

YOUTUBE DATA ANALYSIS USING HADOOP, PIG & HIVE

Name: Sagar Shah

NUID: 001342989

INFO 7250

Engineering of Big-Data Systems Final Project

YouTube Dataset Analysis

Problem Statement

Analyze the Youtube open source API Big dataset using Hadoop, Pig and HIVE based on different column fields to provide various comprehensive insights

Summary

The data-set used in this project is provided by Simon Fraser University. The data-set has five different files based on the data collected by the crawler.

The dataset was available on the following Url:

<http://netsg.cs.sfu.ca/youtubedata/>

We record the following information of a YouTube video in order; they are divided by '\t' in the data file.

video ID	an 11-digit string, which is unique
uploader	a string of the video uploader's username
age	an integer number of days between the date when the video was uploaded and Feb.15, 2007 (YouTube's establishment)
category	a string of the video category chosen by the uploader
length	an integer number of the video length
views	an integer number of the views
rate	a float number of the video rate
ratings	an integer number of the ratings
comments	an integer number of the comments
related IDs	up to 20 strings of the related video IDs

Datasets of User Information

We have collected the information about YouTube users. The crawler retrieves information on the number of uploaded videos and friends of each user from the YouTube API, for a total of more than 1 million users. There is "user.txt", containing the information of number of uploads, watches and friends in order.

- ② Following MapReduce programs and its different design patterns like mentioned below are implemented:
 1. Filtering
 2. Join Patterns
 3. Data Organization
 4. Summarization

- ② Following analysis can be performed on the data-set :
 1. Calculate Max Rating Total Rating and Total Comment Count by Video ID
 2. Moving Rating Average by Video_ID
 3. Best Youtuber based on videos uploaded
 4. Top 50 Favorite YouTube Videos
 5. Total YouTube Videos by Category
 6. Implement Binning based on Categories
 7. Implement Chaining on Binning result to get Top 25 Videos per category
 8. Recommend followers based on connected followers
 9. Total Views based on Video ID

- ② Following Pig analysis is performed on the data-set : (Implemented Visualizations for Pig Analysis found at the bottom in the Pig analysis section)
 1. Calculate top 5 Categories of Youtube Videos
 2. Calculate top 10 Rated of Youtube Videos
 3. Calculate top 10 Rated By Categories of Youtube Videos
 4. Calculate top 10 Viewed of Youtube Videos
 5. Calculate top 10 Viewed By Categories of Youtube Videos

- ② Following Hive analysis is performed on the data-set :
 1. Calculate top 10 channels with maximum number of likes
 2. Calculate top 5 categories with maximum number of comments

