**Project P2: Data Wrangle OpenStreetMaps Data**
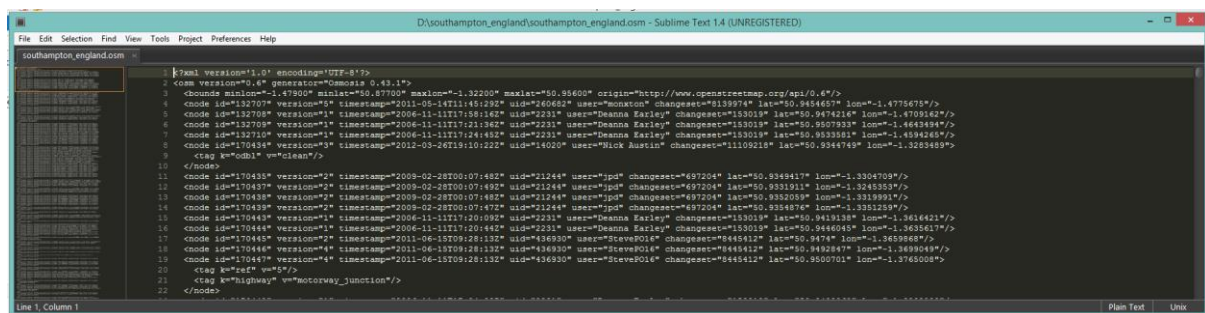
# Olakunle Kuye

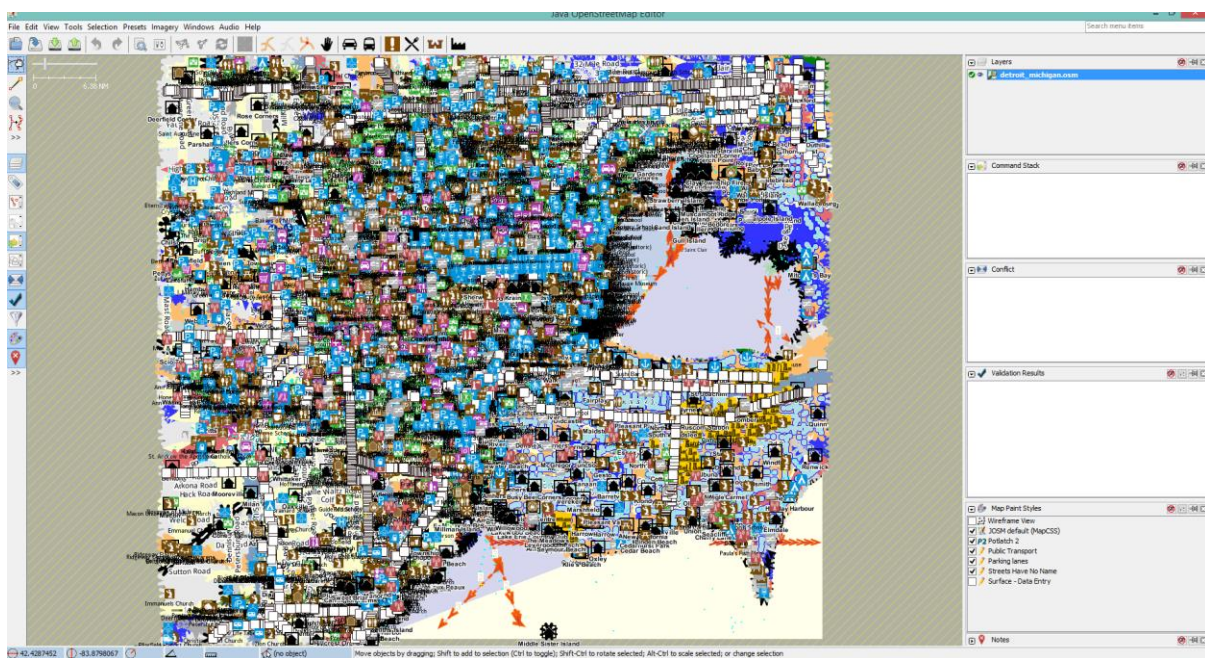# Chosen Map Area

I downloaded the Austin Texas from mapzen metro – extracts:

.

| | | | | | |
|---|---|---|---|---|---|
| Austin, Texas | | | | | |
| OSM PBF | OSM XML | OSM2PGSQL SHP | OSM2PGSQL GEOJSON | IMPOSM SHP | IMPOSM GEOJSON |
| 7.5 MB | 12.1 MB | 15.5 MB | 11.1 MB | 13.4 MB | 15.9 MB |

## Auditing the data

After downloading the OSM XML data, I parsed it using the code I used for the iterative Parsing exercise.

With the result being:

```
{'bounds': 1,
 'member': 13056,
 'nd': 889130,
 'node': 768631,
 'osm': 1,
 'relation': 1283,
 'tag': 535875,
 'way': 79692}
```



I ran the code I used for the Tag Type exercise whereby checks on the dataset were made to determine if they are certain patterns in the tags present in the data.

