

Assignment #4 - Page Rank in Hadoop

Due on August 9th, 2020
Distributed Systems and Cloud Computing
CS549WS—Summer 2020
Professor Dominic Duggan

Daniel Kadyrov
10455680

1 Introduction

This assignment focuses on implementing a Page Ranking algorithm to sort the most popular Wikipedia pages from a provided previously downloaded set. Changes were made to the files shown in the following section. A test using a graph was performed using Hadoop on a local machine was performed. The Page Ranking script was then used on AWS EMR on 5, 10, and 20 clusters.

2 Changes within Provided Files

2.1 DiffMap1

Listing 1: DiffMap1.java

```
1  package edu.stevens.cs549.hadoop.pagerank;
2
3  import java.io.IOException;
4
5  import org.apache.hadoop.mapreduce.*;
6  import org.apache.hadoop.io.*;
7
8  public class DiffMap1 extends Mapper<LongWritable, Text, Text, Text> {
9
10     public void map(LongWritable key, Text value, Context context) throws
        IOException, InterruptedException,
11         IllegalArgumentException {
12         String line = value.toString(); // Converts Line to a String
13         String[] sections = line.split("\t"); // Splits each line
14         if (sections.length > 2) // checks for incorrect data format
15         {
16             throw new IOException("Incorrect data format");
17         }
18         /**
19         * TODO: read node-rank pair and emit: key:node, value:rank
20         */
21
22         // Split into Node and Rank
23         String[] node_rank = sections[0].split("\\+");
24
25         // Emit Node and Rank
26         context.write(new Text(node_rank[0]), new Text(node_rank[1]))
            ;
27     }
28
29 }
```

2.2 DiffMap2

Listing 2: DiffMap2.java

```
1  package edu.stevens.cs549.hadoop.pagerank;
2
3  import java.io.IOException;
4
5  import org.apache.hadoop.mapreduce.*;
6  import org.apache.hadoop.io.*;
7
8  public class DiffMap2 extends Mapper<LongWritable, Text, Text, Text> {
9
10     public void map(LongWritable key, Text value, Context context) throws
        IOException, InterruptedException,
11         IllegalArgumentException {
12         String s = value.toString(); // Converts Line to a String
13
14         /*
15          * TODO: emit: key:"Difference" value: difference calculated in
16          * DiffRed1
17          */
18
19         String[] node_rank = s.split("\\t+");
20
21         // Emit Difference
22         context.write(new Text("Difference"), new Text(node_rank[1]));
23     }
24 }
```

2.3 DiffRed1

Listing 3: DiffRed1.java

```
1  package edu.stevens.cs549.hadoop.pagerank;
2
3  import java.io.*;
4
5  import org.apache.hadoop.mapreduce.*;
6  import org.apache.hadoop.io.*;
7  import java.util.Iterator;
8
9  public class DiffRed1 extends Reducer<Text, Text, Text, Text> {
10
11      public void reduce(Text key, Iterable<Text> values, Context context)
12          throws IOException, InterruptedException {
13          double[] ranks = new double[2];
14          /*
15           * TODO: The list of values should contain two ranks. Compute
16           *       and output their difference.
17           */
18
19          Iterator<Text> iterator = values.iterator();
20
21          // Default Diff set at 0
22          double diff = 0;
23
24          // Rank 1 Calculated
25          if(iterator.hasNext()) {
26              ranks[0] = Double.valueOf(iterator.next().toString());
27          }
28
29          // Rank 2 Calculated
30          if(iterator.hasNext()) {
31              ranks[1] = Double.valueOf(iterator.next().toString());
32          }
33
34          // Calculate Difference Diff
35          diff = Math.abs(ranks[0] - ranks[1]);
36          System.out.println(key.toString() + " " + diff);
37          context.write(key, new Text(String.valueOf(diff)));
38      }
39  }
```

