# 云凡教育大数据学院

## Hive 创建索引

    索引是标准的数据库技术，hive 0.7 版本之后支持索引。Hive 提供有限的索引功能，这不像传统的关系型数据库那样有"键(key)"的概念，用户可以在某些列上创建索引来加速某些操作，给一个表创建的索引数据被保存在另外的表中。Hive 的索引功能现在还相对较晚，提供的选项还较少。但是，索引被设计为可使用内置的可插拔的 java 代码来定制，用户可以扩展这个功能来满足自己的需求。当然不是说有的查询都会受惠于 Hive 索引。用户可以使用 EXPLAIN 语法来分析 HiveQL 语句是否可以使用索引来提升用户查询的性能。像 RDBMS 中的索引一样，需要评估索引创建的是否合理，毕竟，索引需要更多的磁盘空间，并且创建维护索引也会有一定的代价。用户必须要权衡从索引得到的好处和代价。

    下面说说怎么创建索引：

**1、先创建表：**

```
hive> create table user( id int, name string)
    > ROW FORMAT DELIMITED
    > FIELDS TERMINATED BY '\t'
    > STORED AS TEXTFILE;
```

**2、导入数据：**
```
hive> load data local inpath '/export1/tmp/wyp/row.txt'
    > overwrite into table user;
```

**3、创建索引之前测试**
```
hive> select * from user where id =500000;
Total MapReduce jobs = 1
Launching Job 1 out of 1
Number of reduce tasks is set to 0 since there's no reduce operator
Cannot run job locally: Input Size (= 356888890) is larger than
hive.exec.mode.local.auto.inputbytes.max (= 134217728)
Starting Job = job_1384246387966_0247, Tracking URL =
http://l-datalogm1.data.cn1:9981/proxy/application_1384246387966_0247/
```

```
Kill Command=/home/q/hadoop/bin/hadoop job -kill job_1384246387966_0247
Hadoop job information for Stage-1: number of mappers:2; number of reducers:0
2013-11-13 15:09:53,336 Stage-1 map = 0%,    reduce = 0%
2013-11-13 15:09:59,500 Stage-1 map=50%,reduce=0%, Cumulative CPU 2.0 sec
2013-11-13 15:10:00,531 Stage-1 map=100%,reduce=0%, Cumulative CPU 5.63 sec
2013-11-13 15:10:01,560 Stage-1 map=100%,reduce=0%, Cumulative CPU 5.63 sec
MapReduce Total cumulative CPU time: 5 seconds 630 msec
Ended Job = job_1384246387966_0247
```

MapReduce Jobs Launched:

Job 0: Map: 2    Cumulative CPU: 5.63 sec

HDFS Read: 361084006 HDFS Write: 357 SUCCESS

Total MapReduce CPU Time Spent: 5 seconds 630 msec

OK

500000 wyp.

Time taken: 14.107 seconds, Fetched: 1 row(s)

一共用了 14.107s

### 4、对 user 创建索引

```
hive> create index user_index on table user(id)
    > as 'org.apache.hadoop.hive.ql.index.compact.CompactIndexHandler'
    > with deferred rebuild
    > IN TABLE user_index_table;
hive> alter index user_index on user rebuild;
hive> select * from user_index_table limit 5;
```

| 0 | hdfs://mycluster/user/hive/warehouse/table02/000000_0 | [0] |
|---|---|---|
| 1 | hdfs://mycluster/user/hive/warehouse/table02/000000_0 | [352] |
| 2 | hdfs://mycluster/user/hive/warehouse/table02/000000_0 | [704] |
| 3 | hdfs://mycluster/user/hive/warehouse/table02/000000_0 | [1056] |
| 4 | hdfs://mycluster/user/hive/warehouse/table02/000000_0 | [1408] |

Time taken: 0.244 seconds, Fetched: 5 row(s)

这样就对 user 表创建好了一个索引。

### 5、对创建索引后的 user 再进行测试

```
hive> select * from user where id =500000;
```

Total MapReduce jobs = 1

Launching Job 1 out of 1

Number of reduce tasks is set to 0 since there's no reduce operator

Cannot run job locally: Input Size (= 356888890) is larger than

hive.exec.mode.local.auto.inputbytes.max (= 134217728)

Starting Job = job_1384246387966_0247, Tracking URL =

http://l-datalogm1.data.cn1:9981/proxy/application_1384246387966_0247/

Kill Command=/home/q/hadoop/bin/hadoop job -kill job_1384246387966_0247

Hadoop job information for Stage-1: number of mappers:2; number of reducers:0

2013-11-13 15:23:12,336 Stage-1 map = 0%,    reduce = 0%

2013-11-13 15:23:53,240 Stage-1 map=50%,reduce=0%, Cumulative CPU 2.0 sec

2013-11-13 15:24:00,253 Stage-1 map=100%,reduce=0%, Cumulative CPU 5.27 sec

2013-11-13 15:24:01,650 Stage-1 map=100%,reduce=0%, Cumulative CPU 5.27 sec

MapReduce Total cumulative CPU time: 5 seconds 630 msec

Ended Job = job_1384246387966_0247

MapReduce Jobs Launched:

Job 0: Map: 2     Cumulative CPU: 5.63 sec

HDFS Read: 361084006 HDFS Write: 357 SUCCESS

Total MapReduce CPU Time Spent: 5 seconds 630 msec

OK

500000 wyp.

Time taken: 13.042 seconds, Fetched: 1 row(s)

时间用了 13.042s 这和没有创建索引的效果差不多。

　　在 Hive 创建索引还存在 bug：如果表格的模式信息来自 SerDe，Hive 将不能创建索引：云凡教育大数据学院 www.cloudyhadoop.com

01 hive> CREATE INDEX employees_index

02 > ON TABLE employees (country)

03 > AS 'org.apache.hadoop.hive.ql.index.compact.CompactIndexHandler'

04 > WITH DEFERRED REBUILD

05 > IDXPROPERTIES ('creator' = 'me','created_at' = 'some_time')

06 > IN TABLE employees_index_table

07 > COMMENT 'Employees indexed by country and name.';

08 FAILED: Error in metadata: java.lang.RuntimeException:           \

09 Check the index columns, they should appear in the table being indexed.

10　 FAILED: Execution Error, return code 1 from                 \

11 org.apache.hadoop.hive.ql.exec.DDLTask

这个 bug 发生在 Hive0.10.0、0.10.1、0.11.0，在 Hive0.12.0 已经修复了，详情请参见：https://issues.apache.org/jira/browse/HIVE-4251

直击 30 万年薪 业界首播：大数据企业面试+企业大数据实战项目公开课:

1）HDFS、YARN 相关面试题

2）MapReduce 高级编程相关面试题

3）Hbase、Hive、Storm、Spark、Solr 相关面试题

4）项目经验

12 月 16 日，晚 21：00 在 YY：20483828 直线开讲

点击即可进入课堂：http://www.yy.com/20483828

# 云凡教育大数据学院