

Table of Contents

- [Introduction](#)
- [Problems encountered with your Map](#)
 - [File-size > 87MB](#)
 - [Values for cities were redundant](#)
 - [Values for postal codes](#)
 - [Values for street](#)
 - [Necessary adjustments to the original shape element\(\) function](#)
- [Overview of the Data](#)
- [Other ideas about the datasets](#)
 - [Top 5 Contributors' Geo Tag Maps](#)
 - [Mulad](#)
 - [iandees](#)
 - [stucki1](#)
 - [DavidF](#)
 - [sota767](#)
- [Codes](#)
 - [xml2dict\(\)](#) – Code for taking the OSM file and creating dictionary
 - [dict2json\(\)](#) – Code for turning the dictionary into a .json file
 - [audit_xml\(\)](#) – Code for auditing the OSM file
 - [Map Creation Example](#) – Code for creating the Maps for the Top 5 Users
- [Lesson 6 Codes](#)
 - [6.1 Iterative Parsing](#)
 - [6.2 Data Model](#)
 - [6.3 Tag Types](#)
 - [6.4 Exploring Users](#)
 - [6.5 Improving Street Names](#)
 - [6.6 Preparing for Database](#)

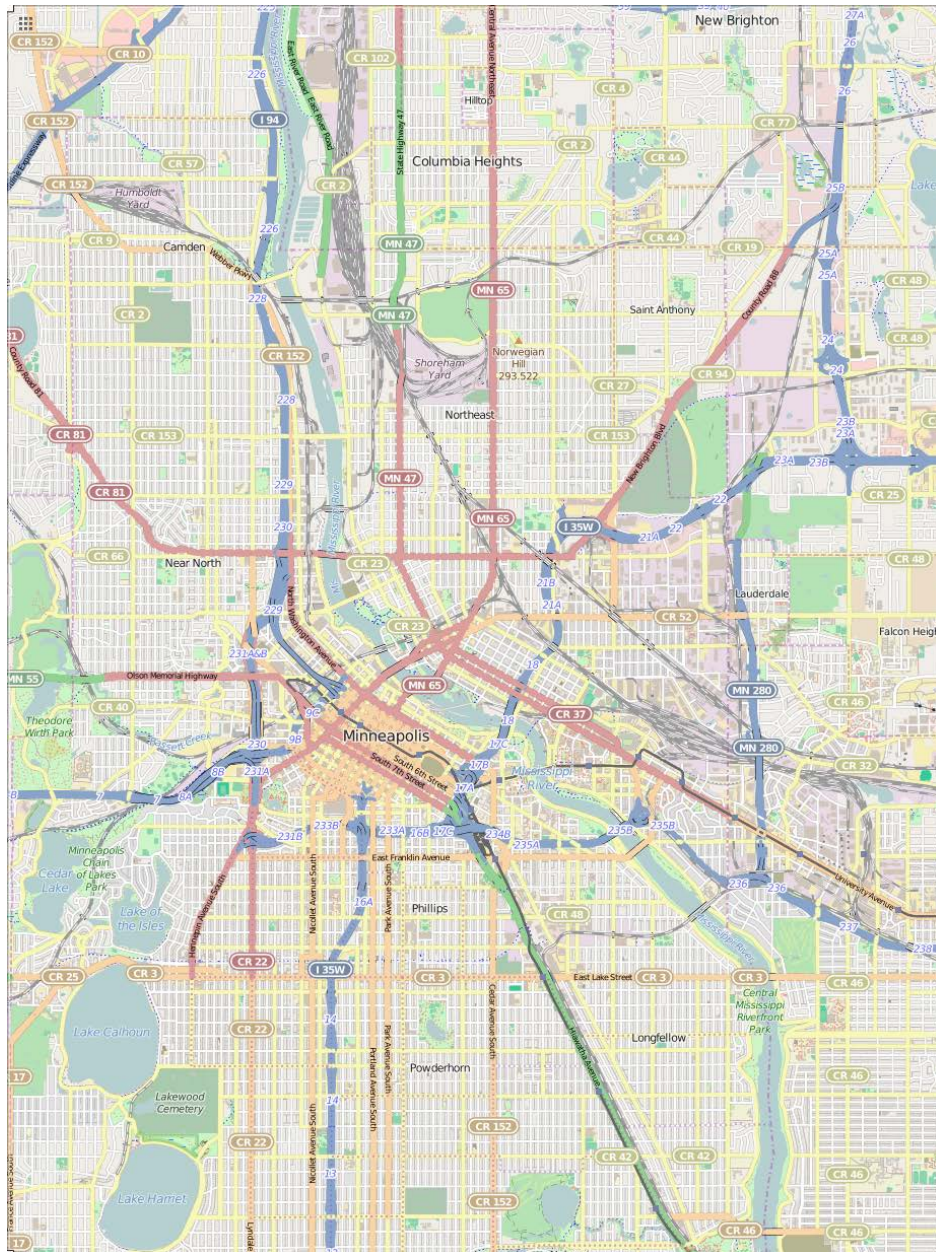
Introduction

For this project I have chosen my home area Minneapolis, MN. The reason is to be better able to validate the data by cross-referencing the data with personal knowledge of the area. The area is mainly Minneapolis, MN with some of the suburbs including New Brighton, Hopkins, and Minneapolis' Twin City, St. Paul. Here is a link to the Overpass API for downloading the OSM for the region.

<http://overpass-api.de/api/map?bbox=-93.3292,44.9150,-93.1713,45.0644>.

All the relevant code will be located in the end of this document with all the links that helped me complete the project. [\[back to top\]](#)

Below is the Map screenshotted from OSM.



1. Problems encountered in your Map

[\[back to top\]](#)

- File-size > 87MB
 - i. Even with iterative parsing, auditing and identifying meaningful patterns took longer than expected
 - ii. Took a smaller sample (<12MB) for creating functions for auditing, parsing, cleaning, etc.
 - iii. Once the functions were complete, they were applied to the 87MB OSM file
- Values for cities were redundant
 - i. For example—there were 4 different representations of Saint Paul. Here are the mappings

```
city_mapping = {"Minneapolis, MN" : "Minneapolis",
               "St. Anthony" : "Saint Anthony",
               "St. Paul" : "Saint Paul",
               "St Anthony" : "Saint Anthony",
               "St Paul" : "Saint Paul"}
```

- Values for postal codes ranged from being 3-digits long to 9-digits, alphabetical, or simply repeated values from different fields (*postcode = street*). One thing to note is out of the 108 unique *postcode* values, once extracting only the 5-digit prefixes, we were left with only 33 unique results which is consistent with the area I've chosen.
 - i. First I chose to limit the postal codes to 5-digit lengths. This is to standardize the data by the 'least common denominator' so the data was more intuitively queryable.
 - ii. The problem postal codes had letters, were too short, or misplaced from a different field. The only way to find the correct postal codes was to search the internet using the other fields for the search criteria and then extrapolate the postal codes. It could be done case-by-case manually but it would be faster if we could do it programmatically.
- iii. Luckily I found [pygeocoder](#)
 - Pygeocoder is a Python Library that utilizes Google's Geocoding Functionality to extrapolate geographic information such as coordinates, postal codes, etc.
 - I decided to use the *reverse_geocode()* function which will take the latitudinal and longitudinal coordinates and return an object from which the postal codes could be called from.

```
> from pygeocoder import Geocoder
> results = Geocoder.reverse_geocode(45.0079126, -93.2473816)
> int(results.postal_code)
> 55418
```

- There were some *postalcode* fields which did not have coordinate systems to use. For this, Pygeocoder also has the *geocode()* function which takes in a string argument, preferably a street address and state, and returns an object which also contains the postal code information.

```
> results = Geocoder.geocode('The Carlyle MN')
> int(results.postal_code)
> 55401
```

- Values for street suffixes were redundant

- i. The original mapping from Lesson 6 only covered a few keys to be normalized. After sifting through the data, I have discovered many more keys. Here is my mapping for street abbreviations.

```
street_mapping = { "St": "Street",
                  "St.": "Street",
                  "Ave": "Avenue",
                  "Rd.": "Road",
                  "Rd": "Road",
                  "SE": "Southeast",
                  "S.E.": "Southeast",
                  "NE": "Northeast",
                  "S": "South",
                  "Dr": "Drive",
                  "Rd/Pkwy": "Road/Parkway",
                  "Ln": "Lane",
                  "Dr.": "Drive",
                  "E": "East",
                  "Pl": "Plain",
                  "ne": "Northeast",
                  "NW": "Northwest",
                  "Ave.": "Avenue",
                  "N.": "North",
                  "W": "West",
                  "Pkwy": "Parkway",
                  "Ter": "Terrace",
                  "Pky": "Parkway",
                  "SW": "Southwest",
                  "N": "North",
                  "Blvd": "Boulevard" }
```

- The original *shape_element* function created in Lesson 6, when applied to this particular OSM file resulted in 447 unique top-level descriptors. A large number of them had similar prefixes meaning they could be grouped under a single dictionary to simplify querying.

