

Open Street Map: A case study on the data

Author: Sagarnil Das

Date: 06/16/2017

Map Area:

City: New York

Link to Dataset: [MapZen New York Dataset](#)

Though I am from India, I have lived a good years of my adult working life when I used to work for Department of Health. For some reason, I fell in love with this place. That's why I decided to work with the dataset of New York and wrangle and clean the data wherever necessary.

Dataset Information:

The original dataset was a huge file of 2.7 GB. So I used the `sorten_osm.py` to iterate through the file and write every 10th top level element. The resulting shorter file was 270 MB which I worked with. The code was provided in the instructors note.

Steps in wrangling, cleaning and auditing the data:

1. **Initial Data Exploration:** Checked in the osm xml file how many unique tags are there to get a feeling of the data I am about to work with.
2. **Check for potential keys:** The 2nd step was to check for potential Keys inside the xml file as these keys would be required in future when I insert the clean data inside MongoDB. Inside the tags, there was an attribute called "k" and another called "v". They acted as a key value pair inside the xml file. So I used regex module to check for patterns inside the keys and look for any problematic characters.
3. **Some more Data Exploration:** I did a little more data exploration again to get a good feel of the data. I coded to find out how many unique users contributed in developing the map in this particular area.

Problems encountered in my Map:

1. **Changing bad Street names:** It was noticed that in the many Street types had inconsistent names. A very high number of unique Street Types were observed (Street, Plaza, Court, Boulevard, Road etc). But the problem was the naming was inconsistent in many places. For example, in multiple places, instead of 'Street', it was written as 'St'. Same goes for 'Boulevard': 'Blvd', 'Parkway': 'Pkw' and 'Avenue': 'Ave'. The types were too large. This is how this problem was handled. I wrote a regex to extract the street type from the address.

```
street_type_re = re.compile(r'\b\s+\.?$', re.IG
```

Then I made a list of some possible street types like 'Street', 'Boulevard', 'Drive', 'Road', 'Court' etc. I also made a mapper dictionary for short abbreviations used in many places, to map them to their full form. This mapping was not built in one go. With an initial mapping, I first built the functions and saw the possible street types and then I modified this mapping to add whatever values I didn't add the first time. Apart from that, these were the functions I wrote for cleaning the street data.

- a) For the auditing purpose these three functions were written:
- b) For updating the name from the abbreviation to the full name, this function was written:

```
def update_name(name, mapping):  
    for search_error in mapping:  
        if search_error in name:  
            name = re.sub(r'\b' + search_error + r'\b\.', mapping[search
```

But now I faced some problems. First I ran the audit function on the OSM file, it worked. After that I tried updating the names. For the 2nd part of the code, that is updating the names, I got a strange error 'KeyError': 'Certain_value' not found. Upon further investigation, I discovered what the problem was. My list of possible city names was not complete. I just gave some initial names. But if while getting the street, it doesn't encounter the same value in the 'expected' list, python thinks that it does not exist. So I went through all the street types in the dictionary I created in python and added all the street types inside the list. Note: Many street names were not incorrect but they came up in my audit as they followed the same pattern. For example: 'East Parkway 4'. So to account for them and the rest, I added them to the final list.

2. **Changing Bad Postal codes:** On investigating the postal codes, my standard was the 5 digit postal code (e.g. 12345). But in many places, I saw there were many incorrect formats like 12345-1234, NY 12345 etc. So just like I cleaned the street names, here also I built a dictionary to replace any of the oddities with the correct 5 digit postal code.

Next Steps:

1. **Data wrangling and making the data shape correctly so that I can insert the data into MongoDB:**

The next task was to wrangle and transform the data into JSON format. All the code is given in project_3_file_1.py. After wrangling and transforming the data, I wrote the data into a JSON file which I will now import into MongoDB.

2. **Data import to MongoDB:**

```
mongoimport --db osm --collection newyorkosm --type json --file sample3.json
```

Data Overview:

1. **File Size:**

New-york_new-york.osm – original OSM file (Size = 2.7 GB)

Sample1.osm – Shortened OSM file (Size = 278 MB)

Sample3.json – Final cleaned JSON file (Size = 514 MB)

2. **MongoDB overview:**

a) Number of documents: 1330111 [db.newyorkosm.find().count()]

b) Number of Nodes: 1149904 [db.newyorkosm.find({"type":"node"}).count()]

c) Number of Ways: 180177 [db.newyorkosm.find({"type":"way"}).count()]

d) Number of Uniques users: 2615 [print len(db.newyorkosm.distinct("created.user"))]

Other ideas about the dataset:

1. Top 10 contributing users for this particular map:

```
pipeline = [{"$group":{"_id":"$created.user",
                        "count":{"$sum":1}}},
             {"$sort":{"count":-1}},
             {"$limit":10}]
```

```
result = db.newyorkosm.aggregate(pipeline)
```

```
for a in result:
    u{'u_id': u'Rub21_nycbuildings', u'count': 488814}
    u{'u_id': u'ingalls_nycbuildings', u'count': 93572}
    u{'u_id': u'MySuffolkNY', u'count': 62646}
    u{'u_id': u'woodpeck_fixbot', u'count': 61621}
    u{'u_id': u'SuffolkNY', u'count': 58062}
```

