

3 DELIVERY: HADOOP

OBJECTIVES

We are introducing Hadoop as a mapReduce platform and the Java classes as the programming framework. Our main objectives are the following:

- Programming a first implementation of word count in python using mapReduce streaming.
- Learning a mapreduce application implementation on Hadoop.
- Launching a job in a Hadoop cluster.

3.1 FIRST IMPLEMENTATION: COUNTING WORDS WITH PYTHON USING HADOOP STREAMING

WORKING WITH HDFS

First of all we create a new data folder in the HDFS file system:

```
1 hdfs dfs -mkdir /hduser/
2 hdfs dfs -mkdir /hduser/input
```

Now we copy an input file (“El Quijote de la Mancha”) from our local folder to HDFS file system:

```
1 wget https://www.gutenberg.org/cache/epub/2000/pg2000.txt
2 hdfs dfs -put ~/pg2000.txt /hduser/input
```

To make sure no problems occurred, we check the existence of the file in HDFS:

```
1 hdfs dfs -ls /hduser/input
```

```
Found 1 items
-rw-r--r--    1 cloudera supergroup    2198927 2018-02-12 07:39
  /hduser/input/pg2000.txt
```

IMPLEMENTATION OF THE MAPPER AND REDUCER

In the snippets 1 and 2 below we can see a simple implementation of a mapper and a reducer to count words:

Snippet 1: word_count_mapper.py implementation

```
1 #!/usr/bin/env python
2 import sys
3
4 for line in sys.stdin:
```

```

5  # Remove leading and trailing characters
6  line=line.strip()
7
8  # Split each line into words
9  words = line.split()
10
11 # Print concat(word, 1) for each word
12 for word in words:
13     print ('%s\t%s' % (word, "1"))

```

One thing we wanted to add to the mapper was to perform a `word.rstrip(',')` for each `word` in `words` to get rid of trailing comma characters, but the mapreduce job gave errors when doing this. So we decided just to keep the mapper untouched, as this is nothing more than a test.

Snippet 2: word_count_reducer.py implementation

```

1  #!/usr/bin/env python
2  import sys
3
4  word2count= {}
5
6  for line in sys.stdin:
7      line = line.strip()
8
9      # Save word and count of each concat(word, 1)
10     word, count = line.split('\t', 1)
11
12     # Check if count is a number
13     try:
14         count = int(count)
15     except ValueError: # Ignore the word
16         continue
17
18     # Add counts of each word in an array
19     try: # Normal behaviour: add + 1 to the counter
20         word2count[word] = word2count[word] + count
21     except: # Add new word to the "database" with value "1"
22         word2count[word] = count
23
24 for word in word2count.keys():
25     print ('%s\t%s' % ( word, word2count[word]))

```

It is always a good idea to make sure the permissions of the scripts are right, so the mapreduce job does not give an error:

```

1  chmod -R a+rx ~/src

```

RUNNING THE MAPREDUCE JOB

Now we need to use the Hadoop runtime with the `hadoop-streaming` tool to use our Python implementation:

```
1 hadoop jar /usr/lib/hadoop-mapreduce/hadoop-streaming.jar -mapper
   ~/src/word_count_mapper.py -reducer ~/src/word_count_reducer.py
   -input /hduser/input -output /hduser/output-python
```

If no errors occurred during the mapreduce job, we should get an output like this:

```
18/02/12 10:55:15 INFO mapreduce.Job: Counters: 50
  File System Counters
    FILE: Number of bytes read=3688605
    ...
    HDFS: Number of write operations=2
  Job Counters
    Killed map tasks=1
    ...
    Total megabyte-milliseconds taken by all reduce tasks=7725056
  Map-Reduce Framework
    Map input records=37861
    ...
    Total committed heap usage (bytes)=391979008
  Shuffle Errors
    BAD_ID=0
    ...
    WRONG_REDUCE=0
  File Input Format Counters
    Bytes Read=2203023
  File Output Format Counters
    Bytes Written=448894
18/02/12 10:55:15 INFO streaming.StreamJob: Output directory:
   /hduser/output-python
```

Now that we have the final output we can check, for instance, the words that appear the most in the text:

```
1 hdfs dfs -cat /hduser/output-python/part-00000 | awk '{print $2, $1}' FS=" " | sort -n | tail -n5
```

```
9575 a
10200 la
15894 y
17988 de
19429 que
```

One question one may ask is if there can be more than one file on the input folder. The answer is yes, but whilst running the mapreduce they are combined in a same output. This is, of course, the intended behaviour of Hadoop.

