

Document for MongoDB Coursework

Part 2

The Architecture of the Data Pipeline and GitHub Link

Shanchuan Wu

Sw9n14@soton.ac.uk

1 The Architecture of Data Pipeline

The architecture of data pipeline consists of several procedures or functions. The procedures or functions are shown as below.

1.1 Importing Data

First of all, data should be imported in the mongodb Database. In the system shell or command prompt, use the script below to insert the documents into the *tweets* collection in the *twitter* database. The *mongoimport* connects to a mongod instance running on localhost on port number 27017.

```
mongoimport -d twitter -c tweets --type csv --file microblog.csv --headerline
```

```
ShanchuanWu@Shanchuan-Wu-MacBook-Pro ~/D/microblogging>
mongoimport -d twitter -c tweets --type csv --file microblog.csv --headerline
connected to: localhost
2015-12-10T15:25:58.770+0000
2015-12-10T15:26:01.768+0000 [##.....] twitter.tweets      19.6 MB/198.5 MB (9.9%)
2015-12-10T15:26:04.770+0000 [####.....] twitter.tweets      39.3 MB/198.5 MB (19.8%)
2015-12-10T15:26:07.770+0000 [#####.....] twitter.tweets      57.8 MB/198.5 MB (29.1%)
2015-12-10T15:26:10.768+0000 [#####.....] twitter.tweets      76.2 MB/198.5 MB (38.4%)
2015-12-10T15:26:13.769+0000 [#####.....] twitter.tweets      89.8 MB/198.5 MB (45.2%)
2015-12-10T15:26:16.767+0000 [#####.....] twitter.tweets     107.0 MB/198.5 MB (53.9%)
2015-12-10T15:26:19.767+0000 [#####.....] twitter.tweets     125.9 MB/198.5 MB (63.4%)
2015-12-10T15:26:22.768+0000 [#####.....] twitter.tweets     145.3 MB/198.5 MB (73.2%)
2015-12-10T15:26:25.768+0000 [#####.....] twitter.tweets     163.8 MB/198.5 MB (82.5%)
2015-12-10T15:26:28.768+0000 [#####.....] twitter.tweets     182.7 MB/198.5 MB (92.0%)
2015-12-10T15:26:31.362+0000 imported 1459861 documents
```

1.2 Preprocessing

When the data are inserted into the database, it is necessary to take some measures to clean the data. Because *id_members* or types of *text* of some messages are not meet the requirements, for instance, their *id_members* are negative. In this way, these messages should be considered as wrong data and need to be removed in this procedure.

The function *isPreProcessDataSuccess* is used to decide whether data are preprocessed or meet the requirements.

```
128 ▼   'isPreProcessDataSuccess': function () {
129 ▼       var count = db.getCollection(this.collectionName).find({
130           '$where': '(this.id_member < 0) || (typeof(this.text) != "string")'
131       }).count();
132 ▼       if (count > 0) {
133           return false;
134 ▼       } else {
135           return true;
136       }
137     },
```

```

139 ▾ 'queryInit': function () { //need to be executed only when the collection's document change or the first time to query
140
141 ▾     if (true == this.isPreProcessDataSuccess()) {
142         print("collection is already preprocessed");
143 ▾     } else {
144         print("begin to preprocess data");
145         this.preProcessData();
146 |
147     };

```

If all data meet the requirements, then function *preProcessData* is executed to analyse every document in the collection.

```

20 ▾ 'preProcessData': function () {
21     // var name = this.collectionName
22     print("begin data parse , invalid data will be removed");
23 ▾     db.getCollection(this.collectionName).find().forEach(function (it) {
24 ▾         if ((it.id_member < 0) || (typeof (it.id_member) != "number") || (typeof (it.text) != "string")) {
25 ▾             db.getCollection(queryByUser.collectionName).remove({
26                 "_id": it._id
27             });
28         }
29     })
30 }
31 },

```

1.3 MapReduce

Map-reduce is a data processing paradigm for condensing large volumes of data into useful aggregated results. In this program, two *mapReduce* functions are defined. The first one is defined to create a new collection named *mrout*. By using this collection, it becomes easy to get the results of query1, query2 and query5. Similarly, another collection created by *mapReduce2* can help to group and count the messages according to different unigrams or bigrams. The detailed functions and queries can be seen in another document.

```

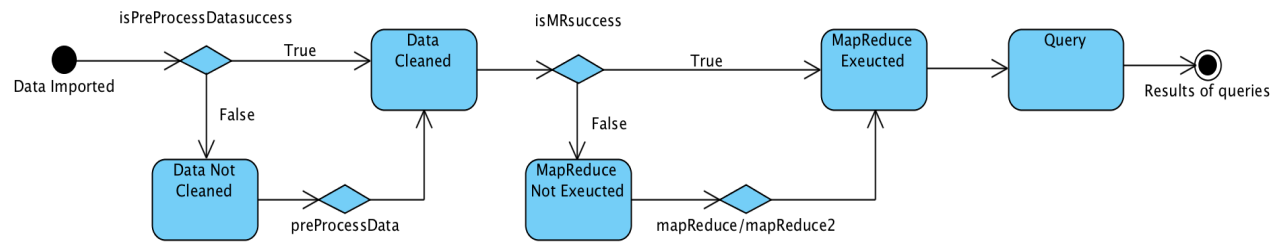
150 ▾     if (true == this.isMRSuccess()) {
151         print("collection is already mapReduced");
152 |
153 ▾     } else {
154         print("begin mapReduce");
155 ▾         db.getCollection(this.collectionName).ensureIndex({
156             'timestamp': 1
157         });
158         this.mapReduce();
159         this.mapReduce2();
160         return this.isMRSuccess();
161     }
162 },

```

1.4 Querying and Output

After two *mapReduce* functions are executed, it is time to use querying function to query the documents in this database. The results of queries and the detailed functions for querying can be seen in another document.

1.5 The State Diagram of Data Pipeline



2 GitHub Link

<https://github.com/358203708/MongodbCoursework>