- ② Implemented Google Data Studio for Visualization purpose

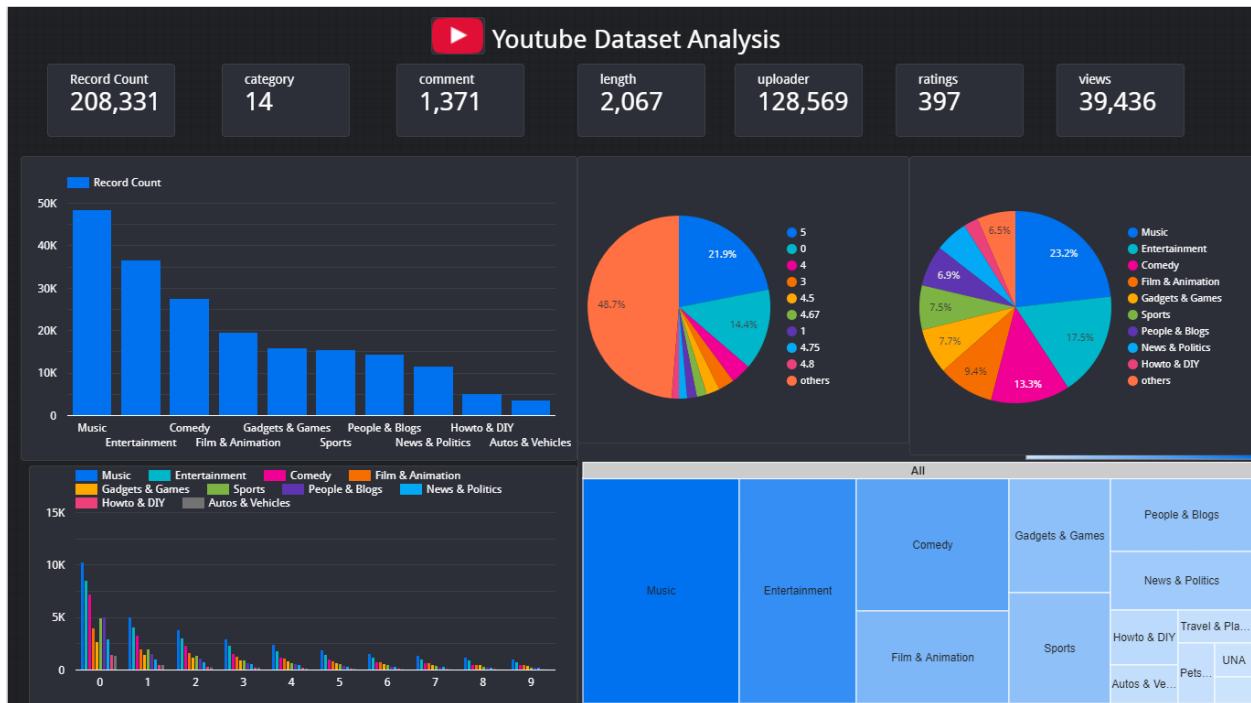
Visualization using Google Data Studio

Link

Link to the dashboard : <https://datastudio.google.com/reporting/97d6fd91-b21f-4576-8f14-5ef154b56f65>

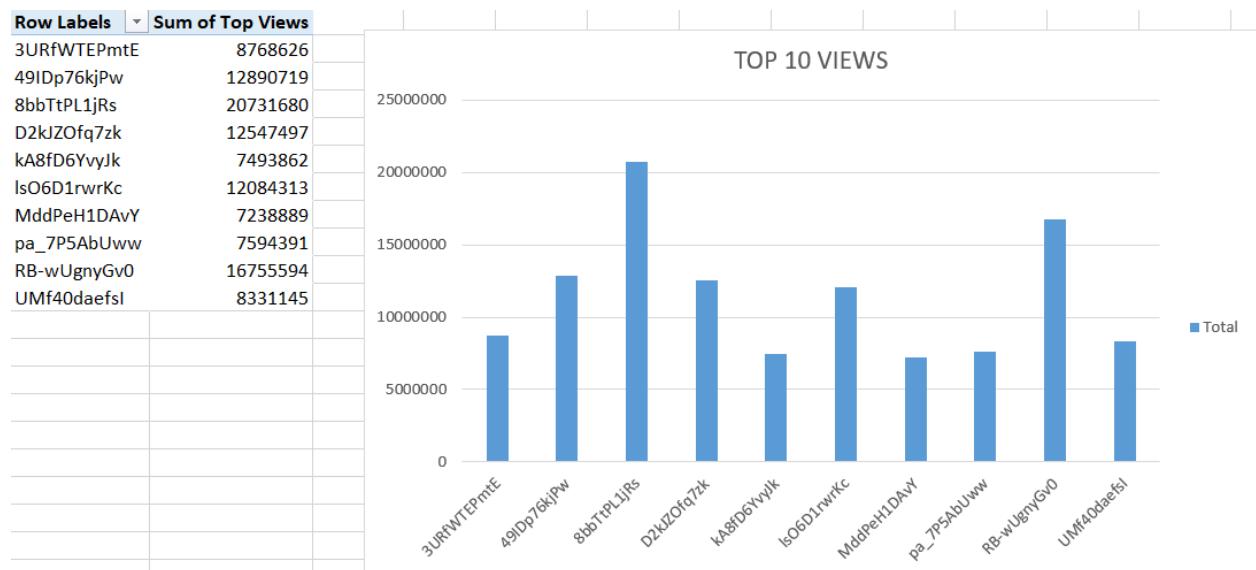
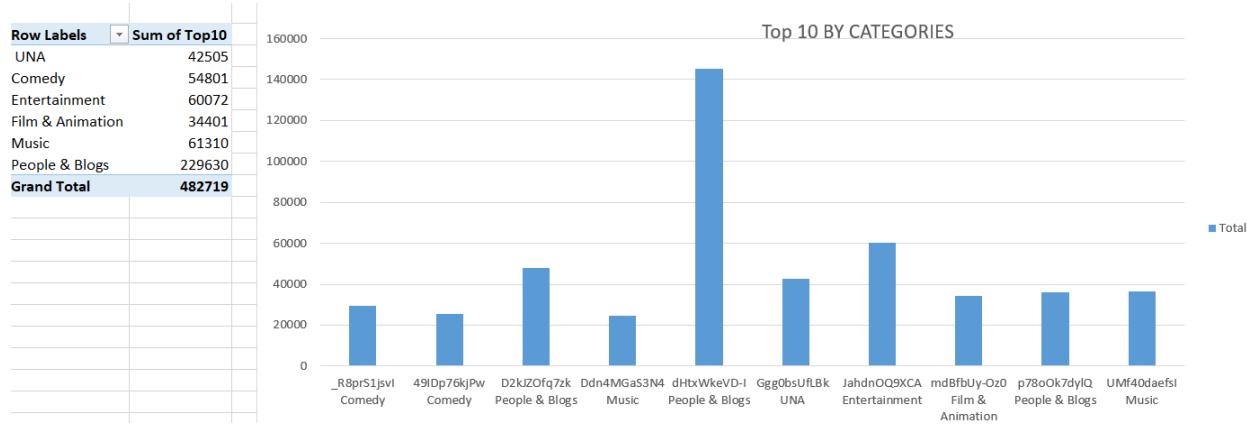
Note : (Have limited the record count for this analysis to 2 Lakh records however, for Map Reduce Analysis records counts is much larger)

