# Curso: Ciência de dados e Big Data

Professor: Cláudio Lúcio
Atividade Prática sobre Map Reduce

Neste caso vamos analisar uma implementação de map reduce no Hadoop.

Vamos utilizar os um exemplo de contagem de palavras já existente na sandbox da HortonWorks. Veja abaixo o programa que será executado (obviamente, implemetanto em Java):

1. Porção relativa ao Map

```java
package org.myorg;

import java.io.IOException;
import java.util.*;

import org.apache.hadoop.fs.Path;
import org.apache.hadoop.conf.*;
import org.apache.hadoop.io.*;
import org.apache.hadoop.mapreduce.*;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.input.TextInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;

public class WordCount {

 public static class Map extends Mapper<LongWritable, Text, Text, IntWritable> {
    private final static IntWritable one = new IntWritable(1);
    private Text word = new Text();

    public void map(LongWritable key, Text value, Context context) throws IOException,
InterruptedException {
        String line = value.toString();
        StringTokenizer tokenizer = new StringTokenizer(line);
        while (tokenizer.hasMoreTokens()) {
            word.set(tokenizer.nextToken());
            context.write(word, one);
        }
    }
 }
}
```

2. Porção relativa ao Reduce

```java
public static class Reduce extends Reducer<Text, IntWritable, Text, IntWritable> {

    public void reduce(Text key, Iterable<IntWritable> values, Context context)
      throws IOException, InterruptedException {
        int sum = 0;
        for (IntWritable val : values) {
            sum += val.get();
        }
        context.write(key, new IntWritable(sum));
    }
}
```

3. Porção relativa ao programa principal

```java
public static void main(String[] args) throws Exception {
    Configuration conf = new Configuration();

    Job job = new Job(conf, "wordcount");

    job.setOutputKeyClass(Text.class);
    job.setOutputValueClass(IntWritable.class);

    job.setMapperClass(Map.class);
    job.setReducerClass(Reduce.class);

    job.setInputFormatClass(TextInputFormat.class);
    job.setOutputFormatClass(TextOutputFormat.class);

    FileInputFormat.addInputPath(job, new Path(args[0]));
    FileOutputFormat.setOutputPath(job, new Path(args[1]));

    job.waitForCompletion(true);
    }

}
```
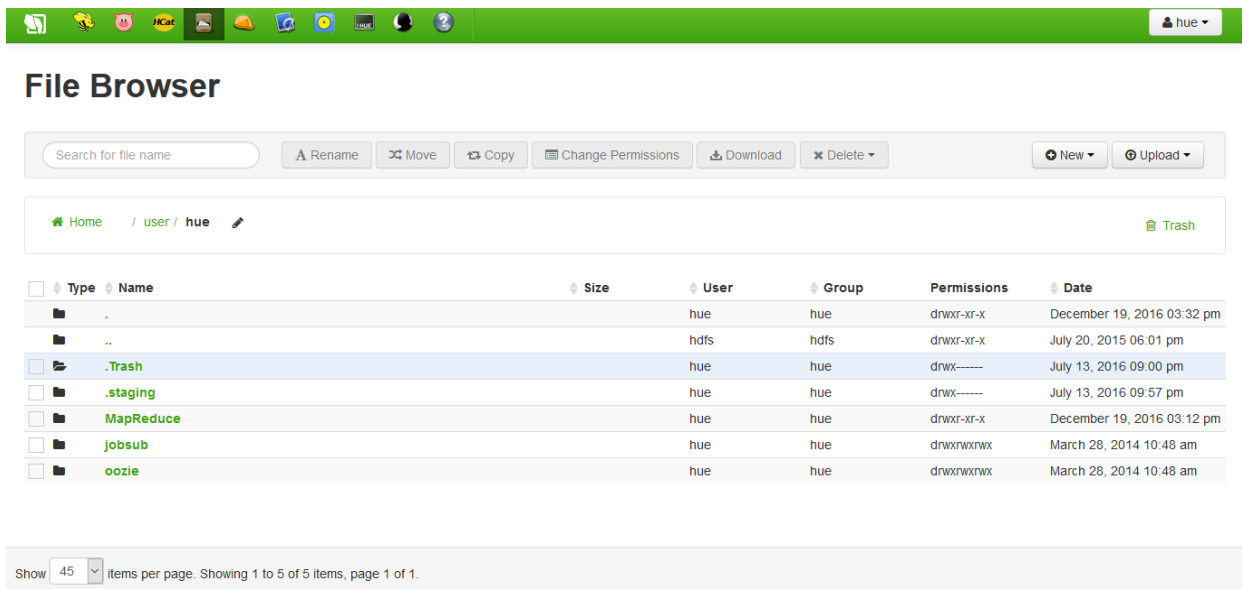
4. Para executar este programa vamos antes criar dois diretórios no usuário Hue, veja a tela a seguir:
   1. MapReduce



