# BIG DATA HADOOP & SPARK TRAINING

Case study I for session 7

Rashmi Krishna

# Case study 1

**Problem Statement:**

**What are the movie titles that the user has rated?**

**How many times a movie has been rated by the user?**

**In question 2 above, what is the average rating given for a movie?**

**\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*RATING MAPPER\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\***

```java
import java.io.IOException;

import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

public class CaseStudyIUseCasesRatingsMapper extends
                            Mapper<LongWritable, Text, Text, Text> {

        public void map(LongWritable key, Text value, Context context)
                                        throws IOException, InterruptedException {

                        try {
                if (key.get() == 0 && value.toString().contains("userId")){
                    return;
                } else {
                                        String record = value.toString();
                                        String[] parts = record.split(",");
                                        context.write(new Text(parts[1]), new Text("ratings\t" +
parts[2]));
                        }
                } catch (Exception e) {
                    e.printStackTrace();
                }

        }
}
```
Explanation:
This code is for mapping the rating:
  ➢ Here we are checking the input received from input and files and bifurcating them accordingly
  ➢ Input values are LongWritable and text formats while outputs are in Text formats.
  ➢ We are taking only *UserID & rating* from this file.
  ➢ We are checking if key and values are null, then return. If not split the inputs by **","** and parts[1] in the parts array is UserID and parts[2] is movierating.
  ➢ This UserID i.e. **Key** and rating i.e. **Value** is sent as output to the reducer from this mapper.
\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*

**\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*MOVIE MAPPER\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\***

```java
import java.io.IOException;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

public class CaseStudyIUseCasesMoviesMapper extends
                            Mapper<LongWritable, Text, Text, Text> {


        public void map(LongWritable key, Text value, Context context)
                                    throws IOException, InterruptedException {


                    try {
        if (key.get() == 0 && value.toString().contains("movieId")){
          return;
        } else {
            String record = value.toString();
                                String[] parts = record.split(",");
                                context.write(new Text(parts[0]), new Text("movies\t" + parts[1]));
        }
    } catch (Exception e) {
      e.printStackTrace();
    }
        }
}
```
Explanation:
This code is for mapping the rating:
  ➢ Here we are checking the input received from input and files and bifurcating them accordingly
  ➢ Input values are LongWritable and text formats while outputs are in Text formats.
  ➢ We are taking only *movieID & moviename* from this file.
  ➢ We are checking if key and values are null, then return. If not split the inputs by **","** and parts[0] in the parts array is movieID and parts[1] is moviename.
  ➢ This movieID i.e. **Key** and moviename i.e. **Value** is sent as output to the reducer from this mapper.

▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪▪

**\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*REDUCER\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\***

```java
import java.io.IOException;

import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;


public class CaseStudyIUseCasesReducer extends
                        Reducer<Text, Text, Text, Text> {

            public void reduce(Text key, Iterable<Text> values, Context context)
                            throws IOException, InterruptedException {
                    String titles = "";
                    double total = 0.0;
                    int count = 0;
                    System.out.println("Text Key    =>"+key.toString());
                    for (Text t : values) {
                            String parts[] = t.toString().split(",");
                            System.out.println("Text values =>"+t.toString());
                            if (parts[0].equals("ratings")) {
                                    count++;
                                    String rating = parts[1].trim();
                                    System.out.println("Rating is =>"+rating);
                                    total += Double.parseDouble(rating);
                            } else if (parts[0].equals("movies")) {
                                    titles = parts[1];
                            }                                 }

                    double average = total / count;
                    String str = String.format("%d\t%f", count, average);
                    context.write(new Text(titles), new Text(str));
            }
    }
```
Explanation:
- ➢ Here outputs of two mappers are inputs to this reducer.
- ➢ Both input and outputs are Text format.
- ➢ Now we check all the inputs and bifurcate them accordingly.
- ➢ UserID and MovieID are the keys, we split the input by "," and check if the part is "*rating*" or not.
    - o If the part is rating then we print the rating and calculate the total.
    - o If the part is not rating then it must moviename, then we pring the moviename and save it in the variable "title"
    - o We calculate the average of the rating for a particular movie title.
- ➢ We print the number of times the movie was rating by the user and the average rating.


