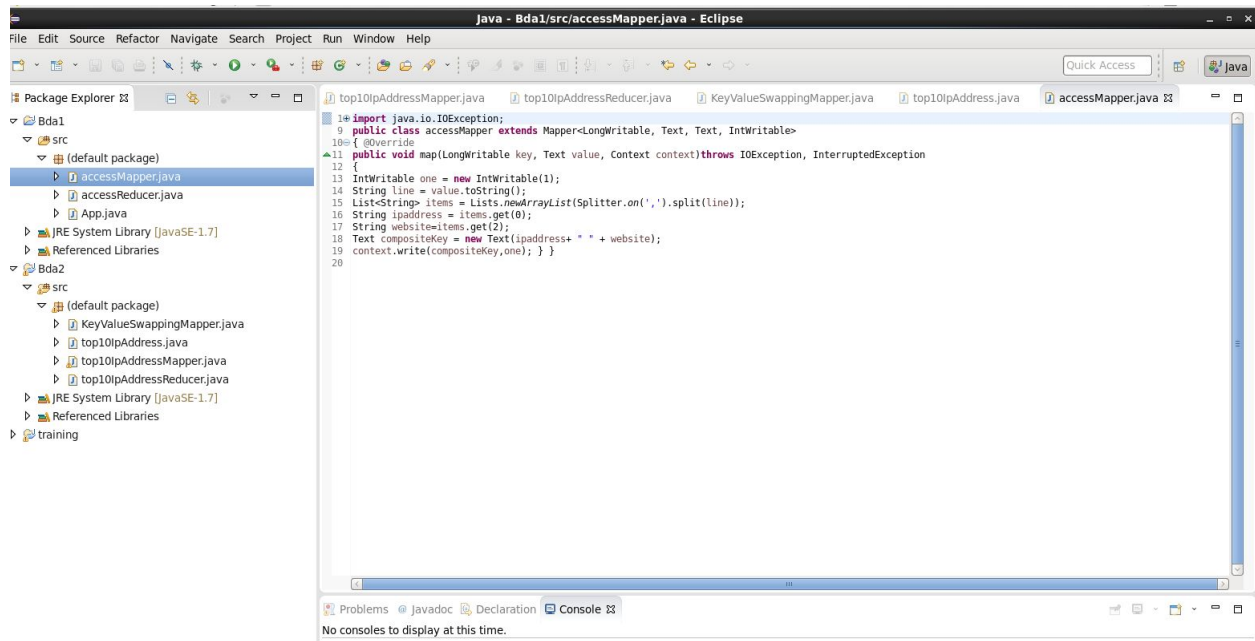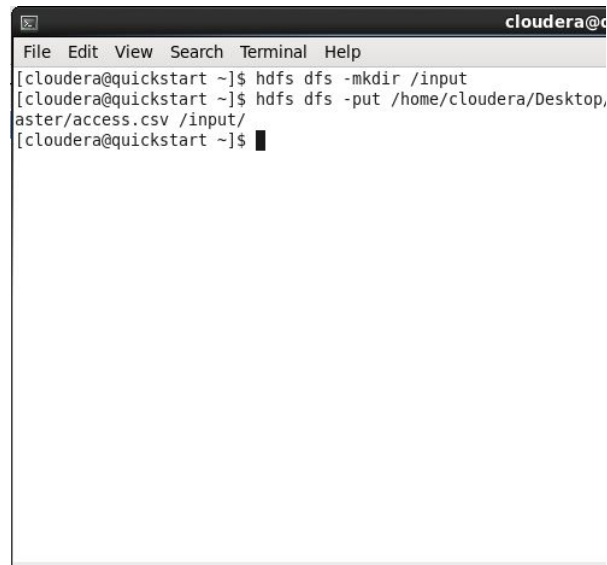# MapReduce:

## Pre requisites:

1) Create .jar files from the collection of java files and save it to /home/cloudera/



2) Two .jar files (Access1.jar and Access2.jar) are saved in /home/cloudera

3) Make a input directory in hdfs dfs and move your dataset into it

```
[cloudera@quickstart ~]$ hdfs dfs -mkdir /input
[cloudera@quickstart ~]$ hdfs dfs -put /home/cloudera/Desktop/
aster/access.csv /input/
[cloudera@quickstart ~]$ █
```

