

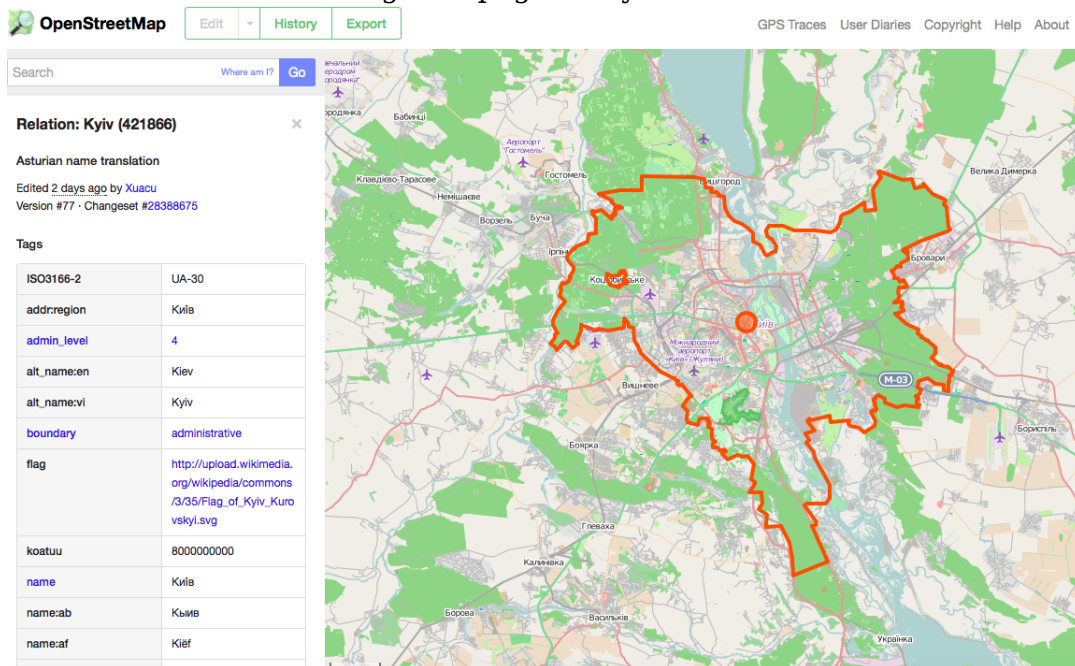
Data Wrangle OpenStreetMaps Data

Oleksii Renov

February 3, 2015

Choosing map area

Choosing map area is very simple for me. I live in Kiev, Ukraine. So i hadn't got another choice. As the very first step I need minimum and maximum longitude and latitude. For this task I'd gone to <http://openstreetmap.org> in the search field I'd typed Kiev and choose the second link. I've the following web page in my browser.



I've done some zooming and choose next latitude and longitude parameters in my case:

- minimum latitude = 29.9048;
- maximum latitude = 50.0986;

- minimum longitude = 31.1600;
- maximum longitude = 50.7052.

Next step is to download my dataset. I've gone in http://overpass-api.de/query_form.html, where i can past some queries to download data what i need.

My query is next:

```
(node(29.9048, 50.0986, 31.1600, 50.7052);<);out;
```

Anyway there is exists numbers of variants loading data. One of them is simply go to this url <http://overpass-api.de/api/map?bbox=29.9048,50.0986,31.1600,50.7052>. The xml file which I've already downloaded is very big, near 250mb. In my usual work i've chance working with xml, and files near 1mb is very hard to analyze. Let's move to processing this dataset.