\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*

******************************DRIVER *********************************************

```java
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;
import org.apache.hadoop.mapreduce.lib.input.TextInputFormat;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.MultipleInputs;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class CaseStudyIUseCasesDriver {

        @SuppressWarnings("deprecation")
        public static void main(String[] args) throws Exception {
   if (args.length != 3) {
     System.err.println("Usage: CaseStudyIUseCase2Driver <input path1> <input path2> <output path>");
      System.exit(-1);
   }

        //Job Related Configurations
        Configuration conf = new Configuration();
        Job job = new Job(conf, "CaseStudyIUseCase2Driver");
        job.setJarByClass(CaseStudyIUseCasesDriver.class);

        //job.setNumReduceTasks(0);

        //Since there are multiple input, there is a slightly different way of specifying input path,
input format and mapper
        MultipleInputs.addInputPath(job, new Path(args[0]),TextInputFormat.class,
CaseStudyIUseCasesMoviesMapper.class);
        MultipleInputs.addInputPath(job, new Path(args[1]),TextInputFormat.class,
CaseStudyIUseCasesRatingsMapper.class);

        //Set the reducer
        job.setReducerClass(CaseStudyIUseCasesReducer.class);

   //set the out path
        Path outputPath = new Path(args[2]);
        FileOutputFormat.setOutputPath(job, outputPath);
        outputPath.getFileSystem(conf).delete(outputPath, true);

   //set up the output key and value classes
   job.setOutputKeyClass(Text.class);
   job.setOutputValueClass(Text.class);

   //execute the job
   System.exit(job.waitForCompletion(true) ? 0 : 1);
 }
}
```

**Explanation:**

➢ Here there are 2 input paths and 1 output path, thereby, we check if all the 3 parameters are entered by the user, if not an error is given saying user has to enter 3 parameters and exits.
➢ Job configuration instance is created and driverclass is set jar by class.
➢ Multiple input path are defined under args[0] and args[1], as we have two csv files. So each csv file is given in two different paths
➢ Output path is defined and also output key and value class

```
                Total vcore-milliseconds taken by all map tasks=1245725
                Total vcore-milliseconds taken by all reduce tasks=477175
                Total megabyte-milliseconds taken by all map tasks=1275622400
                Total megabyte-milliseconds taken by all reduce tasks=488627200
        Map-Reduce Framework
                Map input records=26070134
                Map output records=26070133
                Map output bytes=484518568
                Map output materialized bytes=536658882
                Input split bytes=1988
                Combine input records=0
                Combine output records=0
                Reduce input groups=316067
                Reduce shuffle bytes=536658882
                Reduce input records=26070133
                Reduce output records=316067
                Spilled Records=76775525
                Shuffled Maps =7
                Failed Shuffles=0
                Merged Map outputs=7
                GC time elapsed (ms)=14917
                CPU time spent (ms)=507210
                Physical memory (bytes) snapshot=1639714816
                Virtual memory (bytes) snapshot=16448626688
                Total committed heap usage (bytes)=1188126720
        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=0
        File Output Format Counters
                Bytes Written=2212469
You have new mail in /var/spool/mail/acadgild
```

- Application
  - About
  - Jobs
- Tools

Log Type: stdout
Log Upload Time: 15-Apr-2018 12:46:33
Log Length: 30885947

```
Text Key    =>"""Great Performances"" Cats (1998)"
Text values =>ratings   Musical
Text Key    =>"$1
Text values =>ratings   000 on the Black (1966)"
Text Key    =>"$100
Text values =>ratings   000 for Ringo (1965)"
Text Key    =>"'Human' Factor
Text values =>ratings   The (Human Factor
Text Key    =>"'burbs
Text values =>ratings   The (1989)"
Text Key    =>"1
Text values =>ratings   2
Text Key    =>"10
Text values =>ratings   000 BC (2008)"
Text values =>ratings   000 Black Men Named George (2002)"
Text values =>ratings   000 Dollars for a Massacre (1967)"
Text values =>ratings   000 Days (2014)"
Text Key    =>"1000 Eyes of Dr. Mabuse
Text values =>ratings   The (Die 1000 Augen des Dr. Mabuse) (1960)"
Text Key    =>"10th Judicial Court: Judicial Hearings
Text values =>ratings   The (10e chambre - Instants d'audience) (2004)"
Text Key    =>"10th Kingdom
Text values =>ratings   The (2000)"
Text Key    =>"10th Victim
Text values =>ratings   The (La decima vittima) (1965)"
Text Key    =>"11'09""01 - September 11 (2002)"
Text values =>ratings   Drama
Text Key    =>"11th Hour
Text values =>ratings   The (2007)"
Text Key    =>"12 Dogs of Christmas
Text values =>ratings   The (2005)"
Text Key    =>"13 Frightened Girls! (Candy Web
```