

Wrangle OpenStreetMap Data Project

Project overview

XML OSM dataset of Paris have been downloaded from <https://mapzen.com/data/metro-extracts/> (<https://mapzen.com/data/metro-extracts/>). I chose to analyze Paris, as I recently went there for my holidays and wanted to know if the data concerning this city was consistent. The OSM downloaded is 975Mo. Therefore, it has been parsed with a parameter of $k = 100$, so that codes could be applied faster. Then, munging techniques have been applied to the data, such as assessing the quality of the data for validity, accuracy, completeness, consistency and uniformity, to clean the OpenStreetMap data. Errors in the dataset have been identified and corrected about:

- Amenities
- Street names
- Postal codes

Then, the file has been converted to CSV and finally, SQL has been used to the data schema to complete the project.

The Dataset

OpenStreetMap's data are structured in well-formed XML documents (.osm files) that consist of the following elements:

- **Nodes:** individual dots used to mark specific locations (such as a postal box). Two or more nodes are used to draw line segments or "ways".
- **Ways:** a line of nodes, displayed as connected line segments. "Ways" are used to create roads, paths, rivers, etc.
- **Relations:** When "ways" or areas are linked in some way but do not represent the same physical thing, a "relation" is used to describe the larger entity they are part of. "Relations" are used to create map features, such as cycling routes, turn restrictions, and areas that are not contiguous. The multiple segments of a long way, such as an interstate or a state highway are grouped into a "relation" for that highway. Another example is a national park with several locations that are separated from each other. Those are also grouped into a "relation".
- **Tags:** description of names, types, ... of elements. Like name or type of roads and other attributes.

Numbers of different elements in the dataset

```
In [171]: 1 def count_tags(SAMPLE_FILE):
2         tags = {}
3         for event, elem in ET.iterparse(SAMPLE_FILE):
4             if elem.tag in tags:
5                 tags[elem.tag] += 1
6             else:
7                 tags[elem.tag] = 1
8         return tags
9
10 tag_view = count_tags(SAMPLE_FILE)
11 #pprint.pprint(overview_tags)
12
13 print("In the dataset, there are the following elements:")
14 print(" ")
15 for k, v in tag_view.items():
16     print k,v
```

In the dataset, there are the following elements:

```
node 39321
nd 53162
member 1712
tag 31312
relation 217
way 5515
osm 1
```

List of users

```
In [172]: 1 def process_map(filename):
          2     users = set()
          3     for _, element in ET.iterparse(filename):
          4         if 'user' in element.attrib:
          5             users.add(element.attrib['user'])
          6
          7     return users
          8
          9 paris_users = process_map(SAMPLE_FILE)
         10 print "There are %s users that added information on Openstreetmap a
         11
         12 print(" ")
         13
         14 for idx, line in enumerate(paris_users):
         15     if idx >= 10:
         16         break
         17     print line
```

There are 776 users that added information on Openstreetmap about Paris. These are a sample of 10 users:

OpenBeerMapContributor
Abertzalee
VKaeru
leonremi
Super-Map
JeaRRO
Eric POMMEREAU
Khomille
sevenup
Leonmvd

1. List of amenities

In [173]:

```

1 villages=defaultdict(set)
2 postcodes=defaultdict(set)
3 amenities=defaultdict(set)
4 street=defaultdict(set)
5 country=defaultdict(set)
6
7 for _, elem in ET.iterparse(SAMPLE_FILE, events=("start",)):
8     if elem.tag == "node" or elem.tag == "way":
9         for tag in elem.iter("tag"):
10             if (tag.attrib['k'] == "addr:postcode"):
11                 postcodes[tag.attrib['v']].add(tag.attrib['v'])
12             if (tag.attrib['k'] == "addr:city"):
13                 villages[tag.attrib['v']].add(tag.attrib['v'])
14             if (tag.attrib['k'] == "addr:street"):
15                 street[tag.attrib['v']].add(tag.attrib['v'])
16             if (tag.attrib['k'] == "addr:country"):
17                 country[tag.attrib['v']].add(tag.attrib['v'])
18             if (tag.attrib['k'] == "amenity") :
19                 amenities[tag.attrib['v']].add(tag.attrib['v'])
20
21 print "There are %s amenities in Paris. These are a sample of 10 ar
22
23 print(" ")
24 for idx, line in enumerate(amenities):
25     if idx >= 10:
26         break
27     print line

