# Assignment 1: MapReduce Program Application
COS20028 – Big data Architecture and Applications
Semester 2, 2023

Name: MD NAHID TANJUM
Student ID: 10807068
Date: August 29, 2023

# Driver Code



```
J GenreMovieRating.java ⊠   J *MovieRatingMapper.java      J *MovieRatingReducer.java

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.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;


/*
 * Driver class for the GenreMovieRating MapReduce job.
 * It sets up the job configuration,input and output paths,
 * specifics the Mapper and Reducer classes, and executes the job.
 */
public class GenreMovieRating {

  public static void main(String[] args) throws Exception {

      //Ensure correct usage by expecting exactly 2 command-line arguments:input and output paths.
      if (args.length != 2) {
        System.err.println("Usage: GenreMovieRating <input path> <output path>");
        System.exit(-1);
      }

      /*
       * Set up the job Configuration and name
       */
      Configuration conf = new Configuration();
      Job job = Job.getInstance(conf, "Genre Movie Rating");
```

```java
        /*
         * Set the jar containing the driver,mapper and reducer classes.
         */
        job.setJarByClass(GenreMovieRating.class);

        //Set the Mapper and Reducer classes for the job.
        job.setMapperClass(MovieRatingMapper.class);
        job.setReducerClass(MovieRatingReducer.class);



        /*
         * Specify input and output paths based on command-line arguments.
         */
        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));

        /*
         * Define output key and value classes.
         */
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(LongWritable.class);



        //Metadata:programmer details
        System.out.print("programmer: Md Nahid Tanjum\nStudent ID: 103807068\n");

        /*
         * Run the MapReduce job and exit based on its success.
         */

        System.exit(job.waitForCompletion(true) ? 0 : 1);
```
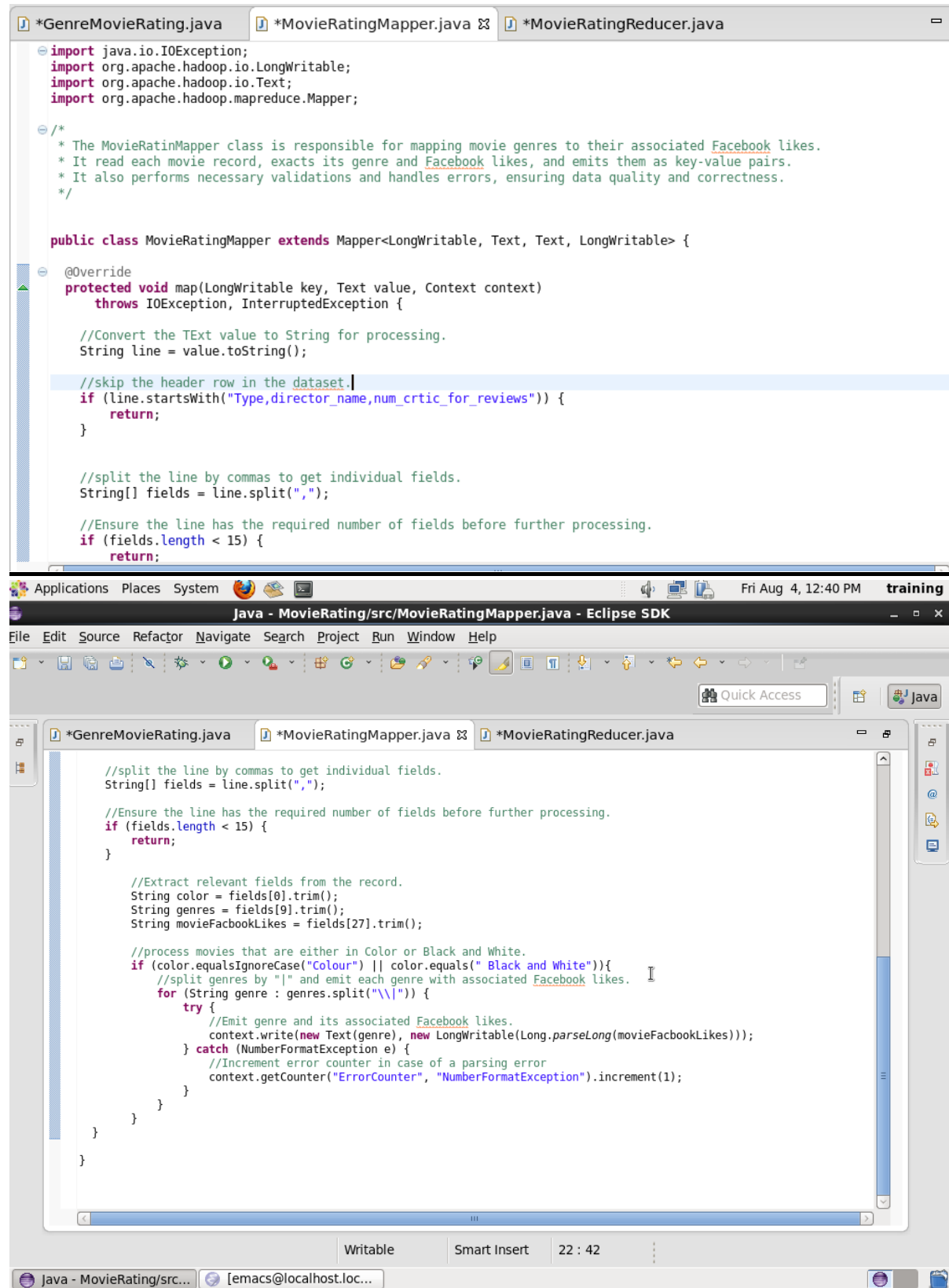
```java
        job.setJarByClass(GenreMovieRating.class);

        //Set the Mapper and Reducer classes for the job.
        job.setMapperClass(MovieRatingMapper.class);
        job.setReducerClass(MovieRatingReducer.class);



        /*
         * Specify input and output paths based on command-line arguments.
         */
        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));

        /*
         * Define output key and value classes.
         */
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(LongWritable.class);



        //Metadata:programmer details
        System.out.print("programmer: Md Nahid Tanjum\nStudent ID: 103807068\n");

        /*
         * Run the MapReduce job and exit based on its success.
         */

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

# Mapper Code

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

/*
 * The MovieRatinMapper class is responsible for mapping movie genres to their associated Facebook likes.
 * It read each movie record, exacts its genre and Facebook likes, and emits them as key-value pairs.
 * It also performs necessary validations and handles errors, ensuring data quality and correctness.
 */


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

    @Override
    protected void map(LongWritable key, Text value, Context context)
        throws IOException, InterruptedException {

        //Convert the TExt value to String for processing.
        String line = value.toString();

        //skip the header row in the dataset.
        if (line.startsWith("Type,director_name,num_crtic_for_reviews")) {
            return;
        }


        //split the line by commas to get individual fields.
        String[] fields = line.split(",");

        //Ensure the line has the required number of fields before further processing.
        if (fields.length < 15) {
            return;
```

