# Lab 02: Map Reduce Programming

## 1. WordCount Program

Step 1: Program's solution

- Import:

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

- Mapper:

```java
  public static class TokenizerMapper
       extends Mapper<Object, Text, Text, IntWritable>{

    private final static IntWritable one = new IntWritable(1);
    private Text word = new Text();

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

- Reducer:

```java
  public static class IntSumReducer
       extends Reducer<Text,IntWritable,Text,IntWritable> {
    private IntWritable result = new IntWritable();

    public void reduce(Text key, Iterable<IntWritable> values,
                       Context context
```

```
                              ) throws IOException, InterruptedException {
      int sum = 0;
      for (IntWritable val : values) {
        sum += val.get();
      }
      result.set(sum);
      context.write(key, result);
    }
  }
```

- Main:

```java
  public static void main(String[] args) throws Exception {
    Configuration conf = new Configuration();
    Job job = Job.getInstance(conf, "word count");
    job.setJarByClass(WordCount.class);
    job.setMapperClass(TokenizerMapper.class);
    job.setCombinerClass(IntSumReducer.class);
    job.setReducerClass(IntSumReducer.class);
    job.setOutputKeyClass(Text.class);
    job.setOutputValueClass(IntWritable.class);
    FileInputFormat.addInputPath(job, new Path(args[0]));
    FileOutputFormat.setOutputPath(job, new Path(args[1]));
    System.exit(job.waitForCompletion(true) ? 0 : 1);
  }
}
```

- Explain:
- In the mapper method: This mapper will take as input an Object (representing the key), and a Text (representing the value) and use the StringTokenizer to separate the words in the value. Then it sends each word to the Reducer with a value of 1.
- In the reducer metho: This reducer will take the words from the Mapper and calculate the total number of occurrences of each word by adding the values of 1s together.

## Step 2: Class Creation
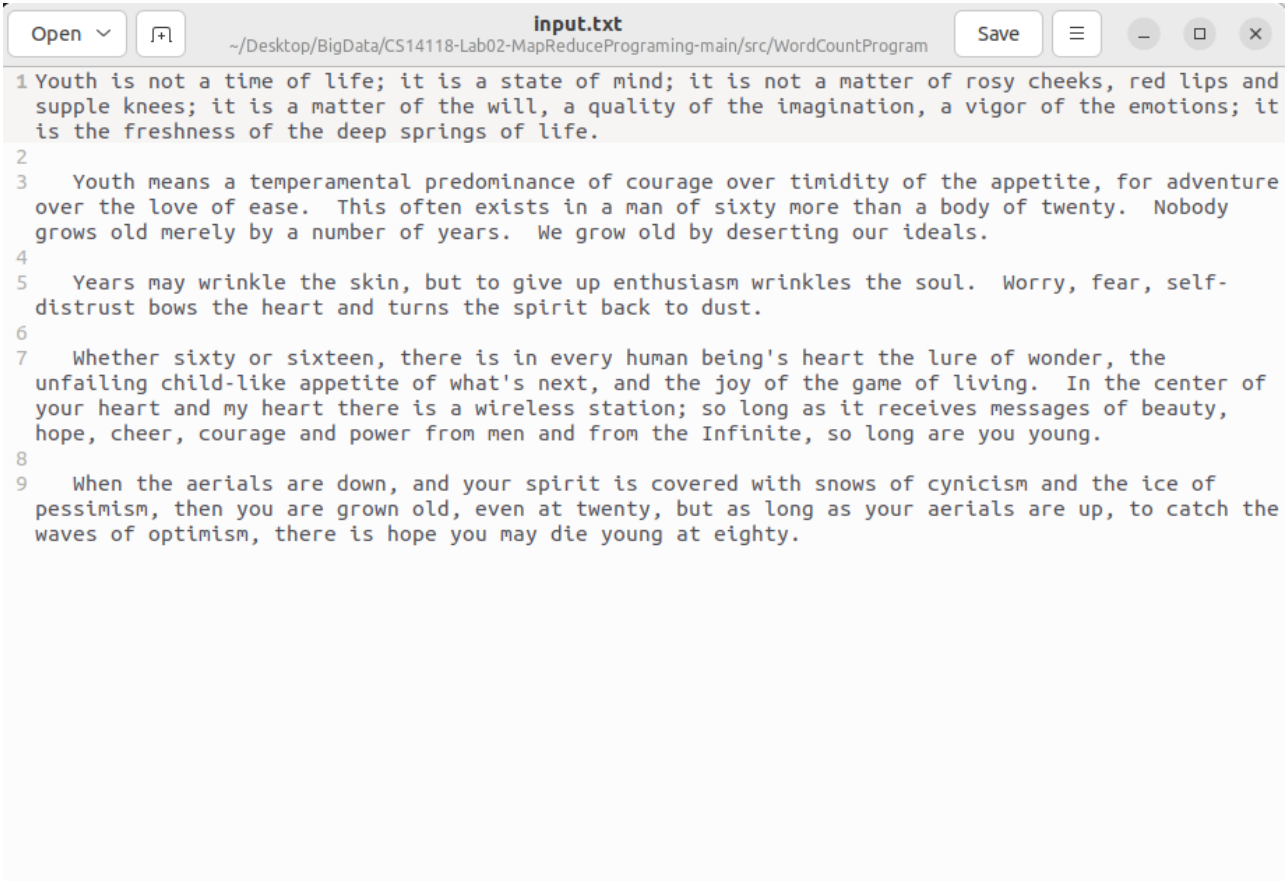
```
jar -cvf WordCount.jar -C classes/ .
```

## Step 3: Create directory structure for program in Hadoop

```
hadoop fs -mkdir /WordCount
hadoop fs -mkdir /WordCount/Input
hadoop fs -put 'local input file's path ' /WordCount/Input
```
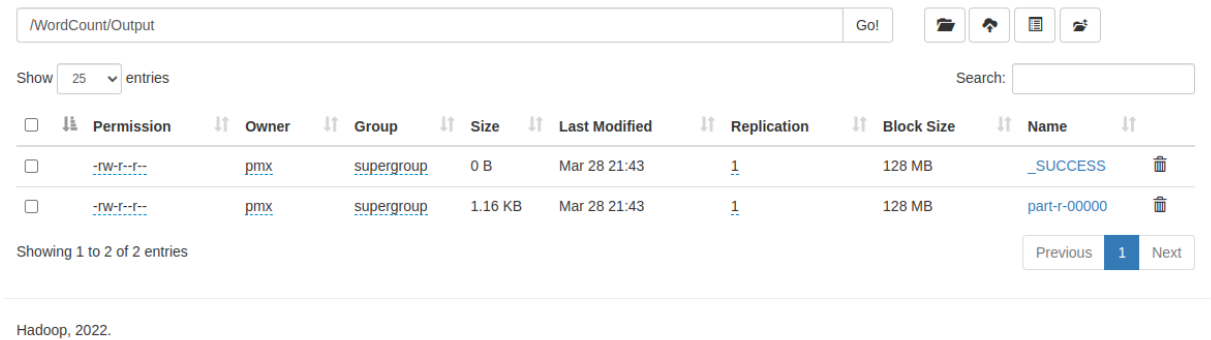
- Example input:

## Step 4: Create Jar File and deploy it to Hadoop

