

## User-defined Counters

While processing information/data using MapReduce job, it is a challenge to monitor the progress of parallel threads running across nodes of distributed clusters. Moreover, it is also complicated to distinguish between the data that has been processed and the data which is yet to be processed. The MapReduce Framework offers a provision of user-defined Counters, which can be effectively utilized to monitor the progress of data across nodes of distributed clusters.

Now, let's use Hadoop Counters to identify the number of complaints pertaining to debt collection, mortgage and other categories in the consumer complaints dataset. The following job code, would help us to accomplish this.

[?](#)

```
1 package com.evoke.bigdata.mr.complaint;
   import org.apache.hadoop.conf.Configured;
2   import org.apache.hadoop.fs.Path;
3   import org.apache.hadoop.io.IntWritable;
4   import org.apache.hadoop.io.Text;
5   import org.apache.hadoop.mapred.FileInputFormat;
6   import org.apache.hadoop.mapred.FileOutputFormat;
7   import org.apache.hadoop.mapred.JobClient;
8   import org.apache.hadoop.mapred.JobConf;
9   import org.apache.hadoop.mapred.RunningJob;
   import org.apache.hadoop.util.Tool;
10  import org.apache.hadoop.util.ToolRunner;
11  public class ComplaintCounter extends Configured implements Tool {
12      @Override
13      public int run(String[] args) throws Exception {
14          String input, output;
15          if(args.length == 2) {
16              input = args[0];
17              output = args[1];
18          } else {
19              input = "your-input-dir";
20              output = "your-output-dir";
          }
```

```

21  JobConf conf = new JobConf(getConf(), ComplaintCounter.class);
22  conf.setJobName(this.getClass().getName());
23  FileInputFormat.setInputPaths(conf, new Path(input));
24  FileOutputFormat.setOutputPath(conf, new Path(output));
25  conf.setMapperClass(ComplaintCounterMapper.class);
26  conf.setMapOutputKeyClass(Text.class);
27  conf.setMapOutputValueClass(IntWritable.class);
28  conf.setOutputKeyClass(Text.class);
29  conf.setOutputValueClass(IntWritable.class);
30  conf.setNumReduceTasks(0);
31  RunningJob job = JobClient.runJob(conf);
32  long debt = job.getCounters().findCounter("Debt-Counter", "debt");
33  long mortgage = job.getCounters().findCounter("Mortgage-Counter", "mortgage");
34  long other = job.getCounters().findCounter("Other-Counter", "other");
35  System.out.println("Debt = " + debt);
36  System.out.println("Mortgage = " + mortgage);
37  System.out.println("OTHER = " + other);
38  return 0;
39  }
40  public static void main(String[] args) throws Exception {
41      int exitCode = ToolRunner.run(new ComplaintCounter(), args);
42      System.exit(exitCode);
43  }
44  }

```

The above Java code contains Map-only job and it displays the values of debt collection, mortgage and other occurrences.

Below is the map code, which is used in the above job code. The below map code scans through each line of the dataset and increments debt collection, mortgage and other occurrences.

[?](#)

```

1
2  package com.evoke.bigdata.mr.complaint;
3  import java.io.IOException;
4  import org.apache.hadoop.io.IntWritable;
5  import org.apache.hadoop.io.LongWritable;
6  import org.apache.hadoop.io.Text;
7  import org.apache.hadoop.mapred.MapReduceBase;
8  import org.apache.hadoop.mapred.Mapper;
9  import org.apache.hadoop.mapred.OutputCollector;
10 import org.apache.hadoop.mapred.Reporter;
11
12 public class ComplaintCounterMapper extends MapReduceBase implements
13     Mapper<LongWritable, Text, Text, IntWritable> {
14     public void map(LongWritable key, Text value,
15         OutputCollector<Text, IntWritable> output, Reporter reporter)
16         throws IOException {
17         String[] fields = value.toString().split(",");
18         if (fields.length > 1) {
19             String fileName = fields[1].toLowerCase();
20             if (fileName.equals("debt collection")) {
21                 reporter.getCounter("Debt-Counter", "debt").increment(1);
22             } else if (fileName.equals("mortgage")) {
23                 reporter.getCounter("Mortgage-Counter", "mortgage").increment(1);
24             } else {
25                 reporter.getCounter("Other-Counter", "other").increment(1);
26             }
27             output.collect(new Text(fileName), new IntWritable(1));
28         }
29     }

