# Analyzing Apache Access Logs in Apache Hive

The following statistics are analyzed:

1. A count of response code's returned from the server.
2. The content size of responses  returned from the server to host.
3. The top ten most popular URL's in the Apache log
4. The average, min, and max content size of responses returned from the server.

**The steps to process data with Apache Hive**

Before proceed the below steps, we have to install the Cloudera Quickstart vm 5.5 and VMwareplayer. The Hadoop 2.6, Java 1.7, Eclipse Luna, Hive, Hbase, Spark, and all required libraries have been included in cloudera.

1. Download the apache log  file from http://www.monitorware.com/en/logsamples/ apache.php and unzip it.

2. Create a loganalyzer/input directory named path in HDFS.

   hadoop fs -mkdir -p /user/cloudera/hive/input

3. Copy the log file from the local file system to directory within the HDFS.

   hadoop fs -put access_log /user/cloudera/hive/input/

4. Create appropriate table for string Apache logs.

```
 1 DROP TABLE IF EXISTS access_log;
 2 CREATE TABLE access_log(
 3     host STRING,
 4     identity STRING,
 5     user STRING,
 6     datetime STRING,
 7     requesturl STRING,
 8     respcode STRING,
 9     size STRING)
10     ROW FORMAT SERDE 'org.apache.hadoop.hive.contrib.serde2.RegexSerDe'
11     WITH SERDEPROPERTIES (
12     "input.regex" = "([^ ]*) ([^ ]*) ([^ ]*) (-|\\[[^\\]]*\\]) ([^ \"]*|\"[^\"]*\") (-|[0-9]*) (-|[0-9]*)",
13     "output.format.string" = "%1$s %2$s %3$s %4$s %5$s %6$s %7$s"
14 )
15 STORED AS TEXTFILE;
```

5. Load access_log file, depending location of file (local file system or  HDFS) do on of followings.

```
 1 LOAD DATA LOCAL INPATH "/home/cloudera/access_log" INTO TABLE access_log;
 2
 3 LOAD DATA INPATH "/user/cloudera/hive/input/access_log" INTO TABLE access_log;
```

6. List the count of response code's returned from the server

```
1 SELECT respcode, count(*) as count
2 FROM access_log
3 GROUP BY respcode;
```

Result

| | respcode | count |
|---|---|---|
| 0 | 200 | 447 |
| 1 | 401 | 98 |
| 2 | 404 | 3 |

7. List the top 10 most popular URL's in the Apache log

```
1  CREATE TABLE  urlsummary(
2      requesturl STRING,
3      numrequest int
4  )
5  STORED AS TEXTFILE;
6
7  INSERT OVERWRITE TABLE urlsummary
8  SELECT requesturl, COUNT(*)
9  FROM access_log
10 WHERE host IS NOT NULL GROUP BY requesturl;
```

```
1 SELECT * FROM urlsummary ORDER BY numrequest DESC LIMIT 10;
```

Result

| | urlsummary.requesturl | urlsummary.numrequest |
|---|---|---|
| 0 | "GET /twiki/bin/view/Main/WebHome HTTP/1.1" | 12 |
| 1 | "GET / HTTP/1.1" | 7 |
| 2 | "GET /twiki/pub/TWiki/TWikiLogos/twikiRobot46x50.gif HTTP/1.1" | 6 |
| 3 | "GET /favicon.ico HTTP/1.1" | 6 |
| 4 | "GET /robots.txt HTTP/1.0" | 5 |
| 5 | "GET /twiki/bin/view/Main/SpamAssassinTaggingOnly HTTP/1.1" | 4 |
| 6 | "GET /twiki/bin/view/Main/SpamAssassinAndPostFix HTTP/1.1" | 4 |
| 7 | "GET /razor.html HTTP/1.1" | 3 |
| 8 | "GET /twiki/bin/view/Main/WebHome HTTP/1.0" | 3 |
| 9 | "GET /twiki/bin/view/Main/DCCAndPostFix HTTP/1.1" | 3 |

8. List the content size of responses returned from the server.

```
1 SELECT host, sum(size) as respsize
2 FROM access_log
3 GROUP BY host;
```

Result

| | host | respsize |
|---|---|---|
| 0 | 10.0.0.153 | 183728 |
| 1 | 128.227.88.79 | 81785 |
| 2 | 200.160.249.68.bmf.com.br | 13269 |
| 3 | 206-15-133-181.dialup.ziplink.net | 0 |
| 4 | 212.92.37.62 | 72981 |
| 5 | 213.181.81.4 | 7649 |
| 6 | 219.95.17.51 | 3169 |
| 7 | 61.9.4.61 | 7936 |

9. List the average, min, and max content size of responses returned from the server

```
1 SELECT
2    max(cast(size as BIGINT)) as max,
3    min(cast(size as BIGINT)) as min,
4    avg(cast(size as BIGINT)) as average
5 FROM access_log;
```

Result

| | max | min | average |
|---|---|---|---|
| 0 | 138789 | 0 | 10190.023722627737 |

