

# OpenStreetMap Data Case Study

---

## Map Area

Charlotte, NC, United States

- <https://www.openstreetmap.org/relation/177415>
- <http://metro.teczno.com/#charlotte>

This map is of my hometown, so I'm more interested to see what database querying reveals, and I'd like an opportunity to contribute to its improvement on OpenStreetMap.org.

## Problems Encountered in the Map

---

After initially downloading a small sample size of the Charlotte area and running it against a provisional data.py file, I noticed five main problems with the data, which I will discuss in the following order:

- Overabbreviated street names ("*S Tryon St Ste 105*")
- Inconsistent postal codes ("*NC28226*", "*282260783*", "*28226*")
- "Incorrect" postal codes (Charlotte area zip codes all begin with "282" however a large portion of all documented zip codes were outside this region.)
- Second level `"k"` tags with the value `"type"` (which overwrites the element's previously processed `node["type"]field`).
- Street names in second level `"k"` tags pulled from Tiger GPS data and divided into segments, in the following format:

```
1 <tag k="tiger:name_base" v="Stonewall"/>
2 <tag k="tiger:name_direction_prefix" v="W"/>
3 <tag k="tiger:name_type" v="St"/>
```

## Overabbreviated Street Names

Once the data was imported to SQL, some basic querying revealed street name abbreviations and postal code inconsistencies. To deal with correcting street names, I opted not use regular expressions, and instead iterated over each word in an address, correcting them to their respective mappings in audit.py using the following function:

```

1 def update(name, mapping):
2     words = name.split()
3     for w in range(len(words)):
4         if words[w] in mapping:
5             if words[w+1].lower() not in ['suite', 'ste.', 'ste']:
6                 # For example, don't update 'Suite E' to 'Suite East'
7                 words[w] = mapping[words[w]] name = " ".join(words)
8     return name

```

This updated all substrings in problematic address strings, such that:

*"S Tryon St Ste 105"*

becomes:

*"South Tryon Street Suite 105"*

## Postal Codes

Postal code strings posed a different sort of problem, forcing a decision to strip all leading and trailing characters before and after the main 5digit zip code. This effectively dropped all leading state characters (as in "NC28226") and 4digit zip code extensions following a hyphen ("28226-0783"). This 5digit restriction allows for more consistent queries.

Regardless, after standardizing inconsistent postal codes, some altogether "incorrect" (or perhaps misplaced?) postal codes surfaced when grouped together with this aggregator:

```

1 SELECT tags.value, COUNT(*) as count
2 FROM (SELECT * FROM nodes_tags
3       UNION ALL
4       SELECT * FROM ways_tags) tags
5 WHERE tags.key='postcode'
6 GROUP BY tags.value
7 ORDER BY count DESC;

```

Here are the top ten results, beginning with the highest count:

1	value count
2	28205 900
3	28208 388
4	28206 268
5	28202 204
6	28204 196
7	28216 174
8	28211 148
9	28203 120
10	28209 104
11	28207 86

These results were taken before accounting for Tiger GPS zip codes residing in second level “k” tags. Considering the relatively few documents that included postal codes, of those, it appears that out of the top ten, seven aren’t even in Charlotte, as marked by a “#”. That struck me as surprisingly high to be a blatant error, and found that the number one postal code and all others starting with “297” lie in Rock Hill, SC. So, I performed another aggregation to verify a certain suspicion...

## Sort cities by count, descending

```
1 | sqlite> SELECT tags.value, COUNT(*) as count
2 | FROM (SELECT * FROM nodes_tags UNION ALL
3 |       SELECT * FROM ways_tags) tags
4 | WHERE tags.key LIKE '%city'
5 | GROUP BY tags.value
6 | ORDER BY count DESC;
```

And, the results, edited for readability:

1	Rock Hill	111
2	Pineville	27
3	Charlotte	26
4	York	24
5	Matthews	10
6	Concord	4
7	3000	3
8	10	2
9	Lake Wylie	2
10	1	1
11	3	1
12	43	1
13	61	1
14	Belmont, N	1
15	Fort Mill,	1

These results confirmed my suspicion that this metro extract would perhaps be more aptly named “Metrolina” or the “Charlotte Metropolitan Area” for its inclusion of surrounding cities in the sprawl. More importantly, three documents need to have their trailing state abbreviations stripped. So, these postal codes aren’t “incorrect,” but simply unexpected. However, one final case proved otherwise.

