

OpenStreetMap Sample Project Data Wrangling with MongoDB

Map area: Vancouver, BC, Canada

<https://www.openstreetmap.org/relation/1852574#map=12/49.2573/-123.1241>

I chose the Vancouver BC area as this is where I reside. The ability to lean on my domain knowledge of this locale proved helpful as I learned how to wrangle data in this project. This area consists of 115 km² and just over 600,000 people. It is the third most populous city in Canada, but the most densely populated Canadian municipality and fourth in North America, behind New York, San Francisco, and Mexico City. Vancouver City is also known for being one of the most ethnically and linguistically diverse city in Canada.

Problems encountered in your map

The data was surprisingly clean. I expected to see a lot more inconsistencies, but with the exception of minor typos, street abbreviations, and field names with colons the data seemed ready to be imported to MongoDB.

Street name typos While very low in occurrence, this type of inconsistency would have gone unnoticed had I not have intimate knowledge of the map area. I added these to my mapping dictionary to correct for. An example of such mistakes is, `Jervis => Jarvis`.

Street name abbreviations While auditing the dataset, I noticed instances where street names were in various formats while contextually being similar. These were also accounted for in the mapping dictionary to standardize the address. The following is a snippet of the mapping dictionary where four variations of `Highway` were corrected for:

```
"Hwy": "Highway", "Hwy.": "Highway", "HWY": "Highway", "HIGHWAY": "Highway"
```

Field names with colon During the data audit, multiple instances of field names with double colons were identified using regular expression. It seemed unlikely that the `lower_colon` regex would catch all colon instances, so I used a slight modification to remove the first colon.

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

Additional note A further point to note about the cleanliness of this dataset was the absence of problem characters.

Overview of the Data

File sizes `vancouver_canada.osm`: 154.5 MB `vancouver_canada.osm.json`: 179.2 MB

This section contains basic statistic about the dataset and the MongoDB queries used on them. The following are snippets of queries found in `query.py`:

Top contributor

```
db.openstreetmap.aggregate([{\n  '$group': {\n    '_id': '$created.user',\n    'count': {\n      '$sum': 1\n    }\n  }\n}, {\n  '$sort': {\n    'count': -1\n  }\n}, {\n  '$limit': 1\n}])\n\n[{'u'_id': u'keithonearth', u'count': 965035}]
```

The user [keithonearth](#) is very active on OpenStreetMap and is a huge contributor to the British Columbia province in general.

Number of users contributing only once

```
db.openstreetmap.aggregate([{\n  '$group': {\n    '_id': '$created.user',\n    'count': {\n      '$sum': 1\n    }\n  }\n}, {\n  '$group': {\n    '_id': '$count',\n    'num_users': {\n      '$sum': 1\n    }\n  }\n}, {\n  '$sort': {\n    '_id': 1\n  }\n}, {\n  '$limit': 1\n}])\n\n[{'u'_id': 5, u'num_users': 156}]
```

Unsurprisingly, there are a large number of single contributors.

Top 10 building types

```

db.openstreetmap.aggregate([
  {
    '$match': {
      'building': {
        '$exists': 1
      }
    }
  }, {
    '$group': {
      '_id': '$building',
      'count': {
        '$sum': 1
      }
    }
  }, {
    '$sort': {
      'count': -1
    }
  }, {
    '$limit': 10
  })
])

[{'_id': 'yes', 'count': 543035},
 {'_id': 'residential', 'count': 9060},
 {'_id': 'apartments', 'count': 5980},
 {'_id': 'house', 'count': 3525},
 {'_id': 'commercial', 'count': 1025},
 {'_id': 'entrance', 'count': 565},
 {'_id': 'roof', 'count': 545},
 {'_id': 'school', 'count': 445},
 {'_id': 'retail', 'count': 400},
 {'_id': 'greenhouse', 'count': 175}]

```

There were an odd number of `yes` entries for buildings. As I investigated in the json file, they seems to be just an *other* category. So I ignored adjusting it to use the amenity name instead.

Top 5 amenities

```

db.openstreetmap.aggregate([
  {
    '$match': {
      'amenity': {
        '$exists': 1
      }
    }
  }, {
    '$group': {
      '_id': '$amenity',
      'count': {
        '$sum': 1
      }
    }
  }, {
    '$sort': {
      'count': -1
    }
  })

```

```

    }, {
      '$limit': 5
    })

[{'_id': 'parking', 'count': 4545},
 {'_id': 'restaurant', 'count': 2460},
 {'_id': 'bench', 'count': 2415},
 {'_id': 'cafe', 'count': 1410},
 {'_id': 'fast_food', 'count': 1035}]

```

This was expected - Vancouverites do love being healthy and the outdoors, which explains the high number of benches and restaurants compared to fast food.

Top 3 cafe

```

db.openstreetmap.aggregate([
  {'$match': {
    'amenity':
      'cafe'
  }},
  {'$group': {
    '_id': '$name',
    'count': {
      '$sum': 1
    }
  }},
  {'$sort': {
    'count': -1
  }},
  {'$limit': 3
  })

[{'_id': 'Starbucks', 'count': 255},
 {'_id': 'Starbucks Coffee', 'count': 115},
 {'_id': 'Tim Hortons', 'count': 65}]

```

The high number of cafes peaked my interest, hence this query. Though it comes to little surprise that Starbucks is everywhere in Vancouver compared to the very Canadian, Tim Hortons. On a side note, the amenities names needs to be scrubbed.

Top 2 fast food

```

db.openstreetmap.aggregate([
  {'$match': {
    'amenity': {
      '$exists': 1
    },
    'amenity':
      'fast_food'
  }}

```

```

    }
  }, {
    '$group': {
      '_id': '$name',
      'count': {
        '$sum': 1
      }
    }
  }, {
    '$sort': {
      'count': -1
    }
  }, {
    '$limit': 2
  })
]]

```

```

[{'_id': 'Subway', 'count': 165}, {'_id': 'McDonald's', 'count': 100}]

```

This query was out of plain curiosity and is inline with the claim that Vancouverites are generally health conscious. Thus the prominence of Subway over McDonald's.

Other ideas about the datasets

The data validation and queries performed were primarily on *node* and *way* tags with a focus on address locations. There were a significant number of *nd* reference tags that were ignored. Depending on the application of this dataset for future analysis, it would be worthwhile to dive into this more and uncover the different tags and see if there are any discrepancies. It would also be interesting to dive into other areas of the data to assess cleanliness, such as postal codes, phone numbers and coordinate data.

Furthermore, there are ample more opportunities to clean this data beyond this cursory exercise as seen in the *Top 3 cafe* query. Nevertheless, the dataset is fairly clean as noted above. OpenStreepMap is crowdsourced and is therefore dependant on users' input, which could lead to incomplete or outdated information. It would be interesting to compare these results with Google Maps'.

References

- [Parsing large XML files](#)
- [Python XML documentation](#)
- [Vancouver wiki](#)
- [Installing MongoDB](#)
- [Importing JSON into MongoDB](#)
- [Discusison Forum](#)