

# hbase批量删除表数据

背景:

hbase的删除功能比较弱，只能单行删除，而且必须指定rowkey。

遇到问题:

今天遇到一个需求，用户导入了大量错误的记录，数据的rowkey开头都是110102，需要删除这些垃圾记录，用hbase shell删除实在不科学。

解决方案:

用hbase的mapreduce工具进行export和import，在export过程中filter掉不需要的数据。

首先说明下表的schema:

```
{NAME => 'freeway.service', FAMILIES => [{NAME => 'service_span_colfam',  
BLOOMFILTER => 'ROW', VERSIONS => '1', MIN_VERSIONS => '0', TTL => '604800',  
IN_MEMORY => 'true'}]}
```

我们使用hbase的export工具在export时filter掉不需要的数据，这边export支持[正则表达式](#)。我们看下export的usage:

```
Usage: Export [-D <property=value>]* <tablename> <outputdir> [<versions>  
[<starttime> [<endtime>]] [^[regex pattern] or [Prefix] to filter]]
```

Note: -D properties will be applied to the conf used.

For example:

```
-D mapred.output.compress=true
```

```
-D
```

```
mapred.output.compression.codec=org.apache.hadoop.io.compress.GzipCodec
```

```
-D mapred.output.compression.type=BLOCK
```

Additionally, the following SCAN properties can be specified to control/limit what is exported..

```
-D hbase.mapreduce.scan.column.family=<familyName>
```

tablename和outputdir是必须的，后面是版本号，starttime, endtime, filter的[正则表达式](#)。

我们这里版本就一个，starttime设为0，endtime设为很大的数，保证把所有数据都拿到。

后面正则表达式要用单引号包住以防Linux的bash解析里面的问号

```
hbase org.apache.hadoop.hbase.mapreduce.Driver export freeway.service  
hdfs://ns/usr/op1/freeway.service 1 0 9999999999999999 '^^(?!110102)'
```

现在这张表的数据就存在hdfs上的一个sequencefile里了。

现在删除原表，再创建一次。

然后import filter后的数据到新的表中：

```
hbase org.apache.hadoop.hbase.mapreduce.Driver import freeway.service  
hdfs://ns/usr/op1/freeway.service/part-m-00000
```