

Spark Fundamentals

Getting started with Spark

Contents

- GETTING START WITH SPARK 3
 - 1.1 SETTING UP THE LAB ENVIRONMENT 4
 - 1.2 STARTING UP THE SPARK SHELL 6
 - 1.2.1 USING SCALA..... 6
 - 1.2.2 USING PYTHON 11
 - 1.3 MORE ON RDD OPERATIONS..... 16
 - 1.3.1 WITH SCALA 16
 - 1.3.2 WITH PYTHON..... 19
 - 1.4 USING SPARK CACHING 23
 - SUMMARY 24

Getting started with Spark

Spark is built around speed and the ease of use. In this section, you will see for yourself how easy it is to get started using Spark with IBM BigInsights. You will be doing some interactive analysis with the Spark Shell, which comes in two flavors. It is available in either Scala or Python. Scala runs on the Java VM and is thus a good way to use existing Java libraries.

Spark's primary abstraction is a distributed collection of items called a Resilient Distributed Dataset or RDD. In a subsequent lab exercise, you will learn more about the details of RDD. In this lab exercise, you'll get to see quickly how it works. RDDs have actions, which return values, and transformations, which return pointers to new RDD. In this lab exercise, you can decide if you wish to do the exercises in the Scala shell or the Python shell – both goes through the same steps, just slightly different syntax.

After completing this hands-on lab, you should be able to:

- Start the Spark shell with Scala and Python
- Perform basic RDD actions and transformations
- Use caching to speed up repeated operations

Allow 30-45 minutes to complete this section of lab.

1.1 Setting up the lab environment

Download the QSE v4 or use the Cloud environment. Review the README file that comes with the QSE for instructions on getting started.

This lab was designed with the Spark cluster on the cloud using an older version of Spark, so the screenshots and texts will be slightly different if you are using the latest QSE. The lab was also designed with Spark 1.1, but tested against the latest QSE (version 4) with Spark 1.2.1.

If you are using the QSE image, set up the following variables in the `~/.bashrc` file to make it easier to work with Spark.

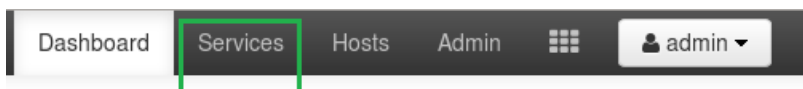
```
export SPARK_HOME=/usr/iop/current/spark-client
```

Be sure you refer to the section to get the labfiles onto the VM image. You will need these files throughout the lab exercise.

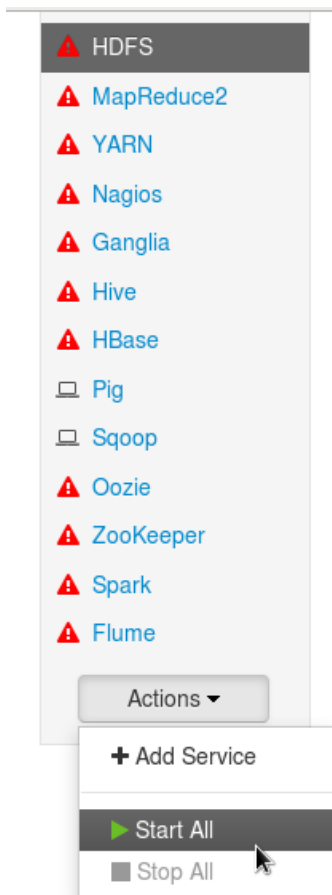
The default user id and password on the QSE image is **virtuser/password**. If you are using the cloud environment, you will use your personal id.

One thing to note: sometimes after you run a command (either Scala or Python), the terminal may seem like got stuck. When that happens, just hit the Enter key to get it back to the prompt.

- __1. Once you have the VM started and you log in to the OS, you must start up the services. The Firefox browser should start up automatically, if not, open it up and point the browser to: `rvm.svl.ibm.com:8080`
- __2. Log in to the Ambari console using the default user id and password: **admin/admin**
- __3. Once logged in, you will be brought to the dashboard. Click on the Services link next to Dashboard:



- __4. On the left side of the window, near the bottom, Click Actions and Start All to start up the services:



It may take a while for all the services to start up. Once it does, you may proceed with the rest of the lab.

Troubleshooting:

Browser cannot load the Ambari user interface

If you have problems with the browser not being able to load the Ambari user interface, then the IP address of the machine might have changed due to a network connectivity reset. Network reset can happen during these conditions:

- Changing from wifi to wired or wired to wifi.
- Shutting down the system, and then re-establishing connection.
- Sending the machine into sleep mode. When it comes back up, the network IP address might be reset.

To resolve this issue, run the following scripts as the **root** user (root/password):

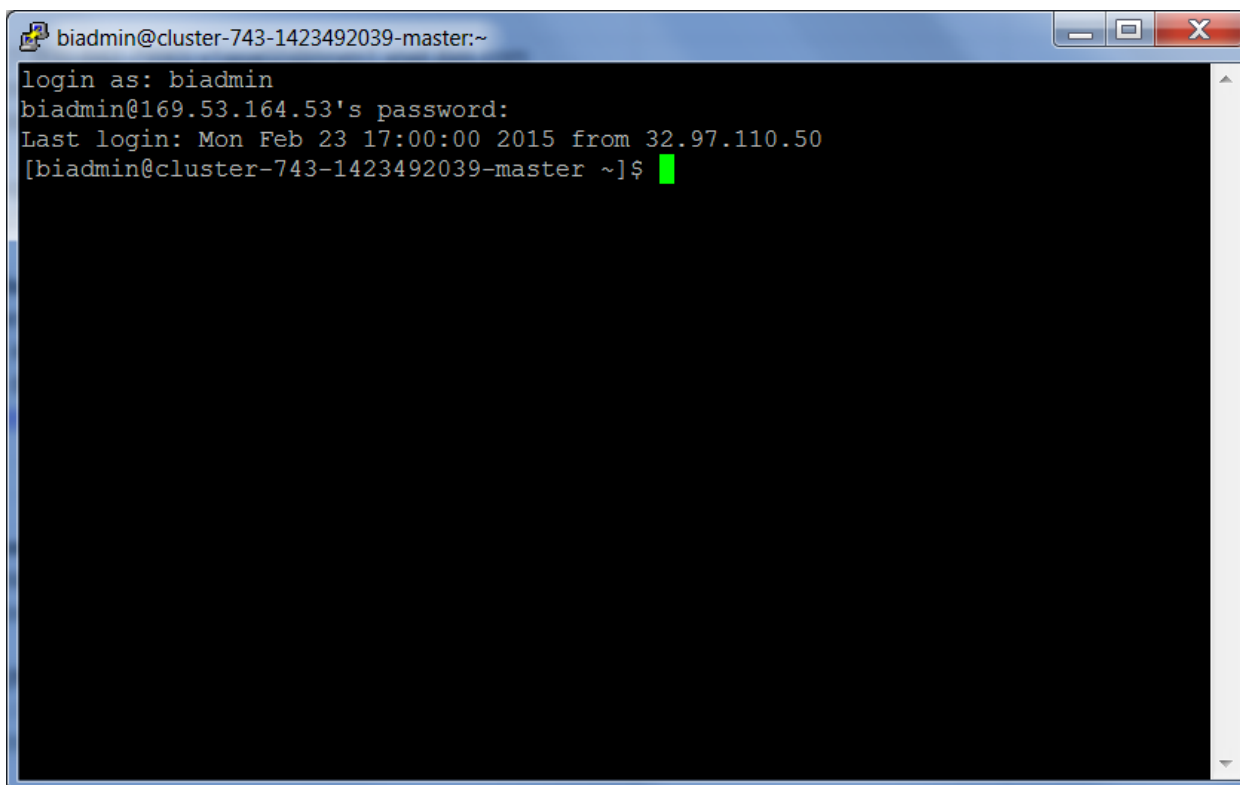
```
/home/virtuser/setHostname.sh  
/home/virtuser/start-all.sh
```

Then, try to connect to `http://rvm.svl.ibm.com:8080`. You can also reboot the VM so that the IP address changes are updated as part of the initialization scripts.

1.2 Starting up the Spark Shell

Spark's shell provides a simple way to learn the APIs. In this section, you will get to use both Scala and Python to work with Spark's API. The first half will cover the shell using Scala. The second half will show the same steps, but using Python instead. To use the Spark shell, you will need to SSH to the cluster. Use your favorite SSH tool to connect to the cluster. This lab exercise will use PuTTY.

