

HDFS Homework.

Task 1.

Here's my screenshot of the converted parquet file:

File Edit Search View Encoding Language Settings Tools Macro Run Plugins Window ?

change.log train_schema.txt destinations_schema.txt train.parquet

```

1 PAR...2014-12-31 21:
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34

```

length: 1,453,677,783 lines: 6,081.56 Ln: 1 Col: 1 Sel: 0 J Macintosh (CR) ANSI INS

And here's the same thing on one of the clusters' nodes:

[illegible]

To run my jar, you need to issue the following command:

```
java -jar program.jar hdfs -schemaPath=schema/path -csvPath=csv/path -newFilePath=file/path -csvSeparator=,
```

Task 2.

HDFS' health can be checked in these 2 ways:

- `sudo -u hdfs hdfs dfsadmin -report`

```

maria_dev@sandbox-hdp/home/maria_dev
[root@sandbox-hdp maria_dev]# sudo -u hdfs hdfs dfsadmin -report
Configured Capacity: 113791799296 (105.98 GB)
Present Capacity: 85713308672 (79.83 GB)
DFS Remaining: 79146093568 (73.71 GB)
DFS Used: 6567215104 (6.12 GB)
DFS Used%: 7.66%
Under replicated blocks: 0
Blocks with corrupt replicas: 0
Missing blocks: 0
Missing blocks (with replication factor 1): 0

-----
Live datanodes (1):

Name: 172.18.0.2:50010 (sandbox-hdp.hortonworks.com)
Hostname: sandbox-hdp.hortonworks.com
Decommission Status : Normal
Configured Capacity: 113791799296 (105.98 GB)
DFS Used: 6567215104 (6.12 GB)
Non DFS Used: 22086615040 (20.57 GB)
DFS Remaining: 79146093568 (73.71 GB)
DFS Used%: 5.77%
DFS Remaining%: 69.55%
Configured Cache Capacity: 0 (0 B)
Cache Used: 0 (0 B)
Cache Remaining: 0 (0 B)
Cache Used%: 100.00%
Cache Remaining%: 0.00%
Xceivers: 4
Last contact: Thu Oct 25 14:51:43 UTC 2018
Last Block Report: Thu Oct 25 14:38:57 UTC 2018

[root@sandbox-hdp maria_dev]#

```

- `hdfs fsck /`

```
[root@sandbox-hdp maria_dev]# hdfs fsck /  
Connecting to namenode via http://sandbox-hdp.hortonworks.com:50070/fsck?ugi=root&path=%2F  
FSCK started by root (auth:SIMPLE) from /172.18.0.2 for path / at Thu Oct 25 14:52:52 UTC 2018  
.....  
.....  
.....  
.....  
.....  
.....  
.....  
.....  
.....Status: HEALTHY  
  
Total size:      6507092582 B  
Total dirs:     245  
Total files:    1093  
Total symlinks:   0 (Files currently being written: 2)  
Total blocks (validated): 1124 (avg. block size 5789228 B) (Total open file blocks (not validated): 2)  
Minimally replicated blocks: 1124 (100.0 %)  
Over-replicated blocks: 0 (0.0 %)  
Under-replicated blocks: 0 (0.0 %)  
Mis-replicated blocks: 0 (0.0 %)  
Default replication factor: 1  
Average block replication: 1.0  
Corrupt blocks: 0  
Missing replicas: 0 (0.0 %)  
Number of data-nodes: 1  
Number of racks: 1  
FSCK ended at Thu Oct 25 14:52:52 UTC 2018 in 91 milliseconds  
  
The filesystem under path '/' is HEALTHY  
[root@sandbox-hdp maria_dev]#
```

I will describe what the output of the first command shows.

- Configured Capacity – is the capacity that is available to HDFS for storage. It's calculated as follows: **Configured Capacity = Total Disk Space – Reserved Space**. Reserved Space is allocated for OS level operations. The parameter responsible for Reserved Space is **dfs.datanode.du.reserved**, it can be changed in **hdfs-site.xml**.
- Present Capacity – is the total amount of storage space which is actually available for storing the files after allocating some space for metadata and open-blocks (Non DFS Used space). So, the difference: **Configured Capacity – Present Capacity** is used for storing metadata and other information.
- DFS Used – is the storage space that has been used up by the HDFS. In order to get the actual size of the files stored in HDFS, you need to divide DFS Used by the replication factor. Replication factor is defined as **dfs.replication** and can be found in **hdfs-site.xml**. So, if DFS Used is 90 GB, and replication factor is 3, then the size of the actual files is 90/3=30GB.
- DFS Remaining – is the amount of storage available to HDFS to store more files. If you have 90GB of remaining space and your replication factor is 3, then you can store 90/3=30GB of files.
Present Capacity = DFS Used + DFS Remaining.
- Non DFS Used is any data in the filesystem of the data nodes that isn't in **dfs.datanode.data.dir**.
Non DFS Used = Configured Capacity – Present Capacity.
- Under-replicated blocks – are blocks whose number of replicas does not meet the target replication number. HDFS Will automatically create new replicas for under-replicated blocks until they meet the target replication. You can view the information about it by issuing the **hdfs dfsadmin -metasave** command.
- Blocks with corrupt replicas – are blocks who have at least one corrupt replica and at least one live replica.
- Missing blocks – blocks that have missing replicas.
- Missing blocks (with replication factor 1) – blocks that are irreparable.
- Xceivers – defines the number of server-side threads (that can be used by sockets for data connections).

And finally, this is how to move or rename a file in hdfs:

hdfs dfs -mv /path/to/file/name /new/path/to/file/new_name

Example:

```
maria_dev@sandbox-hdp:/home/maria_dev
[root@sandbox-hdp maria_dev]# hdfs dfs -mv /user/maria_dev/test.csv /user/maria_dev/test_test_test.csv
[root@sandbox-hdp maria_dev]# hdfs dfs -ls /user/maria_dev
Found 4 items
-rw-r--r--  1 root hdfs  138159416 2018-10-23 10:16 /user/maria_dev/destinations.csv
-rw-r--r--  1 root hdfs   31756066 2018-10-23 10:16 /user/maria_dev/sample_submission.csv
-rw-r--r--  1 root hdfs   276554476 2018-10-23 10:16 /user/maria_dev/test_test_test.csv
-rw-r--r--  1 root hdfs  4070445781 2018-10-23 10:17 /user/maria_dev/train.csv
[root@sandbox-hdp maria_dev]#
```