultiPoint
11	PointZ
13	PolyLineZ
15	PolygonZ
18	MultiPointZ
21	PointM
23	PolyLineM
25	PolygonM
28	MultiPointM
31	MultiPatch

Simple shapefile plot



shapely

<http://toblerity.org/shapely/manual.html>

- Shapely defines geometric objects under `shapely.geometry`:
 - Point
 - Polygon
 - MultiPolygon
- `shape()` Convenience function that creates the appropriate geometric object
- and common operations
 - `.crosses(shape)` - if it partially overlaps `shape`
 - `.contains(shape)` - whether it contains or not the object `shape`
 - `.within(shape)` - whether it is contained by object `shape`
 - `.touches(shape)` - if the boundaries of this object touch `shape`

shapely

<http://toblerity.org/shapely/manual.html>

- **shape** objects provide useful fields to query a shapes properties:
 - **.centroid** - The centroid ("center of mass") of the object
 - **.area** - returns the area of the object
 - **.bounds** - the MBR of the shape in (minx, miny, maxx, maxy) format
 - **.length** - the length of the shape
 - **.geom_type** - the Geometry Type of the object
- **shapely.shape** is also able to easily load **pyshp**'s shape objects to allow for further manipulations.

```
import sys
import shapefile
from shapely.geometry import shape

shp = shapefile.Reader('geofiles/nybb_15c/nybb_wgs84.shp')

recordDict = dict(zip([record[1] for record in shp.iterRecords()], range(shp.numRecords)))

manhattan = shape(shp.shape(recordDict["Manhattan"]))

print("Centroid:", manhattan.centroid)
print("Bounding box:", manhattan.bounds)
print("Geometry type:", manhattan.geom_type)
print("Length:", manhattan.length)
```

Filter points within a Shapefile

```
import sys
import shapefile
from shapely.geometry import shape, Point
import gzip

shp = shapefile.Reader('geofiles/nybb_15c/nybb_wgs84.shp')

recordDict = dict(zip([record[1] for record in shp.iterRecords()], range(shp.numRecords)))

manhattan = shape(shp.shape(recordDict["Manhattan"]))
fp = gzip.open("Manhattan.json.gz", "w")

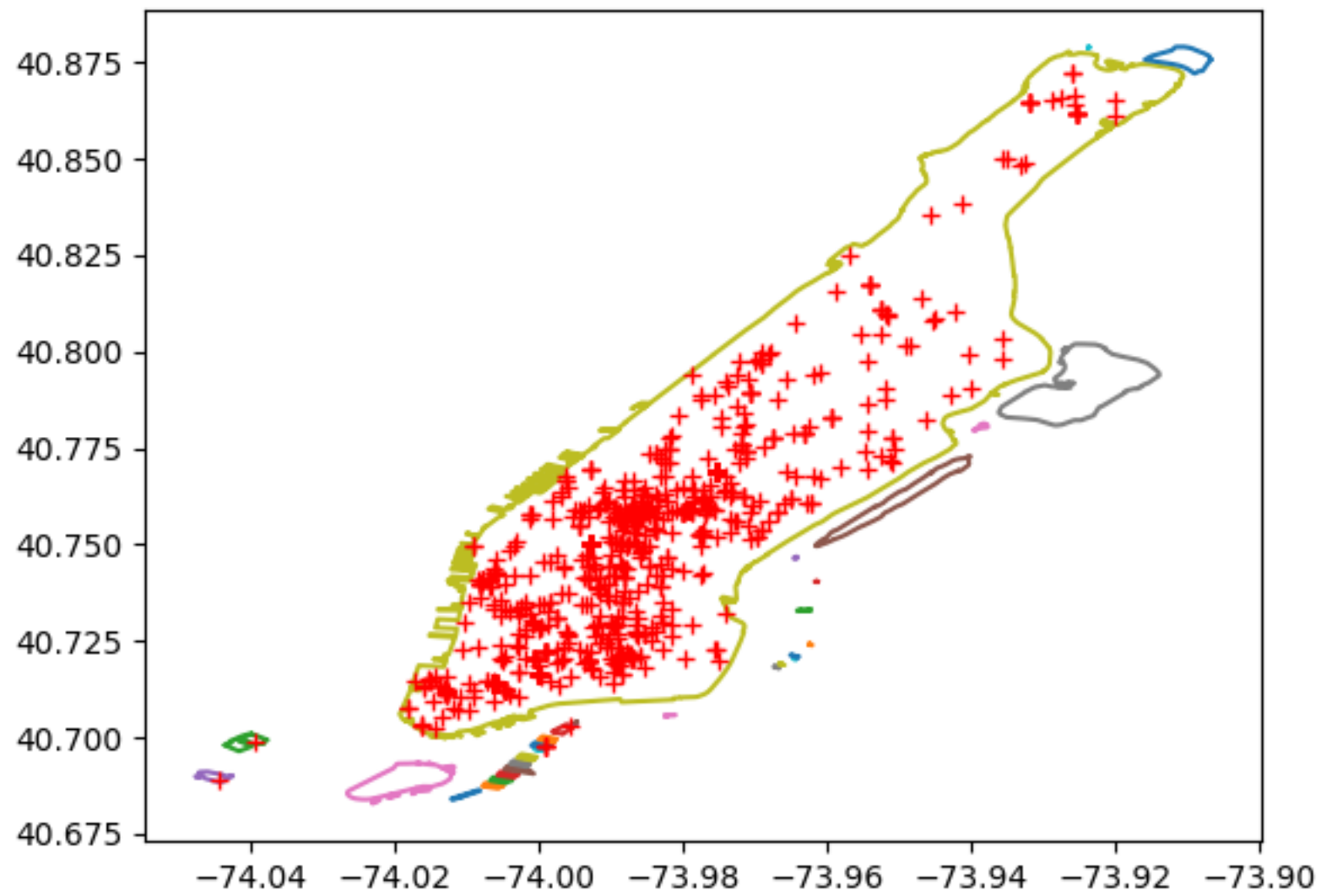
for line in gzip.open("NYC.json.gz"):
    try:
        tweet = eval(line.strip())

        if "coordinates" in tweet and tweet["coordinates"] is not None:
            point = Point(tweet["coordinates"]["coordinates"])

            if manhattan.contains(point):
                fp.write(line)
    except:
        pass

fp.close()
```

