# Michael C. Villadelrey

# **OpenStreetMap Data Case Study**

#### Map Area

Honolulu, Hawaii, United States

- https://mapzen.com/data/metro-extracts/metro/honolulu hawaii/
- <a href="https://www.openstreetmap.org/relation/119231">https://www.openstreetmap.org/relation/119231</a>

Being largest City of Hawaii, Honolulu is the place where you can get the most things offered in Hawaii. Historical landmarks, world class shopping, awesome beaches, parks, falls, great local delicacies, almost everything Hawaii can offer is in Honolulu. This place is one of the destinations I want to have a vacation with my family. So it is a good opportunity to extract data from mapzen and have a good idea on the places on Honolulu.

#### Problems Encountered in the Map

After checking the csv files from the output of the provided data.py file, I have noticed some problems in the data:

- 1. Different street abbreaviations (St., St, Rd, Road, Dr, Drive, etc)
- 2. Inconsistent postal code ("96850", "96734-9998", "HI 96819")

#### **Problems with Street names:**

Both the nodes\_tag and the ways\_tag have key = 'street'. And upon checking the values, I saw inconsistencies in the naming convention of the street names. To deal with this problem, I checked the different conventions in the csv and created list of accepted street names and a dictionary mapping to correct those with unacceptable street names. I have created a function where the street names will be corrected before it is added in the dictionary.

```
def update_name(name, mapping):
# this function will set the street name to an acceptable street name
   badname = street_type_re.search(name).group()
   pos = name.find(badname)

   if badname in mapping:
       goodname = mapping[badname]
       name = name[:pos]+goodname
   return name
```

Sample Corrected Street Names under Ways Tag				
id value (bad street name) value (corrected stre				
268795384	Lusitania St.	Lusitania Street		
280348190	Ala Pumalu St	Ala Pumalu Street		
62541758	Kalakaua Ave	Kalakaua Avenue		

Sample Corrected Street Names under Nodes Tag				
id	value (bad street name)	value (corrected street name)		
367803601	Kipapa Dr	Kipapa Drive		
4198265289	Ala Moana Blvd	Ala Moana Boulevard		
4255639689	Kamehameha Hwy	Kamehameha Highway		
4223222595	Meheula Pkwy	Meheula Parkway		

#### **Problems with the Postal:**

From Wikipedia: "The basic format consists of five numerical digits. An extended **ZIP+4** code, introduced in 1983, includes the five digits of the ZIP Code, a hyphen, and four additional digits that determine a more specific location within a given ZIP Code."

Since most of the values for postal codes in the dataset is five numerical digits, I will adopt this as standard for postal code values.

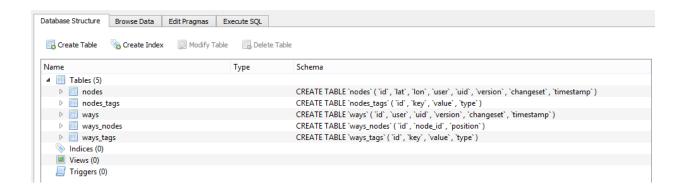
I used regular expression to address this issue:

```
# this function will set the postalcode values to the standard/acceptable values
def update_postalcode(num):
    num = re.findall('[0-9]+',num)[0]
    return num
```

Sample Corrected Postcodes under Ways Tag					
id	value (bad postcode)	value (corrected postcode)			
232136444	96825-9998	96825			
269984254	96826-4427	96826			
302171171	96734-9998	96734			
436080000	96815-2830	96815			
440166239	96817-1713	96817			

Sample Corrected Postcodes under Nodes Tag				
id	value (bad postcode)	value (corrected postcode)		
2609285445	HI 96819	96819		
4263793619	96712-9998	96712		
4338125893	96815-2518	96815		
4339585890	96815-2834	96815		
4339585891	96815-2834	96815		

After the data wrangling stage, I created a database in DB Browser for SQLite and named it as "Honolulu\_Osm.db". I also imported the resulting csv files in SQLite to serve as tables of my Database.



