

# P3: Wrangle OpenStreetMap Data

Edmund Wong

Map Area: Syracuse, NY, United States

## 1. Problems encountered in map

I surveyed the data using Python's ElementTree module. After parsing the entire Syracuse OSM file, I prepared the list of node tags, way tags, and relation tags from the XML data.

```
1. import xml.etree.ElementTree as ET
2. fname = 'syracuse_new-york.osm'
3. tree = ET.parse(fname)
4.
5. nodetags = tree.findall('node/tag')
6. waytags = tree.findall('way/tag')
7. relationtags = tree.findall('relation/tag')
```

To audit this data, I quickly inspected the data to see what kind of common problems there are in fields I suspected to likely have problems. Most phone numbers are in the "XXX XXX XXXX" format. However, I found some phone numbers that were inconsistent to that format. Some of those problematic phone numbers contained dashes and some began with "+1" dialing code. Here is a query I used to detect and clean/fix problematic phone numbers that do not agree with this format. It requires the python phonenumbers module.

```
1. for idx, tag in enumerate(nodetags):
2.     if tag.attrib['k'] == 'phone':
3.         #Get value of phone number
4.         pn = tag.attrib['v']
5.         #Check whether phone number matches XXX XXX XXXX format
6.         if not pn[0:3].isdigit() or pn[3] != ' ' or not pn[4:7].isdigit() or pn[
7.         7] != ' ' or not pn[8:12].isdigit() or len(pn) > 12:
8.             try:
9.                 phone_obj=phonenumbers.parse(pn, 'US')
10.                #replace phone number with correct format
11.                nodetags[idx].attrib['v'] = phonenumbers.format_number(phone_obj
12.                , phonenumbers.PhoneNumberFormat.mro).replace('-', ' ')
13.                print "Fixed:", pn, nodetags[idx].attrib['v']
14.            except NumberParseException:
15.                print pn, "Not valid phone number"
```

Another problem is that the website URLs of organizations and businesses in Syracuse are not standardized. Some top-level-domain URLs have missing trail slashes. Here is an example where I checked and repaired for missing trailing slashes in the website URLs of. Website data exists in nodes, ways, and relations (the example below is waytag specific).

```

1. for idx, tag in enumerate(waytags):
2.     #only top-level domains
3.     if tag.attrib['v'].endswith(('com', '.org', '.edu', '.gov', '.mil', '.net',
        '.int')):
4.         if tag.attrib['k'] == 'website' and not tag.attrib['v'].endswith('/'):
5.             waytags[idx].attrib['v'] = tag.attrib['v'] + '/'

```

## 2. Overview of the Data

Before getting into the overview of the data, the XML data was converted to JSON. Additionally, the data was restructured in a fashion where fields that should belong together are grouped into the same object. For example, here I grouped the “timestamp”, “version”, “uid”, and “user” fields under “created”. It also makes sense that I grouped address fields together. This format was based on James Kao. (<https://jameskao.me>)

An original OSM Node in XML:

```

1. <node id="1694014652" lat="43.0525611" lon="-
    76.1546779" version="6" timestamp="2015-01-
    12T02:07:11Z" changeset="28077497" uid="599436" user="zeromap">
2.   <tag k="name" v="Dinosaur BBQ"/>
3.   <tag k="phone" v="315 476 4937"/>
4.   <tag k="amenity" v="restaurant"/>
5.   <tag k="cuisine" v="american"/>
6.   <tag k="website" v="http://www.dinosaurbarbque.com/bbq-syracuse"/>
7.   <tag k="addr:city" v="Syracuse"/>
8.   <tag k="addr:state" v="NY"/>
9.   <tag k="addr:street" v="West Willow Street"/>
10.  <tag k="addr:country" v="US"/>
11.  <tag k="addr:postcode" v="13202"/>
12.  <tag k="opening_hours" v="Mo-We 11:00-23:00; Th 11:00-24:00; Fr-Sa 11:00-
    1:00; Su 12:00-22:00"/>
13.  <tag k="addr:housenumber" v="246"/>
14. </node>

