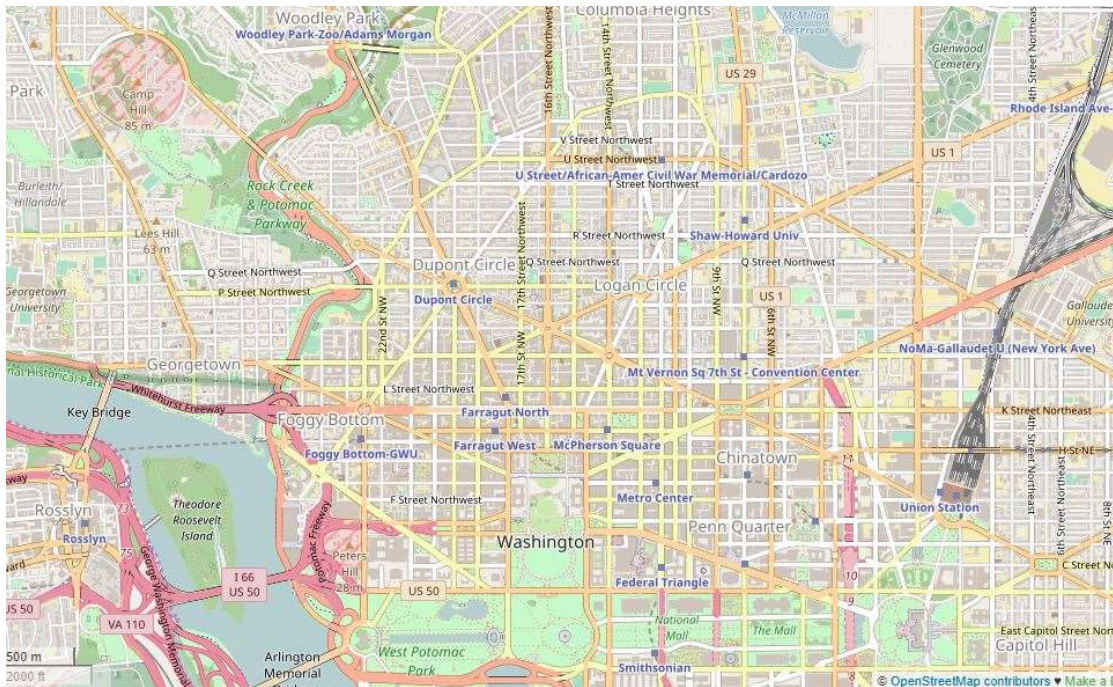


# OpenStreetMap Data Case Study

By Mathias Wainer.

## Map Area

Chosen map was a section of Washington D.C. Most of the area in the chosen section encompasses the Northwest quadrant of DC, this is where the national mall and most government buildings reside.



## Problems Encountered in the Map

After creating a sample file of 5% of the map I ran audits for quadrant names, postal codes and cuisine types. In D.C each address specifies which quadrant it resides. It is possible to have two identical addresses in different quarters (ex: 1900 F street Northwest and 1900 F street Northeast).

## Street

Snippet of street audit code:

```
9 osm_file = open("washington.osm", "r", encoding='utf8')
10
11 street_type_re = re.compile(r'\S+\.?$', re.IGNORECASE)
12 street_types = defaultdict(int)
13
14 def audit_street_type(street_types, street_name):
15     m = street_type_re.search(street_name)
16     if m:
17         street_type = m.group()
18
19         street_types[street_type] += 1
20
```

Output of street audit:

```
Alley 2
Avenue 2
Blvd 1
BN 2
Boulevard 32
Dr 1
Drive 5
Expressway 1
floor 1
Highway 1
Hill 1
Mall 1
n.w. 1
N.W. 1
NE 23
North 1
Northeast 16541
Northwest 20854
northwest 1
NW 51
Road 1
road) 4
S.W. 1
SE 2
```

Many of the quadrants in the addresses were abbreviated, as well as a few street abbreviation errors in neighboring Virginia.

Snippet of street cleaning:

```

56 def is_street_name(elem):
57     return (elem.attrib['k'] == "addr:street")
58
59
60 def audit(osmfile):
61     osm_file = OSMFILE
62     street_types = defaultdict(set)
63     for event, elem in ET.iterparse(osm_file, events=("start",)):
64
65         if elem.tag == "node" or elem.tag == "way":
66             for tag in elem.iter("tag"):
67                 if is_street_name(tag):
68                     audit_street_type(street_types, tag.attrib['v'])
69     osm_file.close()
70     return street_types
71
72 def update_name(name, mapping):
73     m = street_type_re.search(name)
74     if m:
75         street_type = m.group()
76         if street_type in mapping:
77             name = re.sub(street_type_re, mapping[street_type], name)
78     return name

```

Output of street cleaning changes only:

```

NE => Northeast
NW => Northwest
Dr => Drive
St => Street
SW => Southwest
SE => Southeast

```

## Zip-Code

Snippet of zip-code audit:

```

5 post_type_re = re.compile(r'\S+\.?$', re.IGNORECASE)
6 postcodes = defaultdict(int)
7
8 def audit_post_type(postcodes, postcode):
9     m = post_type_re.search(postcode)
10    if m:
11        postcode = m.group()
12
13        postcodes[postcode] += 1
14
15 def print_sorted_dict(d):
16     keys = d.keys()
17     keys = sorted(keys, key=lambda s: s.lower())
18     for k in keys:
19         v = d[k]
20         print(k, v)
21
22 def is_postcode(elem):
23     return (elem.tag == "tag") and (elem.attrib['k'] == "addr:postcode")

```

Partial output of zip-code audit:

```

20001 2492
20002 10011
20003 6692
20004 24
20005 65
20005-1001 1
20005-1009 1
20005-1013 1
20005-1015 2
20005-1019 1
20005-4111 1
20005-5702 1
20005-7700 1
20006 19
20006-5346 1

```

As we can see there are no bad zip code inputs. Some have the 4-digit sub-zip code which subdivides the areas within the zip code.

