

Import OSM Data

Download Data

<https://download.geofabrik.de/index.html>

该网站为各洲、各国家的openstreetmap开源数据库，例如，下载中国地区的数据：

```
1 wget https://download.geofabrik.de/north-america/china-latest.osm.pbf
```

以下过程中，我们使用的数据格式均为osm.pbf。

From osm.pbf To Dataframe

这里我们需要考虑两种情况：Spark3.x和Spark2.x。

Spark3.x

支持直接导入osm.pbf格式数据：

```
1 df = spark.read.format("osm.pbf").load("china-latest.osm.pbf")
```

参考：<https://github.com/simplexspatial/osm4scala/blob/67864d8c9d69039e820c34cee2b0d6aede87fa8f/website/docs/spark-connector.mdx>

Spark 2.x+（这是我现在在用的版本）

由于版本不支持osm.pbf格式的数据，所以需要先将osm.pbf数据转换为该版本可以支持的parquet格式，这里我们需要在终端运行以下代码安装osm-parquetizer工具，利用这个工具转换数据格式。

更加详细的说明请参考：<https://github.com/adrianulbona/osm-parquetizer>

以下为简版说明：

```
1 git clone https://github.com/adrianulbona/osm-parquetizer.git
2 cd osm-parquetizer
3 mvn clean package
4 java -jar target/osm-parquetizer-1.0.1-SNAPSHOT.jar china-latest.osm.pbf
```

上述代码运行时间较长，结果文件（三个）将储存在和china-latest.osm.pbf同一个文件下

```
1 -rw-r--r-- 1 adrianbona adrianbona 145M Apr  3 19:57 romania-
  latest.osm.pbf
2 -rw-r--r-- 1 adrianbona adrianbona 372M Apr  3 19:58 romania-
  latest.osm.pbf.node.parquet
3 -rw-r--r-- 1 adrianbona adrianbona 1.1M Apr  3 19:58 romania-
  latest.osm.pbf.relation.parquet
4 -rw-r--r-- 1 adrianbona adrianbona 123M Apr  3 19:58 romania-
  latest.osm.pbf.way.parquet
```

parquet的储存内容如下：

```

1 node
2   |-- id: long
3   |-- version: integer
4   |-- timestamp: long
5   |-- changeset: long
6   |-- uid: integer
7   |-- user_sid: string
8   |-- tags: array
9     |-- element: struct
10    |   |-- key: string
11    |   |-- value: string
12  |-- latitude: double
13  |-- longitude: double
14
15 way
16   |-- id: long
17   |-- version: integer
18   |-- timestamp: long
19   |-- changeset: long
20   |-- uid: integer
21   |-- user_sid: string
22   |-- tags: array
23     |-- element: struct
24    |   |-- key: string
25    |   |-- value: string
26  |-- nodes: array
27    |-- element: struct
28    |   |-- index: integer
29    |   |-- nodeId: long
30
31 relation
32   |-- id: long
33   |-- version: integer
34   |-- timestamp: long
35   |-- changeset: long
36   |-- uid: integer
37   |-- user_sid: string
38   |-- tags: array
39     |-- element: struct
40    |   |-- key: string
41    |   |-- value: string
42  |-- members: array
43    |-- element: struct
44    |   |-- id: long
45    |   |-- role: string
46    |   |-- type: string

```

在python中即可利用pyspark导入以上osm数据为DataFrame格式（以way数据为例）：

```
1 sqlContext.setConf("spark.sql.parquet.binaryAsString","true")
```

```
2 df = sqlContext.read.parquet("./rawdata/china-latest.osm.pbf.way.parquet")
```

参考 (osm数据分析, DataFrame, 但是不用sedona已经完全可以处理了) :

<https://github.com/bulutenesemre/OsmAnalysis>

<https://github.com/kkahloots/OpenStreetMapsHungary>

From DataFrame to CSV

```
1 def array_to_string(my_list):
2     return '[' + ','.join([str(elem) for elem in my_list]) + ']'
3
4 array_to_string_udf = udf(array_to_string, StringType())
5
6 osmWay = osmWay.withColumn('tags_array',
7     array_to_string_udf(osmWay["tags"]))
8
9 osmWay = osmWay.withColumn('nodes_array',
10     array_to_string_udf(osmWay["nodes"]))
11
12 osmWay.drop("tags").drop("nodes").write.csv('./csvData/osmWay.csv', header=
13     True)
```

这里的osmWay.csv是一个文件夹。用一下代码读入就是正常的一个dataframe, 所以不用担心。但是特别提及的是有一些列是array形式, 需要转换为string的格式才能储存, 导入的时候也需要进行一定转换。

```
1 spark.read.csv("./csvData/osmWay.csv", header=True)
```

但上述保存csv到本地的方法花费时间较长, 且容易报错内存溢出, 下面采用直接将Spark中导入数据到Hbase中去。

Spark-Hbase Connector

<https://stackoverflow.com/questions/35228991/saving-pyspark-rdd-to-hbase-raises-a-ttribute-error>

```
1 host = 'localhost'
2 table = 'transaction_fee_table' #needs to be created before hand in hbase
3     shell
4 conf = {"hbase.zookeeper.quorum": host,
5     "hbase.mapred.outputtable": table,
6     "mapreduce.outputformat.class":
7     "org.apache.hadoop.hbase.mapreduce.TableOutputFormat",
8     "mapreduce.job.output.key.class":
9     "org.apache.hadoop.hbase.io.ImmutableBytesWritable",
10     "mapreduce.job.output.value.class":
11     "org.apache.hadoop.io.Writable"}
12
13 keyConv =
14     "org.apache.spark.examples.pythonconverters.StringToImmutableBytesWritable
15     Converter"
16
17 valueConv =
18     "org.apache.spark.examples.pythonconverters.StringListToPutConverter"
19
20 # data is a rdd
```

```
12 data.saveAsNewAPIHadoopDataset(conf=conf,  
13                               keyConverter=keyConv,  
14                               valueConverter=valueConv)
```

由于hbase环境还没配好，这段代码我还没有试验过