Process Dataset

As very first step, let's check what tags are in our dataset. Also would be helpful to check for attributes. Let's start from tags.

```
1 {'bounds': 1,
2   'member': 75124,
3   'meta': 1,
4   'nd': 1316743,
5   'node': 1079619,
6   'note': 1,
7   'osm': 1,
8   'relation': 5016,
9   'tag': 417122,
10  'way': 153102}
```

Next json is corresponded to tag attributes.

```
1 {'changeset': 1237737,
2   'generator': 1,
3   'id': 1237737,
4   'k': 417122,
5   'lat': 1079619,
6   'lon': 1079619,
```

```

7  'maxlat': 1,
8  'maxlon': 1,
9  'minlat': 1,
10 'minlon': 1,
11 'osm_base': 1,
12 'ref': 1391867,
13 'role': 75124,
14 'timestamp': 1237737,
15 'type': 75124,
16 'uid': 1237737,
17 'user': 1237737,
18 'v': 417122,
19 'version': 1237738}

```

First of all i would like to notice, that 'k' is key and 'v' is value. There are equal number of keys and values what is obvious. But it isn't clear what this keys are holding. I would like to explore more details and understand what keys and values have the most frequency in my dataset. In my opinion top 10 would be perfect for start.

```

1  {('surface', 9780),
2   ('amenity', 9919),
3   ('name:en', 10396),
4   ('name:uk', 10646),
5   ('name:ru', 12823),
6   ('addr:street', 14084),
7   ('name', 23862),
8   ('addr:housenumber', 30646),
9   ('highway', 67770),
10  ('building', 73142)}

```

There are few nice things here. Due to Ukraine history we have official street name in ukrainian, unofficial russian and english names. For convinience I'm going to analyze here only international names. Let's check some values for key 'amenity'. For me it's very interesting.

```

1  {('place_of_worship', 300),
2   ('bench', 331),
3   ('atm', 374),
4   ('kindergarten', 419),
5   ('fuel', 425),

```

```

6 ('restaurant', 514),
7 ('cafe', 520),
8 ('school', 543),
9 ('pharmacy', 566),
10 ('bank', 770),
11 ('parking', 1812)}

```

Let's try to analyze values, we've chosen three regular expressions to check text for lowercase, problem strings and others. Especially we would like to explore data with bad keys.

```

1 {'lower': 258569,
2  'lower_colon': 163,
3  'other': 157676,
4  'problemchars': 714}

```

Let's take a look at keys which values contain problem characters.

```

1 {'building:color': 2,
2  'contact:phone': 40,
3  'phone': 558,
4  'roof:colour': 19,
5  'wheelchair': 1}

```

There are exists more keys, but i want to point at not standartized phone numbers for this dataset. Let's look at them.

```

1 '+3-8-093-207-31-01, +3-8-067-838-47-64, +3-8-067-401-75-61'
2 '+380 44 251-34-34',
3 '+38(098)-109-5001',
4 '+380 44 490 20 30',
5 '+38 067 500 81 89',

```

Oh.. There are lots of phone numbers in different formats. I'm going to create one format for phone number for this dataset. It's can will be a list of phones even list can have only size equal to 1.

From wikipedia http://en.wikipedia.org/wiki/Telephone_numbers_in_Ukraine i've read that typical Ukrainian phone numbers is:

- +380 xx xxx-xx-xx (general phone numbers);

I've written script for this preprocessing, before converting data in JSON and loading data to MongoDB. Next little python script handle different data formats and standartize them.

```

def convertPhone(phone):
    # First step i'm going to remove all non number characters
    phone = re.sub('\D', '', phone)
    if len(phone) == 10:
        phone = "+38" + phone
    elif len(phone) == 7:
        phone = "+38044" + phone
    elif len(phone) == 12:
        phone = "+" + phone
    elif len(phone) == 11:
        phone = "+3" + phone
    elif len(phone) == 9:
        phone = "+380" + phone

    if len(phone) == 13:
        phone = phone[0:4] + '_' + phone[4:6] + '_' + phone[6:9] +
            '-' + phone[9:11] + '-' + phone[11:13]
    return phone

```

Results of transformation:

- '+380(44)593-12-06' to '+380 44 593-12-06';
- '+380-44-599-60-82' to '+380 44 599-60-81';
- '3317730' to '+380 44 331-77-30'.

Nice results. Let's audit streets names. Ukraine has next common street types:

expected = ["вулиця", "бульвар", "проспект", "узвіз", "площа", "провулок", "шосе", "набережна", "тупик", "дорога", "проїзд", "шлях"]

And next mapping:

```

mapping = [ "вул.": "вулиця",
    "пр.": "проспект",
    "пл.": "площа",
    'пров.' : 'провулок'
]

```