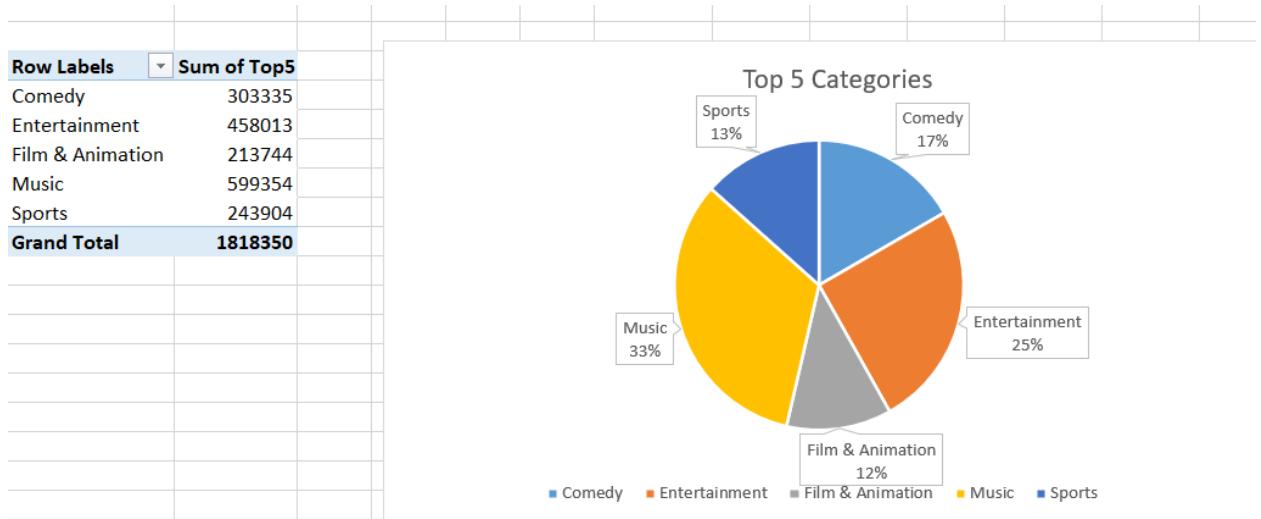
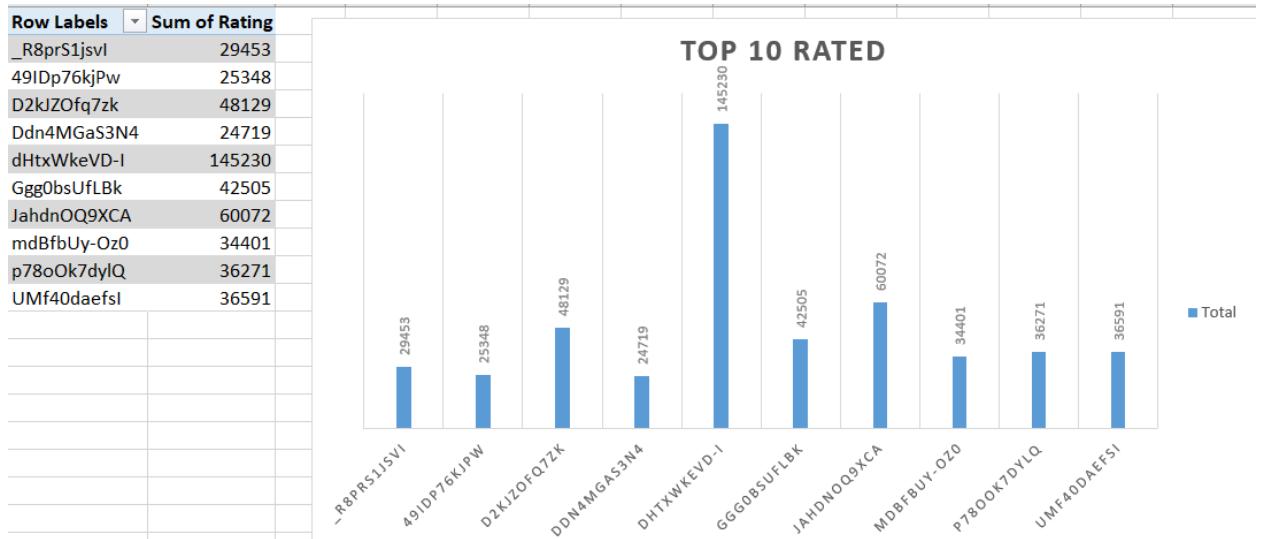
Screenshot