```
hadoop jar "Path to your local file .jar" WordCount /WordCount/Input
/WordCount/Output
```

## Step 5: Final result

- After succesfully calculating, we can check our result in HDFS like below:



Hadoop, 2022.

- 

---

## 2. WordSizeWordCount Program

### Step 1: Program's solution

- Import:

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

- Mapper:

```java
public static class Map extends Mapper<LongWritable, Text, Text,
IntWritable> {
```

```java
        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());
        String length = String.valueOf(word.getLength());
        Text len = new Text(length);
                    context.write(len, one);
                }
            }
        }
```

- Reducer:

```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));
        }
    }
```

- Main:

```java
  public static void main(String[] args) throws Exception {
        Configuration conf = new Configuration();
            Job job = new Job(conf, "WordSizeWordCount");
        job.setJarByClass(WordSizeWordCount.class);
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);
        job.setMapperClass(Map.class);
        job.setCombinerClass(Reduce.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);
    }
```

- Explain:
- In the mapper method: This mapper will take as input an Object (representing the key), and a Text (representing the value) and use the StringTokenizer to separate the words in the value. Then it sends each word's length to the Reducer with a value of 1.
- In the reducer metho: This reducer will take the word's length from the Mapper and calculate the total number of occurrences of each word's length by adding the values of 1s together.
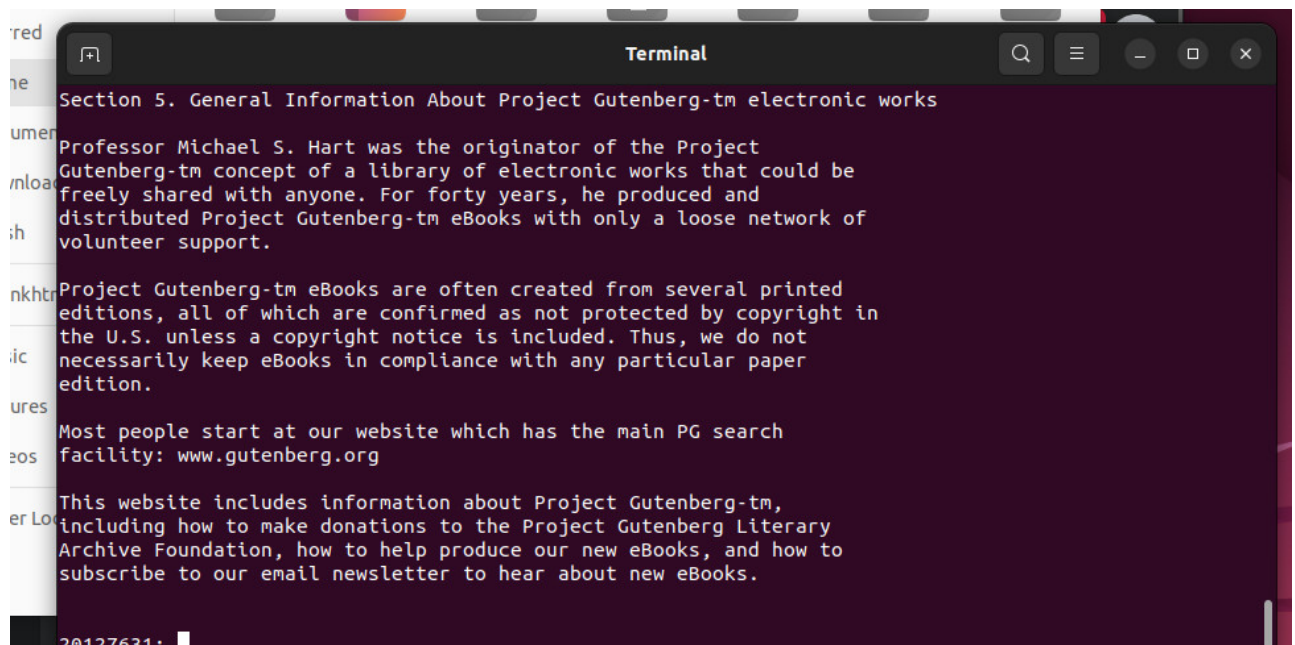
## Step 2: Class Creation

```
jar -cvf WordSizeWordCount.jar -C classes/ .
```

## Step 3: Create directory structure for program in Hadoop

```
hadoop fs -mkdir /WordSizeWordCount
hadoop fs -mkdir /WordSizeWordCount
hadoop fs -put 'local input file's path ' /WordSizeWordCount/Input
```

- Example input:



## Step 4: Create Jar File and deploy it to Hadoop

```
hadoop jar "Path to your local file .jar" WordSizeWordCount
/WordSizeWordCount/Input /WordSizeWordCount/Output
```

## Step 5: Final result

- After succesfully calculating, we can check our result in HDFS like below:



-



-

# 3. WeatherData program

## Step 1: Program's solution

- Import:

```
import java.io.IOException;
import java.util.Iterator;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
```

```
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;
```

- Mapper:

```java
public static class Map
        extends Mapper<Object, Text, Text, Text>{

    public void map(Object key, Text value, Context context
    ) throws IOException, InterruptedException {
        String line = value.toString();
        String[] tokens = line.split("\\s+");
        String date = tokens[1];
        float tempMax = Float.parseFloat(tokens[6].trim());
        float tempMin = Float.parseFloat(tokens[7].trim());

        if(tempMax > 40.0) {
            context.write(new Text("Hot Day " + date), new
Text(String.valueOf(tempMax)));
        }
        if(tempMin < 10.0) {
            context.write(new Text("Cold Day " + date), new
Text(String.valueOf(tempMin)));
        }

    }
}
```

- Reducer:

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

    public void reduce(Text key, Iterator<Text> values,
                        Context context
    ) throws IOException, InterruptedException {
        String temperature = values.next().toString();
        context.write(key, new Text(temperature));
    }
}
```

- Main:

```java
public static void main(String[] args) throws Exception {
        Configuration conf = new Configuration();
        Job job = Job.getInstance(conf, "weather");
        job.setJarByClass(WeatherData.class);
        job.setMapperClass(Map.class);
        job.setCombinerClass(Reduce.class);
        job.setReducerClass(Reduce.class);
        job.setOutputKeyClass(Text.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);
    }
```

- Explain:
- In the mapper method: This mapper will take as input an Object (representing the key), and a Text (representing the value) and use the `split()` method to separate the fields in a line (inside is a regular expression indicate one or more whitespace). It then extract the date, minTemp and maxTemp (at index 1, 6, 7 respectively). Then, it check for the conditions (tempMax > 40 or tempMin < 10) and write the coressponding output to the context.
- In the reducer metho: This reducer will take the words from the Mapper and print out the result.