2. Proportion of top users contribution:

```
pipeline = [{"$group":{"_id": "$created.user",
                        "count": {"$sum": 1}}},
             {"$project": {"proportion": {"$divide": ["$count", db.newyorkosm.find(
                {"_id": "$_id"}, {"$count": 1}).count()]}},
             {"$sort": {"proportion": -1}},
             {"$limit": 3}]
```

```
result = db.newyorkosm.aggregate(pipeline)
```

```
{u_id': u'Rub21_nycbuildings', u'proportion': 0.36749865236811063}
{u_id': u'ingalls_nycbuildings', u'proportion': 0.07034901598438024}
{u_id': u'MySuffolkNY', u'proportion': 0.047098324876645635}
{u_id': u'woodpeck_fixbot', u'proportion': 0.0463277124991824}
```

3. List of Universities:

```
pipeline = [{"$match":{"amenity":{"$exists":1}, "amenity": "university", "name":{"$regex": ".*university.*"}},
             {"$group":{"_id":"$name", "count":{"$sum":1}}},
             {"$sort":{"count":-1}}]
```

```
result = db.newyorkosm.aggregate(pipeline)
```

```
for a in result:
    u{'u_id': u'New Jersey Institute of Technology', u'count': 1}
    u{'u_id': u'Vaughn-Eames Hall', u'count': 1}
```

4. Most Popular Cuisines:

```
pipeline = [{"$match":{"amenity":{"$exists":1}, "amenity":"restaurant", "cuisine":1}},
             {"$group":{"_id":"$cuisine", "count":{"$sum":1}}},
             {"$sort":{"count":-1}},
             {"$limit":10}]
```

```
result = db.newyorkosm.aggregate(pipeline)
```

```
{u'_id': u'italian', u'count': 22}
```

```
{u'_id': u'american', u'count': 22}
```

5. Top 10 amenities:

```
pipeline = [{"$match":{"amenity":{"$exists":1}}},
             {"$group":{"_id":"$amenity", "count":{"$sum":1}}},
             {"$sort":{"count":-1}},
             {"$limit":10}]
```

```
result = db.newyorkosm.aggregate(pipeline)
```

```
for i in result:
```

6. Top 10 postal codes:

```
pipeline = [{"$match":{"address.postcode":{"$exists":1}}},
             {"$group":{"_id":"$address.postcode",
                         "count":{"$sum":1}}},
             {"$sort":{"count":-1}},
             {"$limit": 10}]
```

```
result = db.newyorkosm.aggregate(pipeline)
```

```
for i in result:
```

Checking for other potential problems:

1. Checking for Alphanumeric postal codes: (Result = 0)

```
result = db.newyorkosm.find({"address.postcode": {"$regex": "/^

for a in result:
    pprint.pprint(a)
```

2. Checking for postal codes that start with an uppercase letter followed by one or more lowercase letter: (Result = 0)

```
result = db.newyorkosm.find({"address.postcode": {"$regex": "/^[A-Z] [a-
💡
for a in result:
    pprint.pprint(a)
```

3. Checking for house numbers that has alphabetical characters: (Result = 0)

```
result = db.newyorkosm.find({"address.housenumber": {"$regex": "/[a-zA-

for a in result:
```

4. Checking for House characters that has spaces: (Result = 0)

```
result = db.newyorkosm.find({"address.housenumber": {"$regex"

for a in result:
    pprint.pprint(a)
```

So fortunately, all these additional pattern checks did not bring me any more erroneous values. So right now, the data which resides in my MongoDB is quite clean and has passed my auditing.

Additional Suggestions for Improving Data

From the above queries about user statistics, we see that it is a handful of users who are contributing to the whole map. There is a chance that a couple of them can be Bots also. So I think, if we have the option to check in while I am in a building, this map application can be developed much faster and in an efficient way. Also based on potential gamification approach like Badges, Leader boards and Trophies, an user can be more easily convinced to put his/her own inputs on the map thus increasing the chances of improved Map Datasets.

- **Benefits:**

1. Increase in rate of data thus completing many incomplete territories on the face of the earth improperly plotted in the OpenStreet Map.

2. Worldwide User Involvement.

- **Anticipated Problems:**

1. Chances of incorrect data entry increases as the number of contributor increases maybe even without any background knowledge of effective data auditing and wrangling.
2. Malicious code insertion in the map which can be a potential threat to the other users.

Conclusion:

So, during this whole process of gathering, extracting, cleaning and storing of our data, the biggest problem was fixing the street names with an appropriate street type and incorrect postal codes. The dataset was pretty big and the huge number of street types astounded me. I never would have thought about it. But by doing a cyclical auditing, I was finally able to fix all the names. I did some additional failure checks also all of whose results came back as negative. So even though in my opinion, there are still chances of improvement, I believe the data is sufficiently cleaned for the purpose of this project of Data Wrangling. The additional queries performed on the database documents gave consistent results which supports the theory of clean data. I loved doing this project for the city which was practically my home for 3 years.

References:

1. Maps: <https://www.openstreetmap.org/>
2. Maps: <https://mapzen.com/data/metro-extracts/>
3. Regexes: <https://github.com/SamMorrowDrums/Udacity-OpenStreetmap>
4. Project Layout: [Udacity sample project](#)