__5. Launch a new terminal.



1.2.1 Using Scala

__1. Copy the README.md file from the *labfiles* directory into the /tmp directory on the HDFS. If you don't have the *labfiles* directory, go back to the BDU page and download and set up those files. In the terminal window, execute:

```
hadoop fs -copyFromLocal /home/virtuser/labfiles/README.md /tmp
```

__2. Start the Spark shell with this command:

```
$SPARK_HOME/bin/spark-shell
```

```

biadmin@cluster-743-1423492039-master:~
15/02/24 12:11:22 INFO cluster.YarnClientSchedulerBackend: Application report fr
om ASM:
    appMasterRpcPort: 0
    appStartTime: 1424801478293
    yarnAppState: RUNNING

15/02/24 12:11:24 INFO cluster.YarnClientSchedulerBackend: Registered executor:
Actor[akka.tcp://sparkExecutor@cluster-743-1423492039-master.imdemocloud.com:129
90/user/Executor#546582308] with ID 1
15/02/24 12:11:25 INFO util.RackResolver: Resolved cluster-743-1423492039-master
.imdemocloud.com to /default-rack
15/02/24 12:11:25 INFO storage.BlockManagerMasterActor: Registering block manage
r cluster-743-1423492039-master.imdemocloud.com:10429 with 553.0 MB RAM
15/02/24 12:11:26 INFO cluster.YarnClientSchedulerBackend: Registered executor:
Actor[akka.tcp://sparkExecutor@cluster-743-1423492039-master.imdemocloud.com:278
30/user/Executor#-1740289189] with ID 2
15/02/24 12:11:26 INFO cluster.YarnClientSchedulerBackend: SchedulerBackend is r
eady for scheduling beginning after reached minRegisteredResourcesRatio: 0.8
15/02/24 12:11:26 INFO repl.SparkILoop: Created spark context..
15/02/24 12:11:26 INFO storage.BlockManagerMasterActor: Registering block manage
r cluster-743-1423492039-master.imdemocloud.com:8695 with 553.0 MB RAM
Spark context available as sc.

scala>

```

- ___3. The shell provides code assist. For example, type in "sc." followed by the Tab key to get the list of options associated with the spark context:

```

scala> sc.
accumulableCollection      accumulator              addFile                  addJar
addSparkListener          appName                asInstanceOf             broadcast
cancelJobGroup            clearCallSite         clearFiles               cancelAllJobs
defaultMinPartitions      clearCallSite         clearFiles               clearJobGroup
getAllPools               getCheckpointDir      defaultParallelism      emptyRDD
getLocalProperty          getPoolForName       getExecutorMemoryStatus getExecutorStorageStatus
hadoopConfiguration      getPersistentRDDs    getRDDStorageInfo      getSchedulingMode
isLocal                   jars                 initLocalProperties     getInstanceOf
newAPIHadoopRDD           objectFile            master                  newAPIHadoopFile
sequenceFile              setCallSite           parallelize              runJob
setLocalProperty          sparkUser             setCheckpointDir        setJobDescription
tachyonFolderName        stop                  stop                     submitJob
wholeTextFiles            toString              union                    version

scala> sc.

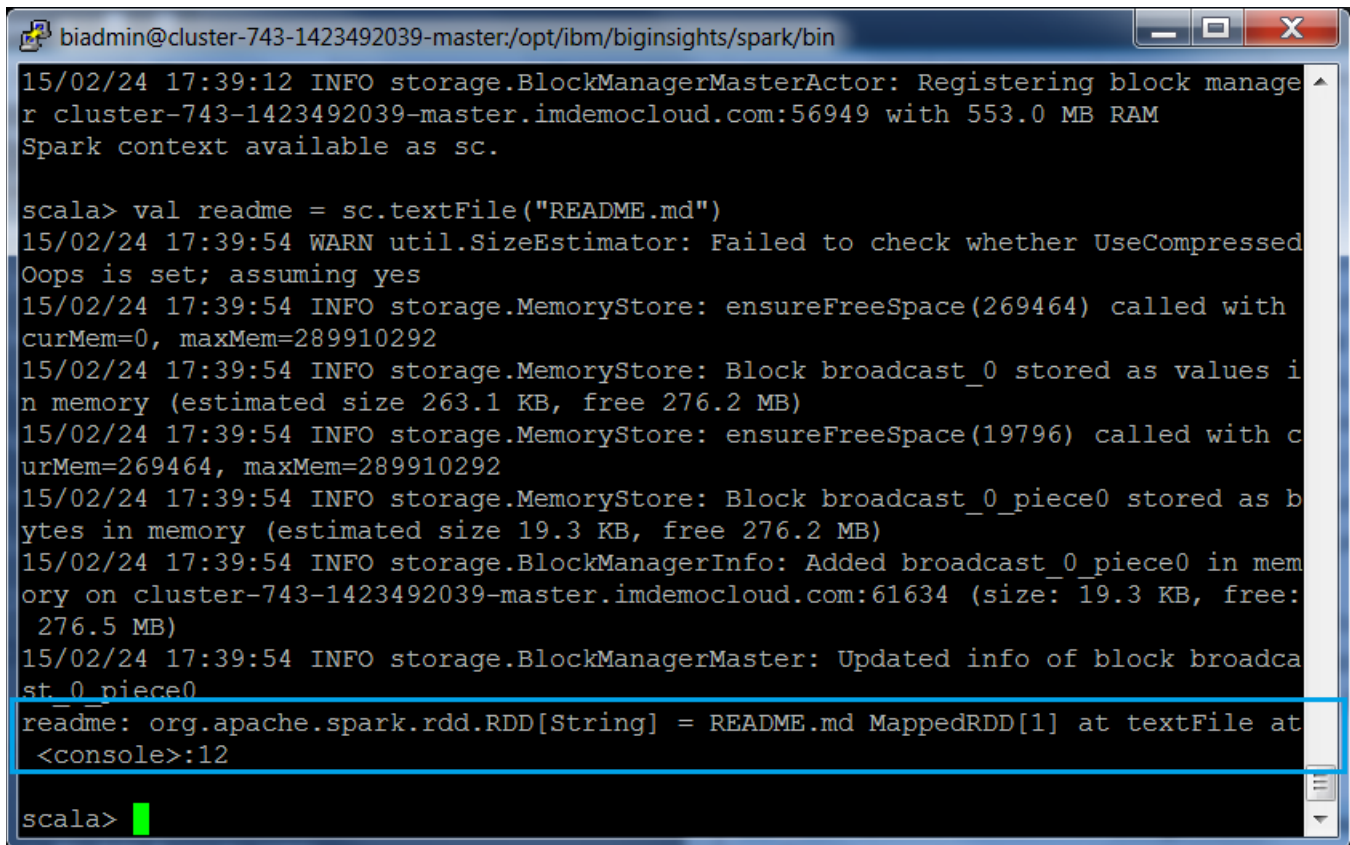
```

- ___4. Note the path of the readme text file is on the HDFS. Type in:

```
val readme = sc.textFile("/tmp/README.md")
```

The parameter of the textFile method is the location within the /tmp directory on the HDFS.

This was a RDD transformation, thus it return a pointer to a RDD, which we have named as *readme*. You can see from the returned line that the readme is a pointer to the RDD.

A screenshot of a terminal window with a blue title bar. The title bar text is 'biadmin@cluster-743-1423492039-master:/opt/ibm/biginsights/spark/bin'. The terminal shows several log messages from Spark, including 'storage.BlockManagerMasterActor: Registering block manager', 'util.SizeEstimator: Failed to check whether UseCompressedOops is set', and 'storage.MemoryStore: ensureFreeSpace'. A Scala command 'scala> val readme = sc.textFile("README.md")' is executed, followed by more logs. A line is highlighted with a blue box: 'readme: org.apache.spark.rdd.RDD[String] = README.md MappedRDD[1] at textFile at <console>:12'. The prompt 'scala>' is followed by a green cursor.

```
biadmin@cluster-743-1423492039-master:/opt/ibm/biginsights/spark/bin
15/02/24 17:39:12 INFO storage.BlockManagerMasterActor: Registering block manager cluster-743-1423492039-master.imdemocloud.com:56949 with 553.0 MB RAM
Spark context available as sc.

scala> val readme = sc.textFile("README.md")
15/02/24 17:39:54 WARN util.SizeEstimator: Failed to check whether UseCompressedOops is set; assuming yes
15/02/24 17:39:54 INFO storage.MemoryStore: ensureFreeSpace(269464) called with curMem=0, maxMem=289910292
15/02/24 17:39:54 INFO storage.MemoryStore: Block broadcast_0 stored as values in memory (estimated size 263.1 KB, free 276.2 MB)
15/02/24 17:39:54 INFO storage.MemoryStore: ensureFreeSpace(19796) called with curMem=269464, maxMem=289910292
15/02/24 17:39:54 INFO storage.MemoryStore: Block broadcast_0_piece0 stored as bytes in memory (estimated size 19.3 KB, free 276.2 MB)
15/02/24 17:39:54 INFO storage.BlockManagerInfo: Added broadcast_0_piece0 in memory on cluster-743-1423492039-master.imdemocloud.com:61634 (size: 19.3 KB, free: 276.5 MB)
15/02/24 17:39:54 INFO storage.BlockManagerMaster: Updated info of block broadcast_0_piece0
readme: org.apache.spark.rdd.RDD[String] = README.md MappedRDD[1] at textFile at <console>:12

scala> █
```

