

Java-Success.com

Prepare to fast-track, choose & go places with 800+ Java & Big Data Q&As with lots of code & diagrams.

[Home](#) [Why? ▾](#) [300+ Java FAQs ▾](#) [300+ Big Data FAQs ▾](#) [Courses ▾](#)

[👤 Membership ▾](#) [Your Career ▾](#)

[Home](#) > [bigdata-success.com](#) > [Tutorials - Big Data](#) > [TUT - Cloudera on Docker](#) > 13:

Docker Tutorial: Apache Spark (spark-shell & pyspark) on Cloudera quickstart

13: Docker Tutorial: Apache Spark (spark-shell & pyspark) on Cloudera quickstart

 Posted on [May 30, 2019](#)

This extends [Docker Tutorial: BigData on Cloudera quickstart via Docker](#).

Step 1: Run the container on a command line.

```
1 ~/projects/docker-hadoop]$ docker run --hostname=qu
2 --privileged=true -t -i -v /Users/arulkumarankumar
3 --publish-all=true -p 8888:8888 -p 80:80 -p 7180:71
4 cloudera/quickstart /usr/bin/docker-quickstart
```

300+ Java Interview FAQs

300+ Java FAQs



16+ Java Key Areas Q&As



150+ Java Architect FAQs



80+ Java Code Quality Q&As



150+ Java Coding Q&As



300+ Big Data Interview FAQs

300+ Big Data FAQs



Tutorials - Big Data



TUT -  Starting Big Data

TUT - Starting Spark & Scala

Step 2: Create a text file “**simple.txt**” in the local file system.

```
1 [root@quickstart /]# touch simple.txt
2 [root@quickstart /]# vi simple.txt
3
1 A big brown fox jumped over a big brown fence.
```

Step 3: Copy this file onto HDFS file system.

```
1 [root@quickstart /]# hdfs dfs -copyFromLocal simple
```

List the file that was copied to HDFS:

```
1 [root@quickstart /]# hdfs dfs -ls /user/root/
2 Found 1 items
3 -rw-r--r-- 1 root supergroup 47 2019-05-2
4 [root@quickstart /]#
```

spark-submit

Step 4: Run the Spark job in the Spark examples.

Cloudera quickstart comes with some Spark examples packaged as a jar file.

```
1 [root@quickstart /]# ls -ltr /usr/lib/spark/examples
2 total 4
3 lrwxrwxrwx 1 root root 64 Apr 6 2016 spark-examp
4 lrwxrwxrwx 1 root root 23 Apr 6 2016 python.tar.g
5
```

Look for the “**JavaWordCount**”.

```
1 [root@quickstart /]# jar -tvf /usr/lib/spark/examp
2 1326 Wed Mar 23 11:54:48 UTC 2016 org/apache/spa
```

TUT - Starting with Python

TUT - Kafka

TUT - Pig

TUT - Apache Storm

TUT - Spark Scala on Zeppelin

TUT - Cloudera

TUT - Cloudera on Docker

TUT - File Formats

TUT - Spark on Docker

TUT - Flume

TUT - Hadoop (HDFS)

TUT - HBase (NoSQL)

TUT - Hive (SQL)

TUT - Hadoop & Spark

TUT - MapReduce

TUT - Spark and Scala

TUT - Spark & Java

TUT - PySpark on Databricks

TUT - Zookeeper

800+ Java Interview Q&As

300+ Core Java Q&As



300+ Enterprise Java Q&As



150+ Java Frameworks Q&As



120+ Companion Tech Q&As



Tutorials - Enterprise Java



```

3 3387 Wed Mar 23 11:54:48 UTC 2016 org/apache/spa
4 1181 Wed Mar 23 11:54:48 UTC 2016 org/apache/spa
5
1 [root@quickstart /]# spark-submit --class org.apac
2 ...
3 ...
4 fox: 1
5 fence.: 1
6 a: 1
7 big: 2
8 A: 1
9 over: 1
10 brown: 2
11 jumped: 1
12