Filter points within a Shapefile





Twitter places

- As we saw before, Twitter defines a “**coordinates**” field in tweets
- There is also a “**place**” field that we glossed over.
- The **place** object contains also geographical information, but at a coarser resolution than the **coordinates** field.
- Each place has a unique **place_id**, a **bounding_box** and some geographical information, such as **country** and **full_name**:

```
{ 'attributes': {},  
  'bounding_box': { 'coordinates': [[[-74.041878, 40.570842],  
    [-74.041878, 40.739434],  
    [-73.855673, 40.739434],  
    [-73.855673, 40.570842]]],  
    'type': 'Polygon'},  
  'country': 'United States',  
  'country_code': 'US',  
  'full_name': 'Brooklyn, NY',  
  'id': '011add077f4d2da3',  
  'name': 'Brooklyn',  
  'place_type': 'city',  
  'url': 'https://api.twitter.com/1.1/geo/id/011add077f4d2da3.json' }
```

The bounding_box field is GeoJSON formatted and compatible with `pyshp.shape`

- places can be of several different types: **'admin'**, **'city'**, **'neighborhood'**, **'poi'**

Twitter places

<https://dev.twitter.com/overview/api/places>

Place Attributes

Place Attributes are metadata about places. An attribute is a key-value pair of arbitrary strings, but with some conventions.

Below are a number of well-known place attributes which may, or may not exist in the returned data. These attributes are provided when the place was created in the Twitter places database.

Key	Description
street_address	
locality	the city the place is in
region	the administrative region the place is in
iso3	the country code
postal_code	in the preferred local format for the place
phone	in the preferred local format for the place, include long distance code
twitter	twitter screen-name, without @
url	official/canonical URL for place
app:id	An ID or comma separated list of IDs representing the place in the applications place database.

Keys can be no longer than 140 characters in length. Values are unicode strings and are restricted to 2000 characters.