- __5. Let's perform some RDD actions on this text file. Count the number of items in the RDD using this command:

```
readme.count()
```

You should see that this RDD action returned a value of 141.


```

biadmin@cluster-743-1423492039-master:/opt/ibm/biginsights/spark/bin
ith 2 tasks
15/02/24 17:42:19 INFO scheduler.TaskSetManager: Starting task 0.0 in stage 2.0
(TID 4, cluster-743-1423492039-master.imdemocloud.com, NODE_LOCAL, 1236 bytes)
15/02/24 17:42:19 INFO scheduler.TaskSetManager: Starting task 1.0 in stage 2.0
(TID 5, cluster-743-1423492039-master.imdemocloud.com, NODE_LOCAL, 1236 bytes)
15/02/24 17:42:19 INFO storage.BlockManagerInfo: Added broadcast_3_piece0 in mem
ory on cluster-743-1423492039-master.imdemocloud.com:56949 (size: 1520.0 B, free
: 552.9 MB)
15/02/24 17:42:19 INFO storage.BlockManagerInfo: Added broadcast_3_piece0 in mem
ory on cluster-743-1423492039-master.imdemocloud.com:7198 (size: 1520.0 B, free:
552.9 MB)
15/02/24 17:42:19 INFO scheduler.TaskSetManager: Finished task 0.0 in stage 2.0
(TID 4) in 111 ms on cluster-743-1423492039-master.imdemocloud.com (1/2)
15/02/24 17:42:20 INFO scheduler.TaskSetManager: Finished task 1.0 in stage 2.0
(TID 5) in 119 ms on cluster-743-1423492039-master.imdemocloud.com (2/2)
15/02/24 17:42:20 INFO cluster.YarnClientClusterScheduler: Removed TaskSet 2.0,
whose tasks have all completed, from pool
15/02/24 17:42:20 INFO scheduler.DAGScheduler: Stage 2 (count at <console>:15) f
inished in 0.120 s
15/02/24 17:42:20 INFO spark.SparkContext: Job finished: count at <console>:15,
took 0.130460308 s
res2: Long = 141

scala>

```

__6. Let's run another action. Run this command to find the first item in the RDD:

```
readme.first()
```

```

biadmin@cluster-743-1423492039-master:/opt/ibm/biginsights/spark/bin
ory on cluster-743-1423492039-master.imdemocloud.com:61634 (size: 1539.0 B, free
: 276.5 MB)
15/02/24 17:46:00 INFO storage.BlockManagerMaster: Updated info of block broadca
st_4_piece0
15/02/24 17:46:00 INFO scheduler.DAGScheduler: Submitting 1 missing tasks from S
tage 3 (README.md MappedRDD[1] at textFile at <console>:12)
15/02/24 17:46:00 INFO cluster.YarnClientClusterScheduler: Adding task set 3.0 w
ith 1 tasks
15/02/24 17:46:00 INFO scheduler.TaskSetManager: Starting task 0.0 in stage 3.0
(TID 6, cluster-743-1423492039-master.imdemocloud.com, NODE_LOCAL, 1236 bytes)
15/02/24 17:46:00 INFO storage.BlockManagerInfo: Added broadcast_4_piece0 in mem
ory on cluster-743-1423492039-master.imdemocloud.com:56949 (size: 1539.0 B, free
: 552.9 MB)
15/02/24 17:46:00 INFO scheduler.DAGScheduler: Stage 3 (first at <console>:15) f
inished in 0.097 s
15/02/24 17:46:00 INFO scheduler.TaskSetManager: Finished task 0.0 in stage 3.0
(TID 6) in 97 ms on cluster-743-1423492039-master.imdemocloud.com (1/1)
15/02/24 17:46:00 INFO cluster.YarnClientClusterScheduler: Removed TaskSet 3.0,
whose tasks have all completed, from pool
15/02/24 17:46:00 INFO spark.SparkContext: Job finished: first at <console>:15,
took 0.107107176 s
res3: String = # Apache Spark
scala>

```

- ___7. Now let's try a transformation. Use the *filter* transformation to return a new RDD with a subset of the items in the file. Type in this command:

```
val linesWithSpark = readme.filter(line => line.contains("Spark"))
```

Again, this returned a pointer to a RDD with the results of the filter transformation.

- ___8. You can even chain together transformations and actions. To find out how many lines contains the word "Spark", type in:

```
readme.filter(line => line.contains("Spark")).count()
```

```

biadmin@cluster-743-1423492039-master:/opt/ibm/biginsights/spark/bin
ith 2 tasks
15/02/24 17:49:47 INFO scheduler.TaskSetManager: Starting task 0.0 in stage 4.0
(TID 7, cluster-743-1423492039-master.imdemocloud.com, NODE_LOCAL, 1236 bytes)
15/02/24 17:49:47 INFO scheduler.TaskSetManager: Starting task 1.0 in stage 4.0
(TID 8, cluster-743-1423492039-master.imdemocloud.com, NODE_LOCAL, 1236 bytes)
15/02/24 17:49:48 INFO storage.BlockManagerInfo: Added broadcast_5_piece0 in mem
ory on cluster-743-1423492039-master.imdemocloud.com:7198 (size: 1642.0 B, free:
552.9 MB)
15/02/24 17:49:48 INFO storage.BlockManagerInfo: Added broadcast_5_piece0 in mem
ory on cluster-743-1423492039-master.imdemocloud.com:56949 (size: 1642.0 B, free
: 552.9 MB)
15/02/24 17:49:48 INFO scheduler.TaskSetManager: Finished task 1.0 in stage 4.0
(TID 8) in 216 ms on cluster-743-1423492039-master.imdemocloud.com (1/2)
15/02/24 17:49:48 INFO scheduler.TaskSetManager: Finished task 0.0 in stage 4.0
(TID 7) in 229 ms on cluster-743-1423492039-master.imdemocloud.com (2/2)
15/02/24 17:49:48 INFO scheduler.DAGScheduler: Stage 4 (count at <console>:15) f
inished in 0.230 s
15/02/24 17:49:48 INFO cluster.YarnClientClusterScheduler: Removed TaskSet 4.0,
whose tasks have all completed, from pool
15/02/24 17:49:48 INFO spark.SparkContext: Job finished: count at <console>:15,
took 0.245317763 s
res4: Long = 21
scala>

```

- __9. Do not close the shell. The rest of the exercises will build upon this one. Open up a new terminal and do the Python exercises on that shell.

1.2.2 Using Python

This section goes over the same tasks as the Scala section, but using the Python shell.

- __1. Load the README file from the local system onto HDFS. If you had done this already, you can skip this step. Execute this command:

```
hadoop fs -put /home/virtuser/README.md /tmp
```

- __2. Start the Python Spark shell with this command:

```
$SPARK_HOME/bin/pyspark
```

```
biadmin@cluster-743-1423492039-master:~  
Actor[akka.tcp://sparkExecutor@cluster-743-1423492039-master.imdemocloud.com:339  
9/user/Executor#-1699729317] with ID 1  
15/02/25 18:17:08 INFO util.RackResolver: Resolved cluster-743-1423492039-master  
.imdemocloud.com to /default-rack  
15/02/25 18:17:08 INFO storage.BlockManagerMasterActor: Registering block manage  
r cluster-743-1423492039-master.imdemocloud.com:39345 with 553.0 MB RAM  
15/02/25 18:17:09 INFO cluster.YarnClientSchedulerBackend: Registered executor:  
Actor[akka.tcp://sparkExecutor@cluster-743-1423492039-master.imdemocloud.com:145  
85/user/Executor#327265516] with ID 2  
15/02/25 18:17:09 INFO cluster.YarnClientSchedulerBackend: SchedulerBackend is r  
eady for scheduling beginning after reached minRegisteredResourcesRatio: 0.8  
Welcome to  
  
      _   _   _   _   _   _  
     / \ / \ / \ / \ / \ / \  
    /___/___/___/___/___/___\  
   /___/___/___/___/___/___/___\  
  /___/___/___/___/___/___/___/___\  
 /___/___/___/___/___/___/___/___/___\  
/_/___/___/___/___/___/___/___/___/___\  
 version 1.1.0  
  
Using Python version 2.6.6 (r266:84292, Nov 21 2013 10:50:32)  
SparkContext available as sc.  
>>> 15/02/25 18:17:09 INFO storage.BlockManagerMasterActor: Registering block ma  
nager cluster-743-1423492039-master.imdemocloud.com:56281 with 553.0 MB RAM  
  
>>>
```