## Step 2: Class Creation

```
jar -cvf WeatherData.jar -C classes/ .
```

## Step 3: Create directory structure for program in Hadoop

```
hadoop fs -mkdir /Weather
hadoop fs -mkdir /Weather/Input
hadoop fs -put 'local input file's path ' /Weather/Input
```

- Example input:

## Step 4: Create Jar File and deploy it to Hadoop

```
hadoop jar "Path to your local file .jar" WeatherData /Weather/Input
/Weather/Output
```

## Step 5: Final result

- After succesfully calculating, we can check our result in HDFS like below:

# 4. Patent Program

Step 1: Program's solution

- Mapper:

```java
public static class PatentMapper
        extends Mapper<Object, Text, Text, Text> {
    Text k = new Text();
    Text v = new Text();

    public void map(Object key, Text value, Context context) throws
```

```java
IOException, InterruptedException {
        String line = value.toString();

        StringTokenizer tokenizer = new StringTokenizer(line, " ");
        while (tokenizer.hasMoreTokens()) {
            String token = tokenizer.nextToken();
            k.set(token);
            String token1 = tokenizer.nextToken();
            v.set(token1);
            context.write(k, v);
        }
    }
}
```

- Reducer:

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

    public void reduce(Text key, Iterable<Text> values,
                        Context context) throws IOException,
InterruptedException {
        int sum = 0;
        for (Text x : values) {
            sum++;
        }
        String result = Integer.toString(sum);
        context.write(key, new Text(result));
    }
}
```

- Main:

```java
public static void main(String[] args) throws Exception {
    Configuration conf = new Configuration();
    Job job = Job.getInstance(conf, "patent program");
    job.setJarByClass(PatentProgram.class);

    job.setMapperClass(PatentMapper.class);
    job.setReducerClass(SumSubPatentReducer.class);

    job.setOutputKeyClass(Text.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);
}
```

- The main idea for this program is that collecting pair of token in Map function, after combining, we count them through the their key and write in output file.

## Step 2: Class Creation

- After complete code in Java, we need to generate file jar from builded classes by below command:

```
jar -cvf PatentProgram.jar -C classes/ .
```

- Notice: Make sure that you export HADOOP_CLASSPATH before buiding file jar

## Step 3: Create directory structure for program in Hadoop

- We need to create folder to store input data in HDFS by below command:

```
hadoop fs -mkdir /PatentProgram
hadoop fs -mkdir /PatentProgram/Input
hadoop fs -put "local input file's path" /PatentProgram/Input
```

- Example input:



-

## Step 4: Create Jar File and deploy it to Hadoop

```
hadoop jar "Path to your local file .jar" PatentProgram
/PatentProgram/Input /PatentProgram/Output
```

## Step 5: Final result

- After succesfully calculating, we can check our result in HDFS like below:



- Hadoop, 2022.



---

# 5. MaxTemp Program

## Step 1: Program's solution

- Import:

```
import java.io.IOException;
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;
```

- Mapper:

```
public static class MaxTemperatureMapper
    extends Mapper<Object, Text, Text, IntWritable>{

  private static final int MISSING = 9999;

  public void map(Object key, Text value, Context context
                  ) throws IOException, InterruptedException {
    String line = value.toString();
    String year = line.substring(0, 4);
    int airTemperature = Integer.parseInt(line.substring(5));
    context.write(new Text(year), new IntWritable(airTemperature));
  }
}
```

- Reducer:

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

  public void reduce(Text key, Iterable<IntWritable> values,
                     Context context
                     ) throws IOException, InterruptedException {
    int maxTemperature = Integer.MIN_VALUE;
    for (IntWritable value : values) {
      maxTemperature = Math.max(maxTemperature, value.get());
    }
    context.write(key, new IntWritable(maxTemperature));
  }
}
```

- Main:

```java
    public static void main(String[] args) throws Exception {
        Configuration conf = new Configuration();
        Job job = Job.getInstance(conf, "max temperature");
        job.setJarByClass(MaxTemp.class);
        job.setMapperClass(MaxTemperatureMapper.class);
        job.setCombinerClass(MaxTemperatureReducer.class);
        job.setReducerClass(MaxTemperatureReducer.class);
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);
        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));
        System.exit(job.waitForCompletion(true) ? 0 : 1);
    }
}
```

- Explain:
- In the mapper method, we extract the year and temperature from each input line and write them to the key/value pair. We do not need to verify the format of the input stream because in this case all the lines have the same format and we can simply use fixed indexes to extract the information.
- In the reduce method, we find the highest temperature for each year by traversing the list of pooled values for the same key.

## Step 2: Class Creation

```
jar -cvf MaxTemp.jar -C classes/ .
```

## Step 3: Create directory structure for program in Hadoop

```
hadoop fs -mkdir /MaxTemp
hadoop fs -mkdir /MaxTemp/Input
hadoop fs -put 'local input file's path ' /MaxTemp/Input
```

- Example input:

## Step 4: Create Jar File and deploy it to Hadoop

```
hadoop jar "Path to your local file .jar" MaxTemp /MaxTemp/Input
/MaxTemp/Output
```

## Step 5: Final result

- After succesfully calculating, we can check our result in HDFS like below:

Block information --  Block 0 ▾

Block ID: 1073741991

Block Pool ID: BP-1926325839-127.0.1.1-1678203211420

Generation Stamp: 1167

Size: 24

Availability:

- pmx

File contents

```
1900    36
1901    48
1902    49
```

-

# 6. AverageSalary Program

## Step 1: Program's solution

- Mapper:

```java
    public static class AvgMapper
            extends Mapper<Object, Text, Text, FloatWritable> {

        private Text id = new Text();
        private FloatWritable salary = new FloatWritable();

        public void map(Object key, Text value, Context context) throws
IOException, InterruptedException {
            String[] values = value.toString().split("\t");
            id.set(values[0]);

            salary.set(Float.parseFloat(values[2]));
            context.write(id, salary);
        }
    }
```

- Reducer

```java
    public static class AvgReducer
            extends Reducer<Text, FloatWritable, Text, FloatWritable> {

        private FloatWritable result = new FloatWritable();

        public void reduce(Text key, Iterable<FloatWritable> values,
                            Context context) throws IOException,
InterruptedException {
            float totalSalary = 0;
            int numberPersons = 0;
            for (FloatWritable salary : values) {
                totalSalary += salary.get();
                numberPersons++;
            }

            result.set(totalSalary/numberPersons);
            context.write(key, result);
        }
    }
```

- Main

```java
    public static void main(String[] args) throws Exception {
        Configuration conf = new Configuration();
        Job job = Job.getInstance(conf, "average salary");
        job.setJarByClass(AverageSalary.class);

        job.setMapperClass(AvgMapper.class);
        job.setReducerClass(AvgReducer.class);

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

        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));
        System.exit(job.waitForCompletion(true) ? 0 : 1);
    }
```

- In Map function, we will collect employee's ID with their salary to make a pair. Then in Reducer, we will take average salary of each employee's ID to write in the output.