Filter points and places

```
import sys
import shapefile
from shapely.geometry import shape, Point
import gzip

shp = shapefile.Reader('geofiles/nybb_15c/nybb_wgs84.shp')

recordDict = dict(zip([record[1] for record in shp.iterRecords()], range(shp.numRecords)))

manhattan = shape(shp.shape(recordDict["Manhattan"]))
fp = gzip.open("Manhattan_places.json.gz", "w")

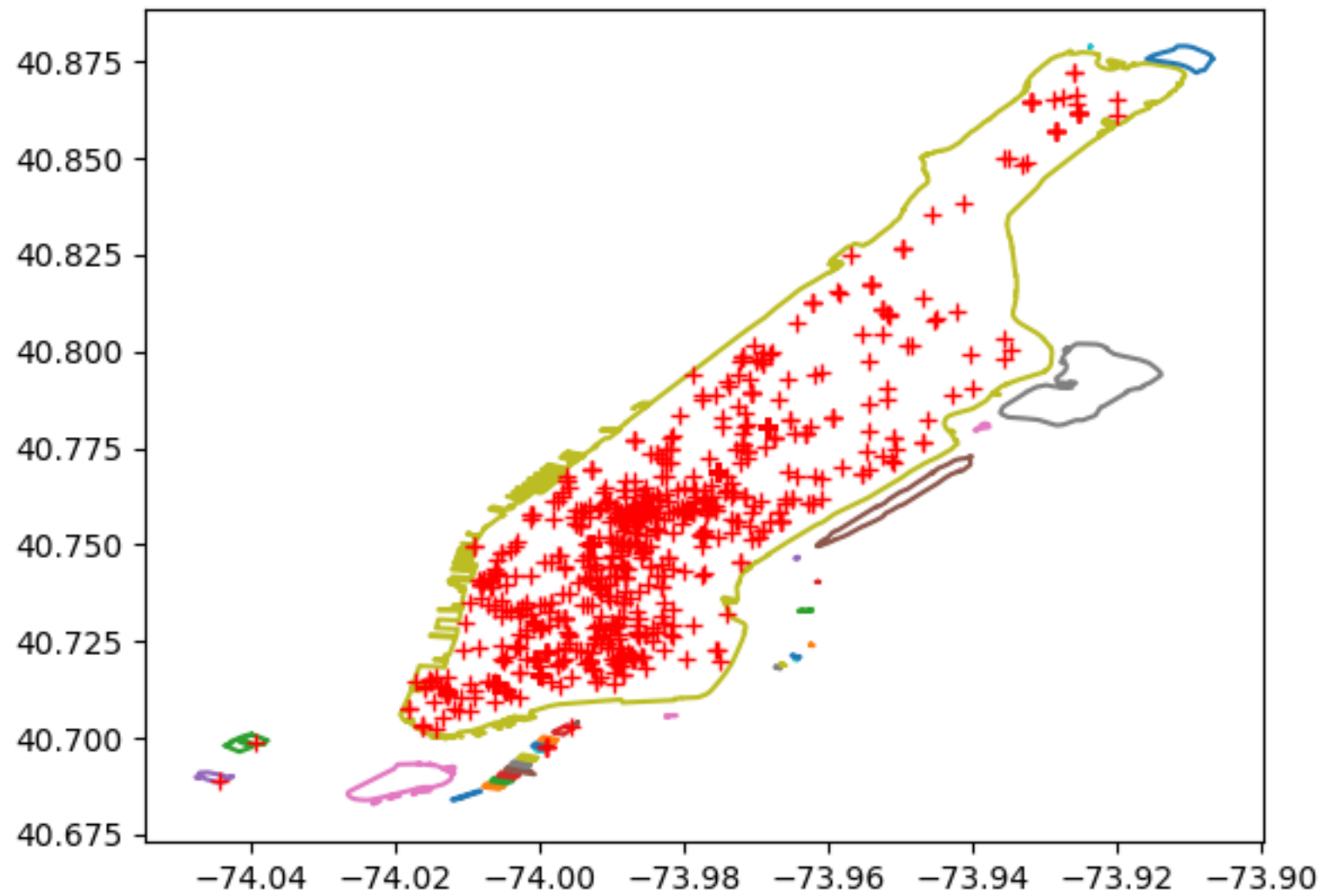
for line in gzip.open("NYC.json.gz"):
    try:
        tweet = eval(line.strip())
        point = None

        if "coordinates" in tweet and tweet["coordinates"] is not None:
            point = Point(tweet["coordinates"]["coordinates"])
        else:
            if "place" in tweet and tweet["place"]["bounding_box"] is not None:
                bbox = shape(tweet["place"]["bounding_box"])
                point = bbox.centroid

        if point is not None and manhattan.contains(point):
            fp.write(line)
    except:
        pass

fp.close()
```