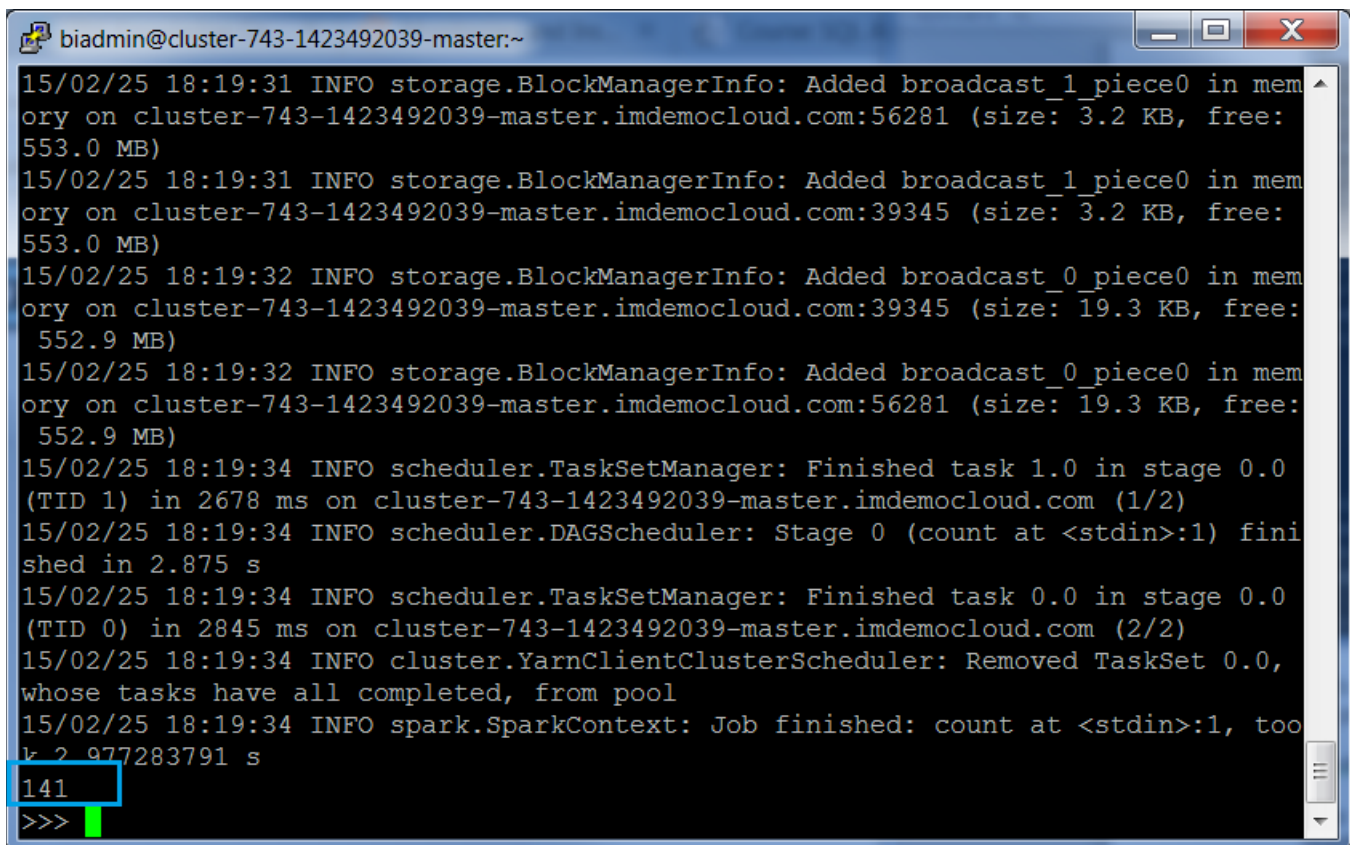
- ___3. Let's start with the example of reading in a text file, and thus converting it to a RDD. Note that the path of the readme text file is on the HDFS. Type in:

```
readme = sc.textFile("/tmp/README.md")
```

This was a RDD transformation, thus it return a pointer to a RDD, which we have named as *readme*.

- ___4. Let's perform some RDD actions on this text file. Count the number of items in the RDD using this command:

```
readme.count()
```

A terminal window titled 'biadmin@cluster-743-1423492039-master:~' displays a series of log messages from the Spark framework. The logs include information about adding broadcast pieces to memory, task completion times, and job status. At the bottom of the terminal, the command '141' is entered and highlighted with a blue box, and the prompt '>>>' is visible with a green cursor.

```
biadmin@cluster-743-1423492039-master:~
15/02/25 18:19:31 INFO storage.BlockManagerInfo: Added broadcast_1_piece0 in mem
ory on cluster-743-1423492039-master.imdemocloud.com:56281 (size: 3.2 KB, free:
553.0 MB)
15/02/25 18:19:31 INFO storage.BlockManagerInfo: Added broadcast_1_piece0 in mem
ory on cluster-743-1423492039-master.imdemocloud.com:39345 (size: 3.2 KB, free:
553.0 MB)
15/02/25 18:19:32 INFO storage.BlockManagerInfo: Added broadcast_0_piece0 in mem
ory on cluster-743-1423492039-master.imdemocloud.com:39345 (size: 19.3 KB, free:
552.9 MB)
15/02/25 18:19:32 INFO storage.BlockManagerInfo: Added broadcast_0_piece0 in mem
ory on cluster-743-1423492039-master.imdemocloud.com:56281 (size: 19.3 KB, free:
552.9 MB)
15/02/25 18:19:34 INFO scheduler.TaskSetManager: Finished task 1.0 in stage 0.0
(TID 1) in 2678 ms on cluster-743-1423492039-master.imdemocloud.com (1/2)
15/02/25 18:19:34 INFO scheduler.DAGScheduler: Stage 0 (count at <stdin>:1) fini
shed in 2.875 s
15/02/25 18:19:34 INFO scheduler.TaskSetManager: Finished task 0.0 in stage 0.0
(TID 0) in 2845 ms on cluster-743-1423492039-master.imdemocloud.com (2/2)
15/02/25 18:19:34 INFO cluster.YarnClientClusterScheduler: Removed TaskSet 0.0,
whose tasks have all completed, from pool
15/02/25 18:19:34 INFO spark.SparkContext: Job finished: count at <stdin>:1, too
k 2 977283791 s
141
>>>
```

You see that this RDD action returned a value of 141.