```

The same node formatted and restructured in JSON:

```

1. {
2.   "website": "http://www.dinosaurbarbque.com/bbq-syracuse",
3.   "cuisine": "american",
4.   "amenity": "restaurant",
5.   "name": "Dinosaur BBQ",
6.   "created": {
7.     "changeset": "28077497",
8.     "user": "zeromap",
9.     "version": "6",
10.    "uid": "599436",
11.    "timestamp": "2015-01-12T02:07:11Z"
12.  },
13.  "opening_hours": "Mo-We 11:00-23:00; Th 11:00-24:00; Fr-Sa 11:00-
    1:00; Su 12:00-22:00",
14.  "pos": [

```

```

15.         43.0525611,
16.         -76.1546779
17.     ],
18.     "phone": "315 476 4937",
19.     "address": {
20.         "city": "Syracuse",
21.         "country": "US",
22.         "state": "NY",
23.         "street": "West Willow Street",
24.         "postcode": "13202",
25.         "housenumber": "246"
26.     },
27.     "type": "node",
28.     "id": "1694014652"
29. }

```

Once the data was imported into mongodb using mongoimport, I performed a couple of queries that will provide a statistical overview of the Syracuse dataset.

*Size of XML:* 64.2 MB

*Size of JSON (after restructuring):* 166.2 MB

*Number of nodes:*

```

> db.syracuse.find({"type":"node"}).count()
279521

```

*Number of ways:*

```

> db.syracuse.find({"type":"way"}).count()
35510

```

*Number of unique users:*

```

> db.syracuse.distinct("created.user").length
234

```

*Number of restaurants:*

```

> db.syracuse.find({"type":"node","amenity":"restaurant"}).count()
77

```

*Number of unique contributing users:*

```

>db.syracuse.distinct("created.user").length
234

```

*Kinds of amenities:*

```

>db.syracuse.distinct("amenity")
[ "parking_entrance", "school", place_of_worship", ... , "driving_range"]

```

*Number of restaurants that are classified by a cuisine type:*

```
> db.syracuse.find({"amenity":"restaurant", "cuisine":{"$exists:1}}).count()  
89
```

### 3. Other ideas about the datasets

#### **Inclusion of elevation data:**

From glancing at the fields of a node, there is missing elevation information that can potentially extend OpenStreetMap to allow bikers and walkers to figure out the difficulty of the commute. Elevation data in node points can also create surface data on the map to visualize hills and other topographical details. In MongoDB, the benefit of its schemaless nature allows us insert and update new kinds of data (such as elevation) into the existing collection easily. Also elevation data can be extended to Ways (Ways that don't share a start and end Node such as highways, railways, and roads) in such a way where overall elevation gain or loss can be calculated.

Potential problems are that elevation data can be problematic to standardize, different units of measurement. It is popularly measured in both feet and meters, so some standard must first be established before incorporating elevation data into OSM. Also, accurate elevation data is not as easily acquirable as the other existing fields and may require expertise of those in the geographical sciences to validate the accuracy of such data.

*Additional exploratory analyses:*

In addition, I have performed other additional exploratory analyses to figure out the count of how often elements have been updated.

```
> db.syracuse.aggregate([{"$match":{"created.version":{"$exists:1}}},  
{"$group":{"_id":"$created.version", "count":{"$sum":1}}}, {"$sort":{"count":-  
1}}])
```

```
{ "_id" : "1", "count" : 202433 }  
{ "_id" : "2", "count" : 93249 }  
{ "_id" : "3", "count" : 12521 }  
{ "_id" : "4", "count" : 3495 }  
{ "_id" : "5", "count" : 1484 }  
{ "_id" : "6", "count" : 703 }  
{ "_id" : "7", "count" : 409 }  
{ "_id" : "8", "count" : 278 }  
{ "_id" : "9", "count" : 186 }  
{ "_id" : "10", "count" : 91 }  
{ "_id" : "11", "count" : 47 }
```

```
{ "_id" : "12", "count" : 33 }  
{ "_id" : "13", "count" : 26 }  
{ "_id" : "14", "count" : 25 }  
{ "_id" : "16", "count" : 13 }  
{ "_id" : "15", "count" : 9 }  
{ "_id" : "17", "count" : 8 }  
{ "_id" : "19", "count" : 8 }  
{ "_id" : "18", "count" : 7 }
```

This output shows the most of these elements (nodes, ways) have been updated only once, and almost less frequently for each additional version change.