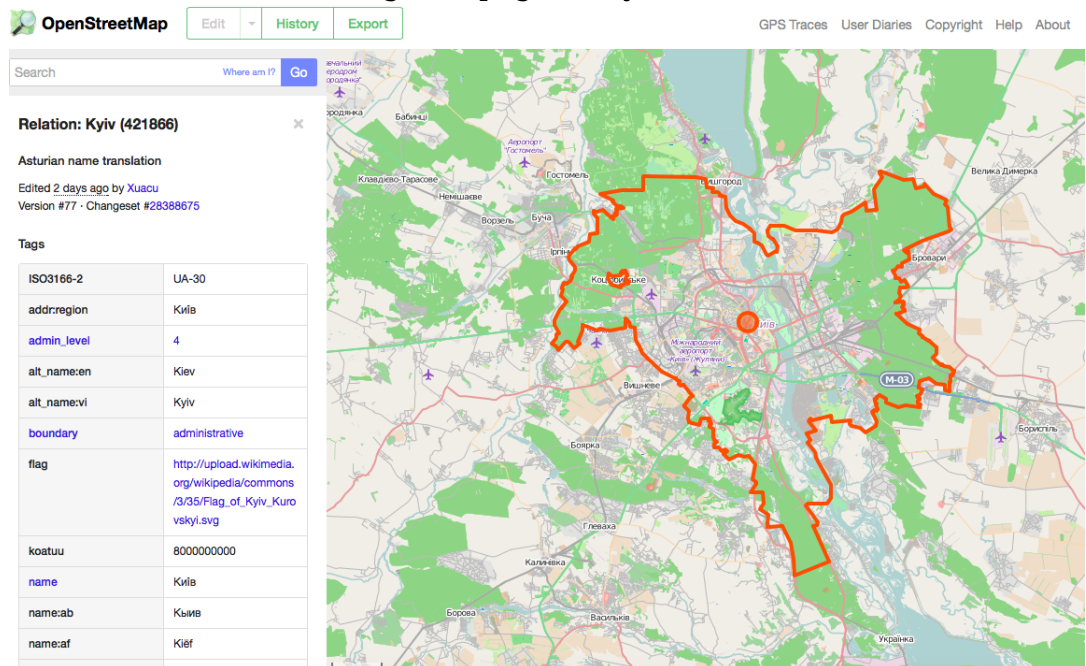# Data Wrangle OpenStreetMaps Data

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 zomming 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)}
```