3.2 SECOND IMPLEMENTATION: COUNTING WORDS WITH JAVA CLASSES

In this part we will reproduce the previous implementation with Python for the mapreduce job, but using Java classes instead.

For the process we will be using the same input file (“El Quijote de la Mancha”), obtained from Project Gutenberg.

IMPLEMENTATION OF THE MAPPER AND REDUCER

The `WordCount` application is quite straight-forward:

- (i) The `Mapper` implementation, via the `map` method, processes one line at a time, as provided by the specified `TextInputFormat`. It then splits the line into tokens separated by white-spaces, via the `StringTokenizer`, and emits a key-value pair of `< <word>, 1>`.
- (ii) `WordCount` also specifies a combiner. Hence, the output of each map is passed through the local combiner (which is same as the `Reducer` as per the job configuration) for local aggregation, after being sorted on the keys.
- (iii) The `Reducer` implementation, via the `reduce` method just sums up the values, which are the occurrence counts for each word.
- (iv) The `main` method specifies various facets of the job, such as the input/output paths (passed via the command line), key/value types, input/output formats etc., in the `Job`. It then calls the `job.waitForCompletion` to submit the job and monitor its progress.

Snippet 3: `WordCount.java` implementation

```
1 import java.io.IOException;
2 import java.util.StringTokenizer;
3 import org.apache.hadoop.conf.Configuration;
4 import org.apache.hadoop.fs.Path;
5 import org.apache.hadoop.io.IntWritable;
6 import org.apache.hadoop.io.Text;
7 import org.apache.hadoop.mapreduce.Job;
8 import org.apache.hadoop.mapreduce.Mapper;
9 import org.apache.hadoop.mapreduce.Reducer;
10 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
11 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
12
13 public class WordCount {
14
15     public static class TokenizerMapper extends Mapper<Object, Text,
16         Text, IntWritable> {
17         private final static IntWritable one = new IntWritable(1);
18         private Text word = new Text();
```

```

19     public void map(Object key, Text value, Context context) throws
        IOException, InterruptedException {
20         StringTokenizer itr = new StringTokenizer(value.toString());
21         while (itr.hasMoreTokens()) {
22             word.set(itr.nextToken());
23             context.write(word, one);
24         }
25     }
26 }
27
28 public static class IntSumReducer extends Reducer<Text,
        IntWritable, Text, IntWritable> {
29     private IntWritable result = new IntWritable();
30
31     public void reduce(Text key, Iterable<IntWritable> values,
        Context context) throws IOException, InterruptedException {
32         int sum = 0;
33         for (IntWritable val : values) {
34             sum += val.get();
35         }
36         result.set(sum);
37         context.write(key, result);
38     }
39 }
40
41 public static void main(String[] args) throws Exception {
42     Configuration conf = new Configuration();
43     Job job = Job.getInstance(conf, "word_count");
44     job.setJarByClass(WordCount.class);
45     job.setMapperClass(TokenizerMapper.class);
46     job.setCombinerClass(IntSumReducer.class);
47     job.setReducerClass(IntSumReducer.class);
48     job.setOutputKeyClass(Text.class);
49     job.setOutputValueClass(IntWritable.class);
50
51     FileInputFormat.addInputPath(job, new Path(args[0]));
52     FileOutputFormat.setOutputPath(job, new Path(args[1]));
53     System.exit(job.waitForCompletion(true) ? 0 : 1);
54 }
55 }

```

CREATING A .JAR BINARY

The first thing to do is to add these lines to the `.bashrc` file to export the required environment variables:

```

1 export PATH=${JAVA_HOME}/bin:${PATH}
2 export HADOOP_CLASSPATH=${JAVA_HOME}/lib/tools.jar

```

We need to get the class files from the `WordCount.java` source code:

```
1 cd ~/src
2 hadoop com.sun.tools.javac.Main WordCount.java
```

The following files will be created:

- WordCount.class
- WordCount\$IntSumReducer.class
- WordCount\$TokenizerMapper.class