## Step 2: Class Creation

- After complete code in Java, we need to generate file jar from builded classes by below command:

```
jar -cvf AverageSalary.jar -C classes/ .
```

- Notice: Make sure that you export HADOOP_CLASSPATH before buiding file jar

## Step 3: Create directory structure for program in Hadoop

- We need to create folder to store input data in HDFS by below command:

```
hadoop fs -mkdir /AverageSalary
hadoop fs -mkdir /AverageSalary/Input
hadoop fs -put "local input file's path" /AverageSalary/Input
```

- Example input:



- 

## Step 4: Create Jar File and deploy it to Hadoop

```
hadoop jar "Path to your local file .jar" AverageSalary
/AverageSalary/Input /AverageSalary/Output
```

## Step 5: Final result

- After succesfully calculating, we can check our result in HDFS like below:



- 



---

# 7. De Identify HealthCare Program

## Step 3: Program's solution

- Mapper

```java
public static Integer[] encryptCol = {2, 3, 4, 5, 6, 7, 8};
private static byte[] key1 = new String("sampleKey1234567").getBytes();
```

```java
    public static class Map
            extends Mapper<Object, Text, NullWritable, Text> {

        public void map(Object key, Text value, Context context) throws
IOException, InterruptedException {
            StringTokenizer tokenizer = new
StringTokenizer(value.toString(), ",");
            List<Integer> list = new ArrayList<>();

            Collections.addAll(list, encryptCol);
            // list = {2, 3, 4, 5, 6, 7, 8}

            System.out.println("Mapper :: one" + value);
            String newStr = "";

            int counter = 1;

            while (tokenizer.hasMoreTokens()) {
                String token = tokenizer.nextToken();
                System.out.println("token" + token);
                System.out.println("i=" + counter);

                if (list.contains(counter)) {
                    if (newStr.length() > 0) {
                        newStr += ",";
                    }
                    newStr += encrypt(token, key1);
                }
                else {
                    if (newStr.length() > 0) {
                        newStr += ",";
                    }
                    newStr += token;
                }
                counter += 1;
            }

            context.write(NullWritable.get(), new Text(newStr.toString()));
        }
    }
```

- Encrypt function

```java
    public static String encrypt(String strToEncrypt, byte[] key)
    {
        try
        {
            Cipher cipher = Cipher.getInstance("AES/ECB/PKCS5Padding");
            SecretKeySpec secretKey = new SecretKeySpec(key, "AES");
            cipher.init(Cipher.ENCRYPT_MODE, secretKey);

            String encryptedString =
```

```
Base64.encodeBase64String(cipher.doFinal(strToEncrypt.getBytes())));

            return encryptedString.trim();
        }
        catch (Exception e)
        {
            logger.error("Error while encrypting", e);
        }
        return null;
    }
```

- Main

```java
    public static void main(String[] args) throws Exception {
        if (args.length != 2) {
            System.out.println("usage: [input] [output]");
            System.exit(-1);
        }

        Configuration conf = new Configuration();
        Job job = Job.getInstance(conf, "de identify data");
        job.setMapperClass(Map.class);

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

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

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

        job.setJarByClass(DeIdentifyData.class);
        job.waitForCompletion(true);
    }
```

- The idea to resolve this question is only using Map function and encrypt function to encrypt data in identified columns which need to be hidden.

## Step 2: Class Creation

- After complete code in Java, we need to generate file jar from builded classes by below command:

```
jar -cvf DeIdentifyData.jar -C classes/ .
```

- Notice: Make sure that you export HADOOP_CLASSPATH before buiding file jar

## Step 3: Create directory structure for program in Hadoop

- We need to create folder to store input data in HDFS by below command:

```
hadoop fs -mkdir /DeIdentifyData
hadoop fs -mkdir /DeIdentifyData/Input
hadoop fs -put "local input file's path" /DeIdentifyData/Input
```

- Example input:



- 

## Step 4: Create Jar File and deploy it to Hadoop

```
hadoop jar "Path to your local file .jar" DeIdentifyData
/DeIdentifyData/Input /DeIdentifyData/Output
```

## Step 5: Final result

- After succesfully calculating, we can check our result in HDFS like below:

Hadoop    Overview    Datanodes    Datanode Volume Failures    Snapshot    Startup Progress    Utilities ▾

## Browse Directory

/DeIdentifyData/Output                                                   Go!

Show 25 ⌄ entries                                                          Search:

| | Permission | Owner | Group | Size | Last Modified | Replication | Block Size | Name | |
|---|---|---|---|---|---|---|---|---|---|
| ☐ | -rw-r--r-- | antrvan | supergroup | 0 B | Mar 23 23:28 | 1 | 128 MB | _SUCCESS | 🗑 |
| ☐ | -rw-r--r-- | antrvan | supergroup | 1.08 KB | Mar 23 23:28 | 1 | 128 MB | part-r-00000 | 🗑 |

Showing 1 to 2 of 2 entries                                   Previous  1  Next

Hadoop, 2022.

Generation Stamp: 1265

Size: 1104

Availability:

- X390

### File contents

11116,YK6bf5BMSnSEeRslpZkTPw==,mrerWhav8woTBBAfVhMCTA==,obp6xFsHvexdv5j4u7
e5ow==,9lN3Gt0BtJBEiv03ARZ12g==,ZEZeGclrL4JQAO86Yocw7w==,Rqc+URijpA/g34bStvD
H8g==,xlGEUUwNfDh9UU5PBLxKMQ==,84
11115,uF0JIaEJwBqzXvFHI2Z/uQ==,9AyNyrJ7W/y0IbEM/adtYg==,obp6xFsHvexdv5j4u7e5o
w==,rOwfjOToY83UsIFJ8u0RFw==,ZEZeGclrL4JQAO86Yocw7w==,JLZyVcPdHSf+7loGsPrH
Ug==,TQGMCgOz1/d1Iol5HCygnQ==,76
11114,dFdQxegqy/HKG8ur1n7wVg==,61JPlFwoYKrcwpSPGx1aMw==,obp6xFsHvexdv5j4u7e
5ow==,zLVYfyf2zPbG/GbGP/ihxQ==,ZEZeGclrL4JQAO86Yocw7w==,Rqc+URijpA/g34bStvDH

Close

---

# 8 Music Track Program

## Step 1: Program's solution

- task1: Number of unique listeners

- Import:

```
import java.util.*;
import java.io.IOException;
import java.util.StringTokenizer;
import org.apache.hadoop.conf.Configuration;
```

```
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapred.JobClient;
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.input.TextInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;
```

- Mapper:

```
  public static class Map extends Mapper<Object, Text, IntWritable,
IntWritable> {
        public void map(Object key, Text value, Context context) throws
IOException, InterruptedException {
            String line = value.toString();
            String[] data = line.split(",");
            IntWritable user = new
IntWritable(Integer.parseInt(data[UserId]));
            IntWritable track = new
IntWritable(Integer.parseInt(data[TrackId]));
            context.write(track, user);


        }
    }
```