2.4 DiffRed2

Listing 4: DiffRed2.java

```
1  package edu.stevens.cs549.hadoop.pagerank;
2
3  import java.io.*;
4  import java.util.Iterator;
5  import org.apache.hadoop.mapreduce.*;
6  import org.apache.hadoop.io.*;
7
8  public class DiffRed2 extends Reducer<Text, Text, Text, Text> {
9
10     public void reduce(Text key, Iterable<Text> values, Context context)
11         throws IOException, InterruptedException {
12         double diff_max = 0.0; // sets diff_max to a default value
13         /*
14          * TODO: Compute and emit the maximum of the differences
15          */
16
17         Iterator<Text> iterator = values.iterator();
18
19         // Calculate Maximum Difference and Emit
20         while(iterator.hasNext()) {
21             diff_max = Math.max(diff_max, Double.valueOf(iterator.
22                 next().toString()));
23         }
24         context.write(new Text(""), new Text(String.valueOf(diff_max))
25             );
26     }
27 }
```

2.5 FinMapper

Listing 5: FinMapper.java

```
1  package edu.stevens.cs549.hadoop.pagerank;
2
3  import java.io.IOException;
4  import org.apache.hadoop.mapreduce.*;
5  import org.apache.hadoop.io.*;
6
7  public class FinMapper extends Mapper < LongWritable, Text, DoubleWritable,
8      Text > {
9      public void map(LongWritable key, Text value, Context context) throws
10         IOException, InterruptedException, IllegalArgumentException {
11         String line = value.toString(); // Converts Line to a String
12         /*
13          * TODO output key:-rank, value: node
14          * See IterMapper for hints on parsing the output of
15            IterReducer.
16         */
17         String[] sections = line.split("\t"); // nodeId+nodeName |
18         rank
19
20         // Check Format of Data
21         if (sections.length > 2) {
22             throw new IOException("Incorrect data format");
23         }
24         if (sections.length != 2) {
25             return;
26         }
27
28         // 0-Rank will Reverse Shuffle the Reducer
29         context.write(new DoubleWritable(0 - Double.valueOf(sections
30             [1])), new Text(sections[0]));
31     }
32 }
```

2.6 FinReducer

Listing 6: FinReducer.java

```
1  package edu.stevens.cs549.hadoop.pagerank;
2
3  import java.io.IOException;
4  import java.util.Iterator;
5  import org.apache.hadoop.io.DoubleWritable;
6  import org.apache.hadoop.io.Text;
7  import org.apache.hadoop.mapreduce.Reducer;
8
9  public class FinReducer extends Reducer < DoubleWritable, Text, Text, Text > {
10
11      public void reduce(DoubleWritable key, Iterable < Text > values,
12          Context context) throws IOException,
13          InterruptedException {
14          /*
15           * TODO: For each value, emit: key:value, value:-rank
16           */
17
18          Iterator <Text> iterator = values.iterator();
19          String node;
20
21          while (iterator.hasNext()) {
22              node = iterator.next().toString();
23
24              // Convert Negative Rank to Rank
25              context.write(new Text(node), new Text(String.valueOf
26                  (0 - key.get())));
27          }
28      }
29  }
```

2.7 InitMapper

Listing 7: InitMapper.java

```
1  package edu.stevens.cs549.hadoop.pagerank;
2
3  import java.io.IOException;
4
5  import org.apache.hadoop.mapreduce.*;
6  import org.apache.hadoop.io.*;
7
8  public class InitMapper extends Mapper < LongWritable, Text, Text, Text > {
9
10     public void map(LongWritable key, Text value, Context context) throws
11         IOException,
12         InterruptedException,
13         IllegalArgumentException {
14         String line = value.toString(); // Converts Line to a String
15         /*
16          * TODO: Just echo the input, since it is already in adjacency
17            list format.
18         */
19
20         // Split the Line using ":" and Emit the Key and Adjacent List
21         String[] pair = line.split(":");
22
23         if (pair != null && pair.length == 2) {
24             context.write(new Text(pair[0].trim()), new Text(pair
25                 [1]));
26         }
27     }
28 }
```