```

Here's the output, when the above map code is executed:

```
node1@ubuntu: ~  
File Edit View Terminal Help  
15/03/12 05:07:56 INFO mapred.JobClient: Map input records=326930  
15/03/12 05:07:56 INFO mapred.JobClient: Map output records=326930  
15/03/12 05:07:56 INFO mapred.JobClient: Input split bytes=214  
15/03/12 05:07:56 INFO mapred.JobClient: Spilled Records=0  
15/03/12 05:07:56 INFO mapred.JobClient: CPU time spent (ms)=5630  
15/03/12 05:07:56 INFO mapred.JobClient: Physical memory (bytes) snapshot=21  
9201536  
15/03/12 05:07:56 INFO mapred.JobClient: Virtual memory (bytes) snapshot=112  
8284160  
15/03/12 05:07:56 INFO mapred.JobClient: Total committed heap usage (bytes)=  
63569920  
15/03/12 05:07:56 INFO mapred.JobClient: Debt-Counter  
15/03/12 05:07:56 INFO mapred.JobClient: debt=47920  
15/03/12 05:07:56 INFO mapred.JobClient: Mortgage-Counter  
15/03/12 05:07:56 INFO mapred.JobClient: mortgage=129548  
15/03/12 05:07:56 INFO mapred.JobClient: Other-Counter  
15/03/12 05:07:56 INFO mapred.JobClient: other=149462  
15/03/12 05:07:56 INFO mapred.JobClient: org.apache.hadoop.mapreduce.lib.input  
.FileInputFormatCounter  
15/03/12 05:07:56 INFO mapred.JobClient: BYTES_READ=53280182  
Debt = 47920  
Mortgage = 129548  
OTHER = 149462  
node1@ubuntu:~$
```

Hadoop offers Job Tracker, an UI tool to determine the status and statistics of all jobs. Using the job tracker UI, developers can view the Counters that have been created (refer the below screenshot for more clarity).

	Counter
File System Counters	FILE: Number of bytes read
	FILE: Number of bytes written
	FILE: Number of read operations
	FILE: Number of large read operations
	FILE: Number of write operations
	HDFS: Number of bytes read
	HDFS: Number of bytes written
	HDFS: Number of read operations
	HDFS: Number of large read operations
	HDFS: Number of write operations
Job Counters	Launched map tasks
	Data-local map tasks
	Total time spent by all maps in occupied slots (ms)
	Total time spent by all reduces in occupied slots (ms)
	Total time spent by all maps waiting after reserving slots (ms)
	Total time spent by all reduces waiting after reserving slots (ms)
Map-Reduce Framework	Map input records
	Map output records
	Input split bytes
	Spilled Records
	CPU time spent (ms)
	Physical memory (bytes) snapshot
	Virtual memory (bytes) snapshot
	Total committed heap usage (bytes)
Debt-Counter	debt
Mortgage-Counter	mortgage
Other-Counter	other
org.apache.hadoop.mapreduce.lib.input.FileInputFormatCounter	BYTES_READ

## Conclusion

Hadoop is considered as one of the best parallel processing architectures around, it offers some really amazing features. Counter is one such feature in Hadoop that allows developers to track the status of data that has been processed. Additionally, it also helps developers

identify junk data and mark it for identification. Hope this post helped developers gain some knowledge on Hadoop Counters, in my next post, I would be talking about Hadoop Combiners.