

## Assignment 5: PageRank

The goal of this assignment was to find the most popular pages in a set of Web pages. You will use your implementation of PageRank to find the 'most popular' pages in a dump of Wikipedia pages.

Decomposing of program, that is breaking program into few small, simple and sequential tasks which can be executed parallel by multiple computers is most important. The basic workflow should be:

***initial mapreduce -> iterate mapreducr -> evaluate difference -> join name -> finish***

I should emphasize on the finish phase. I have broken down the FINISH phase into two parts: ***finishjoin, finishcombine***. Meanwhile, one must change the threshold to adapt to different tasks. Let's move on the next process, DIFF, of which major responsibility is to calculate difference between the intermediate values produced by ITER. If the difference is small enough, we then consider mapreduce is stable and output the result. ITER is a process where we compute the rank. INIT is the first phase we initialize all the value.

Below is the summary of code written in the various files.

In DiffMap1.java file, completed the block of code in Map method to read node-rank pair and emit: key:node, value:rank.

```

18      /**
19       * TODO: read node-rank pair and emit: key:node, value:rank
20       */
21      //Separating Node and Rank
22      String[] noderank = sections[0].split("\\+");
23      //Emits Node and Rank
24      context.write(new Text(noderank[0]), new Text(noderank[1]));
25
26  }
```

In DiffMap2.java, completed the block of code in map method to emit the key, value pair where key is the difference and value is difference calculated in DiffRed1.java.

```

14      /*
15       * TODO: emit: key:"Difference" value: difference calculated in DiffRed1
16       */
17
18      String[] noderank = s.split("\t+");
19      context.write(new Text("Difference"), new Text(noderank[1]));
20  }
21
```

In DiffRed1.java file completed the block of code in reduce method which creates the list that contains the two ranks and out difference between those two ranks.

```

13      /*
14      * TODO: The list of values should contain two ranks. Compute and output their difference.
15      */
16
17      Iterator<Text> iter = values.iterator();
18      double diff = 0;
19      // Rank 1 Calculation
20      if(iter.hasNext()) {
21          ranks[0] = Double.valueOf(iter.next().toString());
22      }
23      // Rank 2 Calculation
24      if(iter.hasNext()) {
25          ranks[1] = Double.valueOf(iter.next().toString());
26      }
27      // Difference Calculation
28      diff = Math.abs(ranks[0] - ranks[1]);
29      System.out.println( key.toString() + " " + diff);
30      context.write(key, new Text(String.valueOf(diff)));
31  }

```

In DiffRed2.java file, completed the block of code in reduce method which computes the emits the maximum differences.

```

12      /*
13      * TODO: Compute and emit the maximum of the differences
14      */
15
16      Iterator<Text> iter = values.iterator();
17      // Find max difference and print in output
18      while(iter.hasNext()) {
19          double diff = Double.valueOf(iter.next().toString());
20          diff_max = diff_max > diff ? diff_max : diff;
21      }
22      context.write(new Text(""), new Text(String.valueOf(diff_max)));
23  }

```

In FindJoinMapper.java and FindJoinReducer.java written the code which joins the final output with linked name. The detailed code of the two files is as below:

### FindJoinMapper.java

```

69 public void map(LongWritable key, Text value, Context context)
70     throws IOException, InterruptedException, IllegalArgumentException {
71     String line = value.toString(); // Converts Line to a String
72     /*
73     * Join final output with link name
74     * input: key: nodeId+rank, text: adjacent list
75     * input: key: nodeId, text: name
76     *
77     * output: key: nodeId, text: rank
78     * output: key: nodeId, text: names
79     */
80     String[] section;
81     if (line.contains(":")) {
82         int ind = line.indexOf(":");
83         section = new String[2];
84         section[0] = line.substring(0, ind);
85         section[1] = line.substring(ind + 1, line.length());
86     } else {
87         section = line.split("\t"); // Splits it into two parts. Part 1: node+rank | Part 2: adj list
88     }
89     if (section.length > 2) // Verify Correct Data Format
90     {
91         throw new IOException("INVALID DATA FORMAT");
92     }
93     String[] noderank = section[0].split("\\+");
94     if (noderank.length == 1) {
95         // nodeID along with name
96         context.write(new Text(noderank[0]), new Text(PageRankDriver.MARKER_NAME + section[1].trim()));
97     }
98     if (noderank.length == 2) {
99         // nodeID along with rank
100        context.write(new Text(noderank[0]), new Text(PageRankDriver.MARKER_RANK + noderank[1]));
101    }
102 }

