

Contents

一 数据源.....	1
二 分析结果及及解决办法.....	1
1) 分析结果	1
2) 解决办法	2
三 数据抽取转换过程.....	2
四 数据统计结果.....	3
1 数据总数.....	3
2 Node 节点的数量.....	3
3 Way 节点的数量	3
4 Suppermarket 的数量	3
5 Highway 的种类以及个数.....	4

一 数据源

选取的是上海地区的 OSM 文件，大小 54M。链接如下：

https://s3.amazonaws.com/metro-extracts.mapzen.com/shanghai_china.osm.bz2

二 分析结果及及解决办法

1) 分析结果

对原始的 OSM 数据进行分析后，发现如下一些问题

- 相同属性的取值不统一
 1. 当 tag 为 “way “时，k=oneway 对应的 v 值不统一，根据官方的解释，应该为 yes 或者 no 但在数据集中有的值为-1。需要统一为 yes 或者 no。官方解释链接如下：
<http://wiki.openstreetmap.org/wiki/Way>
 2. highway 和 Hwy 不一致，需要转换 Hwy 到 highway。
- 当 tag 为 node 或者 way 时，k=name 对应的 v 是中文名字时，有的 node 节点没有 k=name-en 属性。
- 当 tag 为 node 时，user 的值有些是中文。但绝大部分的 user 取值是英文字母。

2) 解决办法

根据以上列出的问题，相应的解决办法如下

- 属性取值不统一的，在抽取数据时，把不规范或者不统一的值修改成规范或者统一值。
- 增加 k=name-en 属性，其 v 的取值采用百度开放的翻译 API 接口翻译成英文。

代码片段如下：

```
def translate_func(text):
    httpClient = None
    myurl = '/api/trans/vip/translate'
    fromLang = 'zh'
    toLang = 'en'
    salt = random.randint(32768, 65536)
    sign = appid + text + str(salt) + secretKey
    m1 = md5.new()
    m1.update(sign)
    sign = m1.hexdigest()
    myurl = myurl + '?appid=' + appid + '&q=' + urllib.quote(
        text) + '&from=' + fromLang + '&to=' + toLang + '&salt=' + str(salt) + '&sign=' + sign

    try:
        httpClient = httplib.HTTPConnection('api.fanyi.baidu.com')
        httpClient.request('GET', myurl)
        response = httpClient.getresponse()
        result = eval(response.read())
        return result["trans_result"][0]["dst"]
    except Exception, e:
        print e
    finally:
        if httpClient:
            httpClient.close()
```

- 将 user 的取值为中文的，用汉语拼音方式改写。

代码片段如下：

```
if u'\u4e00' <= item.attrib["user"] <= u'\u9fff':
    node["user"] = lazy_pinyin(item.attrib["user"])
else:
    node["user"] = item.attrib["user"]
```

三 数据抽取转换过程

- 完成以下内容
 - 1 只处理两种类型的顶级标记：“节点（node）”和“道路（way）”
 - 2 CREATED 数组中的属性应该添加到键“created”下
 - 3 经纬度属性应该添加到“pos”数组中，以用于地理空间索引编制。确保“pos”数组中的值是浮点型，不是字符串。

```
> db.mapdata.find({"shop": "supermarket"}).count()
```

5 Highway 的种类以及个数

```
> db.mapdata.aggregate([{$group:{_id: "$highway",result: { $sum: 1 }}}])
```

```
{ "_id" : "services", "result" : 1 }
```

```
{ "_id" : "tertiary_link", "result" : 9 }
```

```
{ "_id" : "platform", "result" : 1 }
```

```
{ "_id" : "pedestrian", "result" : 21 }
```

```
{ "_id" : "track", "result" : 62 }
```

```
{ "_id" : "living_street", "result" : 10 }
```

```
{ "_id" : "steps", "result" : 21 }
```

```
{ "_id" : "path", "result" : 36 }
```

```
{ "_id" : "trunk_link", "result" : 64 }
```

```
{ "_id" : "service", "result" : 584 }
```

```
{ "_id" : "raceway", "result" : 2 }
```

```
{ "_id" : "construction", "result" : 25 }
```

```
{ "_id" : "cycleway", "result" : 33 }
```

```
{ "_id" : "motorway_link", "result" : 195 }
```

```
{ "_id" : "footway", "result" : 206 }
```

```
{ "_id" : "motorway", "result" : 247 }
```

```
{ "_id" : "primary_link", "result" : 75 }
```

```
{ "_id" : "trunk", "result" : 60 }
```

```
{ "_id" : "secondary", "result" : 454 }
```

```
{ "_id" : "secondary_link", "result" : 25 }
```