
OpenStreetMap data wrangling with MongoDB Documentation

Release

Ignacio Toledo

March 09, 2015

Contents

1	Introduction	1
2	Problems encountered in our map area	2
2.1	Problems related to the parsing of the data	2
2.2	Problems related to the data audition	2
2.3	Validity	2
2.4	Accuracy	2
2.5	Completeness	2
2.6	Consistency	2
2.7	Uniformity	2
3	Overview of the data	2
3.1	File size:	2
3.2	Number of documents	3
3.3	Number of nodes	3
3.4	Number of ways	3
3.5	Number of unique users and top 10 users	3
4	Additional Ideas about the data set	3

1 Introduction

The goal of this project is to use data wrangling (munging) techniques over data obtained from the [OpenStreetMap](#). The data correspond to the area encompassing Santiago de Chile. The data wrangling consisted in :

- Obtaining the data from an API in OSM format.
- Parse the data using an python xml library.
- Audit the data set and determine which problems we will solve
- Once some functions to fix the chosen problems are created, a second parsing is done from the OSM file, but this time fixing some of the issues.
- Transform the parsed data into a valid JSON file to be imported by MongoDB
- Import the JSON file into MongoDB

2 Problems encountered in our map area

2.1 Problems related to the parsing of the data

First we described some basic problems that need to be solve in order to convert the XML file in a JSON file that can be imported by mongodb.

- Problematic characters: the fields names of the JSON document should include only alphabet character. Also, according to the OSM wiki they should be avoided in tags. We removed all *tags* that were matched by the regular expression:

```
/[=\+/\&<>;\'\"`?%#$@\\,\\. \t\r\n]/
```

- Tags with two sets of colons. Colons are present in the OSM's *tags* to allow the use of namespaces. It is not recommend to overuse this feature, so we removed any tags with more than one namespace (i.e., more than one colon)

2.2 Problems related to the data audition

Here a list of the problems we found within the data is presented. The will be classified by the kind of problem defined in the Data Wrangling class.

2.3 Validity

- As in the Lesson 6 exercise, the street's name in the data set do not follow a standard format when dealing with the 2 main kinds of streets used in Santiago's street names: *Avenida* and *Pasaje*. In many cases they are abbreviated, and not in a uniform way.
- Many phone number's do not follow a valid representation for Santiago de Chile. The phone numbers should be in the format: *+(Country_Code) (Area_Code) (Number)*.
 - The *Country_Code* should be 56, and in many cases is not present, nor the plus symbol at the beginning.
 - The *Area_Code* should be always 2 for residential lines, or 9X (with X a number between 5 and 9) for cell phones.
 - The *Number* should consist in an 8 digit number always starting with a 2 for residential lines; if a cell phone, the number should have 7 digits.

2.4 Accuracy

2.5 Completeness

2.6 Consistency

2.7 Uniformity

3 Overview of the data

3.1 File size:

santiago.xlm 216 MB (OSM file) santiago.xlm.json 237 MB (JSON file)

3.2 Number of documents

```
> db.santiago.find().count()
1015506
```

3.3 Number of nodes

```
> db.santiago.find({type: 'node'}).count()
814484
```

3.4 Number of ways

```
> db.santiago.find({type: 'way'}).count()
201022
```

3.5 Number of unique users and top 10 users

```
> db.santiago.distinct("created.user").length
1003

> db.santiago.aggregate(
  [{$group: {'_id': '$created.user',
            'count': {$sum:1}}},
   {$sort: {'count': -1}},
   {$limit: 10}]
)
{ "_id" : "Zambelli Limitada", "count" : 244613 }
{ "_id" : "Fede Borgia", "count" : 204199 }
{ "_id" : "felipeedwards", "count" : 107799 }
{ "_id" : "chesergio", "count" : 61124 }
{ "_id" : "dintrans_g", "count" : 56281 }
{ "_id" : "madek", "count" : 34397 }
{ "_id" : "ALE!", "count" : 30174 }
{ "_id" : "toniello", "count" : 27401 }
{ "_id" : "OttoPilot", "count" : 16241 }
{ "_id" : "Chilestreet", "count" : 15715 }
```

4 Additional Ideas about the data set

There are several tools available to do a quality assurance over an OSM file (JOSM/Validator) or online (Osmose, for selected areas). However they only deal with validity, consistency and uniformity, while accuracy and completeness are left out.

This is understandable, since the accuracy and completeness depend on having *gold standard* databases that we can use to perform these steps in the measure of the data quality. A great project would be to help the local city government to update and maintain the gold standard databases.

Currently, in Chile there is [web page](#) associated to the National Government that aims to provide this information to the users that require it. The information can't only be queried through the web interface, but an API has also been implemented. However many problems remain:

- Data is outdated and no longer maintained: besides some small number of databases, like the database with all the public transportation information for Santiago and the database with all the schools in the Santiago area, most of the geolocation data sets are outdated.
- Several geolocation data sets are missing: a database of the streets of Santiago with valid names, a database related to hospitals and clinics, a database with police stations, etc. They might actually exist, but they are not publicly available in this web page nor through the API.

Another good project would consist in resume the work at the [WikiProject Chile](#) OpenStreetMap page. There is some work already done there, which aims mainly to use a consistent set of tags and values to give consistency and uniformity to chilean map features, but is still clearly incomplete and kind of dead, since no modifications have been done in the last year.