- Reducer:

```
  public static class Reduce extends Reducer<IntWritable, IntWritable,
IntWritable, IntWritable> {
        public void reduce(IntWritable key, Iterable<IntWritable> values,
Context context) throws IOException, InterruptedException {
            Set<Integer> users = new HashSet<Integer>();
            for (IntWritable val : values) {
                users.add(val.get());
            }
            IntWritable result = new IntWritable(users.size());
            context.write(key, result);
        }
    }
```

- Main:

```java
    public static void main(String[] args) throws Exception {
        Configuration conf = new Configuration();
        Job job = new Job(conf, "Listener");
        job.setJarByClass(Listener.class);
        job.setOutputKeyClass(IntWritable.class);
        job.setOutputValueClass(IntWritable.class);
        job.setMapperClass(Map.class);
        job.setCombinerClass(Reduce.class);
        job.setReducerClass(Reduce.class);
        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));
        job.setInputFormatClass(TextInputFormat.class);
        job.setOutputFormatClass(TextOutputFormat.class);
        job.waitForCompletion(true);
    }
```

- Explain:
- In the mapper method: This mapper will take as input an Object (representing the key), and a Text (representing the value) and use the String[] to separate the data in the value. Then it sends each pair<trackId,userId> to the Reducer.
- In the reducer method: This reducer will take the pair<trackId,userId> from the Mapper and add it into Set(HashSet) then return <key,Set.size()>.

- task2: Number of times the track was shared with others

- Import:

```java
import java.util.*;
import java.io.IOException;
import java.util.StringTokenizer;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapred.JobClient;
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.input.TextInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;
```

- Mapper:

```java
    public static class Map extends Mapper<Object, Text, Text, IntWritable> {
        private Text track = new Text();
```

```java
        public void map(Object key, Text value, Context context) throws
IOException, InterruptedException {
            String line = value.toString();
            String[] data = line.split(",");
            track = new Text(data[TrackId]);
            context.write(track, new
IntWritable(Integer.valueOf(data[Shared])));
        }
    }
```

- Reducer:

```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 = sum + val.get();
            }
            context.write(key, new IntWritable(sum));
        }
    }
```

- Main:

```java
  public static void main(String[] args) throws Exception {
        Configuration conf = new Configuration();
        Job job = new Job(conf, "Shared");
        job.setJarByClass(Shared.class);
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);
        job.setMapperClass(Map.class);
        job.setCombinerClass(Reduce.class);
        job.setReducerClass(Reduce.class);
        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));
        job.setInputFormatClass(TextInputFormat.class);
        job.setOutputFormatClass(TextOutputFormat.class);
        job.waitForCompletion(true);
    }
```

- Explain:
- In the mapper method: This mapper will take as input an Object (representing the key), and a Text (representing the value) and use the String[] to separate the data in the value. Then it sends each pair<trackId,Shared> to the Reducer.

- In the reducer method: This reducer will take the pair<trackId,Shared> from the Mapper and calculate the total number of each track was shared with orthers by adding the values of 1s together..

- task3: Number of times the track was listened to on the radio

- Import:

```
import java.util.*;
import java.io.IOException;
import java.util.StringTokenizer;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapred.JobClient;
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.input.TextInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;
```

- Mapper:

```
    public static class Map extends Mapper<Object, Text, Text, IntWritable>
{
        private Text track = new Text();

        public void map(Object key, Text value, Context context) throws
IOException, InterruptedException {
            String line = value.toString();
            String[] data = line.split(",");
            track = new Text(data[TrackId]);
            context.write(track, new
IntWritable(Integer.valueOf(data[Radio])));
        }
    }
```

- Reducer:

```
  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 = sum + val.get();
            }
            context.write(key, new IntWritable(sum));
        }
    }
```

- Main:

```
public static void main(String[] args) throws Exception {
        Configuration conf = new Configuration();
        Job job = new Job(conf, "Radio");
        job.setJarByClass(Radio.class);
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);
        job.setMapperClass(Map.class);
        job.setCombinerClass(Reduce.class);
        job.setReducerClass(Reduce.class);
        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));
        job.setInputFormatClass(TextInputFormat.class);
        job.setOutputFormatClass(TextOutputFormat.class);
        job.waitForCompletion(true);
    }
```

- Explain:
- In the mapper method: This mapper will take as input an Object (representing the key), and a Text (representing the value) and use the String[] to separate the data in the value. Then it sends each pair<trackId,Radio> to the Reducer.
- In the reducer method: This reducer will take the pair<trackId,Radio> from the Mapper and calculate the total number of each track was listened to on radio by adding the values of 1s together..

- task4: Number of times the track was listened to in total

- Import:

```
import java.util.*;
import java.io.IOException;
import java.util.StringTokenizer;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapred.JobClient;
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.input.TextInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;
```

- Mapper:

```java
   public static class Map extends Mapper<Object, Text, Text, IntWritable>
{
        private Text track = new Text();

        public void map(Object key, Text value, Context context) throws
IOException, InterruptedException {
            String line = value.toString();
            String[] data = line.split(",");
            track = new Text(data[TrackId]);
            context.write(track, new
IntWritable(Integer.valueOf(data[Skip])));
        }
    }
```

- Reducer:

```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) {
                if (val.get() == 0)
                    sum++;
            }
            context.write(key, new IntWritable(sum));
        }
    }
```

- Main:

```java
  public static void main(String[] args) throws Exception {
        Configuration conf = new Configuration();
        Job job = new Job(conf, "Total");
        job.setJarByClass(Total.class);
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);
        job.setMapperClass(Map.class);
        job.setCombinerClass(Reduce.class);
```

```
        job.setReducerClass(Reduce.class);
        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));
        job.setInputFormatClass(TextInputFormat.class);
        job.setOutputFormatClass(TextOutputFormat.class);
        job.waitForCompletion(true);
    }
```