2.8 InitReducer

Listing 8: InitReducer.java

```
1  package edu.stevens.cs549.hadoop.pagerank;
2
3  import java.io.*;
4  import java.util.Iterator;
5  import org.apache.hadoop.mapreduce.*;
6  import org.apache.hadoop.io.*;
7
8  public class InitReducer extends Reducer<Text, Text, Text, Text> {
9
10     public void reduce(Text key, Iterable<Text> values, Context context)
11         throws IOException, InterruptedException {
12         /*
13          * TODO: Output key: node+rank, value: adjacency list
14          */
15         // Default Rank is 1
16         int default_rank = 1;
17         Iterator<Text> iterator = values.iterator();
18
19         while(iterator.hasNext()) {
20             // Emit Node and Rank, Value
21             context.write(new Text(key + "+" + default_rank),
22                           iterator.next());
23         }
24     }
25 }
```

2.9 JoinMapper

Listing 9: JoinMapper.java

```
1  package edu.stevens.cs549.hadoop.pagerank;
2
3  import java.io.File;
4  import java.io.IOException;
5  import java.net.URI;
6
7  import org.apache.commons.io.FileUtils;
8  import org.apache.hadoop.fs.Path;
9  import org.apache.hadoop.io.LongWritable;
10 import org.apache.hadoop.io.Text;
11 import org.apache.hadoop.mapreduce.Mapper;
12
13 public class JoinMapper extends Mapper < LongWritable, Text, Text, Text > {
14     @Override
15     protected void setup(Context context) throws IOException,
16         InterruptedException {
17         if (context.getCacheFiles() != null && context.getCacheFiles().length > 0) {
18             URI mappingFileUri = context.getCacheFiles()[0];
19
20             if (mappingFileUri != null) {
21                 System.out.println("Mapping File: " +
22                     FileUtils.readFileToString(new File("./" +
23                         cache)));
24             } else {
25                 System.out.println("no mapping file");
26             }
27         } else {
28             System.out.println("no cache file");
29         }
30     }
31
32     public void map(LongWritable key, Text value, Context context) throws
33         IOException, InterruptedException, IllegalArgumentException {
34         String line = value.toString();
35         String[] sections;
36
37         if (line.contains(":")) {
38             int index = line.indexOf(":");
39             sections = new String[2];
40             sections[0] = line.substring(0, index);
41             sections[1] = line.substring(index + 1, line.length());
42         } else {
43             sections = line.split("\\t");
44         }
45
46         if (sections.length > 2) {
47             throw new IOException("Incorrect data format");
48         }
49
50         String[] node_rank = sections[0].split("\\|+");
51         if (node_rank.length == 1) {
52             context.write(new Text(node_rank[0]), new Text("name:" +
53                 sections[1].trim()));
54         }
55         if (node_rank.length == 2) {
56             context.write(new Text(node_rank[0]), new Text("rank:" +
57                 node_rank[1]));
58         }
59     }
60 }
```

2.10 JoinReducer

Listing 10: JoinReducer.java

```
1  package edu.stevens.cs549.hadoop.pagerank;
2
3  import java.io.IOException;
4  import java.util.Iterator;
5
6  import org.apache.hadoop.io.Text;
7  import org.apache.hadoop.mapreduce.Reducer;
8
9  public class JoinReducer extends Reducer < Text, Text, Text, Text > {
10     public void reduce(Text key, Iterable < Text > values, Context context
        ) throws IOException, InterruptedException,
        IllegalArgumentException {
11         Iterator <Text> iterator = values.iterator();
12         String node_name = "";
13         String rank = "";
14
15         while (iterator.hasNext()) {
16             String tmp = iterator.next().toString();
17             if (tmp.startsWith("name:")) {
18                 node_name = tmp.replaceAll("name:", "");
19             }
20             if (tmp.startsWith("rank:")) {
21                 rank = tmp.replaceAll("rank:", "");
22             }
23         }
24
25         context.write(new Text(key + "+" + node_name), new Text(rank))
                ;
26     }
27 }
```

3 Running on Local Machine

3.1 Graph Files

Listing 11: sample.txt

```

1 1: 2 3
2 2: 4
3 3: 1 4 5
4 4:
5 5: 1 4

```

Listing 12: names.txt

```

1 1: v1
2 2: v2
3 3: v3
4 4: v4
5 5: v5

```

Listing 13: output.txt

```

1 4:0 1.7083333333333333
2 1:0 0.8583333333333333
3 3:0 0.575
4 2:0 0.575
5 5:0 0.4333333333333335

```

3.2 Run Tests

The following command was used to test the local hadoop graph ranking algorithm.