A single zip code stood out as clearly erroneous. Somehow, a “48009” got into the dataset. Let’s display part of its document for closer inspection (for our purposes, only the “address” and “pos” fields are relevant):

```

1  sqlite> SELECT *
2  FROM nodes
3  WHERE id IN (SELECT DISTINCT(id) FROM nodes_tags WHERE key='postcode'
AND value='48009')

```

```

1234706337|35.2134608|-80.8270161|movercash|433196|1|7784874|2011-04-
06T13:16:06Z

```

```

sqlite> SELECT * FROM nodes_tags WHERE id=1234706337 and type='addr';

```

```

1  1234706337|houzenumber|280|addr
2  1234706337|postcode|48009|addr
3  1234706337|street|North Old Woodward Avenue|addr

```

It turns out, “280 North Old Woodward Avenue, 48009” is in Birmingham, Michigan. All data in this document, including those not shown here, are internally consistent and verifiable, except for the latitude and longitude. These coordinates are indeed in Charlotte, NC. I’m not sure about the source of the error, but we can guess it was most likely sitting in front of a computer before this data entered the map. The document can be removed from the database easily enough.

## Data Overview and Additional Ideas

This section contains basic statistics about the dataset, the MongoDB queries used to gather them, and some additional ideas about the data in context.

### File sizes

```

charlotte.osm ..... 294 MB
charlotte.db ..... 129 MB
nodes.csv ..... 144 MB
nodes_tags.csv ..... 0.64 MB
ways.csv ..... 4.7 MB
ways_tags.csv ..... 20 MB
ways_nodes.cv ..... 35 MB

```

### Number of nodes

```

sqlite> SELECT COUNT(*) FROM nodes;

```

1471350

### Number of ways

```

sqlite> SELECT COUNT(*) FROM ways;

```

## Number of unique users

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

337

## Top 10 contributing users

```
1 | sqlite> SELECT e.user, COUNT(*) as num
2 | FROM (SELECT user FROM nodes UNION ALL SELECT user FROM ways) e
3 | GROUP BY e.user
4 | ORDER BY num DESC
5 | LIMIT 10;
```

1	jumbanho	823324
2	woodpeck_f	481549
3	TIGERcnl	44981
4	bot-mode	32033
5	rickmastfa	18875
6	Lightning	16924
7	grossing	15424
8	gopanthers	14988
9	KristenK	11023
10	Lambertus	8066

## Number of users appearing only once (having 1 post)

```
1 | sqlite> SELECT COUNT(*)
2 | FROM
3 |     (SELECT e.user, COUNT(*) as num
4 |      FROM (SELECT user FROM nodes UNION ALL SELECT user FROM ways) e
5 |      GROUP BY e.user
6 |      HAVING num=1) u;
```

56

## Additional Ideas

### Contributor statistics and gamification suggestion

The contributions of users seems incredibly skewed, possibly due to automated versus manual map editing (the word “bot” appears in some usernames). Here are some user percentage statistics:

- Top user contribution percentage (“jumbanho”) 52.92%
- Combined top 2 users' contribution (“jumbanho” and “woodpeck\_fixbot”) 83.87%
- Combined Top 10 users contribution 94.3%
- Combined number of users making up only 1% of posts 287 (about 85% of all users)

Thinking about these user percentages, I'm reminded of “gamification” as a motivating force for contribution. In the context of the OpenStreetMap, if user data were more prominently displayed, perhaps others would take an initiative in submitting more edits to the map. And, if everyone sees that only a handful of power users are creating more than 90% a of given map, that might spur the creation of more efficient bots, especially if certain gamification elements were present, such as rewards, badges, or a leaderboard.

## Additional Data Exploration

---

### Top 10 appearing amenities

```
1 | sqlite> SELECT value, COUNT(*) as num
2 | FROM nodes_tags
3 | WHERE key='amenity'
4 | GROUP BY value
5 | ORDER BY num DESC
6 | LIMIT 10;
```

1	place_of_worship	580
2	school	402
3	restaurant	80
4	grave_yard	75
5	parking	63
6	fast_food	51
7	fire_station	48
8	fuel	31
9	bench	30
10	library	28

### Biggest religion (no surprise here)

```

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

```

christian 571

## Most popular cuisines

```

1  sqlite> SELECT nodes_tags.value, COUNT(*) as num
2  FROM nodes_tags
3      JOIN (SELECT DISTINCT(id) FROM nodes_tags WHERE value='restaurant')
i
4      ON nodes_tags.id=i.id
5  WHERE nodes_tags.key='cuisine'
6  GROUP BY nodes_tags.value
7  ORDER BY num DESC;

```

1	american	9
2	pizza	5
3	steak_hous	4
4	chinese	3
5	japanese	3
6	mexican	3
7	thai	3
8	italian	2
9	sandwich	2
10	barbecue	1

## Conclusion

After this review of the data it's obvious that the Charlotte area is incomplete, though I believe it has been well cleaned for the purposes of this exercise. It interests me to notice a fair amount of GPS data makes it into OpenStreetMap.org on account of users' efforts, whether by scripting a map editing bot or otherwise. With a rough GPS data processor in place and working together with a more robust data processor similar to data.pyl think it would be possible to input a great amount of cleaned data to OpenStreetMap.org.