- Explain:
- In the mapper method: This mapper will take as input an Object (representing the key), and a Text (representing the value) and use the String[] to separate the data in the value. Then it sends each pair<trackId,Skip> to the Reducer.
- In the reducer method: This reducer will take the pair<trackId,Skip> from the Mapper and calculate the total number of each track was'n skipped by adding the values of 1s together..
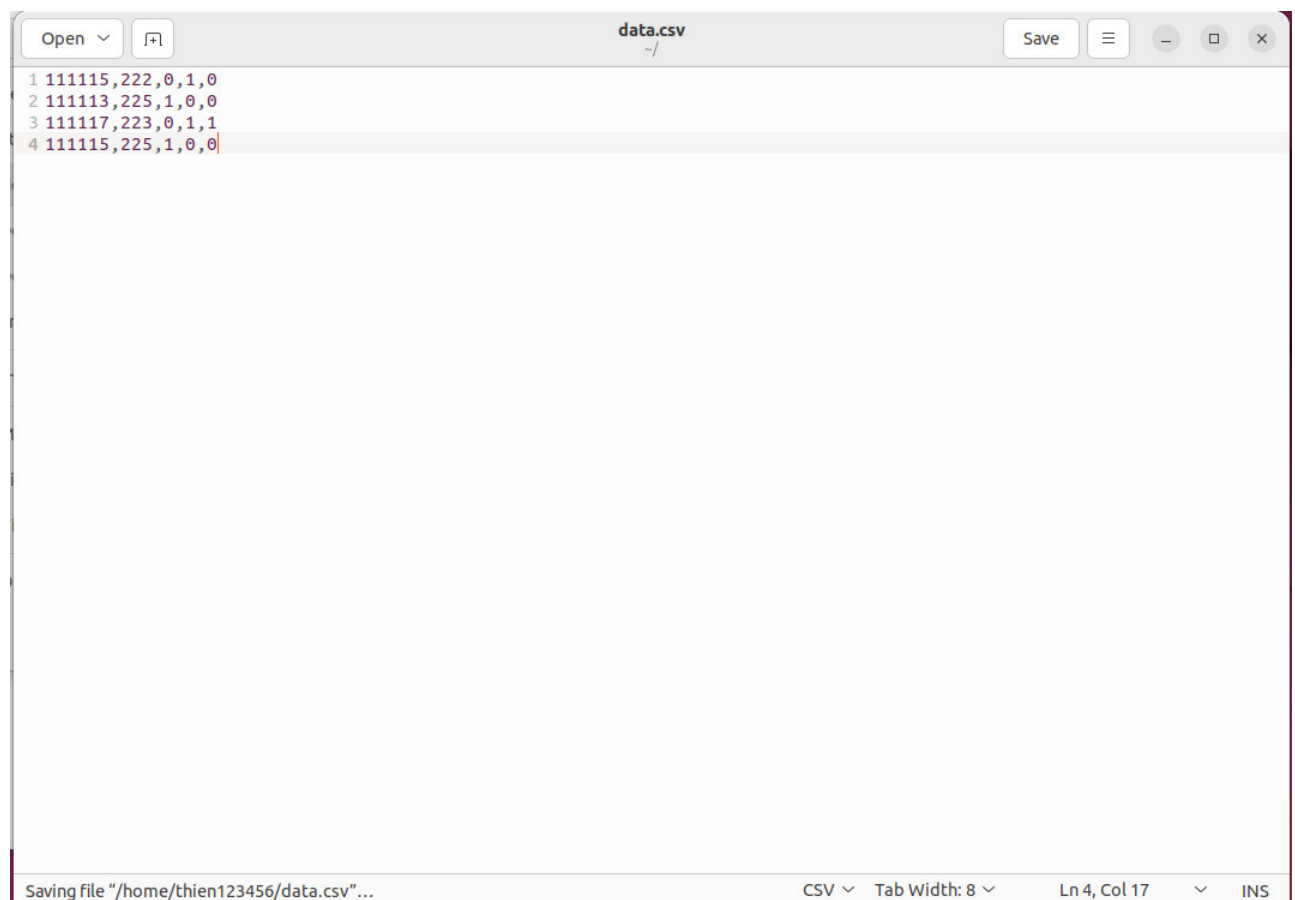
- task5: Number of times the track was skipped on the radio

- Import:

```
import java.util.*;
import java.io.IOException;
import java.util.StringTokenizer;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapred.JobClient;
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.input.TextInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;
```

- Mapper:

```
   public static class Map extends Mapper<Object, Text, Text, IntWritable>
{
        private Text track = new Text();

        public void map(Object key, Text value, Context context) throws
IOException, InterruptedException {
            String line = value.toString();
            String[] data = line.split(",");
            track = new Text(data[TrackId]);
            if (data[Radio].equals("0") || data[Skip].equals("0"))
                context.write(track, new IntWritable(0));
```

```
            else
                context.write(track, new IntWritable(1));

        }
    }
```

- Reducer:

```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 = sum + val.get();
            }
            context.write(key, new IntWritable(sum));
        }
    }
```

- Main:

```java
    public static void main(String[] args) throws Exception {
        Configuration conf = new Configuration();
        Job job = new Job(conf, "Skip_Radio");
        job.setJarByClass(Skip_Radio.class);
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);
        job.setMapperClass(Map.class);
        job.setCombinerClass(Reduce.class);
        job.setReducerClass(Reduce.class);
        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));
        job.setInputFormatClass(TextInputFormat.class);
        job.setOutputFormatClass(TextOutputFormat.class);
        job.waitForCompletion(true);
    }
```

- Explain:
- In the mapper method: This mapper will take as input an Object (representing the key), and a Text (representing the value) and use the String[] to separate the data in the value. Then it sends each pair<trackId,Skip&Radio> to the Reducer.
- In the reducer method: This reducer will take the pair<trackId,Skip&Radio> from the Mapper and calculate the total number of each track was skipped on the radio by adding the values of 1s together..

## Step 2: Class Creation

```
jar -cvf fileName.jar -C classes/ .
```

## Step 3: Create directory structure for program in Hadoop

```
hadoop fs -mkdir /fileName
hadoop fs -mkdir /fileName/Input
hadoop fs -put 'local input file's path ' /fileName/Input
```

- Example input:



## Step 4: Create Jar File and deploy it to Hadoop

```
hadoop jar "Path to your local file .jar" WordCount /WordCount/Input
/WordCount/Output
```

## Step 5: Final result

- After succesfully calculating, we can check our result in HDFS like below:

- task1:

-


-

- task2:

- task3:

- 

- task4:



-

- task5:



-

---

# 9. Telecom Call Data Record Program

## Step 1: Program's solution

- Import:

```java
import java.io.IOException;
import java.text.ParseException;
import java.text.SimpleDateFormat;
import java.util.Date;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.LongWritable;
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;
```

- Mapper:

```java
public static class Map
        extends Mapper<Object, Text, Text, LongWritable>{

    Text phoneNumber = new Text();
    LongWritable minutes = new LongWritable();
```

```java
        public void map(Object key, Text value, Context context
        ) throws IOException, InterruptedException {
            String line = value.toString();
            String[] tokens = line.split("\\|");
            if(tokens[4].equals("1")) {
                phoneNumber.set(tokens[0]);
                minutes.set(calculateTimeInMinutes(tokens[2], tokens[3]));
                context.write(phoneNumber, minutes);
            }
        }

        private long calculateTimeInMinutes(String start, String end) {
            SimpleDateFormat formatter = new SimpleDateFormat(("yyyy-MM-dd
HH:mm:ss"));
            long minutes = -1; // if this value happen then there's an
error
            try {
                // put code in try catch so that java is not angry
                Date startDate = formatter.parse(start);
                Date endDate = formatter.parse(end);
                long duration = endDate.getTime() - startDate.getTime();
                minutes = duration / (1000 * 60);
            } catch (ParseException e) {
                e.printStackTrace();
            }
            return minutes;
        }
    }
```

- Reducer:

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

    public void reduce(Text key, Iterable<LongWritable> values,
                        Context context
    ) throws IOException, InterruptedException {
        long totalMinutes = 0;
        for(LongWritable val: values) {
            totalMinutes += val.get();
        }
        if(totalMinutes > 60) {
            context.write(key, new LongWritable(totalMinutes));
        }


    }
}
```

- Main:

```java
public static void main(String[] args) throws Exception {
        Configuration conf = new Configuration();
        Job job = Job.getInstance(conf, "call data record");
        job.setJarByClass(CallDataRecord.class);
        job.setMapperClass(Map.class);
        job.setCombinerClass(Reduce.class);
        job.setReducerClass(Reduce.class);
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(LongWritable.class);
        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));
        System.exit(job.waitForCompletion(true) ? 0 : 1);
    }
```

- Explain:
- In the mapper method: This mapper will take as input an Object (representing the key), and a Text (representing the value) and use the `split()` method to separate the fields in a line (inside is a regular expression indicate a "|" character). It then check if the std is equal to 1. If true, it will write the FromPhoneNumber and the time in minute (calculate using the calculateTimeInMinutes() utility function) to the context.
- In the reducer metho: This reducer will take the words from the Mapper and add all the values together, it then write the result to the context if the total is greater than 60.
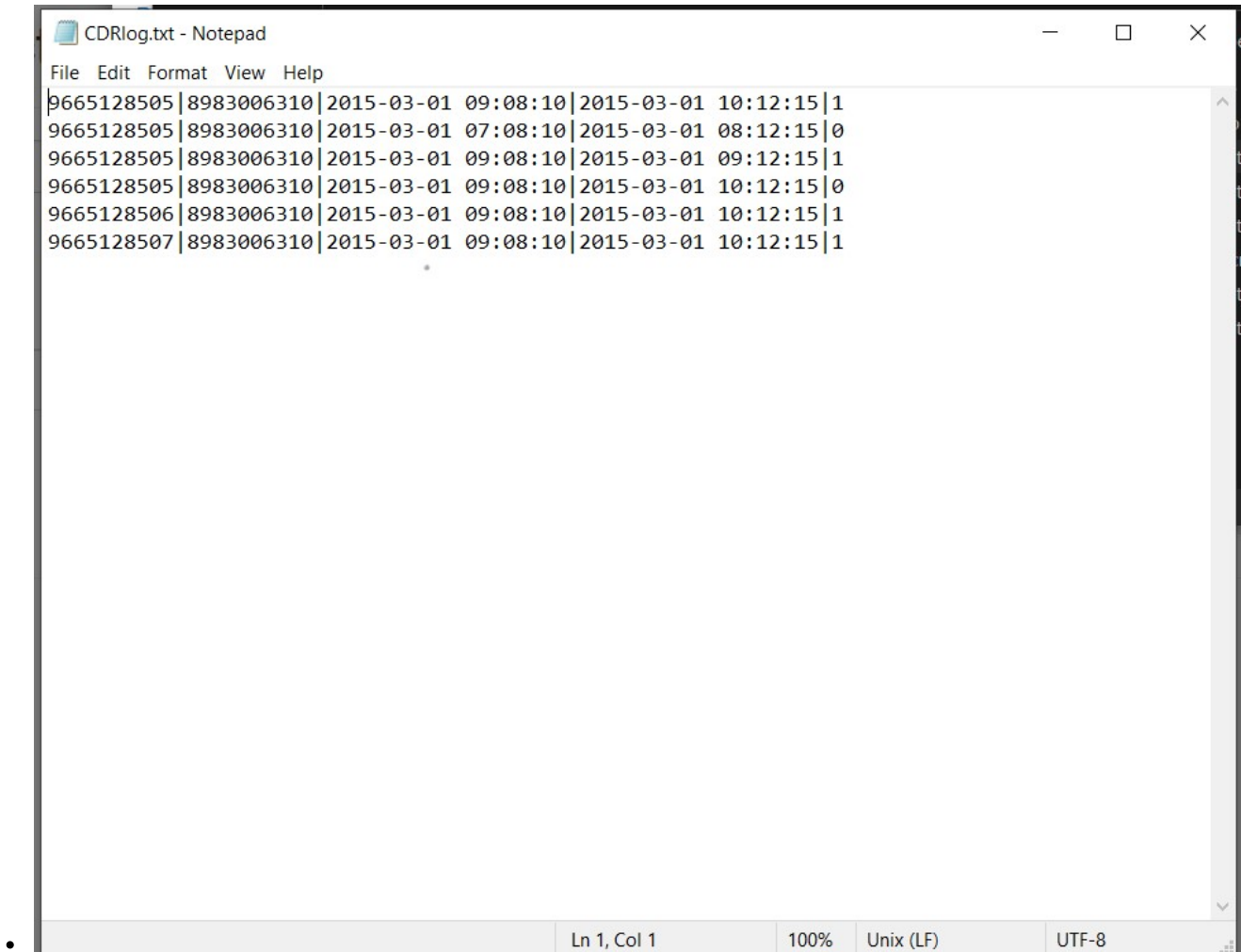
## Step 2: Class Creation

```
jar -cvf CallDataRecord.jar -C classes/ .
```

## Step 3: Create directory structure for program in Hadoop

```
hadoop fs -mkdir /Phone
hadoop fs -mkdir /Phone/Input
hadoop fs -put 'local input file's path ' /Phone/Input
```

- Example input:

## Step 4: Create Jar File and deploy it to Hadoop

```
hadoop jar "Path to your local file .jar" CallDataRecord /Phone/Input
/Phone/Output
```

## Step 5: Final result

- After succesfully calculating, we can check our result in HDFS like below:



Hadoop, 2022.

- 

---

# 10. Count Connected Components Program

Step 1: Program's solution

- Mapper

```
public static class Map
        extends Mapper<Object, Text, Text, Text> {

    public void map(Object key, Text value, Context context
    ) throws IOException, InterruptedException {
        String[] tokens = value.toString().split(" ");
```

```
            String keyValue = tokens[0];
            Arrays.sort(tokens);

            int i = 0;
            while (i < tokens.length) {
                context.write(new Text("map"), new Text(keyValue + "," +
tokens[i]));
                i++;
            }
        }
    }
```

- Reducer

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

        public void reduce(Text key, Iterable<Text> values,
                           Context context
        ) throws IOException, InterruptedException {

            TreeMap<Integer, ArrayList<Integer>> sortedMap = new TreeMap<>
();
            HashMap<Integer, Integer> result = new HashMap<>();
            for (Text value : values) {
                String[] pair = value.toString().split(",");

                int keyItem = Integer.parseInt(pair[0]);
                int valueItem = Integer.parseInt(pair[1]);

                ArrayList<Integer> tmp = sortedMap.getOrDefault(keyItem,
new ArrayList<Integer>());
                tmp.add(valueItem);
                Collections.sort(tmp);
                sortedMap.put(keyItem, tmp);
            }

            for (Integer k : sortedMap.keySet()) {
                Integer start = sortedMap.get(k).get(0);
                if (start.compareTo(k) == 0) {
                    result.put(k, k);
                }
                if (start.compareTo(k) < 0) {
                    result.put(k, result.get(start));
                    for (Integer v : sortedMap.get(k)) {
                        if (v.equals(start)) continue;
                        for (Integer j : result.keySet()) {
                            if (result.get(j).equals(v)) {
                                result.replace(j, v, start);
                            }
                        }
```

```
                }
            }
        }

        HashSet<Integer> components = new HashSet<>();
        boolean b = components.addAll(result.values());

        if (b) {
            context.write(new Text(""), new
Text(String.valueOf(components.size())));
        }
    }
}
```

- Main

```
    public static void main(String[] args) throws Exception {
        Configuration conf = new Configuration();
        Job job = Job.getInstance(conf, "count connected component
program");

        job.setJarByClass(CountConnectedComponentProgram.class);

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

        job.setOutputKeyClass(Text.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);
    }
```

- This task is an intriguing question that calculate numbers of separated components in a graph. To
  resolve this problem, we put pair of source and destination point of every edges in graph to reducer.
  We put all pairs to TreeMap to sort them. Then in each components, we mark all connected vertices
  value to smallest vertex. Finally, the result equals numbers of different values in HashMap.