- __5. Let's run another action. Run this command to find the first item in the RDD:

```
readme.first()
```

```

biadmin@cluster-743-1423492039-master:~
15/02/25 18:20:12 INFO storage.BlockManagerInfo: Added broadcast_2_piece0 in mem
ory on cluster-743-1423492039-master.imdemocloud.com:19868 (size: 2.7 KB, free:
276.5 MB)
15/02/25 18:20:12 INFO storage.BlockManagerMaster: Updated info of block broadca
st_2_piece0
15/02/25 18:20:12 INFO scheduler.DAGScheduler: Submitting 1 missing tasks from S
tage 1 (PythonRDD[3] at RDD at PythonRDD.scala:43)
15/02/25 18:20:12 INFO cluster.YarnClientClusterScheduler: Adding task set 1.0 w
ith 1 tasks
15/02/25 18:20:12 INFO scheduler.TaskSetManager: Starting task 0.0 in stage 1.0
(TID 2, cluster-743-1423492039-master.imdemocloud.com, NODE_LOCAL, 1236 bytes)
15/02/25 18:20:12 INFO storage.BlockManagerInfo: Added broadcast_2_piece0 in mem
ory on cluster-743-1423492039-master.imdemocloud.com:56281 (size: 2.7 KB, free:
552.9 MB)
15/02/25 18:20:13 INFO scheduler.DAGScheduler: Stage 1 (runJob at PythonRDD.sca
la:296) finished in 0.146 s
15/02/25 18:20:13 INFO spark.SparkContext: Job finished: runJob at PythonRDD.sca
la:296, took 0.172304543 s
15/02/25 18:20:13 INFO scheduler.TaskSetManager: Finished task 0.0 in stage 1.0
(TID 2) in 137 ms on cluster-743-1423492039-master.imdemocloud.com (1/1)
15/02/25 18:20:13 INFO cluster.YarnClientClusterScheduler: Removed TaskSet 1.0,
whose tasks have all completed, from pool
u'# Apache Spark'
>>>

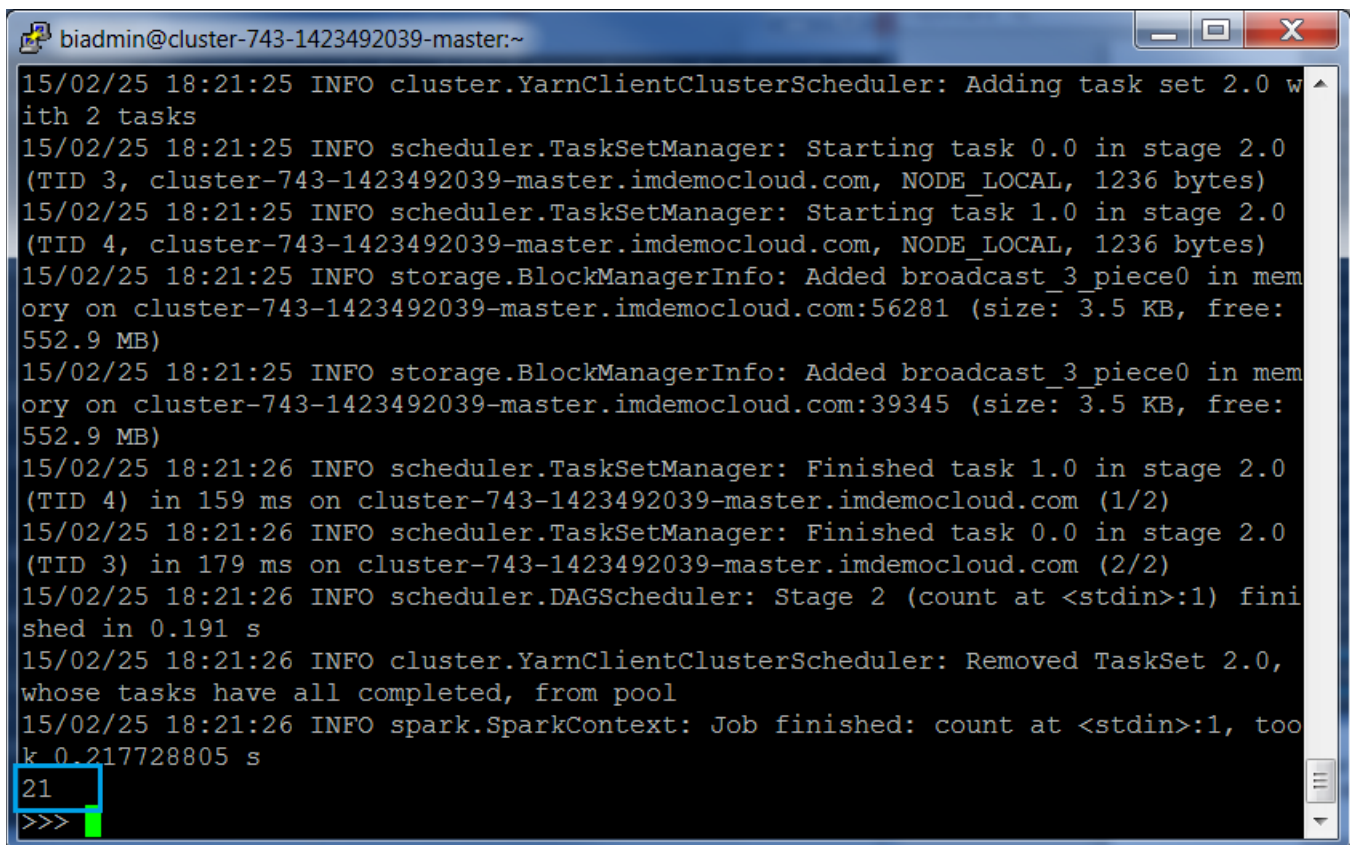
```

- __6. Now let's try a transformation. Use the *filter* transformation to return a new RDD with a subset of the items in the file. Type in this command:

```
linesWithSpark = readme.filter(lambda line: "Spark" in line)
```

- __7. You can even chain together transformations and actions. To find out how many lines contains the word "Spark", type in:

```
readme.filter(lambda line: "Spark" in line).count()
```


A terminal window titled 'biadmin@cluster-743-1423492039-master:~' with standard Windows-style window controls (minimize, maximize, close). The terminal displays a series of log messages from the Spark scheduler and storage components. The messages indicate the addition and execution of tasks in stage 2.0, the completion of tasks 1.0 and 0.0, and the final job completion. The last line shows a prompt '21' followed by '>>>' and a green cursor.

```
15/02/25 18:21:25 INFO cluster.YarnClientClusterScheduler: Adding task set 2.0 with 2 tasks
15/02/25 18:21:25 INFO scheduler.TaskSetManager: Starting task 0.0 in stage 2.0 (TID 3, cluster-743-1423492039-master.imdemocloud.com, NODE_LOCAL, 1236 bytes)
15/02/25 18:21:25 INFO scheduler.TaskSetManager: Starting task 1.0 in stage 2.0 (TID 4, cluster-743-1423492039-master.imdemocloud.com, NODE_LOCAL, 1236 bytes)
15/02/25 18:21:25 INFO storage.BlockManagerInfo: Added broadcast_3_piece0 in memory on cluster-743-1423492039-master.imdemocloud.com:56281 (size: 3.5 KB, free: 552.9 MB)
15/02/25 18:21:25 INFO storage.BlockManagerInfo: Added broadcast_3_piece0 in memory on cluster-743-1423492039-master.imdemocloud.com:39345 (size: 3.5 KB, free: 552.9 MB)
15/02/25 18:21:26 INFO scheduler.TaskSetManager: Finished task 1.0 in stage 2.0 (TID 4) in 159 ms on cluster-743-1423492039-master.imdemocloud.com (1/2)
15/02/25 18:21:26 INFO scheduler.TaskSetManager: Finished task 0.0 in stage 2.0 (TID 3) in 179 ms on cluster-743-1423492039-master.imdemocloud.com (2/2)
15/02/25 18:21:26 INFO scheduler.DAGScheduler: Stage 2 (count at <stdin>:1) finished in 0.191 s
15/02/25 18:21:26 INFO cluster.YarnClientClusterScheduler: Removed TaskSet 2.0, whose tasks have all completed, from pool
15/02/25 18:21:26 INFO spark.SparkContext: Job finished: count at <stdin>:1, took 0.217728805 s
21
>>>
```

1.3 More on RDD Operations

This section builds upon the previous section. You will have to use the same shell from the previous section (either Scala or Python). In this section, you will see that RDD can be used for more complex computations. You will find the line from that readme file with the most words in it.

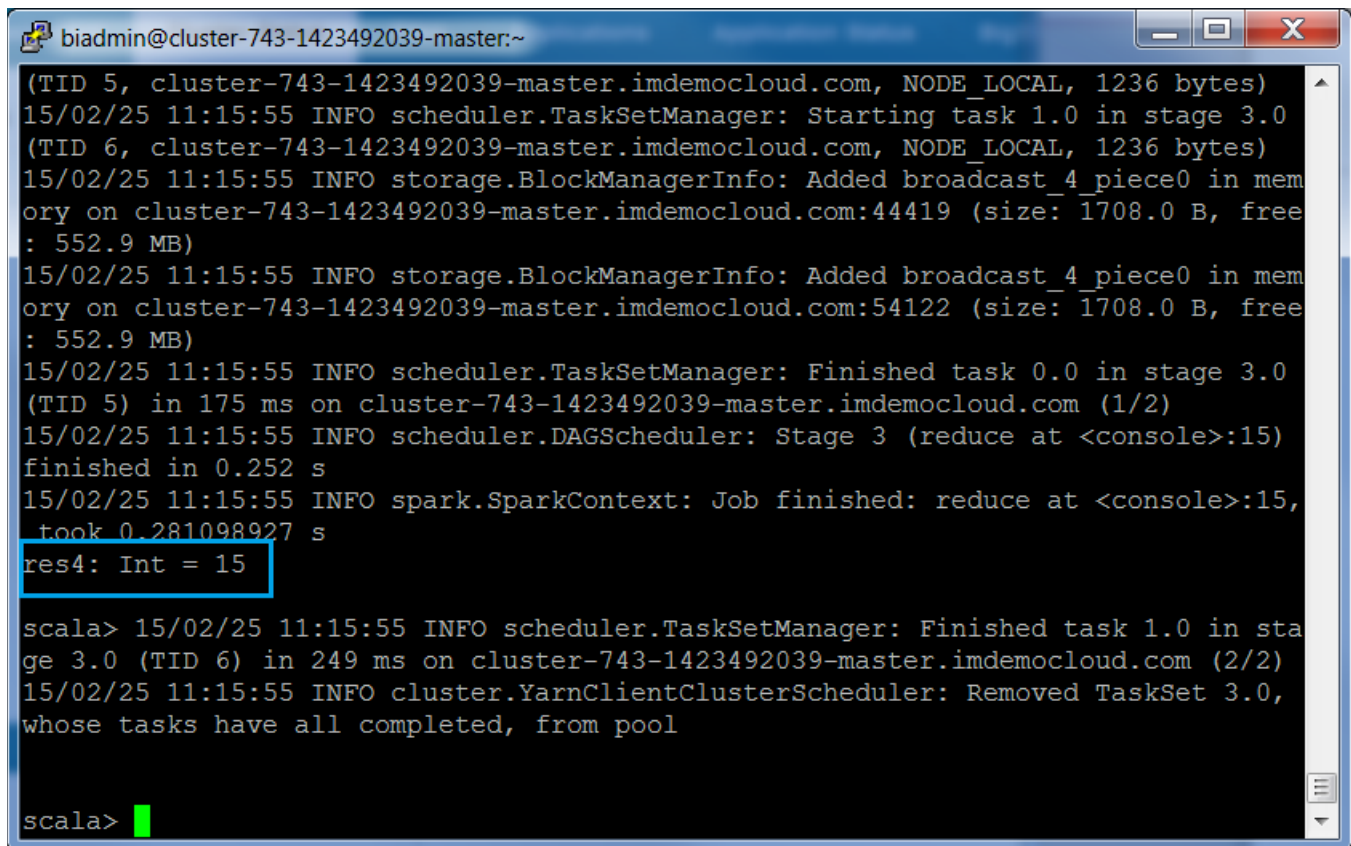
Similar to how the last section was structured, Scala will be shown first, followed by the same steps shown for Python.

