

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("File:///home/cyy292/BDAD/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. | <pre>res1: Long = 23</pre> |
| 4) Provide the collect command you used. | <pre>mydata.collect</pre> |
| 5) Provide the command you used to create the HDFS directory. | <pre>hadoop fs -mkdir loudacre hadoop fs -mkdir loudacre/weblog</pre> |
| 6) Provide the command you used to put the file into HDFS. | <pre>hadoop fs -put 2014-03-15.log loudacre/weblog</pre> |
| 7) Provide the command you used to view the file. | <pre>hadoop fs -cat loudacre/weblog/2014-03-15.log</pre> |

5. Transform a small dataset using RDDs.

| | |
|---|--|
| 8) Initialize <code>logfile</code> . | <pre>val logfile = "/user/cyy292/loudacre/weblog/2014-03-15.log"</pre> |
| 9) Create an RDD from the file. | <pre>val logRDD = sc.textFile(logfile)</pre> |
| 10) View the first 10 lines of the data. | <pre>logRDD.take(10).foreach(println)</pre> |
| 11) Create an RDD containing only lines that are requests for <code>jpg</code> files. | <pre>val logRDD_jpg = logRDD.filter(_ .contains("jpg"))</pre> |
| 12) View the first 10 lines of the data. | <pre>logRDD_jpg.take(10).foreach(println)</pre> |
| 13) Chain the previous commands into a single command that counts the number of JPG requests. | <pre>scala> logRDD.filter(_ .contains("jpg")).count res24: Long = 423</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_len = logRDD.map(_ .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 line_split = logRDD.map(_ .split(' '))</pre> |
| 16) Create an RDD containing only the IP addresses from each line. | <pre>val line_IP = line_split.map(_ (0))</pre> |
| 17) Use <code>foreach(println)</code> to output IP addresses. | <pre>line_IP.foreach(println)</pre> |
| 18) Save the list of IP addresses to an HDFS directory named <code>loudacre/iplist</code> using <code>saveAsTextFile</code> . | <pre>line_IP.saveAsTextFile("loudacre/iplist")</pre> |

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

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

```
[cyy292@login-1-1 BDAD]$ hadoop fs -ls loudacre/iplist
Found 3 items
-rw----- 3 cyy292 users      0 2017-06-12 01:42 loudacre/iplist/_SUCCESS
-rw----- 3 cyy292 users 50653 2017-06-12 01:42 loudacre/iplist/part-00000
-rw----- 3 cyy292 users 50638 2017-06-12 01:42 loudacre/iplist/part-00001
```

6. Transform a large dataset using RDDs.

| | |
|---|--|
| 20) Initialize <code>logfile</code> . | <pre>val logfile = "loudacre/weblogs/FlumeData.1424275921"</pre> |
| 21) Create an RDD from the file. | <pre>var logsRDD: org.apache.spark.rdd.RDD[String] = sc.emptyRDD scala> for (a <- 226 to 536) logsRDD = logsRDD.union(sc.textFile(logfile + a.toString))</pre> |
| 22) View the first 10 lines of the data. | <pre>logsRDD.take(10).foreach(println)</pre> |
| 23) Create an RDD containing only lines that are requests for <code>jpg</code> files. | <pre>val logsRDD_jpg = logsRDD.filter(_.contains("jpg"))</pre> |
| 24) View the first 10 lines of the data. | <pre>logsRDD_jpg.take(10).foreach(println)</pre> |
| 25) Chain the previous commands into a single command that counts the number of JPG requests. | <pre>scala> logsRDD.filter(_ .contains("jpg")).count res53: Long = 64978</pre> |
| 26) Create an RDD using the <code>map</code> function to return the length of each line of the log file | <pre>val logsRDD_length = logsRDD.map(_.length)</pre> |
| 27) Create an RDD using the <code>map</code> and <code>split</code> functions to map an array of words for each line. | <pre>val logsRDD_split = logsRDD.map(_.split(' '))</pre> |
| 28) Create an RDD containing only the IP addresses from each line. | <pre>val logsRDD_IP = logsRDD_split.map(_(0))</pre> |
| 29) Use <code>foreach(println)</code> to output IP addresses. | <pre>logsRDD_IP.foreach(println)</pre> |

30) Save the list of IP addresses to a file in an HDFS directory named `loudacre/bigiplist` - use `saveAsTextFile`.

```
logsRDD_IP.saveAsTextFile("loudacre/bigiplist")
```

