

Big Data Application Development - Summer 2017

Homework 3, Part 2 Answer Sheet

4. Use the REPL to explore Spark RDDs.

1) Provide the command you used to create your RDD.	<pre>val mydata = sc.textFile("/user/lc3397/class3_hw/frostroad.txt")</pre>
2) Provide the command you used to count the elements (lines) in your RDD.	<pre>mydata.count()</pre>
3) Provide the number of elements.	23
4) Provide the collect command you used.	<pre>mydata.collect()</pre>
5) Provide the command you used to create the HDFS directory.	<pre>[lc3397@login-1-1 ~]\$ hdfs dfs -mkdir loudacre [lc3397@login-1-1 ~]\$ hdfs dfs -mkdir loudacre/weblog</pre>
6) Provide the command you used to put the file into HDFS.	<pre>[lc3397@login-1-1 ~]\$ hdfs dfs -put 2014-03-15.log loudacre/weblog</pre>
7) Provide the command you used to view the file.	<pre>hdfs dfs -cat loudacre/weblog/2014-03-15.log</pre>

5. Transform a small dataset using RDDs.

8) Initialize <code>logfile</code> .	<pre>val logfile_name = "/user/lc3397/loudacre/weblog/2014-03-15.log"</pre>
9) Create an RDD from the file.	<pre>val logfile = sc.textFile(logfile_name)</pre>
10) View the first 10 lines of the data.	<pre>logfile.take(10)</pre>
11) Create an RDD containing only lines that are requests for <code>jpg</code> files.	<pre>val pattern = ".jpg .jpeg .JPG .JPEG".r val jpg_data_rdd = logfile.filter(line => pattern.findFirstIn(line).isDefined)</pre>
12) View the first 10 lines of the data.	<pre>jpg_data_rdd.take(10).foreach(println)</pre>
13) Chain the previous commands into a single command that counts the number of JPG requests.	<pre>logfile.map(line=>line.toLowerCase()).filter(line=>line.contains("jpg")).count()</pre>
14) Create an RDD using the <code>map</code> function to return the length of each line of the log file.	<pre>val line_length = logfile.map(s => s.length)</pre>
15) Create an RDD using the <code>map</code> and <code>split</code> functions to map an array of words for each line.	<pre>val splitarray = logfile.map(s => s.split(" "))</pre>
16) Create an RDD containing only the IP addresses from each line.	<pre>val ip_addrs = splitarray.map(s => s(0))</pre>
17) Use <code>foreach(println)</code> to output IP addresses.	<pre>ip_addrs.collect().foreach(println)</pre>
18) Save the list of IP addresses to an HDFS directory named <code>loudacre/iplist</code> using <code>saveAsTextFile</code> .	<pre>ip_addrs.saveAsTextFile("loudacre/iplist_2")</pre>

5. Transform a small dataset using RDDs. (continued)

19) Provide a screenshot of the contents of the `loudacre/iplist` folder. (Paste it below.)

```
[lc3397@login-2-1 ~]$ hdfs dfs -ls loudacre/iplist_2
Found 3 items
-rw-----  3 lc3397 users          0 2017-06-15 20:10 loudacre/iplist_2/_SUCCESS
-rw-----  3 lc3397 users    50653 2017-06-15 20:10 loudacre/iplist_2/part-00000
-rw-----  3 lc3397 users    50638 2017-06-15 20:10 loudacre/iplist_2/part-00001
```

6. Transform a large dataset using RDDs.

20) Initialize <code>logfile</code> .	<code>val logfiles_name = "/user/lc3397/loudacre/weblogs/*"</code>
21) Create an RDD from the file.	<code>val logfiles = sc.textFile(logfiles_name)</code>
22) View the first 10 lines of the data.	<code>logfiles.take(10)</code>
23) Create an RDD containing only lines that are requests for <code>jpg</code> files.	<code>val big_jpg_data = logfiles.map(line=>line.toLowerCase()).filter(line=>line.contains(".jpg"))</code>
24) View the first 10 lines of the data.	<code>big_jpg_data.take(10).foreach(println)</code>
25) Chain the previous commands into a single command that counts the number of JPG requests.	<code>logfiles.map(line=>line.toLowerCase()).filter(line=>line.contains(".jpg")).count()</code>
26) Create an RDD using the <code>map</code> function to return the length of each line of the log file	<code>val line_length = logfiles.map(s => s.length())</code>
27) Create an RDD using the <code>map</code> and <code>split</code> functions to map an array of words for each line.	<code>val splitarray_large = logfiles.map(s => s.split(" "))</code>
28) Create an RDD containing only the IP addresses from each line.	<code>val ip_addrs_large = splitarray_large.map(s => s(0))</code>
29) Use <code>foreach(println)</code> to output IP addresses.	<code>ip_addrs_large.collect().foreach(println)</code>
30) Save the list of IP addresses to a file in an HDFS directory named <code>loudacre/bigiplist</code> - use <code>saveAsTextFile</code> .	<code>ip_addrs_large.saveAsTextFile("loudacre/iplist_large")</code>

6. Transform a large dataset using RDDs. (continued)

31) Provide a screenshot of the contents of the `loudacre/bigiplist` folder. (Paste it below.)

```
[lc3397@login-2-1 ~]$ hdfs dfs -ls loudacre/iplist_large
Found 494 items
-rw----- 3 lc3397 users      0 2017-06-15 22:12 loudacre/iplist_large/_SUCCESS
-rw----- 3 lc3397 users 49265 2017-06-15 22:11 loudacre/iplist_large/part-00000
-rw----- 3 lc3397 users 45854 2017-06-15 22:11 loudacre/iplist_large/part-00001
-rw----- 3 lc3397 users 50031 2017-06-15 22:11 loudacre/iplist_large/part-00002
-rw----- 3 lc3397 users 45898 2017-06-15 22:11 loudacre/iplist_large/part-00003
-rw----- 3 lc3397 users 48070 2017-06-15 22:11 loudacre/iplist_large/part-00004
-rw----- 3 lc3397 users 46430 2017-06-15 22:11 loudacre/iplist_large/part-00005
-rw----- 3 lc3397 users 46177 2017-06-15 22:11 loudacre/iplist_large/part-00006
-rw----- 3 lc3397 users 50720 2017-06-15 22:11 loudacre/iplist_large/part-00007
-rw----- 3 lc3397 users 47314 2017-06-15 22:11 loudacre/iplist_large/part-00008
-rw----- 3 lc3397 users 46282 2017-06-15 22:11 loudacre/iplist_large/part-00009
-rw----- 3 lc3397 users 45998 2017-06-15 22:11 loudacre/iplist_large/part-00010
-rw----- 3 lc3397 users 49261 2017-06-15 22:11 loudacre/iplist_large/part-00011
```

```
[lc3397@login-2-1 ~]$ hdfs dfs -ls loudacre/
Found 5 items
drwx----- - lc3397 users      0 2017-06-15 20:07 loudacre/iplist
drwx----- - lc3397 users      0 2017-06-15 20:10 loudacre/iplist_2
drwx----- - lc3397 users      0 2017-06-15 22:12 loudacre/iplist_large
drwx----- - lc3397 users      0 2017-06-14 23:36 loudacre/weblog
drwx----- - lc3397 users      0 2017-06-15 20:15 loudacre/weblogs
```