## Step 2: Class Creation

- After complete code in Java, we need to generate file jar from builded classes by below command:

```
jar -cvf CountConnectedComponentProgram.jar -C classes/ .
```
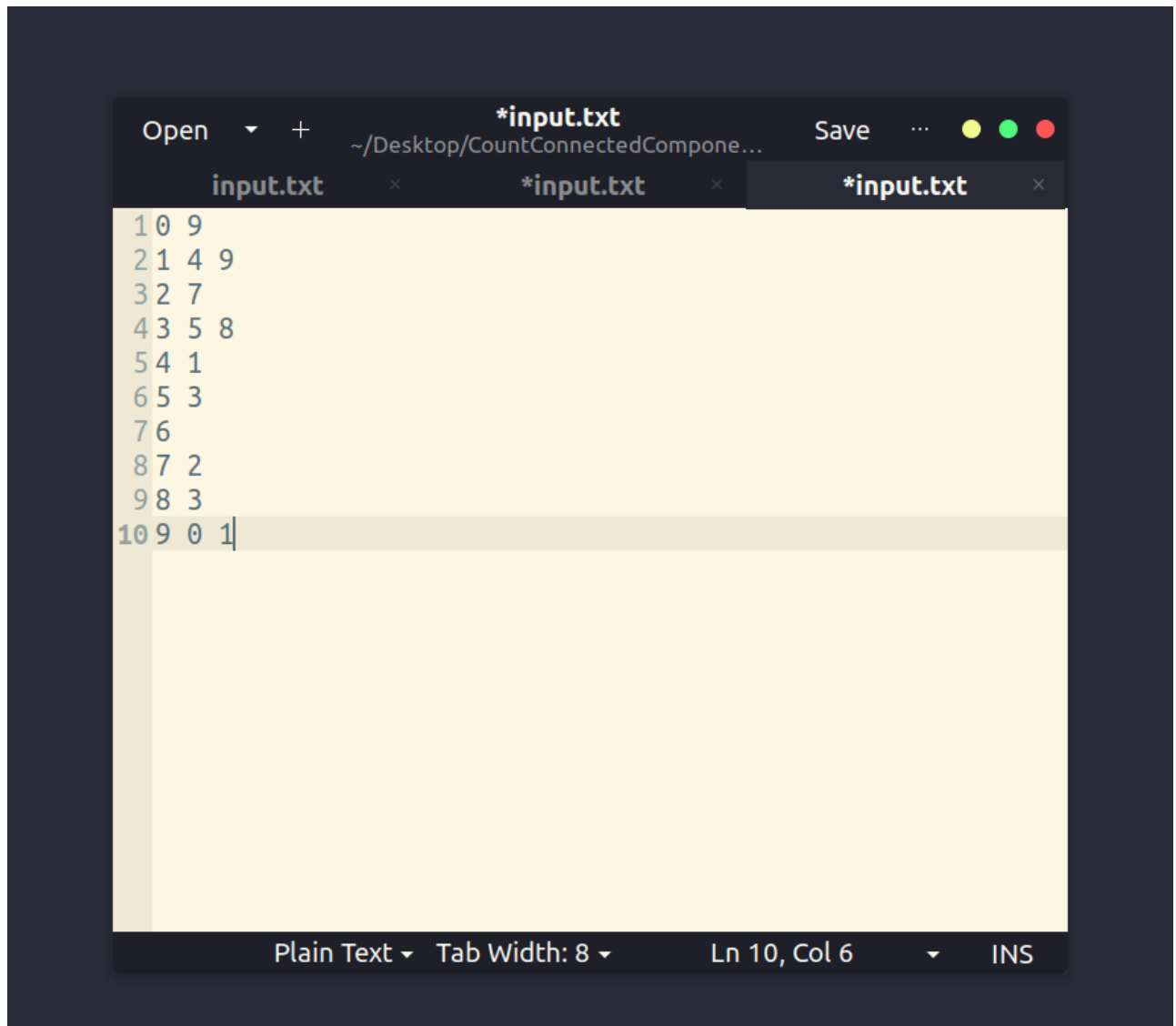
- Notice: Make sure that you export HADOOP_CLASSPATH before buiding file jar

## Step 3: Create directory structure for program in Hadoop

- We need to create folder to store input data in HDFS by below command:

```
hadoop fs -mkdir /CountConnectedComponentProgram
hadoop fs -mkdir /CountConnectedComponentProgram/Input
hadoop fs -put "local input file's path"
/CountConnectedComponentProgram/Input
```

- Example input:



## Step 4: Create Jar File and deploy it to Hadoop

```
hadoop jar "Path to your local file .jar" CountConnectedComponentProgram
/CountConnectedComponentProgram/Input
/CountConnectedComponentProgram/Output
```

## Step 5: Final result

- After succesfully calculating, we can check our result in HDFS like below:

---

# Self-reflection

## 20127435 - Tran Van An

- After completing above tasks, I know more about the useful of MapReduce in real-problems in many aspects as well as get experiences in MapReduce Programing for the midterm test.

## 20127395 - Phan Minh Xuan

- After completing above tasks, I understand how to store, process and manage large data sets, develop skills in the field of big data, especially know more about java language.

20127032 - Bui Gia Huy

- After completing above tasks, I know how to set up and manipulate a basic map reduce program, as well as transforming data using java utility class, as well as familiarize myself with java syntax.

20127631 - Thai Van Thien

- After completing the above tasks, I know how to set up and work with a basic map reduction program and have a preparation for the midterm exam.

## Member's contribution

| Task | Result |
| --- | --- |
| 1.WordCount Program | 100% |
| 2.WordSizeWordCount Program | 100% |
| 3.Weather Data | 100% |
| 4.PatentProgram | 100% |
| 5.MaxTemp Program | 100% |
| 6.Average Salary | 100% |
| 7.De Identify Data | 100% |
| 8.Music Track Program | 100% |
| 9.Telecom Call Data Record Program | 100% |
| 10.Count Connected Components | 100% |

| MSSV | Member | Contribution Percentage |
| --- | --- | --- |
| 20127435 | Tran Van An | 25% |
| 20127395 | Phan Minh Xuan | 25% |
| 20127032 | Bui Gia Huy | 25% |
| 20127631 | Thai Van Thien | 25% |

## References