

# Lab Assignment 1

Author: Vendramin Nicolò

## Imports

```
import java.io.IOException;
import java.util.StringTokenizer;

import java.io.IOException;
import java.util.Map;
import java.util.TreeMap;
import java.util.HashMap;

import org.apache.hadoop.io.NullWritable;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.hadoop.mapreduce.Reducer;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
```

## Main

```
public static void main(String[] args) throws Exception {
    Configuration conf = new Configuration();
    Job job = Job.getInstance(conf, "top ten");
    job.setNumReduceTasks(1);
    job.setJarByClass(TopTen.class);
    job.setMapperClass(TopTenMapper.class);
    job.setCombinerClass(TopTenReducer.class);
    job.setReducerClass(TopTenReducer.class);
    job.setOutputKeyClass(NullWritable.class);
    job.setOutputValueClass(Text.class);
    FileInputFormat.addInputPath(job, new Path(args[0]));
    FileOutputFormat.setOutputPath(job, new Path(args[1]));
    System.exit(job.waitForCompletion(true) ? 0 : 1);
}
```

## Mapper

```
public static Map<String, String> transformXmlToMap(String xml) {
    Map<String, String> map = new HashMap<String, String>();
    try {
        String[] tokens = xml.trim().substring(5, xml.trim().length() - 3)
            .split("\n");

        for (int i = 0; i < tokens.length - 1; i += 2) {
            String key = tokens[i].trim();
            String val = tokens[i + 1].trim();

            map.put(key.substring(0, key.length() - 1), val);
        }
    } catch (StringIndexOutOfBoundsException e) {
        System.err.println(xml);
    }

    return map;
}

public static class TopTenMapper extends Mapper<Object, Text, NullWritable, Text> {
    // Stores a map of user reputation to the record
    private TreeMap<Integer, Text> repToRecordMap = new TreeMap<Integer, Text>();

    public void map(Object key, Text value, Context context) throws IOException, InterruptedException {
        Map<String, String> parsed = transformXmlToMap(value.toString());
        String userId = parsed.get("Id");
        String reputation = parsed.get("Reputation");

        // check that this row contains user data
        if (userId == null || reputation == null) {
            return;
        }

        repToRecordMap.put(Integer.parseInt(reputation), new Text(value));

        // If we have more than ten records, remove the one with the lowest reputation.
        if (repToRecordMap.size() > 10) {
            repToRecordMap.remove(repToRecordMap.firstKey());
        }
    }

    protected void cleanup(Context context) throws IOException, InterruptedException {
        // Output our ten records to the reducers with a null key
        for (Text t : repToRecordMap.values()) {
            context.write(NullWritable.get(), t);
        }
    }
}
```

## Reducer

```
public static class TopTenReducer extends Reducer<NullWritable, Text, NullWritable, Text> {
    // Stores a map of user reputation to the record
    // Overloads the comparator to order the reputations in descending order
    private TreeMap<Integer, Text> repToRecordMap = new TreeMap<Integer, Text>();

    public void reduce(NullWritable key, Iterable<Text> values, Context context) throws IOException, InterruptedException {
        for (Text value : values) {
            Map<String, String> parsed = transformXmlToMap(value.toString());
            repToRecordMap.put(Integer.parseInt(parsed.get("Reputation")), new Text(value));

            // If we have more than ten records, remove the one with the lowest reputation
            if (repToRecordMap.size() > 10) {
                repToRecordMap.remove(repToRecordMap.firstKey());
            }

            // Sort in descending order
            for (Text t : repToRecordMap.descendingMap().values()) {
                // Output our ten records to the file system with a null key
                context.write(NullWritable.get(), t);
            }
        }
    }
}
```

## Output (1)

```
bash-4.1# bin/hadoop jar work/lab1/topten/topten.jar labcode.TopTen /user/root/lab1/topten/input /user/root/lab1/topten/output
17/09/20 11:48:56 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
17/09/20 11:48:57 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
17/09/20 11:48:57 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
17/09/20 11:48:57 INFO input.FileInputFormat: Total input paths to process : 1
17/09/20 11:48:58 INFO mapreduce.JobSubmitter: number of splits:1
17/09/20 11:48:58 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1505902951661_0006
17/09/20 11:48:58 INFO impl.YarnClientImpl: Submitted application application_1505902951661_0006
17/09/20 11:48:58 INFO mapreduce.Job: The url to track the job: http://ca2689501238:8088/proxy/application_1505902951661_0006/
17/09/20 11:49:03 INFO mapreduce.Job: Job job_1505902951661_0006 running in uber mode : false
17/09/20 11:49:03 INFO mapreduce.Job: map 0% reduce 0%
17/09/20 11:49:08 INFO mapreduce.Job: map 100% reduce 0%
17/09/20 11:49:12 INFO mapreduce.Job: map 100% reduce 100%
17/09/20 11:49:13 INFO mapreduce.Job: Job job_1505902951661_0006 completed successfully
17/09/20 11:49:13 INFO mapreduce.Job: Counters: 49
  File System Counters
    FILE: Number of bytes read=7244
    FILE: Number of bytes written=245903
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=5856430
    HDFS: Number of bytes written=7180
    HDFS: Number of read operations=6
    HDFS: Number of large read operations=0
    HDFS: Number of write operations=2
  Job Counters
    Launched map tasks=1
    Launched reduce tasks=1
    Data-local map tasks=1
    Total time spent by all maps in occupied slots (ms)=2501
    Total time spent by all reduces in occupied slots (ms)=2254
    Total time spent by all map tasks (ms)=2501
    Total time spent by all reduce tasks (ms)=2254
    Total vcore-seconds taken by all map tasks=2501
    Total vcore-seconds taken by all reduce tasks=2254
    Total megabyte-seconds taken by all map tasks=2561024
    Total megabyte-seconds taken by all reduce tasks=2308096
  Map-Reduce Framework
    Map input records=13995
    Map output records=10
    Map output bytes=7180
    Map output materialized bytes=7244
    Input split bytes=127
    Combine input records=10
    Combine output records=10
    Reduce input groups=1
    Reduce shuffle bytes=7244
    Reduce input records=10
    Reduce output records=10
    Spilled Records=20
    Shuffled Maps=1
    Failed Shuffles=0
    Merged Map outputs=1
    GC time elapsed (ms)=146
    CPU time spent (ms)=1720
    Physical memory (bytes) snapshot=448331776
    Virtual memory (bytes) snapshot=1466036224
    Total committed heap usage (bytes)=288358400
  Shuffle Errors
    BAD_ID=0
    CONNECTION=0
    IO_ERROR=0
    WRONG_LENGTH=0
    WRONG_MAP=0
    WRONG_REDUCE=0
  File Input Format Counters
    Bytes Read=5856303
  File Output Format Counters
    Bytes Written=7180
```

[illegible]