1.3.1 With Scala

__1. Using the scala shell, copy and paste:

```
readme.map(line => line.split(" ").size).reduce((a, b) => if (a > b) a else b)
```

There are two parts to this. The first maps a line to an integer value, the number of words in that line. In the second part reduce is called to find the line with the most words in it. The arguments to *map* and *reduce* are Scala function literals (closures), but you can use any language feature or Scala/Java library.



```
biadmin@cluster-743-1423492039-master:~
(TID 5, cluster-743-1423492039-master.imdemocloud.com, NODE_LOCAL, 1236 bytes)
15/02/25 11:15:55 INFO scheduler.TaskSetManager: Starting task 1.0 in stage 3.0
(TID 6, cluster-743-1423492039-master.imdemocloud.com, NODE_LOCAL, 1236 bytes)
15/02/25 11:15:55 INFO storage.BlockManagerInfo: Added broadcast_4_piece0 in mem
ory on cluster-743-1423492039-master.imdemocloud.com:44419 (size: 1708.0 B, free
: 552.9 MB)
15/02/25 11:15:55 INFO storage.BlockManagerInfo: Added broadcast_4_piece0 in mem
ory on cluster-743-1423492039-master.imdemocloud.com:54122 (size: 1708.0 B, free
: 552.9 MB)
15/02/25 11:15:55 INFO scheduler.TaskSetManager: Finished task 0.0 in stage 3.0
(TID 5) in 175 ms on cluster-743-1423492039-master.imdemocloud.com (1/2)
15/02/25 11:15:55 INFO scheduler.DAGScheduler: Stage 3 (reduce at <console>:15)
finished in 0.252 s
15/02/25 11:15:55 INFO spark.SparkContext: Job finished: reduce at <console>:15,
took 0.281098927 s
res4: Int = 15

scala> 15/02/25 11:15:55 INFO scheduler.TaskSetManager: Finished task 1.0 in sta
ge 3.0 (TID 6) in 249 ms on cluster-743-1423492039-master.imdemocloud.com (2/2)
15/02/25 11:15:55 INFO cluster.YarnClientClusterScheduler: Removed TaskSet 3.0,
whose tasks have all completed, from pool

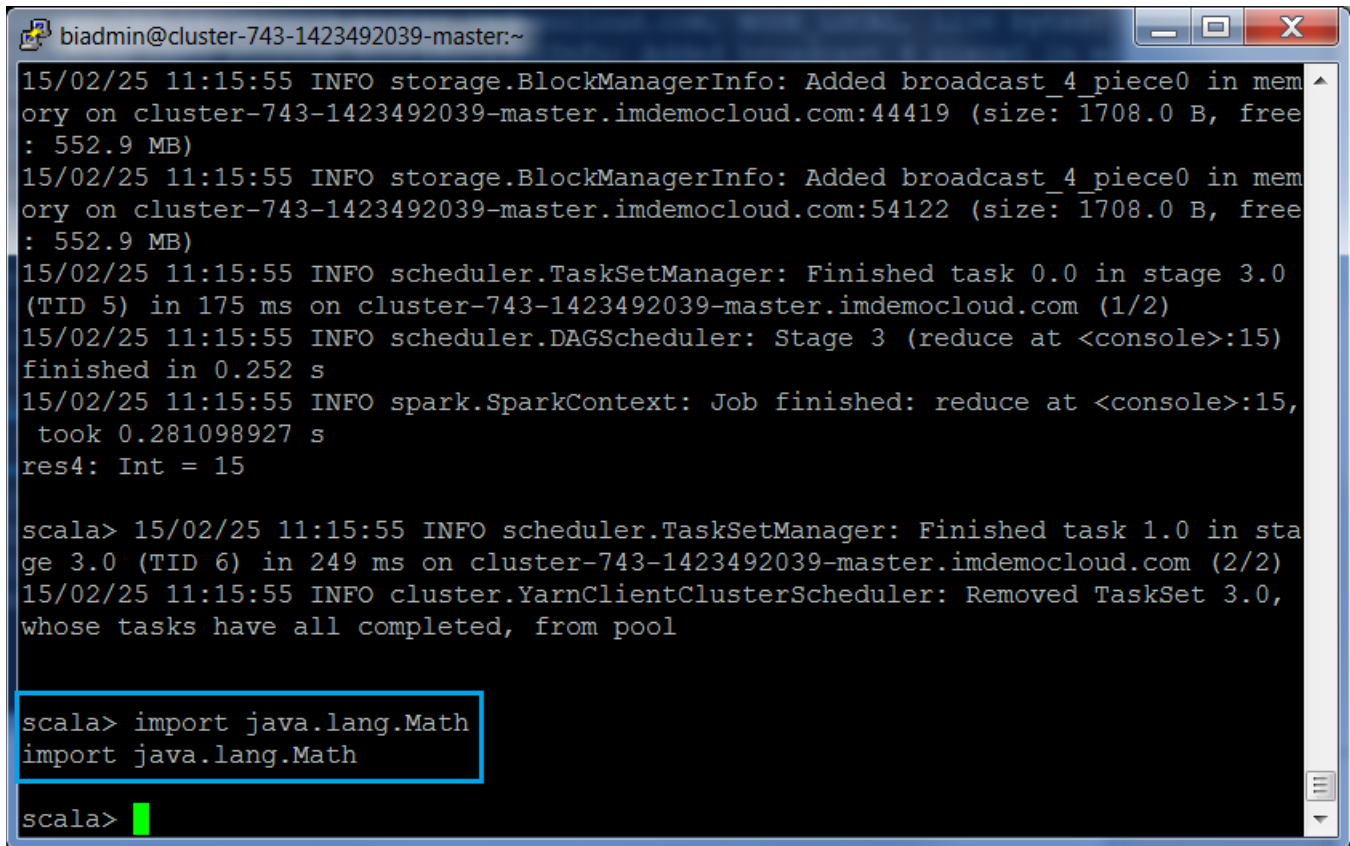
scala>
```

Line 15 contains the most words in it.

In the next step, you use the `Math.max()` function to show that you can indeed use a Java library instead.

- __2. Import in the `java.lang.Math` library. Copy/paste:

```
import java.lang.Math
```

A screenshot of a terminal window with a dark background and light-colored text. The window title is 'biadmin@cluster-743-1423492039-master:~'. The terminal shows several lines of Spark logs, including 'INFO storage.BlockManagerInfo: Added broadcast_4_piece0 in memory on cluster-743-1423492039-master.imdemocloud.com:44419 (size: 1708.0 B, free: 552.9 MB)' and 'INFO scheduler.TaskSetManager: Finished task 0.0 in stage 3.0 (TID 5) in 175 ms on cluster-743-1423492039-master.imdemocloud.com (1/2)'. Below the logs, the output 'res4: Int = 15' is visible. At the bottom, a Scala prompt 'scala>' is followed by the command 'import java.lang.Math', which is highlighted with a red rectangular box. The prompt is followed by a green cursor.

```
biadmin@cluster-743-1423492039-master:~
15/02/25 11:15:55 INFO storage.BlockManagerInfo: Added broadcast_4_piece0 in mem
ory on cluster-743-1423492039-master.imdemocloud.com:44419 (size: 1708.0 B, free
: 552.9 MB)
15/02/25 11:15:55 INFO storage.BlockManagerInfo: Added broadcast_4_piece0 in mem
ory on cluster-743-1423492039-master.imdemocloud.com:54122 (size: 1708.0 B, free
: 552.9 MB)
15/02/25 11:15:55 INFO scheduler.TaskSetManager: Finished task 0.0 in stage 3.0
(TID 5) in 175 ms on cluster-743-1423492039-master.imdemocloud.com (1/2)
15/02/25 11:15:55 INFO scheduler.DAGScheduler: Stage 3 (reduce at <console>:15)
finished in 0.252 s
15/02/25 11:15:55 INFO spark.SparkContext: Job finished: reduce at <console>:15,
took 0.281098927 s
res4: Int = 15

scala> 15/02/25 11:15:55 INFO scheduler.TaskSetManager: Finished task 1.0 in sta
ge 3.0 (TID 6) in 249 ms on cluster-743-1423492039-master.imdemocloud.com (2/2)
15/02/25 11:15:55 INFO cluster.YarnClientClusterScheduler: Removed TaskSet 3.0,
whose tasks have all completed, from pool

scala> import java.lang.Math
import java.lang.Math

scala> █
```