Now we only need to create the WordCount.jar binary:

```
1 jar cf WordCount.jar WordCount*.class
```

RUNNING THE MAPREDUCE JOB

Now we need to use the Hadoop runtime with the `jar` tool to use our Java implementation:

```
1 cd
2 hadoop jar ~/src/WordCount.jar WordCount /hduser/input
   /hduser/output-java
```

If no errors occurred during the mapreduce job, we can check, for instance, the words that appear the most in the text:

```
1 hdfs dfs -cat /hduser/output-java/part-r-00000 | awk '{print $2,
   $1}' FS=" " | sort -n | tail -n5
```

```
9575 a
10200 la
15894 y
17988 de
19429 que
```

As expected, we get the same result we got with the implementation in Python. It is important to notice that this mapper and reducer implementation has the same problem the previous one had: it keeps trailing comma characters.

HOW NOT TO DO THE PREVIOUS STEPS

Notice that in the previous instructions we manually changed the working directory to `~/src` before creating the `.jar` binary; this was a crucial step. Imagine that instead we had done the following:

```
1 hadoop com.sun.tools.javac.Main ~/src/WordCount.java
2 jar cf ~/src/WordCountBad.jar ~/src/WordCount*.class
3 hadoop jar ~/src/WordCountBad.jar WordCount /hduser/input
   /hduser/output-java-bad
```

```
Exception in thread "main" java.lang.ClassNotFoundException:
    WordCount
    at java.net.URLClassLoader$1.run(URLClassLoader.java:366)
    at java.net.URLClassLoader$1.run(URLClassLoader.java:355)
    at java.security.AccessController.doPrivileged(Native Method)
    at java.net.URLClassLoader.findClass(URLClassLoader.java:354)
    at java.lang.ClassLoader.loadClass(ClassLoader.java:425)
    at java.lang.ClassLoader.loadClass(ClassLoader.java:358)
    at java.lang.Class.forName0(Native Method)
    at java.lang.Class.forName(Class.java:270)
    at org.apache.hadoop.util.RunJar.run(RunJar.java:214)
    at org.apache.hadoop.util.RunJar.main(RunJar.java:136)
```

What is going on here? Let's have a look at the two compiled `.jar` files:

```
1 jar -tvf ~/src/WordCount.jar
```

```
0 Sat Feb 17 06:22:42 PST 2018 META-INF/
68 Sat Feb 17 06:22:42 PST 2018 META-INF/MANIFEST.MF
1501 Sat Feb 17 06:22:40 PST 2018 WordCount.class
1739 Sat Feb 17 06:22:40 PST 2018 WordCount$IntSumReducer.class
1736 Sat Feb 17 06:22:40 PST 2018 WordCount$TokenizerMapper.class
```

```
1 jar -tvf ~/src/WordCountBad.jar
```

```
0 Sat Feb 17 15:06:04 PST 2018 META-INF/
68 Sat Feb 17 15:06:04 PST 2018 META-INF/MANIFEST.MF
1501 Sat Feb 17 15:06:04 PST 2018 home/cloudera/src/WordCount.class
1739 Sat Feb 17 15:06:04 PST 2018
    home/cloudera/src/WordCount$IntSumReducer.class
1736 Sat Feb 17 15:06:04 PST 2018
    home/cloudera/src/WordCount$TokenizerMapper.class
```

The issue here is that if we run the `jar cf` command to create the binary from outside the directory where the `*.class` files are located, the binary will include the full *hard* path of the classes, instead of having a *soft* inclusion of them.

It is, of course, possible to run the mapreduce job with the *bad* `.jar` binary, if we call the class by the right name, but it is just simpler to compile the `.jar` binary properly to avoid any mistakes.