```

There are 55 amenities in Paris. These are a sample of 10 amenities:

```

taxi
fast_food
marketplace
conservatoire municipal
car_sharing
public_building
cinema
recycling
social_facility
telephone

```

Corrections to the list of amenities

This is the list of the corrected amenities:

```
1 fast_food => fast_food
2 marketplace => marketplace
3 conservatoire municipal => conservatory
4 car_sharing => car_sharing
5 public_building => public_building
6 cinema => cinema
7 recycling => recycling
8 telephone => telephone
9 college => college
10 parking => parking
```

2. List of street names

This is the list of the different street names in Paris:

```
1 {'1008949616': 'Rue Jasmin',
2  '1789288302': 'Rue des Artistes',
3  '2264901172': 'Rue Desnouettes',
4  '2996097015': 'Rue Maurice Ripoche',
5  '3214241918': 'Rue des Essertes',
6  '3214347988': 'Avenue des Tilleuls',
7  '386235626': 'Rue Anatole France',
8  '714161423': 'Rue de Saintonge',
9  '804245689': 'Rue de Vaugirard',
10 '840887271': 'Rue de Ponthieu'}
```

List of street types and their frequency

```
In [179]: 1 from operator import itemgetter
2
3 def sort_street_types(street_types):
4     result = []
5     for key, value in street_types.iteritems():
6         result.append((key, len(value)))
7         result = sorted(list(result), key=itemgetter(0), reverse=True)
8     return result
9
10 print("These are the 20 most frequent street types and their frequency")
11 print(" ")
12 street_types = sort_street_types(streets)
13 #print a samle of the list
14 street_types[:20]
```

These are the 20 most frequent street types and their frequency:

```
Out[179]: [('rue', 3),
('avenue', 1),
('Voie', 1),
('Villa', 25),
('Square', 14),
('Rue', 786),
('Route', 4),
(u'Rond-', 2),
('Quai', 15),
('Promenade', 1),
('Port', 1),
('Pont', 3),
(u'Place', 25),
('Passage', 21),
(u'Parvis', 1),
('Le', 1),
(u'Impasse', 12),
('Grande', 1),
('Galerie', 1),
('Forum', 1)]
```

List of street types and their abbreviations / typo

This section provides a list of street types and similar words in the document. Therefore, typo can be identified.

In [180]:

```

1 from difflib import get_close_matches
2
3 def populate_expected(street_types, threshold):
4     expected = []
5     for i in street_types[:threshold]:
6         expected.append(i[0])
7
8     return expected
9
10 EXPECTED = populate_expected(street_types, 20)
11 EXPECTED
12
13
14 def find_abbreviations(expected, data):
15     for i in expected:
16         print i, get_close_matches(i, data, 4, 0.5)
17
18 find_abbreviations(EXPECTED, list(streets.keys()))
19

```

```

rue ['rue', 'Rue', 'Route', 'Forum']
avenue ['avenue', 'Avenue', 'Grande']
Voie ['Voie']
Villa ['Villa']
Square ['Square']
Rue ['Rue', 'Route', 'rue']
Route ['Route', 'Rue', 'Autoroute', 'rue']
Rond- [u'Rond-']
Quai ['Quai']
Promenade ['Promenade', 'Grande']
Port ['Port', 'Pont', u'Cour']
Pont ['Pont', 'Port']
Place [u'Place', 'Passage']
Passage ['Passage', u'Impasse', u'Place']
Parvis [u'Parvis']
Le ['Le']
Impasse [u'Impasse', 'Passage']
Grande ['Grande', 'Promenade']
Galerie ['Galerie']
Forum ['Forum', 'rue']