```

spark-shell

```

1 [root@quickstart /]# spark-shell
2 SLF4J: Class path contains multiple SLF4J bindings
3 SLF4J: Found binding in [jar:file:/usr/lib/zookee
4 SLF4J: Found binding in [jar:file:/usr/jars/slf4j
5 SLF4J: See http://www.slf4j.org/codes.html#multip
6 SLF4J: Actual binding is of type [org.slf4j.impl.
7 Setting default log level to "WARN".
8 To adjust logging level use sc.setLogLevel(newLev
9 Welcome to
10
11      ____
12     /  __ \  _  /  _  /  _  /  _  /  _  /  _  /
13    /  /  \ /  \ /  \ /  \ /  \ /  \ /  \ /  \ /  \
14   /  /  \ /  \ /  \ /  \ /  \ /  \ /  \ /  \ /  \
15  /  /  \ /  \ /  \ /  \ /  \ /  \ /  \ /  \ /  \
16 Using Scala version 2.10.5 (Java HotSpot(TM) 64-B
17 Type in expressions to have them evaluated.
18 Type :help for more information.
19 19/05/29 13:45:32 WARN util.NativeCodeLoader: Unab
20 Spark context available as sc (master = local[*],
21 19/05/29 13:45:35 WARN DataNucleus.General: Plugi
22 19/05/29 13:45:35 WARN DataNucleus.General: Plugi
23 19/05/29 13:45:41 WARN metastore.ObjectStore: Vers
24 19/05/29 13:45:41 WARN metastore.ObjectStore: Fai
25 19/05/29 13:45:43 WARN shortcircuit.DomainSocketFo
26 SQL context available as sqlContext.
27
28 scala>
29

```

Interactive Spark code in Scala. The SparkContext will be available as “sc”.

```
1 scala> var map = sc.textFile("hdfs:///user/root/s
2 map: org.apache.spark.rdd.RDD[(String, Int)] = Map
3
4 scala> var counts = map.reduceByKey(_ + _)
5 counts: org.apache.spark.rdd.RDD[(String, Int)] =
6
7 scala> counts.foreach(println)
8 (big,2)
9 (fox,1)
10 (fence.,1)
11 (over,1)
12 (a,1)
13 (brown,2)
14 (A,1)
15 (jumped,1)
16
```

pyspark shell

```
1 [root@quickstart /]# pyspark
2 Python 2.6.6 (r266:84292, Jul 23 2015, 15:22:56)
3 [GCC 4.4.7 20120313 (Red Hat 4.4.7-11)] on linux2
4 Type "help", "copyright", "credits" or "license"
5 SLF4J: Class path contains multiple SLF4J bindings:
6 SLF4J: Found binding in [jar:file:/usr/lib/zookee
7 SLF4J: Found binding in [jar:file:/usr/jars/slf4j-
8 SLF4J: See http://www.slf4j.org/codes.html#multip
9 SLF4J: Actual binding is of type [org.slf4j.impl.L
10 19/05/29 13:57:40 INFO spark.SparkContext: Running
11 19/05/29 13:57:40 WARN util.NativeCodeLoader: Unab
12 19/05/29 13:57:41 INFO spark.SecurityManager: Char
13 19/05/29 13:57:41 INFO spark.SecurityManager: Char
14 19/05/29 13:57:41 INFO spark.SecurityManager: Secu
15 19/05/29 13:57:41 INFO util.Utills: Successfully st
16 19/05/29 13:57:41 INFO slf4j.Slf4jLogger: Slf4jLog
17 19/05/29 13:57:41 INFO Remoting: Starting remoting
18 19/05/29 13:57:41 INFO Remoting: Remoting started
19 19/05/29 13:57:41 INFO Remoting: Remoting now list
20 19/05/29 13:57:41 INFO util.Utills: Successfully st
21 19/05/29 13:57:41 INFO spark.SparkEnv: Registering
22 19/05/29 13:57:41 INFO spark.SparkEnv: Registering
23 19/05/29 13:57:41 INFO storage.DiskBlockManager: C
24 19/05/29 13:57:41 INFO storage.MemoryStore: Memory
25 19/05/29 13:57:41 INFO spark.SparkEnv: Registering
26 19/05/29 13:57:41 INFO server.Server: jetty-8.y.z-
27 19/05/29 13:57:42 INFO server.AbstractConnector: S
```

```

28 19/05/29 13:57:42 INFO util.Utils: Successfully started service 'org.apache.hadoop.yarn.server.NAMER'
29 19/05/29 13:57:42 INFO ui.SparkUI: Started SparkUI
30 19/05/29 13:57:42 INFO executor.Executor: Starting on 'localhost'
31 19/05/29 13:57:42 INFO util.Utils: Successfully started service 'org.apache.hadoop.yarn.server.NAMER'
32 19/05/29 13:57:42 INFO netty.NettyBlockTransferServer: Started
33 19/05/29 13:57:42 INFO storage.BlockManagerMaster: Started
34 19/05/29 13:57:42 INFO storage.BlockManagerMaster: Started
35 19/05/29 13:57:42 INFO storage.BlockManagerMaster: Started
36 Welcome to
37
38      ____
39     /  __ \  _ __| | | |
40    / _ \| | | | | | | |
41   / ___| | | | | | | |
42  /___|_| |_| |_| |_|
43 Using Python version 2.6.6 (r266:84292, Jul 23 2008)
44 SparkContext available as sc, HiveContext available as hc
45 >>>
46