5. Agora utilize o arquivo 2.1 Textos.zip e faça seu upload para o HDFS.

6. Verifique se os arquivos foram todos carregados



7. Com os arquivos vamos abrir uma linha de comando:

```
su hdfs
hadoop  jar  /usr/hdp/2.3.2.0-2950/hadoop-mapreduce/hadoop-mapreduce-examples.jar  wordcount
'/user/hue/MapReduce/2.1 - textos/2.1 - textos/' /user/hue/MapReduceSaida/
```



Ainda em execução:

Finlalização da execução:

```
16/12/19 13:05:03 INFO mapreduce.Job:  map 94% reduce 31%
16/12/19 13:05:45 INFO mapreduce.Job:  map 95% reduce 31%
16/12/19 13:05:46 INFO mapreduce.Job:  map 95% reduce 32%
16/12/19 13:06:27 INFO mapreduce.Job:  map 96% reduce 32%
16/12/19 13:07:05 INFO mapreduce.Job:  map 97% reduce 32%
16/12/19 13:07:33 INFO mapreduce.Job:  map 98% reduce 32%
16/12/19 13:07:34 INFO mapreduce.Job:  map 98% reduce 33%
16/12/19 13:08:22 INFO mapreduce.Job:  map 99% reduce 33%
16/12/19 13:08:50 INFO mapreduce.Job:  map 100% reduce 33%
16/12/19 13:09:01 INFO mapreduce.Job:  map 100% reduce 100%
16/12/19 13:09:04 INFO mapreduce.Job: Job job_1482166846963_0001 completed succe
ssfully
16/12/19 13:09:05 INFO mapreduce.Job: Counters: 51
        File System Counters
                FILE: Number of bytes read=4578315
                FILE: Number of bytes written=108161357
                FILE: Number of read operations=0
                FILE: Number of large read operations=0
                FILE: Number of write operations=0
                HDFS: Number of bytes read=4285076
                HDFS: Number of bytes written=410680
                HDFS: Number of read operations=3003
                HDFS: Number of large read operations=0
```

```
                Failed map tasks=14
                Launched map tasks=1014
                Launched reduce tasks=1
                Other local map tasks=14
                Data-local map tasks=1000
                Total time spent by all maps in occupied slots (ms)=89570481
                Total time spent by all reduces in occupied slots (ms)=12373936
                Total time spent by all map tasks (ms)=89570481
                Total time spent by all reduce tasks (ms)=12373936
                Total vcore-seconds taken by all map tasks=89570481
                Total vcore-seconds taken by all reduce tasks=12373936
                Total megabyte-seconds taken by all map tasks=22392620250
                Total megabyte-seconds taken by all reduce tasks=3093484000
        Map-Reduce Framework
                Map input records=32937
                Map output records=787051
                Map output bytes=7238653
                Map output materialized bytes=4584309
                Input split bytes=160793
                Combine input records=787051
                Combine output records=358530
                Reduce input groups=36805
                Reduce shuffle bytes=4584309
```

```
                Reduce output records=36805
                Spilled Records=717060
                Shuffled Maps =1000
                Failed Shuffles=0
                Merged Map outputs=1000
                GC time elapsed (ms)=102170
                CPU time spent (ms)=594960
                Physical memory (bytes) snapshot=211428573184
                Virtual memory (bytes) snapshot=886697873408
                Total committed heap usage (bytes)=166730924032
        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=4124283
        File Output Format Counters
                Bytes Written=410680
[root@sandbox /]#
```

8. Veja agora os arquivos gerados pela execução do programa map reduce:



9. Veja o conteúdo dos dois arquivos gerados (principalemente part-r-0000):