My goal is to aggregate related descriptors into a single top-level key which can be accessed by calling that key. This will result in a more organized representation of the data. The following describes the grouped descriptors and their sources.

- i. 'address:' fields
 - These were fields dealt with in Lesson 6
- ii. 'metcouncil:' fields
 - These fields with information from the [Metropolitan Council](#) of Minnesota
- iii. 'tiger:' fields
 - Topologically Integrated Geographic Encoding and Referencing system ([TIGER](#)) data was produced by the US Census Bureau and later merged onto OSM database
- iv. 'metrogis:' fields
 - [Regional geographic information](#) system of the seven-county metropolitan area of Minneapolis-St. Paul, Minnesota
- v. 'umn:' fields
 - Fields related to the University of Minnesota, [example](#).
 - This field presented a unique programmatic problem—If it existed, the positional data (Latitude,Longitude) was located in its parent element (in this case *way*), it would be located in a *umn:* k sub-element.
- vi. 'gnis:' fields
 - USGS Geographic Names Information System ([GNIS](#)) contains millions of names for geographic features in the US.

After grouping I reduced the number of descriptors from 447 to 364. The fields mentioned above are important because they determine the sources of information merged onto the OSM user-based database. Cross-referencing these fields could determine accuracy, constituency, and completeness of the data.

2. Overview of the Data

[\[back to top\]](#)

- Statistics

- i. File Sizes

- map.osm – 87MB
 - map.json – 95MB

- ii. Documents

```
> db.minneapolis.find().count()
> 436802
```

- iii. Nodes

```
> db.minneapolis.find({'type' : 'node'}).count()
> 373017
```

- iv. Way

```
> db.minneapolis.find({'type' : 'way'}).count()
> 59143
```

- v. Unique users

```
> len(db.minneapolis.distinct('created.user'))
> 401
```

- vi. Top 5 cuisines in the category of restaurant

```
> pipeline = [{'$match' : {'amenity' : 'restaurant',
                           'cuisine' : {'$exists' : 1}}},
               {'$group' : {'_id' : '$cuisine',
                           'count' : {'$sum' : 1}}},
               {'$sort' : {'count' : -1}},
               {'$limit' : 5}]
> top5 = db.minneapolis.aggregate(pipeline)
> for i in top5['result']:
    print i['_id'], i['count']
> pizza 16
   chinese 15
   american 10
   italian 10
   regional 7
```

3. Other ideas about the Datasets

[\[back to top\]](#)

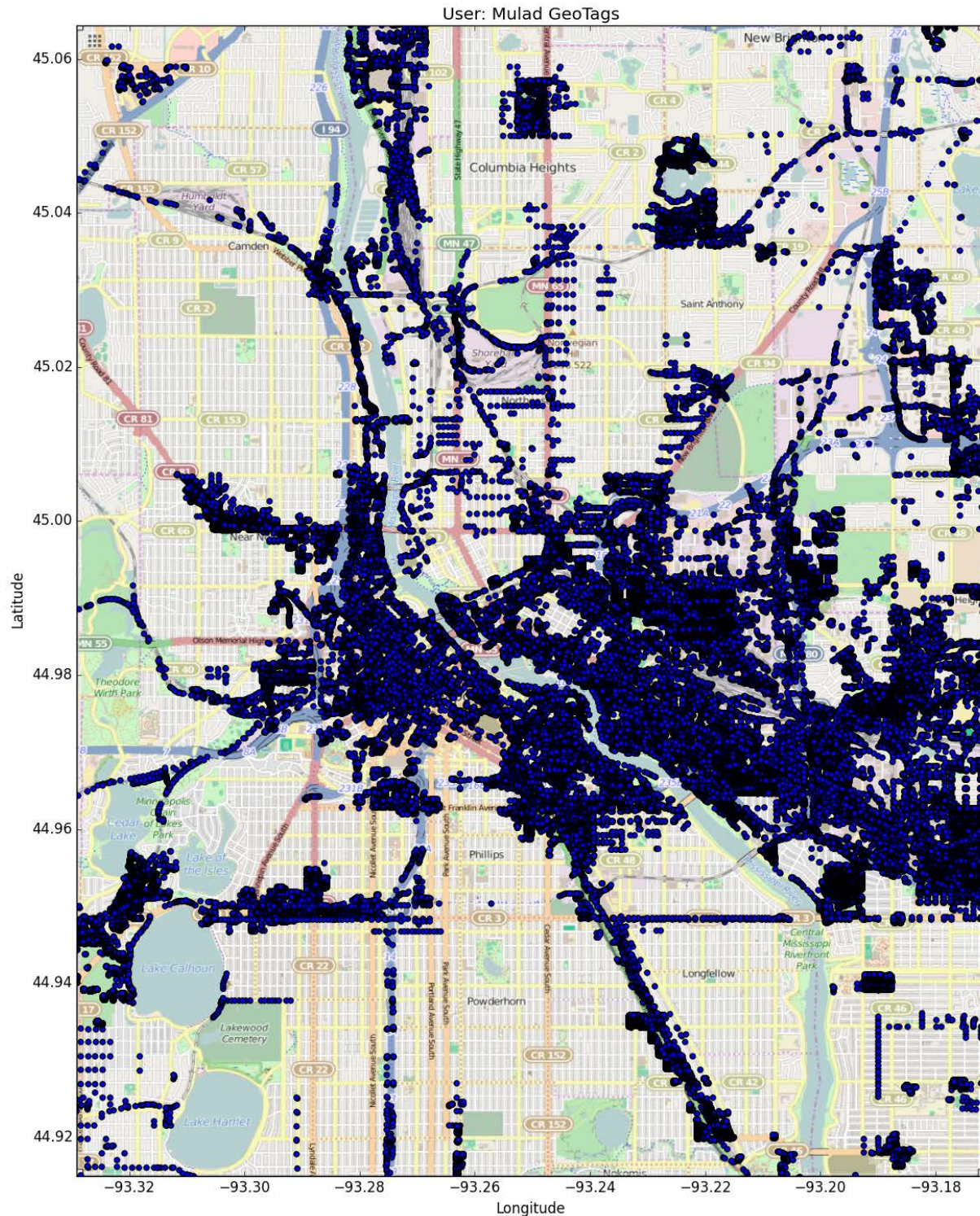
I've decided to take the positional data (latitude, longitude) of the top 5 contributing users and create a graphical display of the places they contributed overlaid on the map area the dataset came from. This should create a sort of 'heat-map' of each users' territory, so to speak. We may be able to see home neighborhoods or traveling patterns. The image of the top 5 contributors on to one map became too confusing to interpret so I split it into 5 individual maps/graphs. The codes for making the maps and an explanation of each step of the code will be available in the [Codes](#) section.

```
> pipeline = [{'$group' : {'_id' : '$created.user', 'count' : {'$sum' : 1}}},
               {'$sort' : {'count' : -1}},
               {'$limit' : 5}]
> top5 = db.minneapolis.aggregate(pipeline)
> for i in top5['result']:
    print i['_id'], i['count']
> Mulad 146082
    iandees 81267
    stucki1 63803
    DavidF 27032
    sota767 25901
```

- From the previous section we have discovered the top 5 contributors to be...
 - i. Mulad – 146082

