

OpenStreetMap Data 案例分析

Map Area:

Seattle, WA, United States

https://mapzen.com/data/metro-extracts/metro/seattle_washington/

数据整理过程中遇到的挑战？

数据整理过程中遇到的最大的挑战是不会应用编程技能来进行相关的数据整理工作，不熟悉如何通过使用正则表达式来匹配相关的字符内容从而筛选出符合要求的数据信息。另外一点挑战是关于数据集中的 Dirty Data 的清理工作，由于 OpenStreetMap 网站提供的数据是由不同的用户自发修改整理的非官方的数据，所以数据的格式有时候并不统一，需要自行进行清理。

Problems Encountered in the Map

- Node tags 中的 access 的 value 值不统一，出现了 yes, no, permissive, 甚至 forestry 等多个值，造成混乱
- Node_tags 中的 Key 值有很多自定义的含义不明确的价值，例如 node id 29546940 存在着 af, am, ar, bg, bs 等多个含义不明显的 key 标签，会对数据的分析造成困难
- 同样含义的 key 标签会有不同的标志 例如 areaway/aerialway electricity/electrical_appliances
- 同样的 value 值可能对应多个 key 值 例如 WA 表示华盛顿州，可能对应 state state_code gnis:ST_num 多个不同的 key 值
- Postcode /Postal_code 存在多个 KEY 值，而且 Value 不规范，例如 V8W, 2L4, V8T4Y3。

Solutions :

```
3 mapping = { "electrical_appliances": "electricity",
4             "aerialway": "aeroway",
5             "state": "state_node",}
6
7 def updata_name(name, mapping):
8     changeword = mapping.keys()
9     for word in changeword:
10         if word in name:
11             name = name.replace(word, mapping.get(word))
12             break
13     return name
14
```

在数据导出之前先对数据进行一定的清洗，将 key 值为 state, aerialway, electrical_appliances, 的标签值分别替换为 electricity, aeroway, state_node。

Data Overview

File sizes

Seattle.osm	1642MB
project.db	994MB
nodes.csv	616MB
nodes_tags.csv	42MB
ways_nodes.csv	126MB
ways_tags.csv	49MB

```
58 #结点的数量
59 SELECT COUNT(*) FROM nodes;
60 7318958
61
62 #道路的数量
63 SELECT COUNT(*) as count FROM ways;
64 452380
```

```
66 #数据contributor 的数量
67 SELECT COUNT(DISTINCT(e.uid))
68 FROM (SELECT uid FROM nodes UNION ALL SELECT uid FROM ways) e;
69 3091
```

查询各个城市出现的次数

```
5 SELECT value, COUNT(*) as count
6 FROM nodes_tags
7 WHERE key LIKE '%city'
8 GROUP BY value
9 ORDER BY count DESC;
```

得到的部分结果，删除了出现次数过小的数据

```
11 Seattle|61314
12 Kirkland|19985
13 Saanich|11511
14 Langford|2802
15 2|2312
16 Oak Bay|2283
17 Colwood|1985
18 Mount Vernon|1964
19 Sooke|1600
20 Esquimalt|1495
21 View Royal|996
22 Metchosin|917
23 Capital H (Part 1)|
24 Victoria|629
25 Highlands|356
26 4|252
27 Woodinville|245
28 Hunts Point|193
29 Issaquah|189
30 Bothell|156
31 Edmonds|124
32 Redmond|107
33 Sedro Woolley|94
34 Bellevue|91
35 Renton|76
36 Everett|65
37 8|60
38 6|59
39 Bainbridge Island|5
40 Gig Harbor|52
41 10|47
42 Lacey|46
43 3|45
44 Olympia|39
45 Freeland|35
46 Burlington|32
47 Kingston|32
48 Becher Bay 1|30
```

查询各种不同类型的 shop 出现次数

```

138 #各种商店的种类名称
139 SELECT nodes_tags.value, COUNT(*) as num
140 FROM nodes_tags
141 WHERE nodes_tags.key = 'shop'
142 GROUP BY nodes_tags.value
143 ORDER BY num DESC;
144
145 convenience|705
146 hairdresser|526
147 clothes|459
148 beauty|425
149 car_repair|304
150 supermarket|273
151 yes|208
152 furniture|180
153 dry_cleaning|155
154 mobile_phone|148
155 car|135
156 pet|131
157 bakery|120
158 gift|116
159 books|108
160 optician|108
161 art|99
162 shoes|92
163 electronics|87
164 tobacco|87
165 jewelry|85
166 sports|84
167 car_parts|82
168 bicycle|79
169 wine|79
170 alcohol|78
171 department_store|66
172 doityourself|64
173 hardware|56
174 massage|56
175 laundry|53
176 garden_centre|50
177 copyshop|47

```

Suggestions about improving data quality

作为开源的项目，应当在每一个新加入的 contributor 做出改动之前给他们提供一个具体的指导文档，来规范他们对数据集的修改，例如可以考虑对于每一个 node 结点都设定固定的 Key 值模板，比如给出 highway, address, country, layer 等固定的 key 值，同时对于不在模板中的 key 值，可以单独保留，定期由人工进行审核复查来判断是否可以导入到新的数据文件中，以此来尽量减少不一致并且有重合的 key 值。预期会出现的问题是模板给出的 Key 值并不能够全面覆盖所有结点的相关属性结点属性，而且随时间变化，会有一些新出现的特殊属性值标签，比如新兴的共享单车专用停车区可能就是一个不同于传统的自行车停车区的新属性值，需要对其进行单独的标注，所以可以考虑正则表达式规则来规范每个 key 值的基本格式并结合人工审查的机制。同时可以考虑给予忠实用户更大的操作权限，而限制新用户的数据编辑权限，这样虽然可能会减少新用户的编辑问题，但是可以使得数据格式更加规范，而且根据统计结果显示，二八原则是使用的，接近 2 成的忠实用户贡献了 8 成左右的数据。

结合使用 excel 和 SQLite，使用 excel 可以在将数据导入到 SQLite 之前对于数据的基本概况有一定了解，比

如缺失值，数据的类型，数据的特征，从而辅助后面使用 SQLite 进行精确的查找。而数据库也是处理大量数据时不可或缺的工具，因为数据库的确比 excel 在查找数据上又更加强大的功能，首先 csv 文件在 excel 中最多显示 1048575 行，数据量过大时,excel 不能显示全部的数据，其次，在 SQLite 中，不同的数据表格之间可以很方便的 join，以此来分析有相关性的数据，最后 SQLite 的强大的查询可以根据要求进行更加精确的查找。

关于 relation 区域的数据不太充足，可以考虑将城区进行更加细致合理的区分和规划，突出每一个区域的功能特点。比如学校，商业区，工业区可以进行特定的表示，同时都可以整合 nodes 和 ways 标签的具体的属性值的特点来进入深入的分析，或许可以指导城市的规划问题，探究解决城市拥堵，安排布局不合理的问题。但是因为功能区的整体特征并没有一个概论，会有功能区的定位模糊的问题存在，可能某个区域既有工业区也有商业区等问题，处理难度较大。