- __3. Now copy and paste in the following to run with the max function:

```
readme.map(line => line.split(" ").size).reduce((a, b) => Math.max(a, b))
```

```

biadmin@cluster-743-1423492039-master:~
(TID 7, cluster-743-1423492039-master.imdemocloud.com, NODE_LOCAL, 1236 bytes)
15/02/25 11:20:11 INFO scheduler.TaskSetManager: Starting task 1.0 in stage 4.0
(TID 8, cluster-743-1423492039-master.imdemocloud.com, NODE_LOCAL, 1236 bytes)
15/02/25 11:20:11 INFO storage.BlockManagerInfo: Added broadcast_5_piece0 in mem
ory on cluster-743-1423492039-master.imdemocloud.com:44419 (size: 1707.0 B, free
: 552.9 MB)
15/02/25 11:20:11 INFO storage.BlockManagerInfo: Added broadcast_5_piece0 in mem
ory on cluster-743-1423492039-master.imdemocloud.com:54122 (size: 1707.0 B, free
: 552.9 MB)
15/02/25 11:20:11 INFO scheduler.TaskSetManager: Finished task 0.0 in stage 4.0
(TID 7) in 126 ms on cluster-743-1423492039-master.imdemocloud.com (1/2)
15/02/25 11:20:11 INFO scheduler.TaskSetManager: Finished task 1.0 in stage 4.0
(TID 8) in 133 ms on cluster-743-1423492039-master.imdemocloud.com (2/2)
15/02/25 11:20:11 INFO scheduler.DAGScheduler: Stage 4 (reduce at <console>:16)
finished in 0.135 s
15/02/25 11:20:11 INFO spark.SparkContext: Job finished: reduce at <console>:16,
took 0.146720021 s
res5: Int = 15

scala> 15/02/25 11:20:11 INFO cluster.YarnClientClusterScheduler: Removed TaskSe
t 4.0, whose tasks have all completed, from pool

scala>

```

- __4. Spark has a MapReduce data flow pattern. We can use this to do a word count on the readme file. Copy and paste in:

```
val wordCounts = readme.flatMap(line => line.split(" ")).map(word => (word, 1)).reduceByKey((a,b) => a + b)
```

Here we combined the *flatMap*, *map*, and the *reduceByKey* functions to do a word count of each word in the readme file.

- __5. To collect the word counts, use the *collect* action.

```
wordCounts.collect()
```

```

biadmin@cluster-743-1423492039-master:~
15/02/25 11:22:31 INFO spark.MapOutputTrackerMasterActor: Asked to send map output locations for shuffle 0 to sparkExecutor@cluster-743-1423492039-master.indemoclou
cloud.com:46920
15/02/25 11:22:32 INFO scheduler.TaskSetManager: Finished task 1.0 in stage 5.0 (TID 12) in 294 ms on cluster-743-1423492039-master.indemoclou
cloud.com (1/2)
15/02/25 11:22:32 INFO scheduler.DAGScheduler: Stage 5 (collect at <console>:18) finished in 0.326 s
15/02/25 11:22:32 INFO spark.SparkContext: Job finished: collect at <console>:18, took 0.866124052 s
15/02/25 11:22:32 INFO scheduler.TaskSetManager: Finished task 0.0 in stage 5.0 (TID 11) in 318 ms on cluster-743-1423492039-master.indemoclou
cloud.com (2/2)
15/02/25 11:22:32 INFO cluster.YarnClientClusterScheduler: Removed TaskSet 5.0, whose tasks have all completed, from pool
res6: Array[(String, Int)] = Array((means,1), (under,2), (this,4), (Because,1), (agree,1), (Python,2), (cluster.,1), (follows.,1), (its,1), (YARN,,3), (have,2), (MRv1,,1), (pre-built,1), (general,2), (locally.,1), (changed,1), (locally,2), (MapReduce,2), (several,1), (only,1), (Configuration,1), (requests,1), (This,2), (sc.parallelize(1,1), (documentation,1), (basic,1), (CLI,1), (learning,,1), (first,1), (graph,1), (without,1), (setting,2), ([params]`,1), (any,2), ("yarn-client",1), (application,1), (prefer,1), (SparkPi,2), (engine,1), (version,3), (file,1), (documentation,,1), (<http://spark.apache.org/>,1), (MASTER,1), (entry,1), (example,3), (are,2), (systems.,1), (<artifactId>hadoop-client</artifactId>,1), (params,1), (scala>,1), (provides,1), (refer,1), (configure,1), (artifact,1)...
scala>

```

You can see the partial results using the collect action.

1.3.2 With Python

__6. Using the python shell, copy in:

```
readme.map(lambda line: len(line.split())).reduce(lambda a, b: a if (a > b) else b)
```

There are two parts to this. The first maps a line to an integer value, the number of words in that line. In the second part reduce is called to find the line with the most words in it. The arguments to *map* and *reduce* are Python anonymous functions (lambdas), but you can use any top level Python functions. In the next step, you'll define a max function to illustrate this feature.

```

biadmin@cluster-743-1423492039-master:~
15/02/25 18:22:19 INFO cluster.YarnClientClusterScheduler: Adding task set 3.0 with 2 tasks
15/02/25 18:22:19 INFO scheduler.TaskSetManager: Starting task 0.0 in stage 3.0 (TID 5, cluster-743-1423492039-master.imdemocloud.com, NODE_LOCAL, 1236 bytes)
15/02/25 18:22:19 INFO scheduler.TaskSetManager: Starting task 1.0 in stage 3.0 (TID 6, cluster-743-1423492039-master.imdemocloud.com, NODE_LOCAL, 1236 bytes)
15/02/25 18:22:19 INFO storage.BlockManagerInfo: Added broadcast_4_piece0 in memory on cluster-743-1423492039-master.imdemocloud.com:56281 (size: 3.1 KB, free: 552.9 MB)
15/02/25 18:22:19 INFO storage.BlockManagerInfo: Added broadcast_4_piece0 in memory on cluster-743-1423492039-master.imdemocloud.com:39345 (size: 3.1 KB, free: 552.9 MB)
15/02/25 18:22:19 INFO scheduler.TaskSetManager: Finished task 0.0 in stage 3.0 (TID 5) in 146 ms on cluster-743-1423492039-master.imdemocloud.com (1/2)
15/02/25 18:22:19 INFO scheduler.TaskSetManager: Finished task 1.0 in stage 3.0 (TID 6) in 159 ms on cluster-743-1423492039-master.imdemocloud.com (2/2)
15/02/25 18:22:19 INFO cluster.YarnClientClusterScheduler: Removed TaskSet 3.0, whose tasks have all completed, from pool
15/02/25 18:22:19 INFO scheduler.DAGScheduler: Stage 3 (reduce at <stdin>:1) finished in 0.171 s
15/02/25 18:22:19 INFO spark.SparkContext: Job finished: reduce at <stdin>:1, took 0.192973301 s
15
>>>

```

__7. Define the max function. You will need to type this in:

```

def max(a, b):
...     if a > b:
...         return a
...     else:
...         return b
...

```

__8. Now copy in the following to run with the max function:

```

readme.map(lambda line: len(line.split())).reduce(max)

```

```

biadmin@cluster-743-1423492039-master:~
15/02/25 18:26:44 INFO cluster.YarnClientClusterScheduler: Adding task set 4.0 with 2 tasks
15/02/25 18:26:44 INFO scheduler.TaskSetManager: Starting task 0.0 in stage 4.0 (TID 7, cluster-743-1423492039-master.imdemocloud.com, NODE_LOCAL, 1236 bytes)
15/02/25 18:26:44 INFO scheduler.TaskSetManager: Starting task 1.0 in stage 4.0 (TID 8, cluster-743-1423492039-master.imdemocloud.com, NODE_LOCAL, 1236 bytes)
15/02/25 18:26:44 INFO storage.BlockManagerInfo: Added broadcast_5_piece0 in memory on cluster-743-1423492039-master.imdemocloud.com:39345 (size: 3.1 KB, free: 552.9 MB)
15/02/25 18:26:44 INFO storage.BlockManagerInfo: Added broadcast_5_piece0 in memory on cluster-743-1423492039-master.imdemocloud.com:56281 (size: 3.1 KB, free: 552.9 MB)
15/02/25 18:26:44 INFO scheduler.TaskSetManager: Finished task 0.0 in stage 4.0 (TID 7) in 139 ms on cluster-743-1423492039-master.imdemocloud.com (1/2)
15/02/25 18:26:44 INFO scheduler.TaskSetManager: Finished task 1.0 in stage 4.0 (TID 8) in 172 ms on cluster-743-1423492039-master.imdemocloud.com (2/2)
15/02/25 18:26:44 INFO scheduler.DAGScheduler: Stage 4 (reduce at <stdin>:1) finished in 0.181 s
15/02/25 18:26:44 INFO cluster.YarnClientClusterScheduler: Removed TaskSet 4.0, whose tasks have all completed, from pool
15/02/25 18:26:44 INFO spark.SparkContext: Job finished: reduce at <stdin>:1, took 0.204416927 s
15
>>>