This is where there are variables declared to search for the following three regular expressions: **lower**, **lower_colon**, and **problemchars**.

**lower**: This variable searches for strings containing lower case characters

**lower_colon**: Searches for strings containing lower case characters and a single colon within the string

**problemchars**: Searches for characters that cannot be used within keys in **MongoDB**.


lower = re.compile(r'^([a-z]|_)*$')

lower_colon = re.compile(r'^([a-z]|_)*:([a-z]|_)*$')

problemchars = re.compile(r'[=\+/&<>;\'"\?%#$@\,\.\t\r\n]')


As the image below shows there are a number of tags with multiple children tags that begin with address .

These will be used with the lower_colon regex to find these keys so that it's possible to build a single address document within a larger json document.

Here is a partial output and a screen shot from running the code used in the Tag Type exercise:

{'lower': 226464, 'lower_colon': 299286, 'other': 10123, 'problemchars': 2}

By modifying the code used in Tag Type exercise to build a set of unique ids found in the xml, outputting the length of the set that represents the number of users making edits in the chosen map area.

```python
def process_map(filename):
    users = set()
    for _, element in ET.iterparse(filename):
        if "uid" in element.attrib and element.tag in ('node', 'way'):
            users.add(element.attrib["uid"])
    return users


def test():
    # You can use another testfile 'map.osm' to look at your solution
    # Note that the assertions will be incorrect then.
    #keys = process_map('D:\\southampton_england\\southampton_england.osm')
    users = process_map('D:\\southampton_england\\southampton_england.osm')

    print len(users)
```

## The output returned was: 910


## Problems with the Data


**Street Names**

The majority of this project will be devoted to auditing and cleaning street names seen within the OSM XML.

Street types used by users in the process of mapping are quite often abbreviated. I will attempt to find these abbreviations and replace them with their full text form. The plan of action is as follows:

The data contained in the OSM XML file will be audited to clean the street types used by users, which are often abbreviated. These abbreviation's will be found and replaced with the full street name.

In order to do this, the following steps will be taken:

- A regex designed to match the last token in a string will be developed, finding the street type in an address.
- I would then proceeded to build a list of expected street types that don't need to be cleaned.
- The XML file will be parsed for tag elements with the **k="addr:street" attributes** .


- A search will be conducted using the regex on the value v attribute of the street name string.
- A dictionary will be built with keys that match the street types in the regex, setting the street names where a specific key was found as a value, this

enables the determination of what needs to be cleaned.

- Another dictionary will be built to contain a map from an offending street type. This along with another regex and a function that will be developed and used to match the offending street types will be used to clean the data. The function will return a cleaned string.

```python
import xml.etree.cElementTree as ET
from collections import defaultdict
import re
import pprint

OSMFILE = "example.osm"
street_type_re = re.compile(r'\b\S+\.?$', re.IGNORECASE)


expected = ["Street", "Avenue", "Boulevard", "Drive",
"Court", "Place", "Square", "Lane", "Road",
            "Trail", "Parkway", "Commons"]
```

The **audit_string** takes in the dictionary of street types, an audit string a regex match and a list of expected street types.

The function will search the string passed to it for a regex and if it finds a match will that is not in the dictionary used to hold expected dictionaries it will be added to it.

```python
def audit_string(match_set_dict, string_to_audit, regex,
expected_matches):
```

```
    m = regex.search(string_to_audit)
    if m:
        match_string = m.group()
        if match_string not in expected_matches:

match_set_dict[match_string].add(string_to_audit)
```

An audit function is defined and it's used to parse and audit street names. The function audits tag elements where **k="addr:street"** and it also matches the tag filter function . The audit function takes in a regex and a list of expected matches.

```
def audit(osmfile, tag_filter, regex, expected_matches =
[]):
    osm_file = open(osmfile, "r")
    match_sets = defaultdict(set)
    # iteratively parse the mapping xml
    for event, elem in ET.iterparse(osm_file,
events=("start",)):
        # node and way tags are of special interest
        if elem.tag == "node" or elem.tag == "way":
            # iterate the "tag" tags within a node or way
            for tag in elem.iter("tag"):
                if tag_filter(tag):
                    audit_string(match_sets,
tag.attrib['v'], regex, expected_matches)
    return match_sets
```

The function **is_street_name** determines if an element contains an attribute **k="addr:street".** I will use **is_street_name** as the **tag_filter** when I call the audit function to audit street names.

```
def is_street_name(elem):
```

```python
        return (elem.attrib['k'] == "addr:street")
```