3.3 THIRD IMPLEMENTATION: PROCESSING WEATHER DATA WITH PYTHON

PROCESSING WEATHER DATA

Many big data applications start with a large size bulk dataset available at a remote server to download. In this case, we are using data from the National Climatic Data Center (NCDC). Each line in the dataset represents a record where a list of meteorological measures is available. The following is an example of all measures that can be found in a record:

```

1 0057
2 332130    # USAF weather station identifier
3 99999     # WBAN weather station identifier
4 19500101  # observation date
5 0300      # observation time
6 4
7 +51317    # latitude (degrees x 1000)
8 +028783   # longitude (degrees x 1000)
9 FM-12
10 +0171     # elevation (meters)
11 99999
12 V020
13 320       # wind direction (degrees)
14 1         # quality code
15 N
16 0072
17 1
18 00450     # sky ceiling height (meters)
19 1         # quality code
20 C
21 N
22 010000    # visibility distance (meters)
23 1         # quality code
24 N
25 9
26 -0128     # air temperature (degrees Celsius x 10)
27 1         # quality code
28 -0139     # dew point temperature (degrees Celsius x 10)
29 1         # quality code
30 10268     # atmospheric pressure (hectopascals x 10)
31 1         # quality code

```

Our problem is to find the maximum and mean temperature of each year found in the dataset. Then, we will be focusing on observation date and air temperature parameters of each record. First we are going to introduce the structure of the application and then we are going to focus on the implementation details of each part of the map reduce application.

Before we start working with Python and Java to construct our mapreduce implementa-

tion, we need to put our data on the HDFS:

```

1 for (( year = 1901; year <= 1920; year ++ )); do
2   hdfs dfs -put ~/data/ncdc-dataset/$year /hduser/input-weather/$year
3 done

```

The dataset for years 1901–1920 can be downloaded from <https://github.com/aldomann/parallel-and-distributed-systems/blob/master/data/ncdc-dataset.zip>.

THE MAP REDUCE JOB

For this process we are interested in three fields from the dataset:

- Year (characters 15–19).
- Air temperature (characters 87–92).
- Quality code for the temperature (characters 92–93).

In the snippet 4 below we can see a simple implementation of a mapper to get the year and the temperature from the NCDC dataset:

Snippet 4: temp_mapper.py implementation

```

1 #!/usr/bin/env python
2
3 import re
4 import sys
5
6 for line in sys.stdin:
7   # Remove heading and tailing characters
8   val = line.strip()
9
10  # String split to get data
11  (year, temp, q) = (val[15:19], val[87:92], val[92:93])
12
13  # Print concat(year, temp) for each year
14  if (temp != "+9999" and re.match("[01459]", q)):
15    print ("%s\t%s" % (year, temp))

```

In the snippet 5, 6, and 7 below we can see three different reducers:

- temp_max_reducer.py looks for the maximum temperature value for each year.
- temp_min_reducer.py looks for the minimum temperature value for each year.
- temp_mean_reducer.py calculates the mean temperature for each year.

Snippet 5: temp_max_reducer.py implementation

```

1 #!/usr/bin/env python
2
3 import sys
4
5 (last_key, max_val) = (None, 0)
6 for line in sys.stdin:

```

```

7   (key, val) = line.strip().split("\t")
8   if last_key and last_key != key:
9       # Print results of non-last key
10      print ("%s\t%s" % (last_key, max_val))
11      # Starts reducing new key
12      (last_key, max_val) = (key, int(val))
13  else:
14      # Process data
15      (last_key, max_val) = (key, max(max_val, int(val)))
16
17  if last_key:
18      # Print results of last key
19      print ("%s\t%s" % (last_key, max_val))

```

Snippet 6: temp_min_reducer.py implementation

```

1  #!/usr/bin/env python
2
3  import sys
4
5  (last_key, min_val) = (None, 0)
6  for line in sys.stdin:
7      (key, val) = line.strip().split("\t")
8      if last_key and last_key != key:
9          # Print results of non-last key
10         print ("%s\t%s" % (last_key, min_val))
11         # Starts reducing new key
12         (last_key, min_val) = (key, int(val))
13     else:
14         # Process data
15         (last_key, min_val) = (key, min(min_val, int(val)))
16
17  if last_key:
18      # Print results of last key
19      print ("%s\t%s" % (last_key, min_val))

```

Snippet 7: temp_mean_reducer.py implementation

```

1  #!/usr/bin/env python
2
3  import sys
4
5  (last_key, max_val, count) = (None, 0, 0)
6  for line in sys.stdin:
7      (key, val) = line.strip().split("\t")
8      if last_key and last_key != key:
9          # Print results of non-last key
10         print ("%s\t%s" % (last_key, max_val/count))
11         # Starts reducing new key
12         (last_key, max_val, count) = (key, int(val), 1)
13     else:
14         # Process data