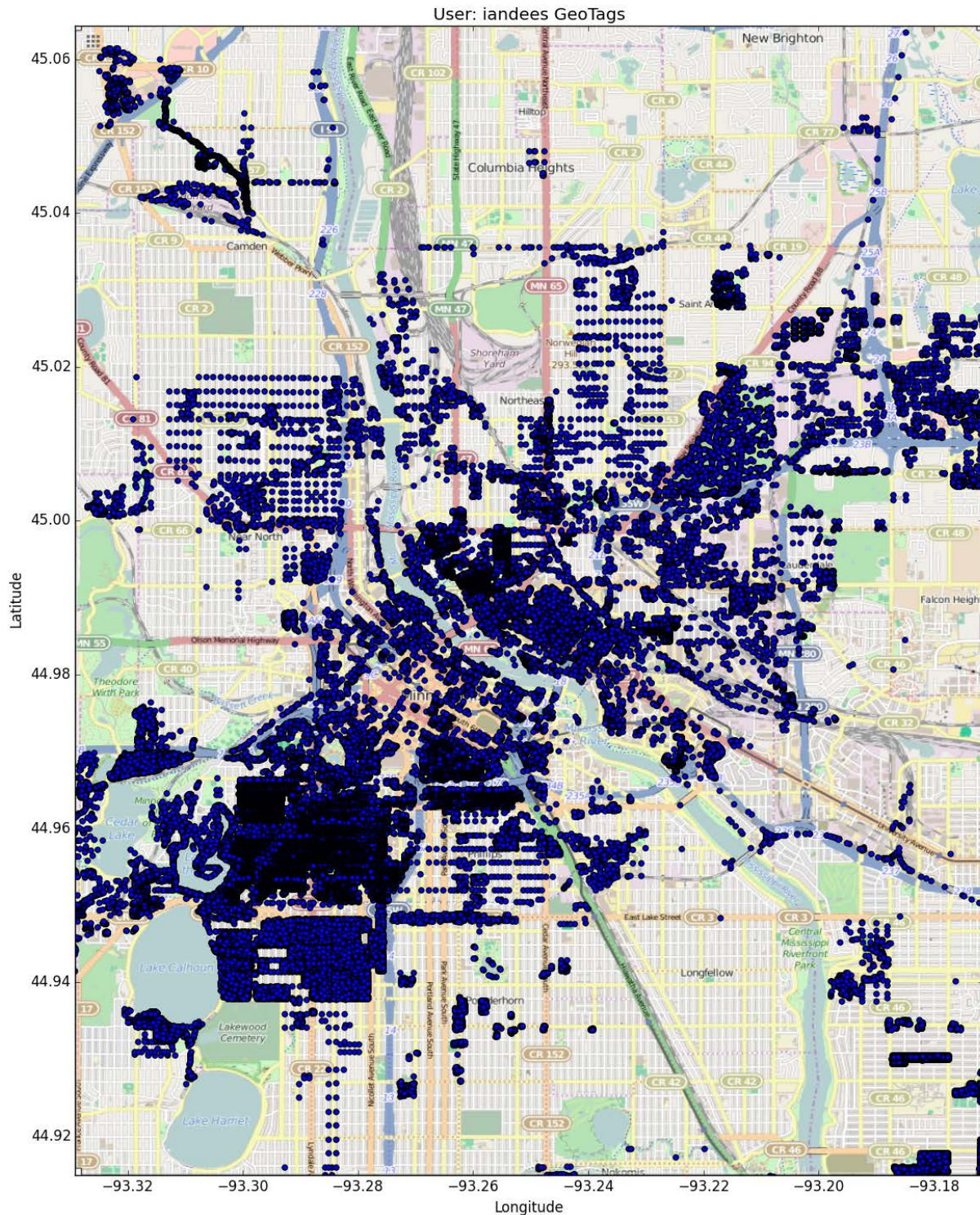
```
## Query to find percentage of total documents user contributed
db.minneapolis.find({'created.user': 'Mulad'}).count() /
float(db.minneapolis.find().count())
```

- ~33.4 % of the total documents
- ii. iandees – 81267
 - ~18.6 % of the total documents
- iii. stucki1 – 63803
 - ~14.6 %
- iv. DavidF – 27032
 - ~6.2 %
- v. sota767 – 25901
 - ~5.9
- There relative map and geotag locations will be shown on the following pages in order of magnitude of contributions.
- I will also include some notes derived from my personal knowledge of this area that could be of some interest.



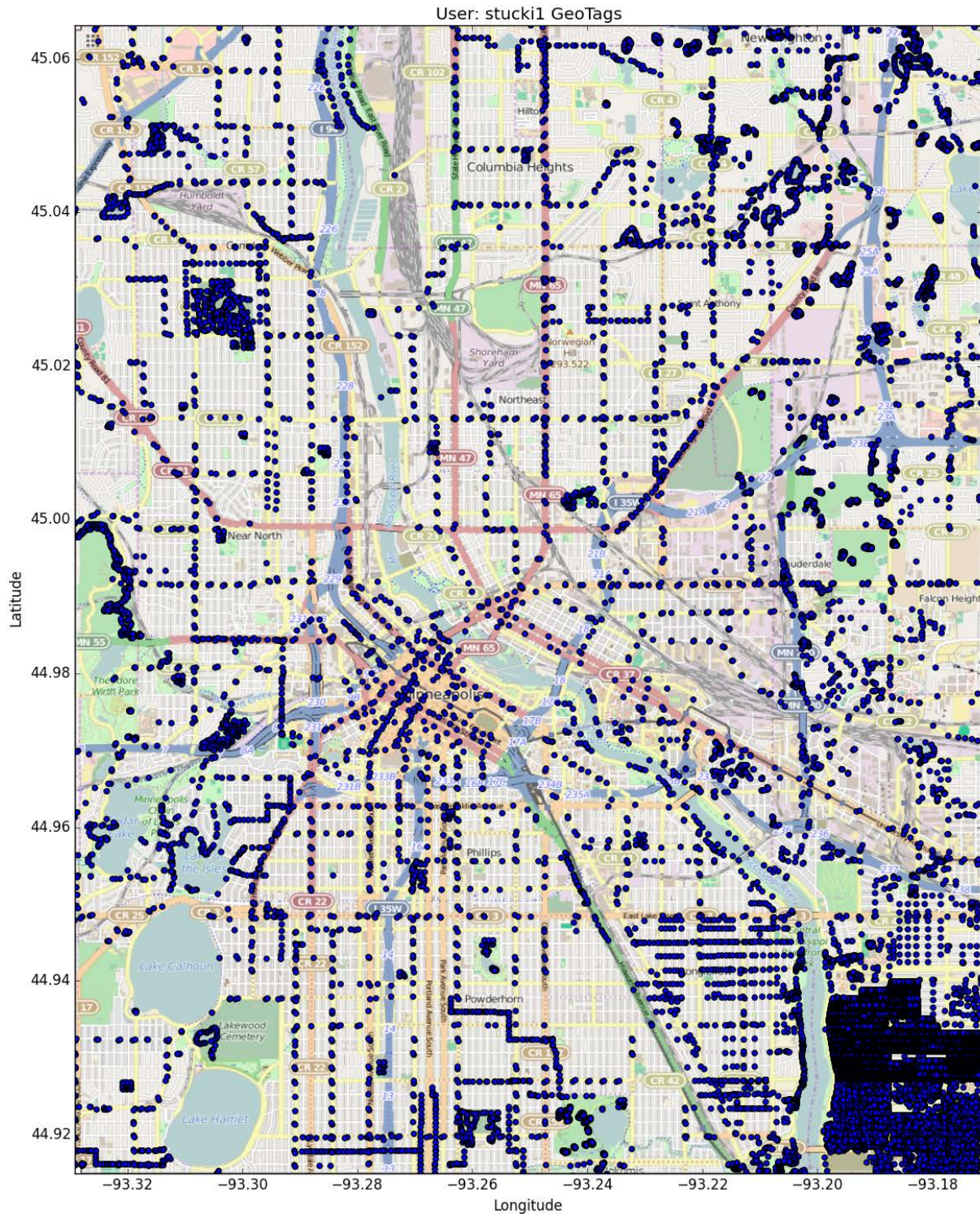
'Mulad' geotagged heavily in downtown Minneapolis, east through UofM and St. Paul, and a south stream through Hiawatha/55. These regions are high-frequency public transportation routes especially the Hiawatha Light Rail which is a primary route to connect the southern suburbs into Minneapolis. This could indicate 'Mulad' depends on public transportation and/or resides in Minneapolis proper.