```

Correction of street names

The street names have been corrected according to the typos earlier identified.

```

1 rue de Tourtille ==> Rue de Tourtille
2 rue Auguste Beau ==> Rue Auguste Beau
3 rue Heloise Isabelle Michaud ==> Rue Heloise Isabelle Michaud (13
  occurrences)
4 r Porte Jaune ==> Rue PoRuede Jaune
5 place de la Défense ==> Place de la Défense
6 rue des Trois Frères ==> Rue des Trois Frères
7 rue du docteur bauer ==> Rue du docteur bauer
8 rue de Rome ==> Rue de Rome (18 occurrences)
9 rue Riquet ==> Rue Riquet (2 occurrences)
10 villa Haussmann ==> Villa Haussmann
11 488 street names were fixed

```

Queries

After having converted the document into a CSV and created SQL tables, queries could have been used to obtain the following statistics about Paris:

In [190]:

```

1 countnodes = "SELECT COUNT(*) FROM nodes;"
2
3 for i in cur.execute(countnodes):
4     print "The number of nodes is:"
5     print i

```

The number of nodes is:
(39321,)

In [191]:

```

1 countways = "SELECT COUNT(*) FROM ways;"
2
3 for i in cur.execute(countways):
4     print "The number of ways is:"
5     print i

```

The number of ways is:
(5515,)

In [192]:

```

1 countusers = "SELECT COUNT(DISTINCT(e.uid)) FROM (SELECT uid FROM i
2
3 for i in cur.execute(countusers):
4     print "The number of ways is:"
5     print i

```

The number of ways is:
(763,)


```
In [193]: 1 # Top 5 amenities
          2
          3 topamenities = "SELECT value, COUNT(*) as num FROM nodes_tags WHERE
          4
          5 for row in cur.execute(topamenities):
          6
          7     print row
          8

('restaurant', 58)
('bench', 43)
('bicycle_parking', 36)
('parking_space', 26)
('waste_basket', 24)
```

```
In [196]: 1 countparking = "SELECT value, COUNT(*) FROM nodes_tags WHERE key='k
          2
          3 for i in cur.execute(countparking):
          4
          5     print i

('parking', 7)
```

Additional ideas

It is suprising to notice that 776 users added information to the map. A lot of elements have been entered in the map. However, if users would have a real motive to do so, even more information could have been added. For instance, each time users would enter an element, they could get points that would be added to their score. With a high score, users could reach higher levels. Therefore, additional data could be implemented, such as grades for the best restaurants, bars and clubs, timetable of shops and their contact.

The benefits of this new procedure are the followings:

- More data will be added (in quantity).
- Users will have to add any type of data, as long as it brings value added. More data will be added in terms of variety.

The potential problems are the followings:

- The data may lack of consistency. The more users, the more there will have different way of writing the data. For instance, some people may write postal codes with the initial of the area and some may write only the postal code numbers.
- The data may not be inaccurate. If users really want to increase their score, some may enter anything just to obtain more points.

Conclusion

As a conclusion, the data in Paris required few corrections and are in general quite consistent. This project focused on street names and amenities. Other elements could have been addressed, such as street numbers or postal codes. This project could have been improved by applying codes not only on the sample created while parsing the document, but on the entire file.

Also, the information given by Mapzen could have been more complete, such as the amount of population in certain area. So that, the concentration of amenities could have been observed according to the population density.

References

- Udacity - <https://www.udacity.com/> (<https://www.udacity.com/>)
- Mapzen - <https://mapzen.com/data/metro-extracts/> (<https://mapzen.com/data/metro-extracts/>)
- OpenStreetMap - <https://www.openstreetmap.org> (<https://www.openstreetmap.org>)
- Overpass API - <http://overpass-api.de/> (<http://overpass-api.de/>)
- Python Software Foundation - <https://www.python.org/> (<https://www.python.org/>)
- Catherine Devlin's Github repository - <https://github.com/catherinedevlin/ipython-sql> (<https://github.com/catherinedevlin/ipython-sql>)
- Python course - <https://www.python-course.eu/> (<https://www.python-course.eu/>)