# Here is a pretty print of the output of the audit function.

```python
street_types = audit(OSMFILE, tag_filter = is_street_name,
regex = street_type_re,
expected_matches = expected)
pprint.pprint(dict(street_types))
```

```
{'100': set(['Jollyville Road Suite 100', 'Old Jollyville Road, Suite 100']),
 '101': set(['4207 James Casey st #101']),
 '104': set(['11410 Century Oaks Terrace Suite #104']),
 '1100': set(['Hwy 290 W, Bldg 1100']),
 '12': set(['Ranch to Market Road 12']),
 '120': set(['Building B Suite 120']),
 '150': set(['Farm-to-Market Road 150']),
 '1626': set(['F.M. 1626', 'FM 1626', 'Farm-to-Market Road 1626']),
 '163': set(['Bee Cave Road Suite 163']),
 '170': set(['County Road 170']),
 '1825': set(['FM 1825']),
 '1826': set(['Farm To Market Road 1826', 'Ranch to Market Road 1826']),
 '183': set(['N HWY 183', 'U.S. 183', 'US 183', 'US Highway 183']),
 '203': set(['West Ben White Boulevard #203',
            'West Ben White Boulevard, #203']),
 '210': set(['Greystone Dr Ste 210']),
 '275': set(['West 11th Street, Suite 275']),
 '290': set(['E Hwy 290',
            'East Hwy 290',
            'Highway 290',
            'U.S. 290',
            'US 290',
            'W HWY 290',
            'W Highway 290',
            'W Hwy 290',
            'W. Highway 290',
            'West Highway 290']),
 '3': set(['Shoal Creek Blvd, Bld 3']),
 '300': set(['321 W Ben White Blvd #300', 'West 5th Street #300']),
 '301': set(['Hwy 290 W. Ste. 301', 'S Bell Blvd., Suite 301']),
 '3177': set(['FM 3177']),
 '35': set(['9600 S IH 35',
            'I 35',
            'Interstate 35',
            'Interstate Highway 35',
            'N I H 35',
            'N Interstate Highway 35',
            'N. IH 35',
            'North IH 35',
            'North Interstate Highway 35',
            'S Interstate 35',
            'S Interstate Highway 35',
```

```
                     'South Interstate 35']),
'4': set(['4201 South Congress Avenue #4']),
'400': set(['North Lakeline Boulevard, Ste. 400']),
'406': set(['North Highway 183 #406']),
'45': set(['TX 45']),
'535': set(['FM 535']),
'6': set(['N IH 35 Bldg. 6']),
'620': set(['FM 620',
            'N FM 620',
            'RM 620',
            'RR 620',
            'Ranch Road 620',
            'Ranch-to-Market Road 620']),
'685': set(['FM 685', 'Farm to Market 685']),
'7': set(['N I H 35 Bldg 7']),
'71': set(['TX 71', 'W Hwy 71', 'West Highway 71']),
'78738': set(['12200 Bee Cave Pkwy, Bee Cave, TX 78738']),
'79': set(['U.S. 79']),
'969': set(['FM 969']),
'A-15': set(['W. Anderson Lane, Suite A-15']),
'A500': set(['Burnet Rd Ste, A500']),
'Austin': set(['Research Blvd, Austin']),
'Ave': set(['Congress Ave',
            'E. University Ave',
            'Round Rock Ave',
            'S Congress Ave',
            'South Congress Ave']),
'Ave.': set(['Woodrow Ave.']),
'Avene': set(['South Congress Avene']),
'B': set(['Avenue B']),
'B100': set(['North Lamar Boulevard suite #B100']),
'Blvd': set(['10000 Research Blvd',
             'Airport Blvd',
             'Escarpment Blvd',
             'Industrial Blvd',
             'N Bell Blvd',
             'North Lamar Blvd',
             'North Research Blvd',
             'Research Blvd',
             'South Bell Blvd',
             'South Lamar Blvd',
             'Stonelake Blvd',
             'Swenson Farms Blvd',
             'West Ben White Blvd',
             'West Louis Henna Blvd']),
'Blvd.': set(['East Old Settlers Blvd.',
              'Research Blvd.',
              'South Exposition Blvd.',
              'South Lamar Blvd.',
              'Tandem Blvd.',
              'W North Loop Blvd.']),
'C-200': set(['9600 IH 35 C-200']),
'CR': set(['Rain Lily CR']),
'Cannon': set(['West William Cannon']),
'Chavez': set(['West Cesar Chavez']),
'Cir': set(['Courtland Cir', 'Dry Run Cir', 'Tanager Cir']),
'Circle': set(['Coronado Circle', 'Morado Circle']),
'Cove': set(['Glacier Cove',
             'Knoll Cove',
```

```
                'Mustang Cove',
                'Rain Dance Cove',
                'Yucca Cove']),
 'Crossing': set(['Esperanza Crossing', 'Kyle Crossing', 'Nuckols Crossing']),
 'Ct': set(['Briarpatch Ct', 'Grapevine Ct']),
 'Cv': set(['Blue Sky Cv',
                'Clear Pond Cv',
                'Enchanted Cv',
