Project 3: OpenStreetMap Data Wrangling with SQL

Name: Shailja Bagri

Map Area: I have chosen Ahmedabad as I have spent a lion share of my childhood there.

- · Location: Ahmedabad, India
- OpenStreetMap URL
- MapZen URL

1. Data Audit

Unique Tags

Looking at the XML file, I found that it uses different types of tags. So, I parse the Ahmedabad,India dataset using ElementTree and count number of the unique tags.

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

- 'bounds': 1
- 'member': 2267,
- 'nd': 640184,
- 'node': 551132
- 'osm': 1,
- 'relation': 504,
- 'tag': 99600,
- 'way': 82383

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": 95729, for tags that contain only lowercase letters and are valid,
- "lower_colon": 2030, for otherwise valid tags with a colon in their names,
- "problemchars": 7, for tags with problematic characters, and
- "other": 34, 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 with the better name. Using audit.py, we updated the names.

- Abbreviations
 - Rd -> Road
- LowerCase
 - o gandhi -> Gandhi
- Misspelling
 - socity -> Society
- Hindi names
 - o rasta -> Road
- UpperCase Words
 - o sbk -> SBK

City name inconsistencies

Using audit.py, we update the names

- LowerCase
 - ahmedabad -> Ahmedabad
- Misspelling

POSTAL CODES

I have writeen the code for postal code inconsistencies, however, no error was detected.

3. Data Overview

File sizes:

- ahmedabad_india.osm: 112.2 MB
- nodes_csv: 28 MB
- nodes_tags.csv: 142 KB
- ways_csv: 0
- MB ways_nodes.csv: 0 MB
- ways_tags.csv: 0 MB
- ahmedabad.db: 5 MB

Number of nodes:

sqlite> SELECT COUNT(*) FROM nodes

Output:

61473

Number of ways:

sqlite> SELECT COUNT(*) FROM ways

Output:

0

Number of unique users:

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

Output:

215

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:

Vkovra 21029

Shravan 15857

Oberaffe 6183

PlaneMad 3392

Heinz_V 3141

ThePatel 1189

Staeublec 768

SRamesh 561

Dhimant 532

Indigome 524

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:

53

4. Additional Data Exploration

Common ammenities:

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

Output:

```
place_of_worship 43
```

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:

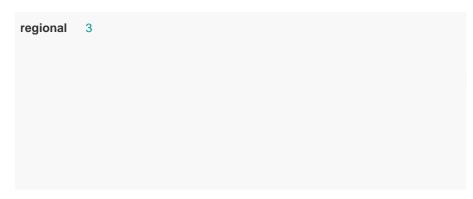
Hindu: 29

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:



5. Conclusion

The OpenStreetMap data of Ahmedabad is of fairly reasonable quality but the typo errors caused by the human inputs are significant. We 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.

Benefits while implementing the improvement

The Ahmedabad Sample file had 2267 members and the contribution by them were good enough. So I could
easily audit the street names.

Problems encountered while implementing the improvement

- The information provided in the map is very limited. This can be proved as while auditing for postal codes, no error were detected and no output for postal codes were generated. Similarly no telephone nos were provided.
- Also for amenities and cuisines. The people staying there and the tourists didn't provide much information so
 that further analysis could be done on that.

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 every time 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.

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 information in the map.
- If we can provide these information then there are more chances to increase views on the map because many people directly enter the famous name on the map.

Files

- Quiz/: scripts completed in lesson Case Study OpenStreetMap
- reportS.doc: this file
 ahmedabad sample.osm: sample data of the OSM file

audit1234.py: audit street, city and update their names

data123.py: build CSV files from OSM and also parse, clean and shape

data database123.py : create database of the CSV files

• mapparser123.py : find unique tags in the data

query123.py: different queries about the database using

SQL report1234.pdf : pdf of this document

• sample.py: extract sample data from the OSM file

• tags123.py : count multiple patterns in the tags