[\[back to top\]](#)



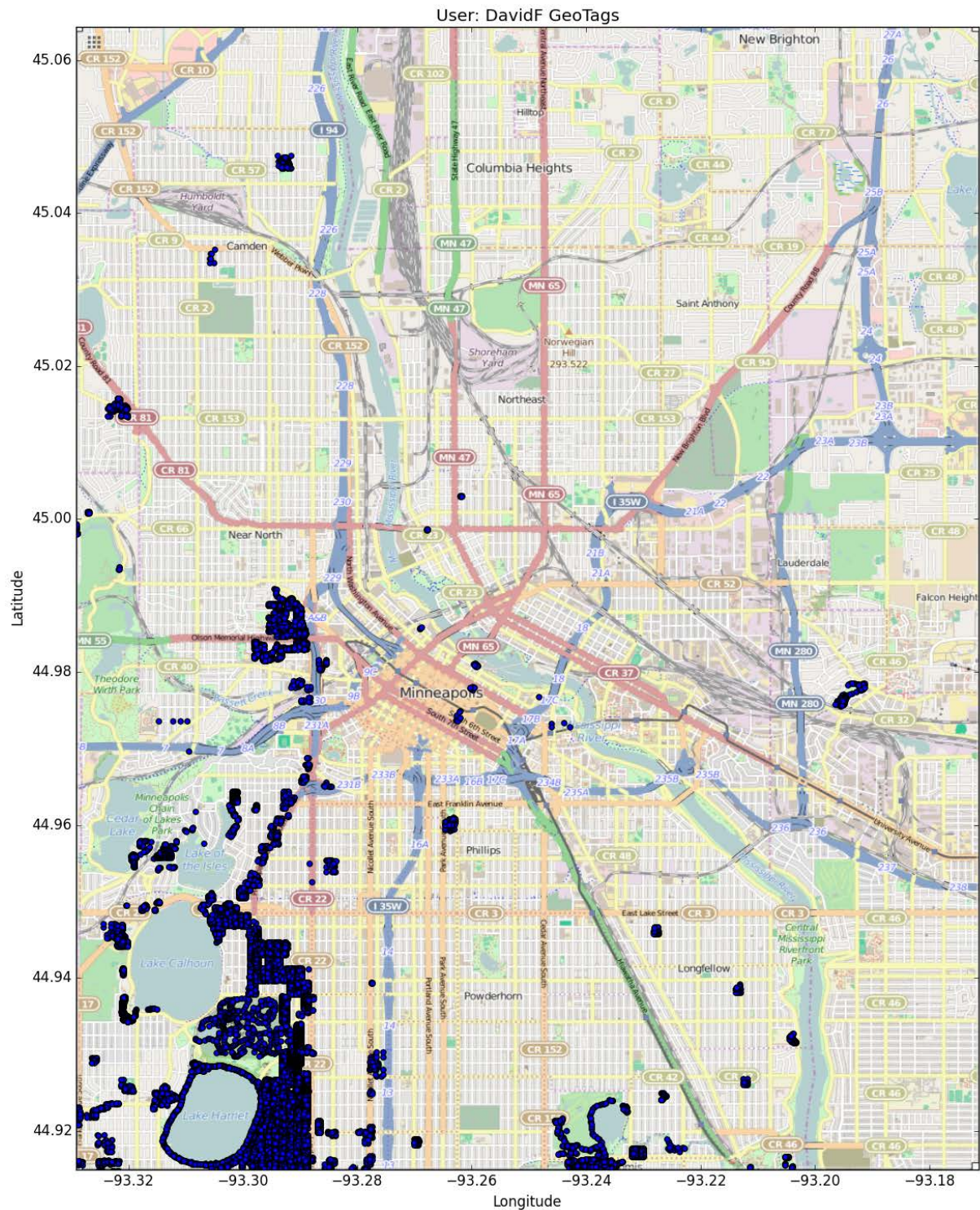
'iandees' geotags are primarily in a region of Minneapolis called 'Uptown' and runs to Northeast Minneapolis. These two regions, especially 'Uptown', are known for a vibrant 'bar scene'. Such a high concentration in the 'Uptown' region could indicate location of home.

[\[back to top\]](#)



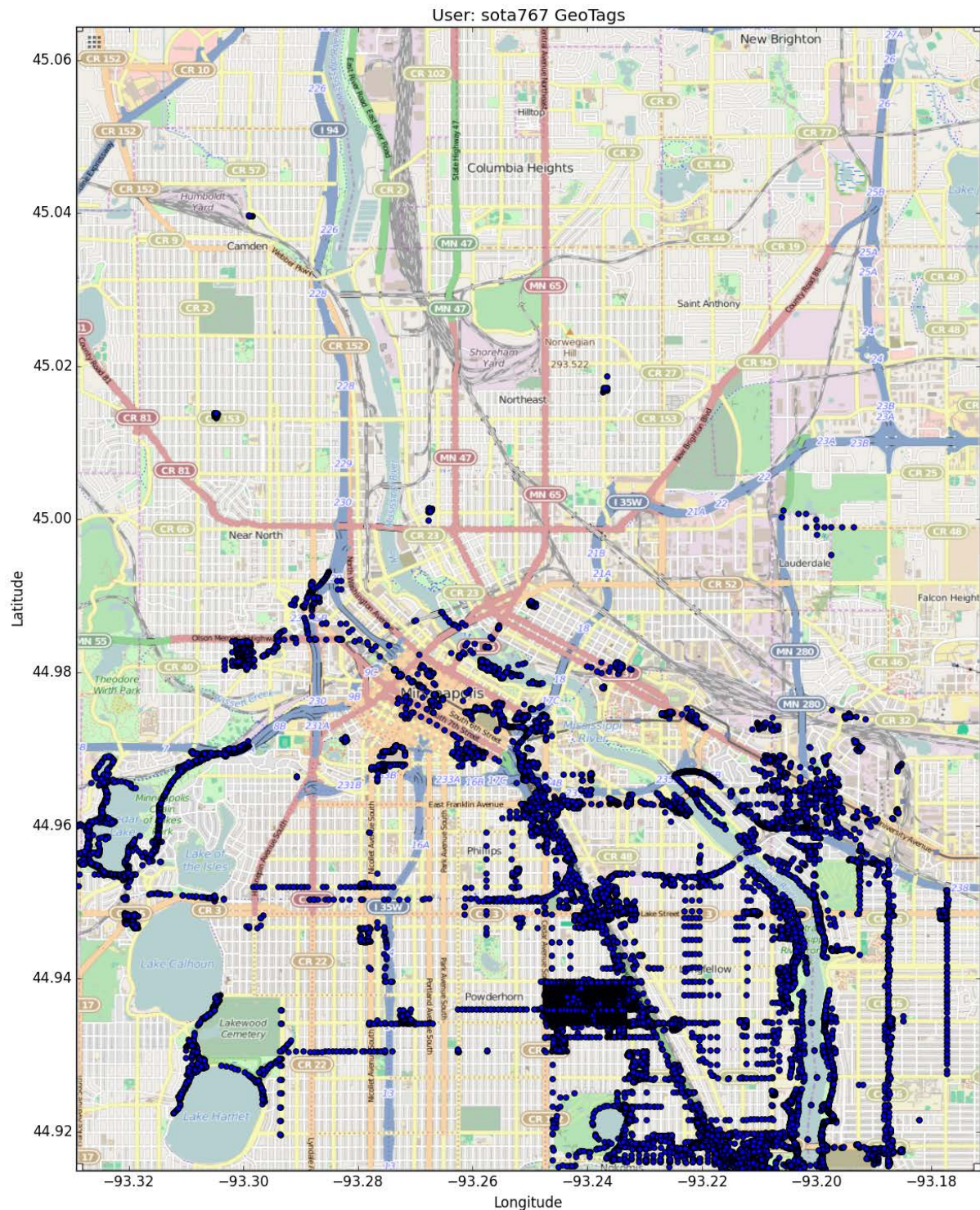
'stucki1's geotag patterns show high concentrations in the St. Paul side and a small concentration Northwest in Minneapolis with sparse tags on various intersections throughout the area. The intersection-tags could indicate 'stucki1' simply tagging while driving and the high-concentrations could indicate a home area and a work area.