# **Data Overview**

# File sizes

Name	Date modified	Type	Size
Monolulu_Osm	10/27/16 5:01 PM	Data Base File	2 KB
Honolulu_Osm.db-journal	10/27/16 5:17 PM	DB-JOURNAL File	4 KB
Final Project 3.ipynb	10/27/16 5:17 PM	IPYNB File	14 KB
nodes nodes	10/27/16 5:16 PM	Microsoft Excel C	20,084 KB
nodes_tags	10/27/16 5:16 PM	Microsoft Excel C	457 KB
ways	10/27/16 5:16 PM	Microsoft Excel C	1,480 KB
ways_nodes	10/27/16 5:16 PM	Microsoft Excel C	6,770 KB
ways_tags	10/27/16 5:16 PM	Microsoft Excel C	3,473 KB
honolulu_hawaii.osm	10/26/16 10:43 PM	OSM File	51,674 KB

#### **Number of Nodes:**

```
import sqlite3
import pandas as pd

conn = sqlite3.connect('Honolulu_Osm.db')

c = conn.cursor()

c.execute('SELECT COUNT(*) FROM nodes')
print c.fetchone()

(242889,)
```

## **Number of Ways:**

```
c.execute('SELECT COUNT(*) FROM ways')
print c.fetchone()

(25540,)
```

# **Number of Unique Users:**

```
c.execute('''SELECT COUNT(DISTINCT(e.uid))
FROM (SELECT uid FROM nodes UNION ALL SELECT uid FROM ways) e''')
print c.fetchone()

(468,)
```

## Number of Users appearing only once:

```
c.execute('''SELECT COUNT(*)
FROM
        (SELECT e.user, COUNT(*) as num
        FROM (SELECT user FROM nodes UNION ALL SELECT user FROM ways) e
        GROUP BY e.user
        HAVING num=1) u''')
print c.fetchone()
(103,)
```

I created a function for queries to display a table:

```
def result(query):
    c.execute(query)
    rows = c.fetchall()
    names = [description[0] for description in c.description]
    df = pd.DataFrame(rows)
    df.columns = names
    print df
```

## **Top 10 Contributing Users:**

```
result('''SELECT e.user, COUNT(*) as num
FROM (SELECT user FROM nodes UNION ALL SELECT user FROM ways) e
GROUP BY e.user
ORDER BY num DESC
LIMIT 10''')
```

```
user
                        num
        Tom Holland
0
                     102158
             cbbaze
1
                      14989
2
              ikiya
                      12807
3
             kr4z33
                      9435
4
     Chris Lawrence
                      9214
5
              pdunn
                      9072
6
             aaront
                       8510
   woodpeck fixbot
7
                      8446
            bdiscoe
8
                       5097
  Mele Sax-Barnett
                      4617
```

## Top 10 appearing amenities:

```
result('''SELECT value, COUNT(*) as num
FROM nodes_tags
WHERE key='amenity'
GROUP BY value
ORDER BY num DESC
LIMIT 10''')
```

```
value num
     restaurant
0
                 196
      fast food 102
1
        parking
2
                  75
3
           cafe
                   62
        toilets
4
                   61
5
  fire station
                   30
  waste basket
6
                   29
        library
7
                   24
           bank
                   22
8
9
          bench
                   22
```

## Most popular cuisines:

```
value num
0
        japanese
                    7
           pizza
1
2
        american
                    5
3
        chinese
                    5
4
        regional
                    5
5 international
                    4
6
           asian
                    3
         italian
7
                    3
            thai
                    3
```

## Top 10 leisure:

```
value count
             pitch
                       338
0
    swimming pool
1
                      138
2
              park
                      137
     picnic table
3
                       65
      golf course
4
                        36
5
    sports centre
                        35
            track
6
                        22
7
           garden
                        19
   nature reserve
8
                        18
       playground
9
                        17
```

### **Religion:**

```
value count
0 christian 25
1 buddhist 7
2 muslim 1
```

## **Suggestions for Improving Data:**

I think it would be best if Mapzen has a "clean up" tool where users can upload corrections for their dataset. In this way, users who notice corrections will be given a chance to correct the dataset themselves. Correction files can be in a form of xml guided by an acceptable schema.

#### **Benefits:**

- It will encourage more users to participate in the objective of having a good and clean dataset for free;
- Cleaning will be faster since more will contribute in the cleaning process.

## **Anticipated Problems:**

- It might result to a more problematic dataset if restrictions on uploading corrections is not properly set;
- Conflicts between users may arise on deciding which correction is better.