Java - MovieRating/src/MovieRatingMapper.java - Eclipse SDK

File  Edit  Source  Refactor  Navigate  Search  Project  Run  Window  Help

*GenreMovieRating.java    *MovieRatingMapper.java    *MovieRatingReducer.java

```java
        //split the line by commas to get individual fields.
        String[] fields = line.split(",");

        //Ensure the line has the required number of fields before further processing.
        if (fields.length < 15) {
            return;
        }

        //Extract relevant fields from the record.
        String color = fields[0].trim();
        String genres = fields[9].trim();
        String movieFacbookLikes = fields[27].trim();

        //process movies that are either in Color or Black and White.
        if (color.equalsIgnoreCase("Colour") || color.equals(" Black and White")){
            //split genres by "|" and emit each genre with associated Facebook likes.
            for (String genre : genres.split("\\|")) {
                try {
                    //Emit genre and its associated Facebook likes.
                    context.write(new Text(genre), new LongWritable(Long.parseLong(movieFacbookLikes)));
                } catch (NumberFormatException e) {
                    //Increment error counter in case of a parsing error
                    context.getCounter("ErrorCounter", "NumberFormatException").increment(1);
                }
            }
        }
    }
}
```

Writable    Smart Insert    22 : 42

# Reducer Code



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

/*
 * The MovieRatingReducer class is responsible for calculating the average Facebook likes for each movie genre.
 * For each genre, it aggregates the total Facebook_likes and computes the average, which is then emitted as a key-value pa
 * key being the genre and value being the average of Facebook likes.
 *
 */
public class MovieRatingReducer extends Reducer<Text, LongWritable, Text, Text> {

    @Override
    protected void reduce(Text genre, Iterable<LongWritable> values, Context context)
            throws IOException, InterruptedException {

        // Initialize variables to keep track of the sum of Facebook likes and the number of movies in each genre.
        long sum = 0;
        int count = 0;

        //Iterate over the Facebook likes for movies in the current genre, aggregating the total likes and counting movies
        for (LongWritable value : values) {

            sum += value.get();
            count++;
        }

        //compute the average number of Facebook likes for the current genre.
        double averageRating = (double) sum / count;

        //Emit the genre and its average Facebook likes as a key-value pair.
        context.write(genre, new Text(String.valueOf(averageRating)));
```



```java
/*
 * The MovieRatingReducer class is responsible for calculating the average Facebook likes for each movie genre.
 * For each genre, it aggregates the total Facebook_likes and computes the average, which is then emitted as a key-value pa
 * key being the genre and value being the average of Facebook likes.
 *
 */
public class MovieRatingReducer extends Reducer<Text, LongWritable, Text, Text> {

    @Override
    protected void reduce(Text genre, Iterable<LongWritable> values, Context context)
            throws IOException, InterruptedException {

        // Initialize variables to keep track of the sum of Facebook likes and the number of movies in each genre.
        long sum = 0;
        int count = 0;

        //Iterate over the Facebook likes for movies in the current genre, aggregating the total likes and counting movies
        for (LongWritable value : values) {

            sum += value.get();
            count++;
        }

        //compute the average number of Facebook likes for the current genre.
        double averageRating = (double) sum / count;

        //Emit the genre and its average Facebook likes as a key-value pair.
        context.write(genre, new Text(String.valueOf(averageRating)));
    }
}
```

# Commands and the Execution Screenshot

# First Page of the Output File Screenshot