[\[back to top\]](#)



'DavidF's geotags are highly concentrated around Lake Harriet/Lake Calhoun area. This area has higher than average housing prices and could indicate personal salary. This could also indicate a frequent area of recreation for the user.

[\[back to top\]](#)



'sota767's geotags are indicative of a home region around the Minnehaha falls area. User could also be a student since there are high concentrations in the UofM campus area. His tagging patterns are concentrated in the Hiawatha Light Rail system and along North-South public transportation routes which could indicate a reliance on public transport.

[\[back to top\]](#)

4. Conclusions

[\[back to top\]](#)

- I am inclined to presume that the top 5 geotaggers were hired by OSM or some other geographical tagging association. What leads me to this hypothesis is that these users tagged in specialized areas of the region that, for the most part, did not overlap and, between them, encompass ~79% of the total tags.
- I would also like to recommend a 'fresh' start with the geotags. Scanning through the documents I found a lot of key value pairs that did not categorize well and most seemed to be in a random format. These outliers are sparsely dispersed throughout the data and do not occur with a great enough frequency to warrant any sort of pattern. They are simply clutter. These values should be removed and a stricter form of geotagging should be implemented so future tags adhere to a certain standard.
- The positional data, once overlaid on top of the data region yielded some intriguing opportunities. If I were to utilize a larger area, I am sure I would be able to extract some meaningful patterns in users' geotagging tendencies relative to known geographical features. Many of the top users tagged in highly populated areas of the region leaving some 'holes'

Codes

[\[back to top\]](#)

```

1. # The xml2dict() is the main function which takes in an OSM file and turns it into a dictionary that can
2. # be accessed by normal dictionary queries. Only specified elements of the file will be used to populate
3. # The dictionary.
4. def xml2dict(filename):
5.     import time
6.     import re
7.     import xml.etree.ElementTree as ET
8.     from pygeocoder import Geocoder
9.     expected = ["Street", "Avenue", "Boulevard", "Drive", "Court", "Place", "Square", "Lane", "Road",
10.                "Trail", "Parkway", "Commons"]
11.     street_type_re = re.compile(r'\b\S+\.?$', re.IGNORECASE)
12.     street_mapping = {"St": "Street",
13.                       "St.": "Street",
14.                       "Ave" : "Avenue",
15.                       "Rd." : "Road",
16.                       "Rd" : "Road",
17.                       "SE" : "Southeast",
18.                       "S.E." : "Southeast",
19.                       "NE" : "Northeast",
20.                       "S" : "South",
21.                       "Dr" : "Drive",
22.                       "Rd/Pkwy" : "Road/Parkway",
23.                       "Ln" : "Lane",
24.                       "Dr." : "Drive",
25.                       "E" : "East",
26.                       "Pl" : "Plain",
27.                       "ne" : "Northeast",
28.                       "NW" : "Northwest",
29.                       "Ave." : "Avenue",
30.                       "N." : "North",
31.                       "W" : "West",
32.                       "Pkwy" : "Parkway",
33.                       "Ter" : "Terrace",
34.                       "Pky" : "Parkway",
35.                       "SW" : "Southwest",
36.                       "N" : "North",
37.                       "Blvd" : "Boulevard"}
38.     city_mapping = {"Minneapolis, MN" : "Minneapolis",
39.                     "St. Anthony" : "Saint Anthony",
40.                     "St. Paul" : "Saint Paul",
41.                     "St Anthony" : "Saint Anthony",
42.                     "St Paul" : "Saint Paul"}
43.     CREATED = [ "version", "changeset", "timestamp", "user", "uid"]
44.     address = ['addr:unit', 'addr:full', 'addr:housenumber', 'addr:postcode', 'addr:street',
45.                'addr:city', 'addr:state', 'addr:country',
46.                'addr:suite', 'addr:housename']
47.     # The shape_element() function takes in iterations of 'ET.iterparse'. Each iteration is a line in the OSM xml file
48.     # and will move forward in the function if the line has a tag that equals 'node' or 'way'.
49.     def shape_element(element):
50.         # First a dictionary is created. This dictionary will receive the cleaned data from the OSM file

```

```

51.         # and then returned at the end.
52.         node = {}
53.         # First we check if the tag is either 'node' or 'way'.
54.         if element.tag == "node" or element.tag == "way" :
55.             # Then we create a key called 'type' which equals either 'node' or 'way'.
56.             node['type'] = element.tag
57.             # Then we check every attribute in the line.
58.             # If an attribute matches any value in CREATED, we create a dictionary call
ed 'created'.
59.             # The loop breaks immedietly after the match and creation of the dictionary
.
60.             for i in element.attrib.keys():
61.                 if i in CREATED:
62.                     node['created'] = {}
63.                     break
64.             # Once the 'created' dictionary is created. The attributes which match the
values in CREATED and their
65.             # respective values are added to the dictionary 'created'. Attributes that
equal 'lat' or 'lon' are skipped.
66.             # The rest of the attributes are made into keys in the node dictionary with
values being their respective
67.             # value pair.
68.             for i in element.attrib.keys():
69.                 if i in CREATED:
70.                     node['created'][i] = element.attrib[i]
71.                 elif i == 'lon' or i == 'lat':
72.                     continue
73.                 else:
74.                     node[i] = element.attrib[i]
75.             # If 'lat' is one of the attributes in the element, it is assumed it has a
'lon' pair as well.
76.             # These two values are combined in a list [Latitude, Longitude]. A key, 'p
os', is created in the
77.             # node dictionary with its value being the [Latitude, Longitude] list.
78.             if 'lat' in element.attrib.keys():
79.                 node['pos'] = [float(element.attrib['lat']), float(element.attrib['lon'
]])
80.             # The following scans the 'k' values of the subelements of 'node' or 'way'
tags and creates a dictionary
81.             # depending on the matches.
82.             # For example, if scanning the 'k' values reveal a value that starts with '
gnis:'
83.             # the code will create a dictionary 'gnis' inside the node dictionary
84.             for i in element:
85.                 if 'k' in i.attrib:
86.                     if i.attrib['k'] in address:
87.                         node['address'] = {}
88.                     elif i.attrib['k'].startswith('metcouncil:'):
89.                         node['metcouncil'] = {}
90.                     elif i.attrib['k'].startswith('tiger:'):
91.                         node['tiger'] = {}
92.                     elif i.attrib['k'].startswith('metrogis:'):
93.                         node['metrogis'] = {}
94.                     elif i.attrib['k'].startswith('umn:'):
95.                         node['umn'] = {}
96.                     elif i.attrib['k'].startswith('gnis:'):
97.                         node['gnis'] = {}
98.             # Some of the Latitudenal and Longitudenal data are not located in the pare
nt levels of 'node' and 'way'
99.             # In particular, if the subelements 'k' values include items that start wi
th 'umn:', the corresponding

