

Project 3: OpenStreetMap Data Wrangling with SQL

Name: Shivam Bhardwaj

Map Area: Pune, Maharashtra, India as i have spent a lot of time in this city.

- Location: Pune, India
- [OpenStreetMap URL](#)
- [MapZen URL](#)

1. Data Audit

Unique Tags

Looking at the XML file, I found that it uses different types of tags. And also my filesize is quite large for myPC so i sampled the data. Then parsed the Pune, India dataset using ElementTree and count number of the unique tags.

mapparser.py is used to count the numbers of unique tags.

- 'member': 195,
- 'nd': 55896,
- 'node': 47133,
- 'osm': 1,
- 'relation': 72,
- 'tag': 10203,
- 'way': 8998

Patterns in the Tags

The "k" value of each tag contain different patterns. Using tags.py, I created 3 regular expressions to check for certain patterns in the tags.

I have counted each of four tag categories.

- "lower" : 9939 , for tags that contain only lowercase letters and are valid,
- "lower_colon" : 260 , for otherwise valid tags with a colon in their names,
- "problemchars" : 0 , for tags with problematic characters, and
- "other" : 4 , for other tags that do not fall into the other three categories.

2. Problems Encountered in the Map

Street address inconsistencies

The main problem we encountered in the dataset is the street name inconsistencies. Below is the old name corrected w better name. Using audit.py, we updated the names.

- **Abbreviations**
 - Rd -> Road
- **LowerCase**
 - pune -> Pune
- **Misspelling**
 - socity -> Society
- **Hindi names**
 - Marg -> Path

City name inconsistencies

Using audit.py, we update the names

- **LowerCase**
 - puna -> Pune
- **Misspelling**

- poona -> Pune

3. Data Overview

File sizes:

- pune_india.osm: 307.3 MB
- nodes_csv: 3.9 MB
- nodes_tags.csv: 16.9 KB
- ways_csv: 549.1 KB
- ways_nodes.csv: 1.3 MB
- ways_tags.csv: 316.8 KB
- pune_sample : 10.3 MB
- pune.db : 7.0 MB

Number of nodes:

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

Output:

```
47133
```

Number of ways:

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

Output:

```
8998
```

Number of unique users:

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

Output:

```
276
```

Top contributing users:

```
sqlite> 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;
```

Output:

```
Sramesh : 1918  
Praveeng : 1904  
Shiva05 : 1740  
Anushapyata : 1616  
Kranthikumar : 1589  
Harishk : 1436  
Saikumar : 1337
```

Number of users contributing only once:

```
sqlite> 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;
```

Output:

103

4. Additional Data Exploration

Biggest religion:

```
sqlite> SELECT nodes_tags.value, COUNT(*) as num  
FROM nodes_tags  
JOIN (SELECT DISTINCT(id) FROM nodes_tags WHERE value='place_of_worship') i  
ON nodes_tags.id=i.id  
WHERE nodes_tags.key='religion'  
GROUP BY nodes_tags.value  
ORDER BY num DESC  
LIMIT 1;
```

Output:

Hindu : 4

Popular cuisines

```
sqlite> SELECT nodes_tags.value, COUNT(*) as num  
FROM nodes_tags  
JOIN (SELECT DISTINCT(id) FROM nodes_tags WHERE value='restaurant') i  
ON nodes_tags.id=i.id  
WHERE nodes_tags.key='cuisine'  
GROUP BY nodes_tags.value  
ORDER BY num DESC;
```

Output:

Barbecue : 1

5. Conclusion

The OpenStreetMap data of Pune is of fairly reasonable quality but the typo errors caused by the human inputs are significant. I have cleaned a significant amount of the data which is required for this project. But, there are lots of improvement needed in the dataset. The dataset contains very less amount of additional information such as amenities, tourist attractions, popular places and other useful interest. The dataset contains very old information which is now incomparable to that of Google Maps or Bing Maps. So, I think there are several opportunities for cleaning and validation of the data in the future.

Additional Suggestion and Ideas

Control typo errors

We can build parser which parse every word input by the users. We can make some rules or patterns to input data which users follow everytime to input their data. This will also restrict users input in their native language. We can develop script or bot to clean the data regularly or certain period. Also we can give user of the maps to authenticate the data on the map.

More information

The tourists or even the city people search map to see the basic amenities provided in the city or what are the popular places and attractions in the city or near outside the city. So, the users must be motivated to also provide these informations in the map. If we can provide these informations then there are more chances to increase views on the map because many people directly enter the famous name on the map.

Files

- `README.md` : description
- `sample.osm` : sample data of the OSM file
- `audit.py` : audit street, city and update their names
- `data.py` : build CSV files from OSM and also parse, clean and shape data
- `database.py` : create database of the CSV files
- `mapparser.py` : find unique tags in the data
- `query.py` : different queries about the database using SQL
- `report.pdf` : pdf of this document
- `sample.py` : extract sample data from the OSM file
- `tags.py` : count multiple patterns in the tags
- `schema.py` : Helper file for data.py