## Cuisine

Code showing snippet of cuisine:

```

1 osm_file = open("washington.osm", "r", encoding='utf8')
2
3 post_type_re = re.compile(r'\S+\.?$', re.IGNORECASE)
4 postcodes = defaultdict(int)
5
6 def audit_post_type(postcodes, postcode):
7     m = post_type_re.search(postcode)
8     if m:
9         postcode = m.group()
10
11         postcodes[postcode] += 1
12
13 def print_sorted_dict(d):
14     keys = d.keys()
15     keys = sorted(keys, key=lambda s: s.lower())
16     for k in keys:
17         v = d[k]
18         print(k, v)
19
20 def is_postcode(elem):
21     return (elem.tag == "tag") and (elem.attrib['k'] == "cuisine")
22
23

```

Partial output of cuisine audit:

```

Afghan 1
afghan 1
american 37
American, Korean 1
american,pizza,beer 1
american;international;fine_dining 1
american;steak_house 1
american_and_geechy 1
asian 12
bagel 2
bagel 2
bagel;coffee_shop;american;breakfast;tea 1
balkan 1
Bar 1
barbecue 3
Bistro 1
Breakfast/Brunch 1
breakfast_bagels,_desserts,_pastries 1
Brunch_Restaurant 1
burger 36

```

Many equal cuisine types are divided due to capitalization. Examples are “Afghan” and “afghan”. Other errors are due to detail such as “coffee” and “coffee\_shop” which should be the same.

Snippet of cuisine cleaning code:

```
33 def is_cuisine_name(elem):
34     return (elem.attrib['k'] == "cuisine")
35
36
37 def audit(osmfile):
38     osm_file = OSMFILE
39     cuisine_types = defaultdict(set)
40     for event, elem in ET.iterparse(osm_file, events=("start",)):
41
42         if elem.tag == "node" or elem.tag == "way":
43             for tag in elem.iter("tag"):
44                 if is_cuisine_name(tag):
45                     audit_cuisine_type(cuisine_types, tag.attrib['v'])
46     osm_file.close()
47     return cuisine_types
48
49 def update_name(name, mapping):
50     m = cuisine_type_re.search(name)
51     if m:
52         cuisine_type = m.group()
53         if cuisine_type in mapping:
54             name = re.sub(cuisine_type_re, mapping[cuisine_type], name)
55     return name
```

Output of cuisine changes only:

```
Afghan => afghan
steak => steak_house
Ethiopian => ethiopian
burger;american => burger
Korean => korean
salads => salad
Salad => salad
latin_america => latin_american
cajun;creole => cajun
coffee => coffee_shop
Diner => diner
```

## Data Overview

File sizes:

washington.osm	116.03 Mb
----------------	-----------



mydb.db	71.86 MB
nodes.csv	38.55 Mb
nodes_tags.csv	3.01 Mb
ways.csv	3.9 Mb
ways_nodes.csv	1.42 Mb
ways_tags.csv	14.59 Mb

Number of nodes:

```
sqlite> SELECT COUNT(*) FROM nodes;
```

482708

Number of ways:

```
sqlite> SELECT COUNT(*) FROM ways;
```

64130

Number of unique users:

```
sqlite> SELECT COUNT(DISTINCT(e.uid)) FROM(SELECT uid FROM nodes UNION ALL
SELECT uid FROM ways) e;
```

442

Number of Chinese restaurants:

```
sqlite> SELECT nodes_tags.value, count(*) as num
```

```
FROM nodes_tags JOIN (SELECT DISTINCT(id) FROM nodes_tags WHERE
value='restaurant') ON nodes_tags.id=i.id
```

```
WHERE nodes_tags.key='cuisine' and nodes_tags.value='chinese';
```

24

Looking at restaurants and fastfood:

```
SELECT nodes_tags.value, COUNT(*) as num
```

```
FROM nodes_tags
```

```
JOIN(SELECT DISTINCT(id) FROM nodes_tags WHERE value='fast_food') i
```

```
ON nodes_tags.id=i.id
```

```
WHERE nodes_tags.key='name'
```

```
GROUP BY nodes_tags.value
```

ORDER BY num DESC

LIMIT 4;

Fast Food	Restaurant
Subway 26	Chipotle 3
McDonald's 15	Sweet Green 3
Chipotle 8	&Pizza 2
Taylor Gourmet 6	Banana Leaves 2

(restaurants would have the same query except for replacing fast\_food with restaurant)

### Future Problem to Fix

Some restaurants like “chipotle” are counted in both or only some braches are. This is important because DC is considered a food desert and a lot of research is being done to figure out the situation and possible solutions for a food desert. I would remove the fast\_food tag completely and group everything under key: “restaurant”. Fast food is a somewhat subjective term since Taylor Gourmet is an upscale sandwich chain and there are a lot of restaurants where you pick up food at a counter that would not be considered fast foods.

The benefits of removing the fast\_food key is that it becomes easier to find where Washingtonians eat. The issue with removing the fast\_food tag is that we lose a little bit of detail in our map. If a researcher studying DC’s food desert wanted to look at fast food only he/she would have to query for specific restaurants such as “Subway” and “McDonalds”. With increasing levels of detail errors such as these are more likely to occur.

### Conclusion

I was pleasantly surprised how clean the Washington DC dataset was. After the SQL queries ran above I learned some interesting facts about my college town. There could be some improvement in the classification of fast food and restaurants for example.