```

```
15     (last_key, max_val, count) = (key, int(max_val) + int(val),  
16         int(count) + 1)  
17 if last_key:  
18     # Print results of last key  
19     print ("%s\t%s" % (last_key, max_val/count))
```

3.3.1 RUNNING THE MAXIMUM TEMPERATURE MAPREDUCE JOB

As we already did before, we need to use the Hadoop runtime with the `hadoop-streaming` tool to use our Python implementation:

```
1 hadoop jar /usr/lib/hadoop-mapreduce/hadoop-streaming.jar -mapper  
  ~/src/temp_mapper.py -reducer ~/src/temp_max_reducer.py -input  
  /hduser/input-weather -output /hduser/output-max-temp
```

After the mapreduce job finishes successfully, we just need to check the output:

```
1 hdfs dfs -cat /hduser/output-max-temp/part-00000
```

```
1901  317  
1902  244  
1903  289  
1904  256  
1905  283  
1906  294  
1907  283  
1908  289  
1909  278  
1910  294  
1911  306  
1912  322  
1913  300  
1914  333  
1915  294  
1916  278  
1917  317  
1918  322  
1919  378  
1920  294
```

3.3.2 RUNNING THE MINIMUM TEMPERATURE MAPREDUCE JOB

Now we just need to use the Hadoop runtime with the `hadoop-streaming` tool to use our Python implementation:

```
1 hadoop jar /usr/lib/hadoop-mapreduce/hadoop-streaming.jar -mapper  
  ~/src/temp_mapper.py -reducer ~/src/temp_min_reducer.py -input  
  /hduser/input-weather -output /hduser/output-min-temp
```

After the mapreduce job finishes successfully, we just need to check the output:

```
1 hdfs dfs -cat /hduser/output-min-temp/part-00000
```

```
1901 -333
1902 -328
1903 -306
1904 -294
1905 -328
1906 -250
1907 -350
1908 -378
1909 -378
1910 -372
1911 -378
1912 -411
1913 -372
1914 -378
1915 -411
1916 -289
1917 -478
1918 -450
1919 -428
1920 -344
```

3.3.3 RUNNING THE MEAN TEMPERATURE MAPREDUCE JOB

Now we just need to use the Hadoop runtime with the `hadoop-streaming` tool to use our Python implementation:

```
1 hadoop jar /usr/lib/hadoop-mapreduce/hadoop-streaming.jar -mapper
  ~/src/temp_mapper.py -reducer ~/src/temp_mean_reducer.py -input
  /hduser/input-weather -output /hduser/output-mean-temp
```

After the mapreduce job finishes successfully, we just need to check the output:

```
1 hdfs dfs -cat /hduser/output-mean-temp/part-00000
```

```
1901 46.698507007922
1902 21.659558263518658
1903 48.241744739671326
1904 33.32224247948952
1905 43.3322664228014
1906 47.0834855681403
1907 31.76414576084966
1908 28.836573511543136
1909 26.565303955402175
1910 35.558665794637015
1911 30.719045120671563
1912 16.801145236855803
1913 29.958786491127647
```

1914	29.817932296431838
1915	5.098548073625243
1916	21.42393787117405
1917	22.91685727355901
1918	31.36519845111326
1919	27.605149653640048
1920	43.508667830133795

3.4 FOURTH IMPLEMENTATION: PROCESSING WEATHER DATA

The idea here is the same, to implement a mapper and a reducer to find the maximum, mean, and minimum temperature for each year in the dataset.

First we need to define our *main* class `MaxTemperature`, where we define the *main* function that calls to the mapper and the reducer:

Snippet 8: `MaxTemperature.java` implementation

```

1 import org.apache.hadoop.fs.Path;
2 import org.apache.hadoop.io.IntWritable;
3 import org.apache.hadoop.io.Text;
4 import org.apache.hadoop.mapreduce.Job;
5 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
6 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
7
8 public class MaxTemperature {
9
10     public static void main(String[] args) throws Exception {
11         if (args.length != 2) {
12             System.err.println("Usage: MaxTemperature <input path> <output <
13                 path>");
14             System.exit(-1);
15         }
16
17         @SuppressWarnings("deprecation")
18         Job job = new Job();
19         job.setJarByClass(MaxTemperature.class);
20         job.setJobName("Max temperature");
21
22         FileInputFormat.addInputPath(job, new Path(args[0]));
23         FileOutputFormat.setOutputPath(job, new Path(args[1]));
24
25         job.setMapperClass(TemperatureMapper.class);
26         job.setReducerClass(MaxTemperatureReducer.class);
27
28         job.setOutputKeyClass(Text.class);
29         job.setOutputValueClass(IntWritable.class);
30
31         System.exit(job.waitForCompletion(true) ? 0 : 1);
32     }
33 }

```

The `MinTemperature.java` and `MeanTemperature.java` source codes are essentially the same, so we will not include them here explicitly. One just needs to be careful to call the classes by the right name.

For all of the mapreduce jobs we will use the same mapper, as the map job just extracts the relevant information from the dataset. The implementation can be seen in snippet 9 below:

 Snippet 9: TemperatureMapper.java implementation

```

1 import java.io.IOException;
2
3 import org.apache.hadoop.io.IntWritable;
4 import org.apache.hadoop.io.LongWritable;
5 import org.apache.hadoop.io.Text;
6 import org.apache.hadoop.mapreduce.Mapper;
7
8 public class TemperatureMapper
9     extends Mapper<LongWritable, Text, Text, IntWritable> {
10
11     private static final int MISSING = 9999;
12
13     @Override
14     public void map(LongWritable key, Text value, Context context)
15         throws IOException, InterruptedException {
16
17         String line = value.toString();
18         System.out.println(line);
19
20         String year = line.substring(15, 19);
21         int airTemperature;
22         if (line.charAt(87) == '+') { // parseInt doesn't like leading
23             plus signs
24             airTemperature = Integer.parseInt(line.substring(88, 92));
25         } else {
26             airTemperature = Integer.parseInt(line.substring(87, 92));
27         }
28         String quality = line.substring(92, 93);
29         if (airTemperature != MISSING && quality.matches("[01459]")) {
30             context.write(new Text(year), new IntWritable(airTemperature));
31         }
32     }
  
```

The reducer is the usual maximum calculation algorithm, the main difference with the implementation in Python is that `maxValue` is initialised to `Integer.MIN_VALUE = -2147483648`. The implementation for `MaxTemperatureReducer.java` can be seen in snippet 10 below:

 Snippet 10: MaxTemperatureReducer.java implementation

```

1 import java.io.IOException;
2
3 import org.apache.hadoop.io.IntWritable;
4 import org.apache.hadoop.io.Text;
5 import org.apache.hadoop.mapreduce.Reducer;
6
7 public class MaxTemperatureReducer
8     extends Reducer<Text, IntWritable, Text, IntWritable> {
9
  
```

```

10  @Override
11  public void reduce(Text key, Iterable<IntWritable> values,
12      Context context)
13      throws IOException, InterruptedException {
14
15      int maxValue = Integer.MIN_VALUE;
16      for (IntWritable value : values) {
17          maxValue = Math.max(maxValue, value.get());
18      }
19      context.write(key, new IntWritable(maxValue));
20  }
21  }

```

For `MinTemperatureReducer.java` the main difference is in the initialisation of `minValue` using `Integer.MAX_VALUE = 2147483647`, and the calculation of the minimum. The changes in the implementation can be seen in snippet 11 below:

Snippet 11: `MinTemperatureReducer.java` implementation

```

15  int minValue = Integer.MAX_VALUE;
16  for (IntWritable value : values) {
17      minValue = Math.min(minValue, value.get());
18  }

```

For `MeanTemperatureReducer.java` we need to work with the `DoubleWritable` class, as we need to calculate the ratio between the total sum of temperature values (`totalValue`) and the amount of records per year (`count`). The implementation can be seen in snippet 12 below:

Snippet 12: `MeanTemperatureReducer.java` implementation

```

1  import java.io.IOException;
2
3  import org.apache.hadoop.io.IntWritable;
4  import org.apache.hadoop.io.DoubleWritable;
5  import org.apache.hadoop.io.Text;
6  import org.apache.hadoop.mapreduce.Reducer;
7
8  public class MeanTemperatureReducer
9      extends Reducer<Text, IntWritable, Text, DoubleWritable> {
10
11      @Override
12      public void reduce(Text key, Iterable<IntWritable> values,
13          Context context)
14          throws IOException, InterruptedException {
15
16          double totalValue = 0.0;
17          int count = 0;
18          for (IntWritable value : values) {
19              totalValue += value.get();
20              count += 1;