```



```

100.         # locational data is located in 'umn:BuildingCenterXYLatitude' and '
umn:BuildingCenterXYLongitude'
101.         # The following code extracts those values and includes them in the
'pos' list.
102.         for i in element:
103.             if 'k' in i.attrib:
104.                 if i.attrib['k'] in ['umn:BuildingCenterXYLatitude', 'umn:Bu
ildingCenterXYLongitude']:
105.                     node['pos'] = []
106.                     break
107.         for i in element:
108.             if 'k' in i.attrib:
109.                 if i.attrib['k'] == 'umn:BuildingCenterXYLatitude':
110.                     node['pos'].append(float(i.attrib['v']))
111.                     break
112.         for i in element:
113.             if 'k' in i.attrib:
114.                 if i.attrib['k'] == 'umn:BuildingCenterXYLongitude':
115.                     node['pos'].append(float(i.attrib['v']))
116.                     break
117.         # As instructed in lesson 6, some of the subelements whose tags are
'ref' need to be grouped
118.         # under the list 'node_refs'. The following creates that list.
119.         for i in element:
120.             if 'ref' in i.attrib:
121.                 node['node_refs'] = []
122.         # The following code populates the previously created groups in the
node dictionary
123.         # by scanning the subelements and tests the existence of the necessa
ry value-pairs.
124.         # If the necessary value-
pairs exist, the 'k' value prefix is stripped and added
125.         # as a key under the respective group. Then the 'v' is added as teh
value/
126.         # For example. 'k' = 'addr:street', 'v' = 'dorland' is added as
127.         # 'street' : 'dorland'
128.         for i in element:
129.             if 'k' in i.attrib:
130.                 if i.attrib['k'] in address:
131.                     if i.attrib['k'] == 'addr:city':
132.                         if i.attrib['v'] in city_mapping.keys():
133.                             node['address'][re.sub('addr:', '', i.attrib['k'
])] = city_mapping[i.attrib['v']]
134.                     else:
135.                         node['address'][re.sub('addr:', '', i.attrib['k'
])] = i.attrib['v']
136.                     else:
137.                         node['address'][re.sub('addr:', '', i.attrib['k'])]
= i.attrib['v']
138.                     elif i.attrib['k'].startswith('metcouncil:'):
139.                         node['metcouncil'][re.sub('metcouncil:', '', i.attrib['k
'])] = i.attrib['v']
140.                     elif i.attrib['k'].startswith('tiger:'):
141.                         node['tiger'][re.sub('tiger:', '', i.attrib['k'])] = i.a
ttrib['v']
142.                     elif i.attrib['k'].startswith('metrogis:'):
143.                         node['metrogis'][re.sub('metrogis:', '', i.attrib['k'])]
= i.attrib['v']
144.                     elif (i.attrib['k'].startswith('umn:') and
145.                           i.attrib['k'] not in ['umn:BuildingCenterXYLatitude'
, 'umn:BuildingCenterXYLongitude']):

```

```

146.         node['umn'][re.sub('umn:', '', i.attrib['k'])] = i.attrib['v']
147.         elif i.attrib['k'].startswith('gnis:'):
148.             node['gnis'][re.sub('gnis:', '', i.attrib['k'])] = i.attrib['v']
149.         elif ('addr:street:' in i.attrib['k'] or
150.              i.attrib['k'] in ['umn:BuildingCenterXYLatitude', 'umn:BuildingCenterXYLongitude']):
151.             continue
152.         else:
153.             # All the remaining value pairs are then added with the
154.             # 'k' value being the key
155.             # and the 'v' value being the value of that key
156.             node[i.attrib['k']] = i.attrib['v']
157.             # Earlier we created the 'node_refs' list inside the node dictionary
158.             # The following will add the values to this list.
159.             if 'ref' in i.attrib:
160.                 node['node_refs'].append(i.attrib['ref'])
161.             # Finally after the node dictionary is created from the particular iteration of ET.iterparse(),
162.             # it is returned
163.             return node
164.         # 2 empty lists are created. We will append each dictionary returned by feeding the iterations of the OSM
165.         # file into shape_element() into the 'temp' list. Since some of the iterations will not contain the tags
166.         # 'node' and 'way', the shape_element() will return 'None' for these instances.
167.         # We will then run through each iteration of temp and only append the iterations that contains information
168.         # to the data list.
169.         temp = []
170.         data = []
171.         for _, element in ET.iterparse(filename):
172.             temp.append(shape_element(element))
173.         for i in temp:
174.             if i != None:
175.                 data.append(i)
176.             # The following corrects the street values for the data. We run through each iteration of the data list.
177.             # If the iteration contains a key 'address', we move forward and check if the key 'address' contains a key
178.             # 'street'. If these conditions are met, we check the street value's suffix is in the keys of 'street_mapping'.
179.             # If it is, the code replaces that word for the value of the key it matches.
180.             for i in data:
181.                 if 'address' in i:
182.                     if 'street' in i['address']:
183.                         search = street_type_re.search(i['address']['street'])
184.                         if search:
185.                             if street_type_re.search(i['address']['street']).group() in street_mapping.keys():
186.                                 data[data.index(i)]['address']['street'] = re.sub(street_type_re.search(i['address']['street']).group(),
187.                                                                                     street_mapping[street_type_re.search(i['address']['street']).group()],
188.                                                                                     data[data.index(i)]['address']['street'])

```

```

188.         # The following corrects the incorrect postal codes. We run through each it
    eration of the data list.
189.         # The first 'if' statement checks if the iteration has the keys 'address' an
    d 'pos'
190.         # If so, we proceed to check if the postal code is incorrect.
191.         # If it is, we use the coordinates in the 'pos' list and apply it to the Geo
    coder.reverse_geocode()
192.         # The postal code from that query will replace the incorrect postal code.
193.         #
194.         # The second 'elif' statement is for iterations that have an 'address' but d
    o not have a 'pos' list
195.         # to reverse_geocode() the postal code.
196.         # It instead will use the Geocoder.geocode() function and search specific el
    ements of the 'address'
197.         # key to return the correct postal code.
198.         for i in data:
199.             if 'address' in i and 'pos' in i:
200.                 if 'postcode' in i['address']:
201.                     if len(i['address']['postcode']) < 5 or re.search('[a-zA-
    Z]', i['address']['postcode']):
202.                         results = Geocoder.reverse_geocode(i['pos'][0], i['pos'][1])
203.                         i['address']['postcode'] = str(results.postal_code)
204.                 elif 'address' in i and 'pos' not in i:
205.                     if 'postcode' in i['address']:
206.                         if len(i['address']['postcode']) < 5 or re.search('[a-zA-
    Z]', i['address']['postcode']):
207.                             q = ''
208.                             if 'houseName' in i['address']:
209.                                 q = i['address']['houseName'] + ' MN'
210.                             elif 'houseNumber' in i['address'] and 'street' in i['addres
    s']:
211.                                 q = i['address']['houseNumber'] + ' ' + i['address']['st
    reet'] + ' MN'
212.                             results = Geocoder.geocode(q)
213.                             i['address']['postcode'] = str(results.postal_code)
214.         # The following will standardize 2 part postal codes to 1 part postal codes.
215.         # For example, 55404-1234 will be turned to 55404
216.         for i in data:
217.             if 'address' in i:
218.                 if 'postcode' in i['address']:
219.                     if len(i['address']['postcode']) > 5:
220.                         i['address']['postcode'] = i['address']['postcode'][0:5]
221.         # The cleaned and sorted data is ready to be returned.
222.         return data
223.
224.         # The dict2json() simply takes the dictionary produced by xml2dict() and saves i
    t as the specified
225.         # file in .json format
226.         def dict2json(dict, output_file):
227.             import codecs
228.             import json
229.             with codecs.open(output_file, 'w') as fo:
230.                 for i in dict:
231.                     fo.write(json.dumps(i) + '\n')
232.             fo.close()
233.
234.
235.

```



```

236.     # The audit_xml() scans through the original OSM file, which is formatted in XML
237.     # and returns a specified dictionary of audits
238.     def audit_xml(filename, form = 'all', value = None):
239.         import xml.etree.ElementTree as ET
240.         # count_attrib() will count the unique values of each element in the file
241.         def count_attrib(filename):
242.             attrib_count = {}
243.             for _, element in ET.iterparse(filename):
244.                 if element.tag in ['node', 'way', 'tag', 'nd']:
245.                     for i in element.attrib.keys():
246.                         if i not in attrib_count.keys():
247.                             attrib_count[i] = {'count': 1}
248.                         else:
249.                             attrib_count[i]['count'] += 1
250.             return attrib_count
251.         # count_val() will count a specific value within the file
252.         def count_val(x, filename):
253.             k = {}
254.             for _, element in ET.iterparse(filename):
255.                 if element.tag in ['node', 'way', 'tag', 'nd']:
256.                     if x in element.attrib:
257.                         if element.attrib[x] not in k:
258.                             k[element.attrib[x]] = 1
259.                         else:
260.                             k[element.attrib[x]] += 1
261.             return k
262.         # tag_count simply counts the unique tags located in the file
263.         def tag_count(filename):
264.             tags = {}
265.             for _, element in ET.iterparse(filename):
266.                 if element.tag in ['node', 'way', 'tag', 'nd']:
267.                     for i in element:
268.                         if i.tag not in tags:
269.                             tags[i.tag] = 1
270.                         else:
271.                             tags[i.tag] += 1
272.             return tags
273.         # depending on the 'form' and/or 'value' argument, audit_xml() will return t
274.         # audit type
275.         if form.lower() == 'all':
276.             data = {}
277.             attrib = count_attrib(filename)
278.             for i in attrib:
279.                 data[i] = count_val(i, filename)
280.             for i in data:
281.                 data[i]['TOTAL'] = attrib[i]['count']
282.             data['TAGS'] = tag_count(filename)
283.             return data
284.         elif form.lower() == 'tags':
285.             return tag_count(filename)
286.         elif form.lower() == 'attributes':
287.             return count_attrib(filename)
288.         elif form.lower() == 'values':
289.             return count_val(v, filename)
290.         else:
291.             if not form or form not in ['all', 'tags', 'attributes', 'values']:
292.                 print "Invalid Audit Type"

```

The following code created the map which represented the individual user's geotag. I basically used 'cbook' and 'imread' to load the image I screenshotted from OSM of my region as the graph background.