```

### FindJoinReducer.java

```

1 package edu.stevens.cs549.hadoop.pagerank;
2
3 import java.io.IOException;
4
5 public class FindJoinReducer extends Reducer<Text, Text, Text, Text> {
6
7     public void reduce(Text key, Iterable<Text> values, Context context)
8         throws IOException, InterruptedException, IllegalArgumentException {
9
10        /*
11        * Join final output with linked name
12        * input: key: nodeId, val: (mark_rank)rank
13        * input: key: nodeId, val: (mark_name)name
14        *
15        * output: key: nodeId+names, text: rank
16        */
17        Iterator<Text> iter = values.iterator();
18        String nodeName = "";
19        String rank = "";
20        while (iter.hasNext()) {
21            String temp = iter.next().toString();
22            if (temp.startsWith(PageRankDriver.MARKER_NAME)) {
23                nodeName = temp.replaceAll(PageRankDriver.MARKER_NAME, "");
24            }
25            if (temp.startsWith(PageRankDriver.MARKER_RANK)) {
26                rank = temp.replaceAll(PageRankDriver.MARKER_RANK, "");
27            }
28        }
29
30        context.write(new Text(key + "+" + nodeName), new Text(rank));
31    }
32 }

```

In FinMapper.jav file, completed the block of code which output the rank and node as key, value.

```

8 public class FinMapper extends Mapper<LongWritable, Text, DoubleWritable, Text> {
9
10 public void map(LongWritable key, Text value, Context context)
11     throws IOException, InterruptedException, IllegalArgumentException {
12     String line = value.toString(); // Converts Line to a String
13     /*
14      * TODO output key:-rank, value: node
15      * See IterMapper for hints on parsing the output of IterReducer.
16      */
17
18     String[] section = line.split("\t"); // nodeId+nodeName | rank
19
20     if (section.length > 2) //Verify correct data format
21     {
22         throw new IOException("INVALID DATA FORMAT");
23     }
24     if (section.length != 2) {
25         return;
26     }
27     // Reverse shuffle redce logic
28     context.write(new DoubleWritable(0 - Double.valueOf(section[1])), new Text(section[0]));
29
30 }

```

In FinReducer.java completed the block of code to emit each key value pair of value and rank.

```

9 public class FinReducer extends Reducer<DoubleWritable, Text, Text, Text> {
10
11 public void reduce(DoubleWritable key, Iterable<Text> values, Context context) throws IOException,
12     InterruptedException {
13     /*
14      * TODO: For each value, emit: key:value, value:-rank
15      */
16
17     Iterator<Text> iter = values.iterator();
18     String node;
19     while(iter.hasNext()) {
20         node = iter.next().toString();
21         // convert -rank back to rank
22         context.write(new Text(node), new Text(String.valueOf(0 - key.get())));
23     }
24 }
25 }
26

```

In InitMapper.java file, completed the block of code to echo the input which is already in the adjacency list.

```

10 public void map(LongWritable key, Text value, Context context) throws IOException, In
11     IllegalArgumentException {
12     String line = value.toString(); // Converts Line to a String
13     /*
14     * TODO: Just echo the input, since it is already in adjacency list format.
15     */
16
17     /*
18     * split the line by symbol ":", and output key, adjacent list to reducer
19     */
20     String[] p = line.split(":");
21     if(p != null && p.length == 2) {
22         context.write(new Text(p[0].trim()), new Text(p[1]));
23     }
24 }
25 }
26 }
27 }
28 }

```

In InitReducer.java file, completed the block of code which outputs the node rank and value in the adjacency list.

```

3 import java.io.*;
7
8 public class InitReducer extends Reducer<Text, Text, Text, Text> {
9
10 public void reduce(Text key, Iterable<Text> values, Context context) throws IOException,
11     /*
12     * TODO: Output key: node+rank, value: adjacency list
13     */
14
15     /*
16     * Since default rank is 1, so we need only output node+rank and adjacency list
17     */
18     int default_rank = 1;
19     Iterator<Text> val = values.iterator();
20     while(val.hasNext()) {
21         context.write(new Text(key + "+" + default_rank), val.next());
22     }
23 }
24 }
25 }
26 }
27 }
28 }