# Hive and HBase Integration

**Steps for hive and hbase integration**

1. Create a table access_log_hbase and columnfamily as m in hbase
   Create 'access_log_hbase', 'm'

2. List the table in hbase

```
hbase(main):001:0> list
TABLE
access_log_hbase
1 row(s) in 1.7030 seconds
```

3. View the data in **access_log_hbase**

```
hbase(main):002:0> scan 'access_log_hbase'
ROW                    COLUMN+CELL
0 row(s) in 0.7910 seconds

hbase(main):003:0>
```

4. Create an external table in hive as **access_log_hive**

```
1  CREATE EXTERNAL TABLE access_log_hive(
2    datetime STRING,
3    host STRING,
4    identity STRING,
5    user STRING,
6    requesturl STRING,
7    respcode STRING,
8    size STRING
9  )
10 STORED BY 'org.apache.hadoop.hive.hbase.HBaseStorageHandler'
11 WITH SERDEPROPERTIES (
12   "hbase.columns.mapping" = ":key,m:host,m:identity,m:user,m:requesturl,m:respcode,m:size"
13 )
14 TBLPROPERTIES ("hbase.table.name" = "access_log_hbase");
```

5. Overwrite the hive as **access_log_hive** with existing hive table **access_log**

```
1  INSERT OVERWRITE TABLE access_log_hive
2  SELECT datetime,
3  host,
4  identity,
5  user,
6  requesturl,
7  respcode,
8  size
9  FROM access_log;
```

## 6. View the data in access_log_hbase

```
hbase(main):001:0> list
TABLE
access_log_hbase
1 row(s) in 3.1710 seconds

=> ["access_log_hbase"]
hbase(main):002:0> scan 'access_log_hbase'
ROW                              COLUMN+CELL
 [07/Mar/2004:16:05:49 -0800]    column=m:host, timestamp=1450057257703, value=64.242.88.10
 [07/Mar/2004:16:05:49 -0800]    column=m:identity, timestamp=1450057257703, value=-
 [07/Mar/2004:16:05:49 -0800]    column=m:requesturl, timestamp=1450057257703, value="GET /twiki/bin/edit/Main/Double_bounce_sender?topicparent=Main.
                                 ConfigurationVariables HTTP/1.1"
 [07/Mar/2004:16:05:49 -0800]    column=m:respcode, timestamp=1450057257703, value=401
 [07/Mar/2004:16:05:49 -0800]    column=m:size, timestamp=1450057257703, value=12846
 [07/Mar/2004:16:05:49 -0800]    column=m:user, timestamp=1450057257703, value=-
 [07/Mar/2004:16:06:51 -0800]    column=m:host, timestamp=1450057257703, value=64.242.88.10
 [07/Mar/2004:16:06:51 -0800]    column=m:identity, timestamp=1450057257703, value=-
 [07/Mar/2004:16:06:51 -0800]    column=m:requesturl, timestamp=1450057257703, value="GET /twiki/bin/rdiff/TWiki/NewUserTemplate?rev1=1.3&rev2=1.2 HT
                                 TP/1.1"
 [07/Mar/2004:16:06:51 -0800]    column=m:respcode, timestamp=1450057257703, value=200
 [07/Mar/2004:16:06:51 -0800]    column=m:size, timestamp=1450057257703, value=4523
 [07/Mar/2004:16:06:51 -0800]    column=m:user, timestamp=1450057257703, value=-
 [07/Mar/2004:16:10:02 -0800]    column=m:host, timestamp=1450057257703, value=64.242.88.10
 [07/Mar/2004:16:10:02 -0800]    column=m:identity, timestamp=1450057257703, value=-
 [07/Mar/2004:16:10:02 -0800]    column=m:requesturl, timestamp=1450057257703, value="GET /mailman/listinfo/hsdivision HTTP/1.1"
 [07/Mar/2004:16:10:02 -0800]    column=m:respcode, timestamp=1450057257703, value=200
 [07/Mar/2004:16:10:02 -0800]    column=m:size, timestamp=1450057257703, value=6291
 [07/Mar/2004:16:10:02 -0800]    column=m:user, timestamp=1450057257703, value=-
 [07/Mar/2004:16:11:58 -0800]    column=m:host, timestamp=1450057257703, value=64.242.88.10
 [07/Mar/2004:16:11:58 -0800]    column=m:identity, timestamp=1450057257703, value=-
 [07/Mar/2004:16:11:58 -0800]    column=m:requesturl, timestamp=1450057257703, value="GET /twiki/bin/view/TWiki/WikiSyntax HTTP/1.1"
 [07/Mar/2004:16:11:58 -0800]    column=m:respcode, timestamp=1450057257703, value=200
 [07/Mar/2004:16:11:58 -0800]    column=m:size, timestamp=1450057257703, value=7352
 [07/Mar/2004:16:11:58 -0800]    column=m:user, timestamp=1450057257703, value=-
 [07/Mar/2004:16:20:55 -0800]    column=m:host, timestamp=1450057257703, value=64.242.88.10
 [07/Mar/2004:16:20:55 -0800]    column=m:identity, timestamp=1450057257703, value=-
 [07/Mar/2004:16:20:55 -0800]    column=m:requesturl, timestamp=1450057257703, value="GET /twiki/bin/view/Main/DCCAndPostFix HTTP/1.1"
 [07/Mar/2004:16:20:55 -0800]    column=m:respcode, timestamp=1450057257703, value=200
```