Listing 14: Local Command

```

1 edu.stevens.cs549.hadoop.pagerank.PageRankDriver composite s3://cs549/input s3
  ://cs549/output inter1 inter2 diff 10

```

Table 1: Spread of COVID-19 by Continent

Number of Reducers	Processing Time
5	13s
10	18s
20	40s

Increasing the number of reducers increases the processing time of the algorithm.

4 Running on EMR

Listing 15: Custom JAR Command

```
1 edu.stevens.cs549.hadoop.pagerank.PageRankDriver composite s3://cs549/input s3
  ://cs549/output inter1 inter2 diff 10
```

The algorithm took 23 minutes to run on a 5 machine cluster with 10 reducers. The following ranking was outputted:

Listing 16: output.txt

```
1      5302153:0      12810.69631667126
2      84707:0       8159.113728956683
3      88822:0       7886.093686178045
4      1921890:0     7275.127835822817
5      5300058:0     5826.635099830321
6      81615:0       5014.404767519376
7      1804986:0     4197.8888148045
8      5535280:0     4164.226675136804
9      896161:0     3769.910788409694
10     5535664:0     3714.2571317056295
11     79583:0       3583.7680798360816
12     1601519:0     3488.3225964390367
13     687324:0     3455.841896902911
14     1948883:0     3364.344625704658
15     5308545:0     3190.9223377001085
16     505135:0     3064.802824618247
17     1603276:0     2941.357165895597
18     5596267:0     2914.369814878134
19     2497500:0     2825.4558209203
20     2995510:0     2730.416154609054
21     1650573:0     2659.3656421983587
22     2370447:0     2565.2363930365195
23     77935:0      2528.8948762668074
24     4141787:0     2407.816159087741
25     3492254:0     2356.5909651517313
26     2437900:0     2304.3565862362984
27     2401294:0     2295.6349478790767
28     4189168:0     2293.5300406199854
29     687618:0     2224.083020041053
30     3072654:0     2102.6401724186544
31     434174:0     2102.335180035135
32     3988566:0     2096.664388916008
33     4015997:0     2018.9134163444628
34     686242:0     2005.6101902332678
35     4696900:0     1961.1958496038294
36     76573:0      1960.9973826414164
37     4351989:0     1959.317983917727
38     478879:0     1904.298745200261
39     3603437:0     1898.3236554574517
40     5115901:0     1879.5875947888485
41     2876077:0     1846.0542169601665
42     1386743:0     1841.197756971188
43     5492723:0     1780.7479241121116
44     3997849:0     1712.6427307459228
45     75323:0      1630.2723881913055
46     3587465:0     1575.4630614468424
47     4089591:0     1571.4228562480328
48     181909:0     1553.487944478633
49     4015913:0     1537.6157012671888
50     3013310:0     1487.1735404348779
```

When trying to run more than 5 reducers or more than 5 clusters on the Amazon EMR service resulted in the following error:

```
Terminated with errorsThe request to create the EMR cluster or add EC2 instances to  
it failed. The number of vCPUs for instance type m5.xlarge exceeds the EC2 service  
quota for that type. Request a service quota increase for your AWS account or choose  
a different instance type and retry the request.  
For more information, see https://docs.aws.amazon.com/console/elasticmapreduce/vcpu-limit
```

I cannot find any information on Canvas about setting up EMR or addressing this issue. I asked a question on Piazza and awaiting for the reply at the time of submitting this assignment.