```

In IterMapper.java completed the block of code to emit adjacent vertex and computed weight as key, value pair. While doing so I have taken care to emit input adjacency list of the node.

```

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 it into two parts. Part 1: node+rank | Part 2: adj list
14
15     if (sections.length > 2) // Checks if the data is in the incorrect format
16     {
17         throw new IOException("Incorrect data format");
18     }
19     if (sections.length != 2) {
20         return;
21     }
22
23     /*
24     * TODO: emit key: adj vertex, value: computed weight.
25     *
26     * Remember to also emit the input adjacency list for this node!
27     * Put a marker on the string value to indicate it is an adjacency list.
28     */
29
30     // split node+rank
31     String[] node_rank = sections[0].split("\t");
32     String node = String.valueOf(node_rank[0]);
33     double rank = Double.valueOf(node_rank[1]);
34     String adj_list = sections[1].toString().trim();
35
36     String[] adj_nodes = adj_list.split(" ");
37     int N = adj_nodes.length; // outgoing links number
38     // 1/n * rank
39     double weightOfPage = (double)1/N * rank; // calculate current page weight if outgoing to other links
40     for(String adj_node : adj_nodes) {
41         context.write(new Text(adj_node), new Text(String.valueOf(weightOfPage)));
42     }
43     // at the same time, emit current node's adj_list list with marker "ADJ:"
44     context.write(new Text(node), new Text(PageRankDriver.MARKER + sections[1]));
45 }
46 }
47 }
48 }
49 }
50 }

```

In IterReducer.java file, completed block of code in reduce method which outputs the node+rank and adjacency list as key, value pair. Here as suggested by professor, I have used PageRank algorithm to compute the rank from weights contributed by incoming edges.

```
public void reduce(Text key, Iterable<Text> values, Context context) throws
    IOException, InterruptedException {
    double d = PageRankDriver.DECAY; // Decay factor
    /*
     * TODO: emit key:node+rank, value: adjacency list
     * Use PageRank algorithm to compute rank from weights contributed by
     * Remember that one of the values will be marked as the adjacency li
     */
    Iterator<Text> iter = values.iterator();
    double curr_Rank = 0; // default rank is 1 - d
    String adj_list = "";
    while(iter.hasNext()) {
        String line = iter.next().toString();
        if(!line.startsWith(PageRankDriver.MARKER)) {
            curr_Rank += Double.valueOf(line);
        } else {
            adj_list = line.replaceAll(PageRankDriver.MARKER, "");
        }
    }
    //sum calculations
    curr_Rank = 1 - d + curr_Rank * d;
    context.write(new Text(key + "+" + curr_Rank), new Text(adj_list));
}
```

In PageRankDriver.java file, written FinishJoin method through which operation is implemented to perform the join function. Through which you can set the Mapper and Reducer.

```
214@ static void finishJoin(String input, String output, int reducers)
215     throws Exception {
216     String temp_Join_Dir = "tempJoin";
217     System.out.println("Finish Job Started");
218     Job j = Job.getInstance(); // Creates a new Job
219     j.setJarByClass(PageRankDriver.class); // Sets the Driver class
220     j.setNumReduceTasks(reducers); // Sets the number of reducers
221
222     FileInputFormat.addInputPath(j, new Path(input)); // Adds input and output paths
223     //Input
224     FileInputFormat.addInputPath(j, new Path(CACHE_DIR));
225     FileOutputFormat.setOutputPath(j, new Path(temp_Join_Dir));
226     // Sets Mapper and Reducer Classes
227     j.setMapperClass(FindJoinMapper.class);
228     j.setReducerClass(FindJoinReducer.class);
229     // Sets Mapper and Reducer output types
230     j.setMapOutputKeyClass(Text.class);
231     j.setMapOutputValueClass(Text.class);
232
233     j.setOutputKeyClass(Text.class);
234     j.setOutputValueClass(Text.class);
235
236     // Prints message on successful completion or error.
237     if(j.waitForCompletion(true)) {
238         finishCombine(temp_Join_Dir, output, reducers);
239     }
240 }
```