# Analysis:

## 1) MapReduce in Hadoop to find the number of times each IP accessed the website.

## Output:

Command: hadoop jar /home/cloudera/Access1.jar App /input/access.csv /out1

```
                                            cloudera@quickstart:~                                    _ □ □
File  Edit  View  Search  Terminal  Help
33.187.141.154 " ""GET /ftp/plugins/access.ssh/checkInstall.php?destServer=%7C%7Cecho%2088663 HTT
P/1.1"" 404 237"          1
33.187.141.154 " ""GET /gestion_documentos/plugins/access.ssh/checkInstall.php?destServer=%7C%7Ce
:ho%2093559 HTTP/1.1"" 404 252"   1
33.187.141.154 " ""GET /intranet/plugins/access.ssh/checkInstall.php?destServer=%7C%7Cecho%207922
l HTTP/1.1"" 404 242"   1
33.187.141.154 " ""GET /lab/ajaxplorer/plugins/access.ssh/checkInstall.php?destServer=%7C%7Cecho%
2063186 HTTP/1.1"" 404 248"       1
33.187.141.154 " ""GET /login/plugins/access.ssh/checkInstall.php?destServer=%7C%7Cecho%2099436 H
ITP/1.1"" 404 239"        1
33.187.141.154 " ""GET /manager/ajaxplorer-core-3.1.1/plugins/access.ssh/checkInstall.php?destSer
ver=%7C%7Cecho%2049443 HTTP/1.1"" 404 263"        1
33.187.141.154 " ""GET /neos/plugins/access.ssh/checkInstall.php?destServer=%7C%7Cecho%2098 HTTP/
l.1"" 404 238"   1
33.187.141.154 " ""GET /newsdm/plugins/access.ssh/checkInstall.php?destServer=%7C%7Cecho%2094530
ITTP/1.1"" 404 240"       1
33.187.141.154 " ""GET /partners/plugins/access.ssh/checkInstall.php?destServer=%7C%7Cecho%205941
3 HTTP/1.1"" 404 242"   1
33.187.141.154 " ""GET /pdf_and_image_library/plugins/access.ssh/checkInstall.php?destServer=%7C%
7Cecho%207553 HTTP/1.1"" 404 255"       1
33.187.141.154 " ""GET /plugins/access.ssh/checkInstall.php?destServer=%7C%7Cecho%2067708 HTTP/1.
l"" 404 233"    1
33.187.141.154 " ""GET /pool/plugins/access.ssh/checkInstall.php?destServer=%7C%7Cecho%2057510 HT
TP/1.1"" 404 238"   1
33.187.141.154 " ""GET /prints/plugins/access.ssh/checkInstall.php?destServer=%7C%7Cecho%2018300
ITTP/1.1"" 404 240"       1
33.187.141.154 " ""GET /repo/plugins/access.ssh/checkInstall.php?destServer=%7C%7Cecho%2061484 HT
TP/1.1"" 404 238"   1
33.187.141.154 " ""GET /repository/plugins/access.ssh/checkInstall.php?destServer=%7C%7Cecho%2029
047 HTTP/1.1"" 404 244" 1
33.187.141.154 " ""GET /share/plugins/access.ssh/checkInstall.php?destServer=%7C%7Cecho%2051302 H
ITP/1.1"" 404 239"        1
33.187.141.154 " ""GET /test/plugins/access.ssh/checkInstall.php?destServer=%7C%7Cecho%207365 HTT
P/1.1"" 404 238"          1
33.187.141.154 " ""GET /transfer/plugins/access.ssh/checkInstall.php?destServer=%7C%7Cecho%208010
3 HTTP/1.1"" 404 242"   1
33.187.141.154 " ""GET /upload/plugins/access.ssh/checkInstall.php?destServer=%7C%7Cecho%208916 H
ITP/1.1"" 404 240"        1
33.187.141.154 " ""GET /uploader/plugins/access.ssh/checkInstall.php?destServer=%7C%7Cecho%203781
l HTTP/1.1"" 404 242"   1
33.187.141.154 " ""GET /uploads/plugins/access.ssh/checkInstall.php?destServer=%7C%7Cecho%2069800
 HTTP/1.1"" 404 241"      1
33.187.141.154 " ""GET /web/plugins/access.ssh/checkInstall.php?destServer=%7C%7Cecho%2040027 HTT
P/1.1"" 404 237"          1
33.187.141.154 " ""GET /webdav/plugins/access.ssh/checkInstall.php?destServer=%7C%7Cecho%2069212
ITTP/1.1"" 404 240"       1
 🖴 [Java - Bda1/src/acces...  🔲 [Hadoop-Mapreduce-...  🔲 [HadoopMapReduce1....  🖳 cloudera@quickstart:~  🖳 cloudera@quickstart:~
```