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

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

```
[cyy292@login-1-1 BDAD]$ hadoop fs -ls loudacre/bigiplist
Found 623 items
-rw----- 3 cyy292 users          0 2017-06-12 02:42 loudacre/bigiplist/_SUCCESS
-rw----- 3 cyy292 users    24900 2017-06-12 02:41 loudacre/bigiplist/part-00000
-rw----- 3 cyy292 users    25004 2017-06-12 02:41 loudacre/bigiplist/part-00001
-rw----- 3 cyy292 users    24931 2017-06-12 02:41 loudacre/bigiplist/part-00002
-rw----- 3 cyy292 users    24973 2017-06-12 02:41 loudacre/bigiplist/part-00003
-rw----- 3 cyy292 users    25046 2017-06-12 02:40 loudacre/bigiplist/part-00004
-rw----- 3 cyy292 users    24765 2017-06-12 02:40 loudacre/bigiplist/part-00005
-rw----- 3 cyy292 users    24907 2017-06-12 02:41 loudacre/bigiplist/part-00006
-rw----- 3 cyy292 users    25103 2017-06-12 02:41 loudacre/bigiplist/part-00007
-rw----- 3 cyy292 users    24954 2017-06-12 02:41 loudacre/bigiplist/part-00008
-rw----- 3 cyy292 users    24886 2017-06-12 02:41 loudacre/bigiplist/part-00009
-rw----- 3 cyy292 users    24832 2017-06-12 02:41 loudacre/bigiplist/part-00010
-rw----- 3 cyy292 users    24831 2017-06-12 02:41 loudacre/bigiplist/part-00011
-rw----- 3 cyy292 users    25049 2017-06-12 02:41 loudacre/bigiplist/part-00012
-rw----- 3 cyy292 users    24986 2017-06-12 02:41 loudacre/bigiplist/part-00013
-rw----- 3 cyy292 users    24917 2017-06-12 02:41 loudacre/bigiplist/part-00014
-rw----- 3 cyy292 users    24950 2017-06-12 02:41 loudacre/bigiplist/part-00015
-rw----- 3 cyy292 users    24965 2017-06-12 02:40 loudacre/bigiplist/part-00016
-rw----- 3 cyy292 users    24950 2017-06-12 02:40 loudacre/bigiplist/part-00017
-rw----- 3 cyy292 users    24911 2017-06-12 02:40 loudacre/bigiplist/part-00018
-rw----- 3 cyy292 users    25053 2017-06-12 02:40 loudacre/bigiplist/part-00019
-rw----- 3 cyy292 users    24948 2017-06-12 02:40 loudacre/bigiplist/part-00020
-rw----- 3 cyy292 users    24914 2017-06-12 02:40 loudacre/bigiplist/part-00021
-rw----- 3 cyy292 users    24971 2017-06-12 02:41 loudacre/bigiplist/part-00022
-rw----- 3 cyy292 users    25045 2017-06-12 02:41 loudacre/bigiplist/part-00023
-rw----- 3 cyy292 users    24944 2017-06-12 02:41 loudacre/bigiplist/part-00024
-rw----- 3 cyy292 users    24956 2017-06-12 02:41 loudacre/bigiplist/part-00025
-rw----- 3 cyy292 users    24911 2017-06-12 02:41 loudacre/bigiplist/part-00026
-rw----- 3 cyy292 users    24929 2017-06-12 02:41 loudacre/bigiplist/part-00027
-rw----- 3 cyy292 users    24830 2017-06-12 02:40 loudacre/bigiplist/part-00028
-rw----- 3 cyy292 users    25024 2017-06-12 02:40 loudacre/bigiplist/part-00029
-rw----- 3 cyy292 users    24917 2017-06-12 02:41 loudacre/bigiplist/part-00030
-rw----- 3 cyy292 users    24929 2017-06-12 02:41 loudacre/bigiplist/part-00031
-rw----- 3 cyy292 users    24971 2017-06-12 02:41 loudacre/bigiplist/part-00032
-rw----- 3 cyy292 users    25070 2017-06-12 02:41 loudacre/bigiplist/part-00033
-rw----- 3 cyy292 users    24823 2017-06-12 02:41 loudacre/bigiplist/part-00034
-rw----- 3 cyy292 users    24788 2017-06-12 02:41 loudacre/bigiplist/part-00035
-rw----- 3 cyy292 users    24958 2017-06-12 02:41 loudacre/bigiplist/part-00036
-rw----- 3 cyy292 users    24870 2017-06-12 02:41 loudacre/bigiplist/part-00037
-rw----- 3 cyy292 users    25002 2017-06-12 02:41 loudacre/bigiplist/part-00038
```