Within the same file, written a function named finishCombine which will finish the job end and make the job started.

```
static void finishCombine(String input, String output, int reducers)
    throws Exception {
    System.out.println("Finish Join Job end, and Finish Job Started");
    // Creates a new Job
    Job j = Job.getInstance();
    // Sets the Driver class
    j.setJarByClass(PageRankDriver.class);
    // Sets the number of reducers
    j.setNumReduceTasks(reducers);
    // Adds input and output paths
    FileInputFormat.addInputPath(j, new Path(input));
    FileOutputFormat.setOutputPath(j, new Path(output));

    j.setMapperClass(FinMapper.class); // Sets Mapper and Reducer Classes
    j.setReducerClass(FinReducer.class);

    j.setMapOutputKeyClass(DoubleWritable.class); // Sets Mapper and Reducer output types
    j.setMapOutputValueClass(Text.class);

    j.setOutputKeyClass(Text.class);
    j.setOutputValueClass(Text.class);

    // Prints message on successful completion or error.
    System.out.println(j.waitForCompletion(true) ? "Finish Job Completed" : "Finish Job Error");
    // Exits once job finishes.

    deleteDirectory("tempJoin"); // Deletes the input directory
}
```

In Composite method, as suggested I have added my name and CWID.

```
273 public static void composite(String input, String output, String interim1,
274     String interim2, String diff, int reducers) throws Exception {
275     /*
276     * TODO
277     */
278     System.out.println("Abhilash Ugaonkar (10415787)");
279 }
```

## Experimenting with different number of Reducers

On EMR, I first ran the job on test data given in the assignment specification. While running that particular job I had used Reducers= 2

The output of the above job is as below:

```
D 1.7083333333333333
A 0.8583333333333333
C 0.575
B 0.575
E 0.43333333333333335
```

The job was completed in two minutes. Here are the details of job run.

Resource Groups

Clone Terminate AWS CLI export

Cluster: CS-549-EMR Terminated Terminated by user request

Summary Application history Monitoring Hardware Events Steps Configurations Bootstrap actions

Add step Clone step Cancel step

Steps

Filter: Completed steps  6 steps (all loaded)

ID	Name	Status	Start time (UTC-5)	Elapsed time	Log files	Actions
s-3M311WVRPK36E	Custom JAR	Completed	2017-12-10 13:31 (UTC-5)	41 minutes	<a href="#">View logs</a>	<a href="#">View jobs</a>
s-3160R3B31CMNR	Custom JAR	Completed	2017-12-10 11:23 (UTC-5)	41 minutes	<a href="#">View logs</a>	<a href="#">View jobs</a>
s-3MJY6EF2AJ4QL	Custom JAR	Completed	2017-12-10 01:45 (UTC-5)	2 minutes	<a href="#">View logs</a>	<a href="#">View jobs</a>

JAR location : s3://cs-549-assignment5/my\_test\_jar/PageRank-1.0.0.jar  
 Main class : None  
 Arguments : edu.stevens.cs549.hadoop.pagerank.PageRankDriver composite s3://cs-549-assignment5/links1 s3://cs-549-assignment5/output inter1 inter2 diffout 2  
 Action on failure: Continue

Thereafter, I ran the job on Wikipedia data downloaded from the links given in the assignment specification.

Reducers = 4

***edu.stevens.cs549.hadoop.pagerank.PageRankDriver composite s3://cs-549-assignment5/links s3://cs-549-assignment5/output inter1 inter2 diffout 4***

Resource Groups

Clone Terminate AWS CLI export

Cluster: CS-549-EMR Terminated Terminated by user request

Summary Application history Monitoring Hardware Events Steps Configurations Bootstrap actions

Add step Clone step Cancel step

Steps

Filter: Completed steps  6 steps (all loaded)

ID	Name	Status	Start time (UTC-5)	Elapsed time	Log files	Actions
s-3M311WVRPK36E	Custom JAR	Completed	2017-12-10 13:31 (UTC-5)	41 minutes	<a href="#">View logs</a>	<a href="#">View jobs</a>
s-3160R3B31CMNR	Custom JAR	Completed	2017-12-10 11:23 (UTC-5)	41 minutes	<a href="#">View logs</a>	<a href="#">View jobs</a>
s-3MJY6EF2AJ4QL	Custom JAR	Completed	2017-12-10 01:45 (UTC-5)	2 minutes	<a href="#">View logs</a>	<a href="#">View jobs</a>
s-2ETYGY1SXIWGR	Custom JAR	Completed	2017-12-10 01:33 (UTC-5)	2 minutes	<a href="#">View logs</a>	<a href="#">View jobs</a>
s-2M4HWZ2LYLZJS	Custom JAR	Completed	2017-12-09 19:54 (UTC-5)	47 minutes	<a href="#">View logs</a>	<a href="#">View jobs</a>