## 2) Top 10 most visited IP addresses

## Output:

**Command:** hadoop jar /home/cloudera/Access2.jar top10IpAddress /input/access.csv /out2

```
[cloudera@quickstart ~]$ hadoop jar /home/cloudera/Access2.jar top10IpAddress /input/access.csv /out2
20/12/04 04:48:51 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
20/12/04 04:48:52 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy
20/12/04 04:48:53 INFO input.FileInputFormat: Total input paths to process : 1
20/12/04 04:48:53 WARN hdfs.DFSClient: Caught exception
java.lang.InterruptedException
        at java.lang.Object.wait(Native Method)
        at java.lang.Thread.join(Thread.java:1281)
        at java.lang.Thread.join(Thread.java:1355)
        at org.apache.hadoop.hdfs.DFSOutputStream$DataStreamer.closeResponder(DFSOutputStream.java:967)
        at org.apache.hadoop.hdfs.DFSOutputStream$DataStreamer.endBlock(DFSOutputStream.java:705)
        at org.apache.hadoop.hdfs.DFSOutputStream$DataStreamer.run(DFSOutputStream.java:894)
20/12/04 04:48:53 WARN hdfs.DFSClient: Caught exception
java.lang.InterruptedException
        at java.lang.Object.wait(Native Method)
        at java.lang.Thread.join(Thread.java:1281)
        at java.lang.Thread.join(Thread.java:1355)
        at org.apache.hadoop.hdfs.DFSOutputStream$DataStreamer.closeResponder(DFSOutputStream.java:967)
        at org.apache.hadoop.hdfs.DFSOutputStream$DataStreamer.endBlock(DFSOutputStream.java:705)
        at org.apache.hadoop.hdfs.DFSOutputStream$DataStreamer.run(DFSOutputStream.java:894)
20/12/04 04:48:53 INFO mapreduce.JobSubmitter: number of splits:1
20/12/04 04:48:53 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1607081095019_0002
20/12/04 04:48:54 INFO impl.YarnClientImpl: Submitted application application_1607081095019_0002
20/12/04 04:48:54 INFO mapreduce.Job: The url to track the job: http://quickstart.cloudera:8088/proxy/application_1607081095019_0002/
20/12/04 04:48:54 INFO mapreduce.Job: Running job: job 1607081095019 0002
```

**Top 10 most visited IP addresses**

```
[cloudera@quickstart ~]$ hdfs dfs -cat /out2/Out2/part-r-00000 | head -n10
4958    155.33.18.236
3724    207.248.55.246
2812    10.15.10.129
2108    10.15.10.135
1501    129.10.65.240
1279    107.20.213.124
765     168.144.67.144
667     50.63.154.43
643     72.158.153.33
642     118.102.182.196
```
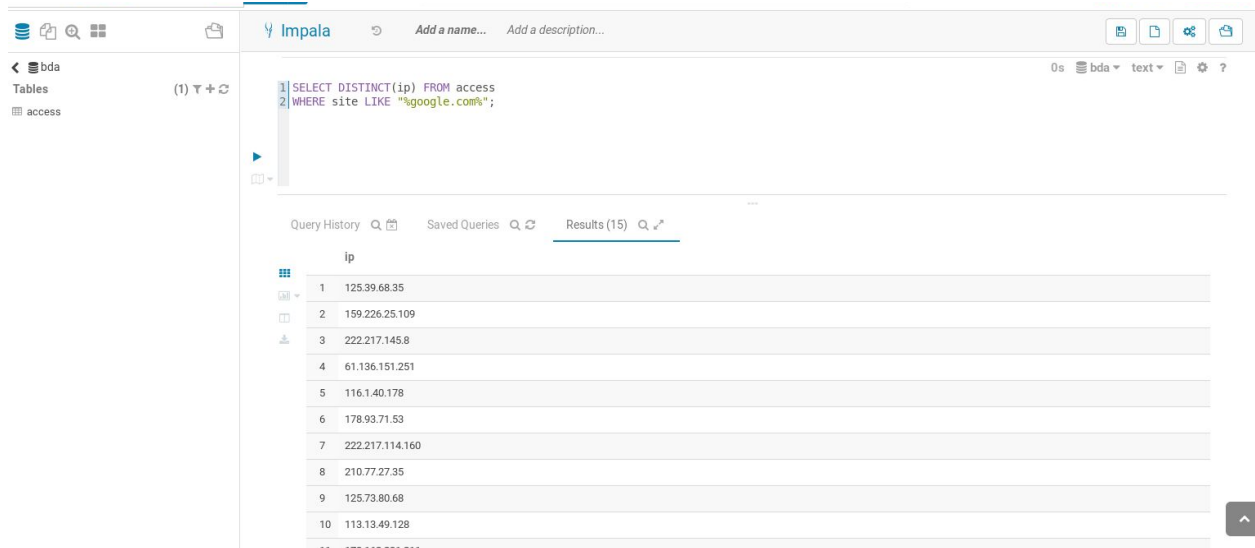
# HIVE:

# Prerequisites :

1)  Upload the dataset into hive by specifying the data types and column names.

## Analysis:

### 1) Ip users that visited "google.com"



### 2) Number of users that visited "google.com"

## 3) Number of times GET method was used



## 1) Sites that were visited in October



## 5) Number of sites visited in October

## 6) Most visited sites

7) Number of users having ips between 127.0.0.0 and 129.0.0.0



8) Ips between 127.0.0.0 and 129.0.0.0 ordered in descending order.