                'Mallard Cv',
                'Red River Cv']),
 'D1': set(['W Hwy 71, STE D1']),
 'Dr': set(['Blazing Star Dr',
                'Carranzo Dr',
                'Country Creek Dr',
                'Executive Center Dr',
                'Forest Creek Dr',
                'Jack Nicklaus Dr',
                'Jimmy Clay Dr',
                'Kirkglen Dr',
                'Kyle Center Dr',
                'Round Rock West Dr',
                'W Lake Dr',
                'W William Cannon Dr',
                'Walter Seaholm Dr',
                'Wild Rose Dr',
                'Wild Turkey Dr']),
 'Dr.': set(['11928 Stonehollow Dr.', 'Sportsplex Dr.', 'Stonebridge Dr.']),
 'East': set(['Hwy 290 East', 'US 290 East']),
 'Expressway': set(['North Mo-Pac Expressway',
                    'North MoPac Expressway',
                    'North Mopac Expressway',
                    'South MoPac Expressway']),
 'Expwy': set(['N Mopac Expwy']),
 'FM1431': set(['FM1431']),
 'Fork': set(['Roaring Fork']),
 'Fortuna': set(['Via Fortuna']),
 'Gonzales': set(['Gonzales']),
 'Grande': set(['Rio Grande']),
 'Guadalupe': set(['Guadalupe']),
 'H': set(['Avenue H']),
 'Highway': set(['Capital of Texas Highway',
                'North Capitol of Texas Highway',
                'S Capital of Texas Highway',
                'South Capital of Texas Highway']),
 'Hollow': set(['Lost Oasis Hollow']),
 'Hwy': set(['S Capital of Texas Hwy']),
 'I35': set(['I35']),
 'IH-35': set(['N IH-35', 'South IH-35']),
 'IH35,': set(['N. IH35,']),
 'J': set(['Business Interstate Highway 35 J']),
 'Jr': set(['Robert Martinez Jr']),
 'Liberty': set(['W. Liberty']),
 'Ln': set(['Honeybee Ln', 'Rockwood Ln', 'W Anderson Ln']),
 'Loop': set(['Amberwood Loop', 'E. North Loop', 'Mirasol Loop']),
 'N': set(['Capital of TX Hwy N', 'Farm to Market 620 N']),
 'North': set(['Ashwood North',
                'Cuernavaca Drive North',
                'FM 620 North',
                'IH 35 North',
```

```
                    'Interstate Highway 35 Service Road North',
                    'Ranch Road 620 North']),
'Northwest': set(['Carlos G Parker Boulevard Northwest']),
'Oaks': set(['Sierra Oaks']),
'Oltorf': set(['East Oltorf']),
'Pass': set(['Continental Pass', 'Turks Cap Pass']),
'Path': set(['Daisy Path']),
'Pflugerville': set(['N. IH 35 Pflugerville']),
'Picadilly': set(['Picadilly']),
'Pkwy': set(['Riata Trace Pkwy', 'Roger Hanks Pkwy']),
'Quarry': set(['Quarry']),
'RD': set(['Grassy Field RD']),
'Rd': set(['Barley Rd',
           'Barton Springs Rd',
           'Big Meadow Rd',
           'Burnet Rd',
           'Commons Rd',
           'Dry Creek Rd',
           'Elderberry Rd',
           'Grassy Field Rd',
           'N Interstate 35 Frontage Rd',
           'Open Sky Rd',
           'Saddleback Rd',
           'Spicewood Springs Rd',
           'W US Highway 290 Service Rd']),
'Real': set(['El Camino Real']),
'River': set(['Red River']),
'SB': set(['South I-35 Service SB']),
'Seco': set(['Arroyo Seco']),
'South': set(['Venture Boulevard South']),
'Speedway': set(['Speedway']),
'Spicewood': set(['State Highway 71 W Spicewood']),
'St': set(['Duval St',
           'E 43rd St',
           'E Live Oak St',
           'E Oltorf St',
           'Guadalupe St',
           'N Main St',
           'Pecan St',
           'Rio Grande St',
           'S 1st St',
           'W 6th St',
           'W Annie St',
           'West Lynn St']),
'St.': set(['E 38th 1/2 St.', 'E. 43rd St.', 'Guadalupe St.', 'Pecan St.']),
'Sutton': set(['Sutton']),
'Terrace': set(['Norseman Terrace']),
'W': set(['East HWY 290 W', 'Highway 71 W', 'Hwy 290 W']),
'Way': set(['Brista Way',
            'Dell Way',
            'Desert Willow Way',
            'Mountain Laurel Way',
            'Pure Brook Way',
            'Running Water Way',
            'Technology Way',
            'Whispering Wind Way',
            'Wild Plum Way',
            'Winecap Way']),
'West': set(['Old 2243 West', 'US Highway 290 West']),
```

```
'Willo': set(['Switch Willo']),
'Wow': set(['Pow Wow']),
'lane': set(['brigadoon lane']),
'street': set(['East main street', 'South 1st street']),
'suite#L131': set(['N. Lamar suite#L131']),
'texas': set(['9901 N Capital Of texas'])}
```