```

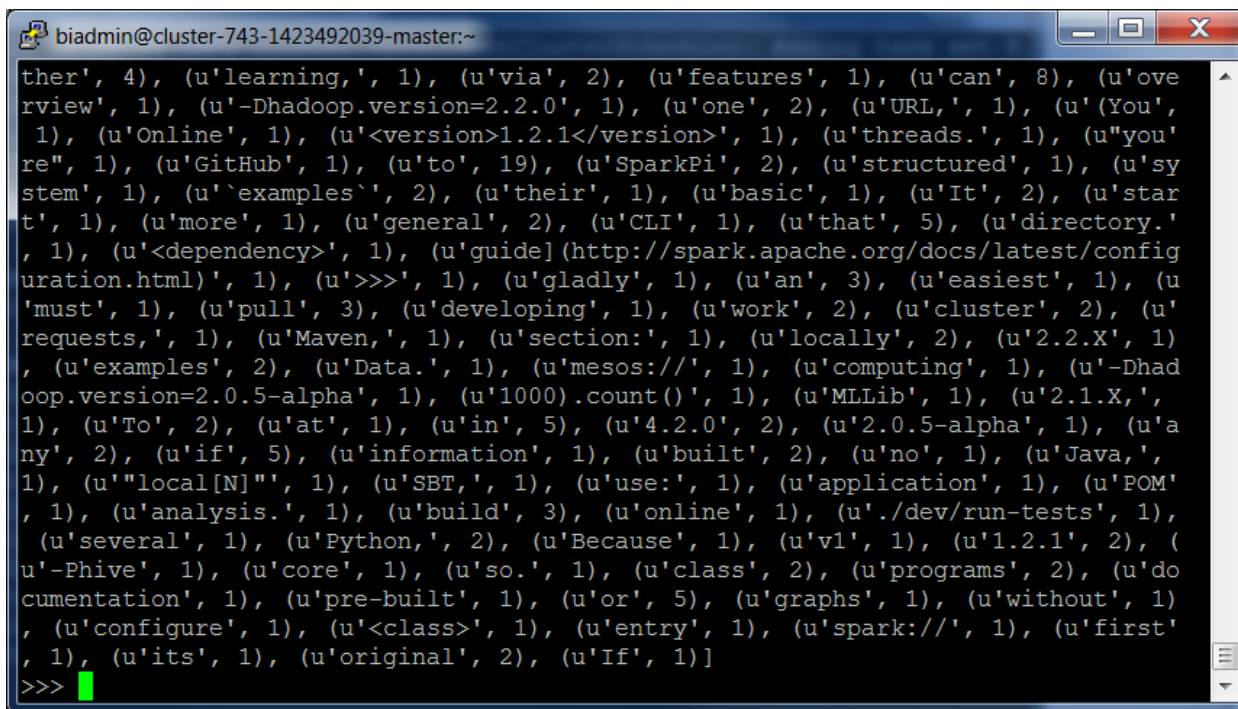
- __9. Spark has a MapReduce data flow pattern. We can use this to do a word count on the readme file. Copy in:

```
wordCounts = readme.flatMap(lambda line: line.split()).map(lambda word: (word, 1)).reduceByKey(lambda a, b: a+b)
```

Here we combined the *flatMap*, *map*, and the *reduceByKey* functions to do a word count of each word in the readme file.

- __10. To collect the word counts, use the *collect* action.

```
wordCounts.collect()
```



```

biadmin@cluster-743-1423492039-master:~
ther', 4), (u'learning', 1), (u'via', 2), (u'features', 1), (u'can', 8), (u'ove
rview', 1), (u'-Dhadoop.version=2.2.0', 1), (u'one', 2), (u'URL', 1), (u'(You',
1), (u'Online', 1), (u'<version>1.2.1</version>', 1), (u'threads.', 1), (u"you'
re", 1), (u'GitHub', 1), (u'to', 19), (u'SparkPi', 2), (u'structured', 1), (u'sy
stem', 1), (u'`examples`', 2), (u'their', 1), (u'basic', 1), (u'It', 2), (u'star
t', 1), (u'more', 1), (u'general', 2), (u'CLI', 1), (u'that', 5), (u'directory.'
, 1), (u'<dependency>', 1), (u'guide[http://spark.apache.org/docs/latest/config
uration.html]', 1), (u'>>>', 1), (u'gladly', 1), (u'an', 3), (u'easiest', 1), (u
'must', 1), (u'pull', 3), (u'developing', 1), (u'work', 2), (u'cluster', 2), (u'
requests', 1), (u'Maven', 1), (u'section:', 1), (u'locally', 2), (u'2.2.X', 1)
, (u'examples', 2), (u'Data.', 1), (u'mesos://', 1), (u'computing', 1), (u'-Dhad
oop.version=2.0.5-alpha', 1), (u'1000).count()', 1), (u'MLLib', 1), (u'2.1.X',
1), (u'To', 2), (u'at', 1), (u'in', 5), (u'4.2.0', 2), (u'2.0.5-alpha', 1), (u'a
ny', 2), (u'if', 5), (u'information', 1), (u'built', 2), (u'no', 1), (u'Java',
1), (u'"local[N]"', 1), (u'SBT', 1), (u'use:', 1), (u'application', 1), (u'POM'
, 1), (u'analysis.', 1), (u'build', 3), (u'online', 1), (u'./dev/run-tests', 1),
(u'several', 1), (u'Python', 2), (u'Because', 1), (u'v1', 1), (u'1.2.1', 2), (
u'-Phive', 1), (u'core', 1), (u'so.', 1), (u'class', 2), (u'programs', 2), (u'do
cumentation', 1), (u'pre-built', 1), (u'or', 5), (u'graphs', 1), (u'without', 1)
, (u'configure', 1), (u'<class>', 1), (u'entry', 1), (u'spark://', 1), (u'first'
, 1), (u'its', 1), (u'original', 2), (u'If', 1)]
>>>

```

1.4 Using Spark caching

In this short section, you'll see how Spark caching can be used to pull data sets into a cluster-wide in-memory cache. This is very useful for accessing repeated data, such as querying a small "hot" dataset or when running an iterative algorithm. Both Python and Scala use the same commands, so you can input these into any of the two shells.

- __1. As a simple example, let's mark our *linesWithSpark* dataset to be cached and then invoke the first count operation to tell Spark to cache it. Remember that transformation operations such as `cache` does not get processed until some action like `count()` is called. Once you run the second `count()` operation, you should notice a small increase in speed.

```
linesWithSpark.cache()  
  
linesWithSpark.count()  
  
linesWithSpark.count()
```

It may seem silly to cache such a small file, but for larger data sets across tens or hundreds of nodes, this would still work. The second *linesWithSpark.count()* action runs against the cache and would perform significantly better for large datasets.

- __2. At this time, you can terminate the shells. Use the key combination: CTRL + D

Summary

Having completed this exercise, you should now be able to log in to your environment and use the Spark shell to run simple actions and transformations for Scala and/or Python. You understand that Spark caching can be used to cache large datasets and subsequent operations on it will utilize the data in the cache rather than re-fetching it from HDFS.

NOTES

[illegible]

NOTES

[illegible]



© Copyright IBM Corporation 2015.

The information contained in these materials is provided for informational purposes only, and is provided AS IS without warranty of any kind, express or implied. IBM shall not be responsible for any damages arising out of the use of, or otherwise related to, these materials. Nothing contained in these materials is intended to, nor shall have the effect of, creating any warranties or representations from IBM or its suppliers or licensors, or altering the terms and conditions of the applicable license agreement governing the use of IBM software. References in these materials to IBM products, programs, or services do not imply that they will be available in all countries in which IBM operates. This information is based on current IBM product plans and strategy, which are subject to change by IBM without notice. Product release dates and/or capabilities referenced in these materials may change at any time at IBM's sole discretion based on market opportunities or other factors, and are not intended to be a commitment to future product or feature availability in any way.

IBM, the IBM logo and ibm.com are trademarks of International Business Machines Corp., registered in many jurisdictions worldwide. Other product and service names might be trademarks of IBM or other companies. A current list of IBM trademarks is available on the Web at "Copyright and trademark information" at www.ibm.com/legal/copytrade.shtml.



Please Recycle
