

Hadoop Commands (2.4.0)

CSE Workshop VM • Ubuntu 14.04 Trusty Tahr

User names: `cse-user` (for general exercises in MATLAB, Python, R); `hduser` (for Hadoop)

Password: `cse-user1` (all users)

This VM has a dedicated Hadoop user account, a recommended practice to separate the Hadoop installation from other software applications and user accounts on the same machine. To log in for Hadoop exercises:

```
cse-user@csevm$ sudo login hduser
Enter password for sudo: cse-user1
Enter password: cse-user1
hduser@csevm$
```

FS Shell Commands

`hadoop fs <args> or hdfs dfs <args>`

All FS shell commands take paths as arguments in the format `scheme://authority/path`. For HDFS the scheme is `hdfs`, and for the local filesystem the scheme is `file`. `scheme` and `authority` are optional. An HDFS file or directory such as `/parent/child` can be specified as `hdfs://namenodehost/parent/child` or simply as `/parent/child`. Most commands in FS shell behave like corresponding Unix commands.

<code>cat</code>	Copies source path to stdout. <code>hadoop fs -cat /user/script.sh</code>	<code>mkdir</code>	Takes path as argument and creates directories. <code>hadoop fs -mkdir /user/userdir</code>
<code>copyFromLocal</code>	Like <code>put</code> , except source restricted to local file reference. <code>hadoop fs -copyFromLocal ./data/file01.csv /user/data</code>	<code>mv</code>	Moves files from source to destination (but not across filesystems). <code>hadoop fs -mv /user/hadoop/file1 /user/hadoop/file2</code>
<code>copyToLocal</code>	Like <code>get</code> , except that destination is restricted to local file reference. <code>hadoop fs -copyToLocal [-ignorecrc] [-crc] URI <localdst></code>	<code>put</code>	Copy single or multiple sources from local file system to destination filesystem. <code>hadoop fs -put localfile /user/hdfsfile</code>
<code>cp</code>	Copy files from source to destination. <code>hadoop fs -cp /user/oldfile /user/newfile</code>	<code>rm</code>	Delete files specified as args. Only deletes non empty directory and files. <code>hadoop fs -rm /user/emptydir</code>
<code>du</code>	Displays aggregate length of files contained in directory or length of file. <code>hadoop fs -du /user/data</code>	<code>rmdir</code>	Recursive version of delete. <code>hadoop fs -rmdir /user/nonemptydir</code>
<code>expunge</code>	Empty the Trash. (See <i>HDFS Design</i> for more information.) <code>hadoop fs -expunge</code>	<code>tail</code>	Displays last kilobyte of the file to stdout. <code>hadoop fs -tail /user/data/01.dat</code>
<code>get</code>	Copy files to the local file system. <code>hadoop fs -get /user/hdfsfile localfile</code>	<code>test -e</code>	Check if file exists. (0 if true)
<code>ls</code>	For file, returns stat on the file; for directory, returns list of direct children <code>hadoop fs -ls /user/file</code>	<code>-z</code>	Check if zero length file. (0 if true)
		<code>-d</code>	Check if directory (return 1) or not. <code>hadoop fs -test -e filename</code>

Running Hadoop Locally

1. Format the filesystem:
`$ hdfs namenode -format`
2. Start NameNode daemon and DataNode daemon:
`$ start-all.sh`
3. Browse the web interface for the NameNode, available at <http://localhost:50070>.
4. Make the HDFS directories required to execute MapReduce jobs:
`$ hdfs dfs -mkdir /user`
`$ hdfs dfs -mkdir /user/hduser`
`$ hdfs dfs -mkdir /user/hduser/input`
5. Copy the input files into the distributed filesystem:
`$ hdfs dfs -copyFromLocal $HADOOP_HOME/etc/hadoop /user/hduser/input`
6. Run an example, one which here counts the number of instances of each string matching the regular expression:
`$ hadoop jar $HADOOP_HOME/share/hadoop/mapreduce/hadoop-mapreduce-examples-2.4.0.jar`
`grep /user/hduser/input/hadoop /user/hduser/output '<[a-z.]+>'`
7. Examine the output files.
 - a. Copy output files from the distributed filesystem to the local filesystem and examine them:
`$ hdfs dfs -get output output`
`$ cat output/*`
 - b. View output files on the distributed filesystem:
`$ hdfs dfs -cat output/*`
8. When done, stop the daemons with:
`$ stop-dfs.sh`

Hadoop Daemon Ports (<http://localhost:#####>)

MAP/REDUCE

50030	job tracker	50060	task tracker
-------	-------------	-------	--------------

HDFS

8020	name node	50090	secondary name node
50070			
50010	data nodes	50100	backup/checkpoint node
50020		50105	
50075			