There were a number unexpected abbreviated street types as well as cardinal directions and highway directions. I created a function to replace the abbreviated street types taking a regex to search for, updating a mapping directory.

```python
def update(string_to_update, mapping, regex):
    m = regex.search(string_to_update)
    if m:
        match = m.group()
        if match in mapping:
            string_to_update = re.sub(regex,
mapping[match], string_to_update)
    return string_to_update
```

The results of the audit function were used to build a dictionary to map abbreviations to their full and clean representation.

```python
map_street_types = \
{
"Ave" : "Avenue",
"BLVD" : "Boulevard",
"Blvd" : "Boulevard",
"Blvd." : "Boulevard",
"Cir" : "Circle",
"Dr" : "Drive",
"Ln" : "Lane",
"Ln" : "Lane",
"Pkwy" : "Parkway",
```

```
"Rd" : "Road",
"Rd." : "Road",
"St" : "Street",
"St." : "Street"
}
```

# The keys are replaced where they appear in a string using:

```
bad_streets =
"|".join(map_street_types.keys()).replace('.', '')

street_type_updater_re = re.compile(r'\b(' + bad_streets +
r')\b\.?', re.IGNORECASE)
```

# Here are the results showing the abbreviated street types have been updated as expected:

```
W US Highway 290 Service Rd => W US Highway 290 Service Road
Saddleback Rd => Saddleback Road
Barley Rd => Barley Road
Open Sky Rd => Open Sky Road
Elderberry Rd => Elderberry Road
N Interstate 35 Frontage Rd => N Interstate 35 Frontage Road
Spicewood Springs Rd => Spicewood Springs Road
Big Meadow Rd => Big Meadow Road
Commons Rd => Commons Road
Burnet Rd => Burnet Road
Barton Springs Rd => Barton Springs Road
Dry Creek Rd => Dry Creek Road
Grassy Field Rd => Grassy Field Road
Roger Hanks Pkwy => Roger Hanks Parkway
Riata Trace Pkwy => Riata Trace Parkway
Tanager Cir => Tanager Circle
Dry Run Cir => Dry Run Circle
Courtland Cir => Courtland Circle
Congress Ave => Congress Avenue
E. University Ave => E. University Avenue
South Congress Ave => South Congress Avenue
Round Rock Ave => Round Rock Avenue
S Congress Ave => S Congress Avenue
Tandem Blvd. => Tandem Boulevard
South Lamar Blvd. => South Lamar Boulevard
South Exposition Blvd. => South Exposition Boulevard
```

```
W North Loop Blvd. => W North Loop Boulevard
East Old Settlers Blvd. => East Old Settlers Boulevard
Research Blvd. => Research Boulevard
Pecan St. => Pecan Street
Guadalupe St. => Guadalupe Street
E. 43rd St. => E. 43rd Street
E 38th 1/2 St. => E 38th 1/2 Street
Honeybee Ln => Honeybee Lane
W Anderson Ln => W Anderson Lane
Rockwood Ln => Rockwood Lane
Kyle Center Dr => Kyle Center Drive
Kirkglen Dr => Kirkglen Drive
Country Creek Dr => Country Creek Drive
Blazing Star Dr => Blazing Star Drive
Walter Seaholm Dr => Walter Seaholm Drive
Wild Turkey Dr => Wild Turkey Drive
W Lake Dr => W Lake Drive
W William Cannon Dr => W William Cannon Drive
Executive Center Dr => Executive Center Drive
Jack Nicklaus Dr => Jack Nicklaus Drive
Wild Rose Dr => Wild Rose Drive
Round Rock West Dr => Round Rock West Drive
Carranzo Dr => Carranzo Drive
Forest Creek Dr => Forest Creek Drive
Jimmy Clay Dr => Jimmy Clay Drive
Pecan St => Pecan Street
N Main St => N Main Street
S 1st St => S 1st Street
E Oltorf St => E Oltorf Street
E 43rd St => E 43rd Street
E Live Oak St => E Live Oak Street
W Annie St => W Annie Street
Rio Grande St => Rio Grande Street
W 6th St => W 6th Street
West Lynn St => West Lynn Street
Guadalupe St => Guadalupe Street
Duval St => Duval Street
West Louis Henna Blvd => West Louis Henna Boulevard
Airport Blvd => Airport Boulevard
Swenson Farms Blvd => Swenson Farms Boulevard
West Ben White Blvd => West Ben White Boulevard
N Bell Blvd => N Bell Boulevard
North Lamar Blvd => North Lamar Boulevard
Industrial Blvd => Industrial Boulevard
Stonelake Blvd => Stonelake Boulevard
Escarpment Blvd => Escarpment Boulevard
South Lamar Blvd => South Lamar Boulevard
South Bell Blvd => South Bell Boulevard
North Research Blvd => North Research Boulevard
Research Blvd => Research Boulevard
10000 Research Blvd => 10000 Research Boulevard
```

# The cardinal directions North, South, East and West appear to be abbreviated in

the data, so using a similar technique to set the correct the abbreviated cardinals by creating a regex to match NSEW at the begging of a string, followed by an optional period.