[\[back to top\]](#)

```

2. import numpy as np
3. import matplotlib.pyplot as plt
4. from scipy.misc import imread
5. import matplotlib.cbook as cbook
6. from pymongo import MongoClient
7. client = MongoClient()
8. db = client['maps']
9. ## Top 5 users found from the data statistics are of this document
10. users = ['Mulad', 'iandees', 'stuckii', 'DavidF', 'sota767']
11. ## Loading only the documents where 'created.user' match the users list
12. q = db.minneapolis.find({'created.user' : {'$in' : users}})
13. temp = []
14. for i in q:
15.     if 'pos' in i:
16.         temp.append(i)
17. ## Creating a dictionary consisting of users as the key and their
18. ## Latitude and Longitude in a list for values of that key
19. data = {}
20. for i in temp:
21.     if i['created']['user'] not in data:
22.         data[i['created']['user']] = []
23.         data[i['created']['user']].append(i['pos'])
24. ## Defining which user to create the graph for
25. user = 'DavidF'
26. lon = []
27. lat = []
28. ## Loading that users positional data for latitude and longitude into the
29. ## 'lat' list and 'lon' list
30. for i in data[user]:
31.     lon.append(i[1])
32.     lat.append(i[0])
33. ## Creating the boundaries of the graph from the original XML max/min lat/lon
34. extent = [-93.3292, -93.1713, 44.915, 45.0644]
35. ## loading the image from the screen shot used in the beginning of this report
36. datafile = cbook.get_sample_data('map.png')
37. img = imread(datafile)
38. ## Using 'lat' values as the y-axis and 'lon' values as the x-axis, we will
39. ## create a scatter plot
40. plt.scatter(lon,lat,zorder=0.2)
41. ## Labeling the graph
42. plt.ylabel('Latitude')
43. plt.xlabel('Longitude')
44. plt.title('User: {} GeoTags'.format(user))
45. ## Producing the graph with the aspect ratio being the original screenshot's
46. ## pixel ratio and the extent of its limits the boundaries stated in the list
47. ## 'extent'
48. plt.imshow(img, zorder = 0, extent = extent, aspect = (1235/float(924)))

```

Lesson 6 Codes

6.1 Iterative Parsing [\[back to top\]](#)

```
1. #!/usr/bin/env python
2. # -*- coding: utf-8 -*-
3. """
4. Your task is to use the iterative parsing to process the map file and
5. find out not only what tags are there, but also how many, to get the
6. feeling on how much of which data you can expect to have in the map.
7. The output should be a dictionary with the tag name as the key
8. and number of times this tag can be encountered in the map as value.
9.
10. Note that your code will be tested with a different data file than the 'example.osm'
11. """
12. import xml.etree.ElementTree as ET
13. import pprint
14.
15. def count_tags(filename):
16.     # YOUR CODE HERE
17.     tree = ET.parse(filename)
18.     data = []
19.     dic = {}
20.     tags = tree.getiterator()
21.     for i in range(len(tags)):
22.         data.append(tags[i].tag)
23.
24.     utags = list(set(data))
25.     for i in utags:
26.         dic[i] = 0
27.
28.     for i in range(len(tags)):
29.         dic[tags[i].tag] += 1
30.     return dic
31.
32.
33.
34. def test():
35.
36.     tags = count_tags('example.osm')
37.     pprint.pprint(tags)
38.     assert tags == {'bounds': 1,
39.                    'member': 3,
40.                    'nd': 4,
41.                    'node': 20,
42.                    'osm': 1,
43.                    'relation': 1,
44.                    'tag': 7,
45.                    'way': 1}
46.
47.
48.
49. if __name__ == "__main__":
50.     test()
```


6.2 Data Model [\[back to top\]](#)

We would like to get the data into a database.
But we need to decide on a data model!
Which of the following models would you prefer
for the following data node?

```
{
  "_id": "2406124091",
  "visible": "true",
  "version": "2",
  "changeset": "17206049",
  "timestamp": "2013-08-03T16:43:42Z",
  "user": "linuxUser16",
  "uid": "1219059",
  "pos": [41.9757030, -87.6921867],
  "tags": [
    { "k": "addr:housenumber", "v": "5157" },
    { "k": "addr:postcode", "v": "60625" },
    { "k": "addr:street", "v": "North Lincoln Ave" },
    { "k": "amenity", "v": "restaurant" }
  ]
}
```

```
{
  "_id": "2406124091",
  "visible": "true",
  "created": {
    "version": "2",
    "changeset": "17206049",
    "timestamp": "2013-08-03T16:43:42Z",
    "user": "linuxUser16",
    "uid": "1219059"
  },
  "pos": [41.9757030, -87.6921867],
  "address": {
    "houenumber": "5157",
    "postcode": "60625",
    "street": "North Lincoln Ave"
  },
  "amenity": "restaurant",
  "cuisine": "mexican",
  "name": "La Cabana De Don Luis",
  "phone": "1 (773)-271-5176"
}
```

```
<node id="2406124091" visible="true" version="2" changeset="17206049"
  timestamp="2013-08-03T16:43:42Z" user="linuxUser16" uid="1219059"
  lat="41.9757030" lon="-87.6921867">
  <tag k="addr:city" v="Chicago"/>
  <tag k="addr:housenumber" v="5157"/>
  <tag k="addr:postcode" v="60625"/>
  <tag k="addr:street" v="North Lincoln Ave"/>
  <tag k="amenity" v="restaurant"/>
  <tag k="cuisine" v="mexican"/>
  <tag k="name" v="La Cabana De Don Luis"/>
  <tag k="outdoor_seating" v="no"/>
  <tag k="phone" v="1 (773)-271-5176"/>
  <tag k="seating" v="no"/>
  <tag k="takeaway" v="yes"/>
</node>
```

6.3 Tag Types [\[back to top\]](#)

```
1. import xml.etree.ElementTree as ET
2. import pprint
3. import re
4.
5. lower = re.compile(r'^([a-z]|_)*$')
6. lower_colon = re.compile(r'^([a-z]|_)*:([a-z]|_)*$')
7. problemchars = re.compile(r'[\+/\&<>;\'\"\\?%#$@,\.\ \t\r\n]')
8.
9.
10. def key_type(element, keys):
11.     if element.tag == 'tag':
12.         k = element.attrib['k']
13.         search = lower.search(k)
14.         if search:
15.             keys['lower'] += 1
16.         else:
17.             search = lower_colon.search(k)
18.             if search:
19.                 keys['lower_colon'] += 1
20.             else:
21.                 search = problemchars.search(k)
22.                 if search:
23.                     keys['problemchars'] += 1
24.                 else:
25.                     keys['other'] += 1
26.         # YOUR CODE HERE
27.         return keys
28.         pass
29.
30.     return keys
31.
32.
33.
34. def process_map(filename):
35.     keys = {"lower": 0, "lower_colon": 0, "problemchars": 0, "other": 0}
36.     for _, element in ET.iterparse(filename):
37.         keys = key_type(element, keys)
38.
39.     return keys
40.
41.
42.
43. def test():
44.     # You can use another testfile 'map.osm' to look at your solution
45.     # Note that the assertions will be incorrect then.
46.     keys = process_map('example.osm')
47.     pprint.pprint(keys)
48.     assert keys == {'lower': 5, 'lower_colon': 0, 'other': 1, 'problemchars': 1}
49.
50.
51. if __name__ == "__main__":
52.     test()
```