```

```

1 Using Python version 2.6.6 (r266:84292, Jul 23 2008)
2 SparkContext available as sc, HiveContext available as hc
3 >>> input_file = sc.textFile("hdfs:///user/root/simple-out.txt")
4 19/05/29 14:00:11 INFO storage.MemoryStore: Block written to memory
5 19/05/29 14:00:11 INFO storage.MemoryStore: Block written to memory
6 19/05/29 14:00:11 INFO storage.BlockManagerInfo: 1.0 MB
7 19/05/29 14:00:11 INFO spark.SparkContext: Created mapred job
8 >>> map = input_file.flatMap(lambda line: line.split(" "))
9 >>> counts = map.reduceByKey(lambda a, b: a + b)
10 19/05/29 14:00:54 WARN shortcircuit.DomainSocketFactory: No domain sockets found for 'localhost'
11 19/05/29 14:00:54 INFO mapred.FileInputFormat: Total input files: 1
12 >>> for x in counts.collect():
13 ...     print(x)
14 ...
15 (u'A', 1)
16 (u'a', 1)
17 (u'over', 1)
18 (u'fox', 1)
19 (u'brown', 2)
20 (u'jumped', 1)
21 (u'big', 2)
22 (u'fence.', 1)
23 >>>
24

```

Writing the output to another file “simple-out.txt”.

```

1 >>> counts.saveAsTextFile("hdfs:///user/root/simple-out.txt")

```

Quit the interactive shell with “quit()”.

List the “simple-out.txt” folder.

```
1 root@quickstart [/]# hdfs dfs -ls hdfs:///user/root/
2 Found 3 items
3 -rw-r--r--    1 root supergroup      0 2019-05-20 16:00 hdfs:///user/root/simple-out.txt
4 -rw-r--r--    1 root supergroup    45 2019-05-20 16:00 hdfs:///user/root/simple-out.txt
5 -rw-r--r--    1 root supergroup    56 2019-05-20 16:00 hdfs:///user/root/simple-out.txt
6 [root@quickstart [/]#
```

View the “simple-out.txt”

```
1 [root@quickstart [/]# hdfs dfs -cat hdfs:///user/root/simple-out.txt
2 (u'A', 1)
3 (u'a', 1)
4 (u'over', 1)
5 (u'fox', 1)
6 [root@quickstart [/]# hdfs dfs -cat hdfs:///user/root/simple-out.txt
7 (u'brown', 2)
8 (u'jumped', 1)
9 (u'big', 2)
10 (u'fence.', 1)
11 [root@quickstart [/]#
```

◀ 12: Docker Tutorial: Hadoop Big Data configuration files on Cloudera quickstart

14: Docker Tutorial: Hive (via beeline) on Cloudera quickstart ▶

Disclaimer

The contents in this Java-Success are copyrighted and from EmpoweringTech pty ltd. The EmpoweringTech pty ltd has the right to correct or enhance the current content without any prior notice. These are general advice only, and one needs to take his/her own circumstances into consideration. The EmpoweringTech pty ltd will not be held liable for any damages caused or alleged to be caused either directly or indirectly by these materials and resources. Any trademarked names or labels used in this blog remain the property of their respective trademark owners. Links to external sites do not imply endorsement of the linked-to sites. [Privacy Policy](#)