Here is the output of auditing using

**cardinal_directions = audit(OSMFILE, is_street_name, cardinal_dir_re)**

**outputs :**

```
{'E': set(['E 38th 1/2 St.',
          'E 43rd St',
          'E Hwy 290',
          'E Live Oak St',
          'E Main Street',
          'E Oltorf St',
          'E Palm Valley Boulevard',
          'E Riverside Drive']),
 'E.': set(['E. 43rd St.', 'E. North Loop', 'E. University Ave']),
 'N': set(['N Bell Blvd',
          'N FM 620',
          'N HWY 183',
          'N I H 35',
          'N I H 35 Bldg 7',
          'N IH 35 Bldg. 6',
          'N IH-35',
          'N Interstate 35 Frontage Rd',
          'N Interstate Highway 35',
          'N Lake Creek Drive',
          'N Main St',
          'N Mopac Expwy']),
 'N.': set(['N. IH 35',
          'N. IH 35 Pflugerville',
          'N. IH35,',
          'N. Lamar suite#L131']),
 'S': set(['S 1st St',
          'S Bell Blvd., Suite 301',
          'S Capital of Texas Highway',
          'S Capital of Texas Hwy',
          'S Congress Ave',
          'S Interstate 35',
          'S Interstate Highway 35']),
 'W': set(['W 35th Street',
          'W 6th St',
```

```
        'W Anderson Ln',
        'W Annie St',
        'W HWY 290',
        'W Highway 290',
        'W Hwy 290',
        'W Hwy 71',
        'W Hwy 71, STE D1',
        'W Lake Dr',
        'W North Loop Blvd.',
        'W Parmer Lane',
        'W US Highway 290 Service Rd',
        'W William Cannon Dr']),
 'W.': set(['W. Anderson Lane, Suite A-15',
        'W. Highway 290',
        'W. Liberty',
        'W. Slaughter Lane'])}
```

The output shows the cardinals at the beginning street names, but creating an exhaustive mapping will resolve this using:

```
map_cardinal_directions = \
    {
        "E" : "East",
        "E." : "East",
        "N" : "North",
        "N." : "North",
        "S" : "South",
        "S." : "South",
        "W" : "West",
        "W." : "West"
    }
```

Whereby the keys will be replaced anywhere in the string using a regex: exhaustive mapping

```
bad_directions =
"|".join(map_cardinal_directions.keys()).replace('.', '')
```

```python
cardinal_dir_updater_re = re.compile(r'\b(' +
bad_directions + r')\b\.?', re.IGNORECASE)
```

# I modified the code used in the audit.py for exercise 6, traversing the dictionary applying updates for street types and the appropriate cardinal direction.

## Using:

```python
    cardinal_directions = audit(OSMFILE, is_street_name,
cardinal_dir_re)
    for cardinal_direction, ways in
cardinal_directions.iteritems():
        if cardinal_direction in map_cardinal_directions:
            for name in ways:
                better_name = update(name,
map_street_types, street_type_updater_re)
                best_name = update(better_name,
map_cardinal_directions, cardinal_dir_updater_re)
                print name, "=>", better_name, "=>",
best_name
```