Filter points and places



Filter

```
import sys
import gzip
import numpy as np
import shapefile
from shapely.geometry import shape, Point
import matplotlib.pyplot as plt

shp = shapefile.Reader('geofiles/nybb_15c/nybb_wgs84.shp')
recordDict = dict(zip([record[1] for record in shp.iterRecords()],
range(shp.numRecords)))

manhattan = shp.shape(recordDict["Manhattan"])

points = np.array(manhattan.points)
parts = manhattan.parts
parts.append(len(manhattan.points))

for i in range(len(parts)-1):
    plt.plot(points.T[0][parts[i]:parts[i+1]], points.T[1][parts[i]:parts[i+1]])

points_X = []
points_Y = []

manhattan = shape(shp.shape(recordDict["Manhattan"]))

for line in gzip.open(sys.argv[1]):
    try:
        tweet = eval(line.strip())
        point = None

        if "coordinates" in tweet and tweet["coordinates"] is not None:
            point = Point(tweet["coordinates"]["coordinates"])
        else:
            if "place" in tweet and tweet["place"]["bounding_box"] is not None:
                bbox = shape(tweet["place"]["bounding_box"])
                point = bbox.centroid

        if point is not None and manhattan.contains(point):
            points_X.append(point.x)
            points_Y.append(point.y)
    except:
        pass

plt.plot(points_X, points_Y, 'r+')
plt.savefig(sys.argv[1] + '.png')
```

@bgonca

plot_shapefile_points.py

Aggregation

```
import sys
import numpy as np
import shapefile
from shapely.geometry import shape, Point
import matplotlib.pyplot as plt
import gzip

def map_points(xllcorner, yllcorner, cellsize, nrows, x, y):
    x = int((x-xllcorner)/cellsize)
    y = (nrows-1)-int((y-yllcorner)/cellsize)

    return x, y

def save_asc(data, xllcorner, yllcorner, cellsize, filename):
    fp = open(filename, "w")

    nrows, ncols = data.shape

    print("ncols", ncols, file=fp)
    print("nrows", nrows, file=fp)
    print("xllcorner", xllcorner, file=fp)
    print("yllcorner", yllcorner, file=fp)
    print("cellsize", cellsize, file=fp)
    print("NODATA_value", "-9999", file=fp)

    for i in range(nrows):
        for j in range(ncols):
            print((" %u " % data[i, j]), end="", file=fp)

        print("\n", end="", file=fp)

    fp.close()
```

```

shp = shapefile.Reader('geofiles/nybb_15c/nybb_wgs84.shp')
recordDict = dict(zip([record[1] for record in shp.iterRecords()], range(shp.numRecords)))
manhattan = shape(shp.shape(recordDict["Manhattan"]))

xllcorner, yllcorner, xurcorner, yurcorner = manhattan.bounds
cellsize = 0.01

ncols = int((xurcorner-xllcorner)/cellsize)
nrows = int((yurcorner-yllcorner)/cellsize)

data = np.zeros((nrows, ncols), dtype='int')

for line in gzip.open("NYC.json.gz"):
    try:
        tweet = eval(line.strip())
        point = None

        if "coordinates" in tweet and tweet["coordinates"] is not None:
            point = Point(tweet["coordinates"]["coordinates"])
        else:
            if "place" in tweet and tweet["place"]["bounding_box"] is not None:
                bbox = shape(tweet["place"]["bounding_box"])
                point = bbox.centroid

        if point is not None and manhattan.contains(point):
            coord_x, coord_y = map_points(xllcorner, yllcorner, cellsize, nrows, point.x, point.y)
            data[coord_y, coord_x] += 1

    except:
        pass

save_asc(data, xllcorner, yllcorner, cellsize, "Manhattan.asc")

plt.imshow(np.log(data+1))
plt.colorbar()
plt.savefig('Manhattan_cells.png')

```

Aggregation