6.4 Exploring Users [\[back to top\]](#)

```
1. #!/usr/bin/env python
2. # -*- coding: utf-8 -*-
3. import xml.etree.ElementTree as ET
4. import pprint
5. import re
6. """
7. Your task is to explore the data a bit more.
8. The first task is a fun one - find out how many unique users
9. have contributed to the map in this particular area!
10.
11. The function process_map should return a set of unique user IDs ("uid")
12. """
13.
14. def get_user(element):
15.     return
16.
17.
18. def process_map(filename):
19.     users = set()
20.     for _, element in ET.iterparse(filename):
21.         att = element.attrib
22.         if 'uid' in att:
23.             users.add(att['uid'])
24.
25.
26.     return users
27.
28.
29. def test():
30.
31.     users = process_map('example.osm')
32.     pprint.pprint(users)
33.     assert len(users) == 6
34.
35.
36.
37. if __name__ == "__main__":
38.     test()
```

[\[back to top\]](#)

6.5 Improving Street Names [\[back to top\]](#)

```
1. """
2. Your task in this exercise has two steps:
3.
4. - audit the OSMFILE and change the variable 'mapping' to reflect the changes needed to
   fix
5.   the unexpected street types to the appropriate ones in the expected list.
6.   You have to add mappings only for the actual problems you find in this OSMFILE,
7.   not a generalized solution, since that may and will depend on the particular area y
   ou are auditing.
8. - write the update_name function, to actually fix the street name.
9.   The function takes a string with street name as an argument and should return the f
   ixed name
10.  We have provided a simple test so that you see what exactly is expected
11. """
12. import xml.etree.cElementTree as ET
13. from collections import defaultdict
14. import re
15. import pprint
16.
17. OSMFILE = "example.osm"
18. street_type_re = re.compile(r'\b\S+\.?$', re.IGNORECASE)
19.
20.
21. expected = ["Street", "Avenue", "Boulevard", "Drive", "Court", "Place", "Square", "Lane",
   "Road",
22.             "Trail", "Parkway", "Commons"]
23.
24. # UPDATE THIS VARIABLE
25. mapping = { "St": "Street",
26.             "St.": "Street",
27.             "Ave" : "Avenue",
28.             "Rd." : "Road"}
29.
30.
31. def audit_street_type(street_types, street_name):
32.     m = street_type_re.search(street_name)
33.     if m:
34.         street_type = m.group()
35.         if street_type not in expected:
36.             street_types[street_type].add(street_name)
37.
38.
39. def is_street_name(elem):
40.     return (elem.attrib['k'] == "addr:street")
41.
42.
43. def audit(osmfile):
44.     osm_file = open(osmfile, "r")
45.     street_types = defaultdict(set)
46.     for event, elem in ET.iterparse(osm_file, events=("start",)):
47.
48.         if elem.tag == "node" or elem.tag == "way":
49.             for tag in elem.iter("tag"):
50.                 if is_street_name(tag):
51.                     audit_street_type(street_types, tag.attrib['v'])
52.
```



```
53.     return street_types
54.
55.
56. def update_name(name, mapping):
57.
58.     # YOUR CODE HERE
59.     if (street_type_re.search(name).group() not in expected and
60.         street_type_re.search(name).group() in mapping.keys()):
61.         name = re.sub(street_type_re.search(name).group(),
62.                       mapping[street_type_re.search(name).group()],
63.                       name)
64.
65.     return name
66.
67.
68. def test():
69.     st_types = audit(OSMFILE)
70.     assert len(st_types) == 3
71.     pprint.pprint(dict(st_types))
72.
73.     for st_type, ways in st_types.iteritems():
74.         for name in ways:
75.             better_name = update_name(name, mapping)
76.             print name, "=>", better_name
77.             if name == "West Lexington St.":
78.                 assert better_name == "West Lexington Street"
79.             if name == "Baldwin Rd.":
80.                 assert better_name == "Baldwin Road"
81.
82.
83. if __name__ == '__main__':
84.     test()
```

[\[back to top\]](#)

6.6 Preparing for Database [\[back to top\]](#)

```

1. #!/usr/bin/env python
2. # -*- coding: utf-8 -*-
3. import xml.etree.ElementTree as ET
4. import pprint
5. import re
6. import codecs
7. import json
8.
9. lower = re.compile(r'^([a-z]|_)*$')
10. lower_colon = re.compile(r'^([a-z]|_)*:([a-z]|_)*$')
11. problemchars = re.compile(r'[=+/&<>;\'\"?%$@,\.\ \t\r\n]')
12.
13. CREATED = [ "version", "changeset", "timestamp", "user", "uid"]
14. address = ['addr:housenumber', 'addr:postcode', 'addr:street', 'addr:city']
15.
16. def shape_element(element):
17.     node = {}
18.     if element.tag == "node" or element.tag == "way" :
19.         # YOUR CODE HERE
20.         node['type'] = element.tag
21.         for i in element.attrib.keys():
22.             if i in CREATED:
23.                 node['created'] = {}
24.                 break
25.         for i in element.attrib.keys():
26.             if i in CREATED:
27.                 node['created'][i] = element.attrib[i]
28.             elif i == 'lon' or i == 'lat':
29.                 continue
30.             else:
31.                 node[i] = element.attrib[i]
32.         if 'lat' in element.attrib.keys():
33.             node['pos'] = [float(element.attrib['lat']), float(element.attrib['lon'])]
34.
35.         for i in element:
36.             if 'k' in i.attrib:
37.                 if i.attrib['k'] in address:
38.                     node['address'] = {}
39.                     break
40.         for i in element:
41.             if 'ref' in i.attrib:
42.                 node['node_refs'] = []
43.         for i in element:
44.             if 'k' in i.attrib:
45.                 if i.attrib['k'] in address:
46.                     node['address'][re.sub('addr:', '', i.attrib['k'])] = i.attrib['v']
47.
48.                 elif 'addr:street:' in i.attrib['k']:
49.                     continue
50.                 else:
51.                     node[i.attrib['k']] = i.attrib['v']
52.         if 'ref' in i.attrib:
53.             node['node_refs'].append(i.attrib['ref'])
54.
55.     return node
56. else:

```

```
55.         return None
56.
57.
58. def process_map(file_in, pretty = False):
59.     # You do not need to change this file
60.     file_out = "{0}.json".format(file_in)
61.     data = []
62.     with codecs.open(file_out, "w") as fo:
63.         for _, element in ET.iterparse(file_in):
64.             el = shape_element(element)
65.             if el:
66.                 data.append(el)
67.                 if pretty:
68.                     fo.write(json.dumps(el, indent=2)+"\n")
69.                 else:
70.                     fo.write(json.dumps(el) + "\n")
71.     return data
72.
73. def test():
74.     # NOTE: if you are running this code on your computer, with a larger dataset,
75.     # call the process_map procedure with pretty=False. The pretty=True option adds
76.     # additional spaces to the output, making it significantly larger.
77.     data = process_map('example.osm', True)
78.     #pprint.pprint(data)
79.
80.     correct_first_elem = {
81.         "id": "261114295",
82.         "visible": "true",
83.         "type": "node",
84.         "pos": [41.9730791, -87.6866303],
85.         "created": {
86.             "changeset": "11129782",
87.             "user": "bbmiller",
88.             "version": "7",
89.             "uid": "451048",
90.             "timestamp": "2012-03-28T18:31:23Z"
91.         }
92.     }
93.     assert data[0] == correct_first_elem
94.     assert data[-1]["address"] == {
95.         "street": "West Lexington St.",
96.         "housenumber": "1412"
97.     }
98.     assert data[-
99.         1]["node_refs"] == [ "2199822281", "2199822390", "2199822392", "2199822369",
100.                             "2199822370", "2199822284", "2199822281"]
101.
102.     if __name__ == "__main__":
103.         test()
```

[\[back to top\]](#)