## This produces the output:

```
E. University Ave => E. University Avenue => East University Avenue
E. 43rd St. => E. 43rd Street => East 43rd Street
E. North Loop => E. North Loop => East North Loop
E Oltorf St => E Oltorf Street => East Oltorf Street
E Live Oak St => E Live Oak Street => East Live Oak Street
E Main Street => E Main Street => East Main Street
E 43rd St => E 43rd Street => East 43rd Street
E Riverside Drive => E Riverside Drive => East Riverside Drive
E Hwy 290 => E Hwy 290 => East Hwy 290
E 38th 1/2 St. => E 38th 1/2 Street => East 38th 1/2 Street
E Palm Valley Boulevard => E Palm Valley Boulevard => East Palm Valley Boulevard
N. IH35, => N. IH35, => North IH35,
N. Lamar suite#L131 => N. Lamar suite#L131 => North Lamar suite#L131
```

```
N. IH 35 => N. IH 35 => North IH 35
N. IH 35 Pflugerville => N. IH 35 Pflugerville => North IH 35 Pflugerville
W. Liberty => W. Liberty => West Liberty
W. Slaughter Lane => W. Slaughter Lane => West Slaughter Lane
W. Highway 290 => W. Highway 290 => West Highway 290
W. Anderson Lane, Suite A-15 => W. Anderson Lane, Suite A-15 => West Anderson
Lane, Suite A-15
N Interstate Highway 35 => N Interstate Highway 35 => North Interstate Highway 35
N Bell Blvd => N Bell Boulevard => North Bell Boulevard
N Interstate 35 Frontage Rd => N Interstate 35 Frontage Road => North Interstate
35 Frontage Road
N IH 35 Bldg. 6 => N IH 35 Bldg. 6 => North IH 35 Bldg. 6
N FM 620 => N FM 620 => North FM 620
N IH-35 => N IH-35 => North IH-35
N I H 35 Bldg 7 => N I H 35 Bldg 7 => North I H 35 Bldg 7
N Lake Creek Drive => N Lake Creek Drive => North Lake Creek Drive
N I H 35 => N I H 35 => North I H 35
N Main St => N Main Street => North Main Street
N HWY 183 => N HWY 183 => North HWY 183
N Mopac Expwy => N Mopac Expwy => North Mopac Expwy
S Bell Blvd., Suite 301 => S Bell Boulevard, Suite 301 => South Bell Boulevard,
Suite 301
S 1st St => S 1st Street => South 1st Street
S Interstate 35 => S Interstate 35 => South Interstate 35
S Capital of Texas Highway => S Capital of Texas Highway => South Capital of Texas
Highway
S Congress Ave => S Congress Avenue => South Congress Avenue
S Interstate Highway 35 => S Interstate Highway 35 => South Interstate Highway 35
S Capital of Texas Hwy => S Capital of Texas Hwy => South Capital of Texas Hwy
W US Highway 290 Service Rd => W US Highway 290 Service Road => West US Highway
290 Service Road
W Hwy 290 => W Hwy 290 => West Hwy 290
W 35th Street => W 35th Street => West 35th Street
W Hwy 71, STE D1 => W Hwy 71, STE D1 => West Hwy 71, STE D1
W Parmer Lane => W Parmer Lane => West Parmer Lane
W Annie St => W Annie Street => West Annie Street
W North Loop Blvd. => W North Loop Boulevard => West North Loop Boulevard
W Anderson Ln => W Anderson Lane => West Anderson Lane
W Lake Dr => W Lake Drive => West Lake Drive
W 6th St => W 6th Street => West 6th Street
W William Cannon Dr => W William Cannon Drive => West William Cannon Drive
W Hwy 71 => W Hwy 71 => West Hwy 71
W Highway 290 => W Highway 290 => West Highway 290
W HWY 290 => W HWY 290 => West HWY 290
```

# To load data into MongoDB I transformed the data.

# The rules I chose included the following:

To only process two types of top level tags "node" and "way".

Attributes except from created: "version, changeset, timestamp, user, uid" should be added a key.
The longitudes and latitudes will be added to an array "pos", to be used in a geospatial indexing, ensuring the values inside the "pos" are floats and not strings.
Problematics characters in the second level of the tag "k" would be ignored.
The dictionary in the second tag level with values that starts with "addr" but containing ":" should be added.
The presence of a second tag containing character ":" should be ignored .the character separates the type and direction of a street.

For the transformation will make use of a function that processes an element. Within the function there is an update function that uses regex and mapping dictionaries to clean addresses.

**Overview Of the Data**

The downloaded file is 176.584186 MB
The json file is 256.771823 MB

The diagram below shows the json file
being imported into MongoDB.



The diagram below shows the json file
being imported into MongoDB.



Using the following bit of code to
connect to the MongoDB database, I was
able to determine the number of
documents.

```python
db_name = "osm"

client =
MongoClient('192.168.255.133:27017')
db = client[db_name]

collection = db.austin_texas

print collection.count()
```

**Result : 848323**

The number of unique users can be determined by modifying the previous code to include:

```python
Print
len(collection.distinct('created.user'))
```

**Result: 910**
The top contributing user was known by running the following code :

```python
db.austin_texas.aggregate([{"$group" :
{"_id" : "$created.user", "count" :
{"$sum" : 1}}}, {"$sort" : {"count" : -
1}},{"$limit" : 1}])
```

{ "_id" : "woodpeck_fixbot", "count" :
241116 }



The number of nodes and ways can be known
by running the following code in the
database (within MongoDB) :

```
aggregate({"$group" : {"_id" : "$type",
"count" : {"$sum" : 1}}})['result']
```

To give :
{ "_id" : "way", "count" : 79692 }
{ "_id" : "node", "count" : 768631 }

To determine the number of containing a street address, I modified the previous code (within MongoDB) to use:
db.austin_texas.find({"address.street" : {"$exists" : 1}}).count()

**Results: 2014**



To display a list of cleaned zip codes I used the following code (within MongoDB) :
db.austin_texas.aggregate([{"$match" : {"address.postcode" : {"$exists" :

```
1}}},{"$group" : {"_id" :
"$address.postcode", "count" : {"$sum" :
1}}}, {"$sort" : {"count" : -1}}])
```

**Results :**
```
{ "_id" : "78704", "count" : 63 }
{ "_id" : "78705", "count" : 54 }
{ "_id" : "78757", "count" : 48 }
{ "_id" : "78681", "count" : 47 }
{ "_id" : "78640", "count" : 46 }
{ "_id" : "78759", "count" : 46 }
{ "_id" : "78702", "count" : 40 }
{ "_id" : "78701", "count" : 39 }
{ "_id" : "78746", "count" : 38 }
{ "_id" : "78745", "count" : 32 }
```