JAR location : s3://cs-549-assignment5/jars/PageRank-1.0.0.jar  
 Main class : None  
 Arguments : edu.stevens.cs549.hadoop.pagerank.PageRankDriver composite s3://cs-549-assignment5/links s3://cs-549-assignment5/output inter1 inter2 diffout 4  
 Action on failure: Continue

It took 47 minutes to finish the job. Output of the job is attached in my submission.



Going further, I increased the number of reducers from 4 to 9.  
Reducers = 9

***edu.stevens.cs549.hadoop.pagerank.PageRankDriver composite s3://cs-549-assignment5/links s3://cs-549-assignment5/output inter1 inter2 diffout 9***

The screenshot shows the AWS EMR console for cluster CS-549-EMR. The 'Steps' tab is selected, displaying a table of 6 steps. The second step, 's-3160R3B31CMNR', is highlighted. Its details show a 'Completed' status, a start time of 2017-12-10 11:23 (UTC-5), and an elapsed time of 41 minutes. The arguments for the job are: `edu.stevens.cs549.hadoop.pagerank.PageRankDriver composite s3://cs-549-assignment5/links s3://cs-549-assignment5/output inter1 inter2 diffout 9`.

ID	Name	Status	Start time (UTC-5)	Elapsed time	Log files	Actions
s-3M311WVRPK36E	Custom JAR	Completed	2017-12-10 13:31 (UTC-5)	41 minutes	<a href="#">View logs</a>	<a href="#">View jobs</a>
s-3160R3B31CMNR	Custom JAR	Completed	2017-12-10 11:23 (UTC-5)	41 minutes	<a href="#">View logs</a>	<a href="#">View jobs</a>

JAR location : s3://cs-549-assignment5/jars/PageRank-1.0.0.jar  
Main class : None  
Arguments : edu.stevens.cs549.hadoop.pagerank.PageRankDriver composite s3://cs-549-assignment5/links s3://cs-549-assignment5/output inter1 inter2 diffout 9  
Action on failure: Continue

Surprisingly, even though the number of reducers were increased it took only 41 minutes to complete the job. I have recorded the execution steps in one of the demo videos.

After this I thought of increasing number of reducers further to 11 and below were the results.

Reducers = 11

***edu.stevens.cs549.hadoop.pagerank.PageRankDriver composite s3://cs-549-assignment5/links s3://cs-549-assignment5/output inter1 inter2 diffout 11***

The screenshot shows the AWS EMR console for cluster CS-549-EMR. The 'Steps' tab is selected, displaying a table of 6 steps. The first step, 's-3M311WVRPK36E', is highlighted. Its details show a 'Completed' status, a start time of 2017-12-10 13:31 (UTC-5), and an elapsed time of 41 minutes. The arguments for the job are: `edu.stevens.cs549.hadoop.pagerank.PageRankDriver composite s3://cs-549-assignment5/links s3://cs-549-assignment5/output inter1 inter2 diffout 11`.

ID	Name	Status	Start time (UTC-5)	Elapsed time	Log files	Actions
s-3M311WVRPK36E	Custom JAR	Completed	2017-12-10 13:31 (UTC-5)	41 minutes	<a href="#">View logs</a>	<a href="#">View jobs</a>

JAR location : s3://cs-549-assignment5/jars/PageRank-1.0.0.jar  
Main class : None  
Arguments : edu.stevens.cs549.hadoop.pagerank.PageRankDriver composite s3://cs-549-assignment5/links s3://cs-549-assignment5/output inter1 inter2 diffout 11  
Action on failure: Continue

Surprisingly, it again took 41 minutes to finish the job. I have attached the output of the job in the submission.

I had thought that as the number of reducers will increase time to execute will also increase, however, in my case when I ran the above jobs the taken decreased and remain same equal for reducers = 9 and reducers = 11.

Ideally, my experience is as the number of reducers increases it takes longer to complete job.

I have attached the demonstration video, code, jars and report in the submission.