```



```
21     }
22     context.write(key, new DoubleWritable(totalValue/count));
23 }
24 }
```

CREATING THE .JAR BINARIES

We need to get the class files from the different .java source codes:

```
1 cd ~/src
2 hadoop com.sun.tools.javac.Main MaxTemperature.java
   TemperatureMapper.java MaxTemperatureReducer.java
3 hadoop com.sun.tools.javac.Main MinTemperature.java
   TemperatureMapper.java MinTemperatureReducer.java
4 hadoop com.sun.tools.javac.Main MeanTemperature.java
   TemperatureMapper.java MeanTemperatureReducer.java
```

The following files will be created:

- MaxTemperature.class
- MinTemperature.class
- MeanTemperature.class
- TemperatureMapper.class
- MaxTemperatureReducer.class
- MinTemperatureReducer.class
- MeanTemperatureReducer.class

To create the .jar binaries, we use:

```
1 jar cf MaxTemperature.jar MaxTemperature.class
   TemperatureMapper.class MaxTemperatureReducer.class
2 jar cf MinTemperature.jar MinTemperature.class
   TemperatureMapper.class MinTemperatureReducer.class
3 jar cf MeanTemperature.jar MeanTemperature.class
   TemperatureMapper.class MeanTemperatureReducer.class
```

3.4.1 RUNNING THE MAXIMUM TEMPERATURE MAPREDUCE JOB

As we already did before, we need to use the Hadoop runtime with the `jar` tool to use our Java implementation:

```
1 hadoop jar ~/src/MaxTemperature.jar MaxTemperature
   /hduser/input-weather /hduser/output-max-temp-java
```

After the mapreduce job finishes successfully, we just need to check the output:

```
1 hdfs dfs -cat /hduser/output-max-temp-java/part-r-00000
```

```
1901 317
1902 244
1903 289
1904 256
1905 283
1906 294
1907 283
1908 289
1909 278
1910 294
1911 306
1912 322
1913 300
1914 333
1915 294
1916 278
1917 317
1918 322
1919 378
1920 294
```

3.4.2 RUNNING THE MINIMUM TEMPERATURE MAPREDUCE JOB

We just need to use the Hadoop runtime with the `jar` tool to use our Java implementation:

```
1 hadoop jar ~/src/MinTemperature.jar MinTemperature
   /hduser/input-weather /hduser/output-min-temp-java
```

After the mapreduce job finishes successfully, we just need to check the output:

```
1 hdfs dfs -cat /hduser/output-min-temp-java/part-r-00000
```

```
1901 -333
1902 -328
1903 -306
1904 -294
1905 -328
1906 -250
1907 -350
1908 -378
1909 -378
1910 -372
1911 -378
1912 -411
1913 -372
1914 -378
1915 -411
1916 -289
1917 -478
1918 -450
```

1919	-428
1920	-344

3.4.3 RUNNING THE MEAN TEMPERATURE MAPREDUCE JOB

We just need to use the Hadoop runtime with the `jar` tool to use our Java implementation:

```
1 hadoop jar ~/src/MeanTemperature.jar MeanTemperature  
   /hduser/input-weather /hduser/output-mean-temp-java
```

After the mapreduce job finishes successfully, we just need to check the output:

```
1 hdfs dfs -cat /hduser/output-mean-temp-java/part-r-00000
```

1901	46.698507007922
1902	21.659558263518658
1903	48.241744739671326
1904	33.32224247948952
1905	43.3322664228014
1906	47.0834855681403
1907	31.76414576084966
1908	28.836573511543136
1909	26.565303955402175
1910	35.558665794637015
1911	30.719045120671563
1912	16.801145236855803
1913	29.958786491127647
1914	29.817932296431838
1915	5.098548073625243
1916	21.42393787117405
1917	22.91685727355901
1918	31.36519845111326
1919	27.605149653640048
1920	43.508667830133795