To determine the top five cites most
common using the following code (within
MongoDB) :

```
db.austin_texas.aggregate([{"$match" :
{"address.city" : {"$exists" : 1}}},
{"$group" : {"_id" : "$address.city",
"count" : {"$sum" : 1}}}, {"$sort" :
{"count" : -1}}, {"$limit" : 5}])
```

**Results :**

```
{ "_id" : "Austin", "count" : 541 }
{ "_id" : "Round Rock", "count" : 65 }
{ "_id" : "Austin, TX", "count" : 48 }
{ "_id" : "Kyle", "count" : 47 }
{ "_id" : "Buda", "count" : 21 }
```

To determine the top ten amenities using
: db.austin_texas.aggregate([{"$match" :
{"amenity" : {"$exists" : 1}}}, {"$group"
: {"_id" : "$amenity", "count" : {"$sum"
: 1}}}, {"$sort" : {"count" : -1}},
{"$limit" : 10}])

**Results :**

```
{ "_id" : "parking", "count" : 1845 }
{ "_id" : "restaurant", "count" : 684 }
{ "_id" : "waste_basket", "count" : 591 }
{ "_id" : "school", "count" : 562 }
{ "_id" : "fast_food", "count" : 503 }
```

```
{ "_id" : "place_of_worship", "count" :
488 }
{ "_id" : "fuel", "count" : 369 }
{ "_id" : "bench", "count" : 349 }
{ "_id" : "shelter", "count" : 231 }
{ "_id" : "bank", "count" : 150 }
```

Top religions with denominations can be determined by running the following code in mongodb to give :

```
db.austin_texas.aggregate([{"$match" :
{"amenity" :
"place_of_worship"}},{"$group" : {"_id" :
{"religion" : "$religion", "denomination"
: "$denomination"}, "count" : {"$sum" :
1}}},{"$sort" : {"count" : -1}}])
```

**Results :**

```
{ "_id" : { "religion" : "christian" },
"count" : 196 }
{ "_id" : { "religion" : "christian",
"denomination" : "baptist" }, "count" :
108 }
{ "_id" : { "religion" : "christian",
"denomination" : "lutheran" }, "count" :
30 }
```

```
{ "_id" : { "religion" : "christian",
"denomination" : "methodist" }, "count" :
29 }
{ "_id" : {   }, "count" : 25 }
{ "_id" : { "religion" : "christian",
"denomination" : "catholic" }, "count" :
22 }
{ "_id" : { "religion" : "christian",
"denomination" : "presbyterian" },
"count" : 21 }
{ "_id" : { "religion" : "christian",
"denomination" : "pentecostal" }, "count"
: 12 }
{ "_id" : { "religion" : "christian",
"denomination" : "mormon" }, "count" : 6
}
{ "_id" : { "religion" : "buddhist" },
"count" : 5 }
{ "_id" : { "religion" : "christian",
"denomination" : "roman_catholic" },
"count" : 5 }
{ "_id" : { "religion" : "christian",
"denomination" : "seventh_day_adventist"
}, "count" : 4 }
{ "_id" : { "religion" : "jewish" },
"count" : 3 }
{ "_id" : { "religion" : "christian",
"denomination" : "jehovahs_witness" },
"count" : 2 }
```

```
{ "_id" : { "religion" : "christian",
"denomination" : "anglican" }, "count" :
2 }
{ "_id" : { "religion" : "muslim" },
"count" : 2 }
{ "_id" : { "religion" :
"unitarian_universalist" }, "count" : 2 }
{ "_id" : { "religion" : "christian",
"denomination" : "episcopal" }, "count" :
2 }
{ "_id" : { "religion" : "christian",
"denomination" : "mormon3" }, "count" : 1
}
{ "_id" : { "denomination" : "baptist" },
"count" : 1 }
```

The top 10 leisure within the data can be
determined using

```
db.austin_texas.aggregate([{"$match" :
{"leisure" : {"$exists" : 1}}}, {"$group"
: {"_id" : "$leisure", "count" : {"$sum"
: 1}}}, {"$sort" : {"count" : -1}},
{"$limit" : 10}])
```

**Results :**

```
{ "_id" : "pitch", "count" : 828 }
```

```
{ "_id" : "park", "count" : 421 }
{ "_id" : "playground", "count" : 104 }
{ "_id" : "sports_centre", "count" : 70 }
{ "_id" : "golf_course", "count" : 70 }
{ "_id" : "track", "count" : 64 }
{ "_id" : "swimming_pool", "count" : 55 }
{ "_id" : "garden", "count" : 21 }
{ "_id" : "stadium", "count" : 20 }
{ "_id" : "slipway", "count" : 12 }
```

## About the Dataset:

I believe the data is well formed and structured making it suitable for extension whereby user will be able to add reviews with regards to bicycle routes, restaurants and schools etc.