Visualizations for Pig Analysis

The detailed analysis can be found in the bottom, these are just the screenshots





Map Reduce Analysis

Merge Files : Merges all the CSV Files into one and stores it into HDFS

Code

```
package mergedataset;
import java.io.IOException;
public class MergeCSV {
    public static void main(String[] args) throws IOException {
        Configuration conf = new Configuration();
        FileSystem hdfs = FileSystem.get(conf);
        FileSystem local = FileSystem.getLocal(conf);
        Path inputDir = new Path(args[0]);
        Path hdfsFile = new Path(args[1]);
        try {
            FileStatus[] inputFiles = local.listStatus(inputDir);
            FSDataOutputStream out = hdfs.create(hdfsFile);
            for(int i = 0; i < inputFiles.length; i++) {
                System.out.println(inputFiles[i].getPath().getName());
                FSDataInputStream in = local.open(inputFiles[i].getPath());
                byte buffer[] = new byte[256];
                int bytesRead = 0;
                while((bytesRead = in.read(buffer)) > 0) {
                    out.write(buffer, 0, bytesRead);
                }
                in.close();
            }
            out.close();
        }
        catch (IOException e) {
            e.printStackTrace();
        }
    }
}
```

Execution

```
sagarshah95@ubuntu: /usr/local/bin/hadoop-3.3.0/bin$ ./hadoop jar /home/sagarshah95/putmergeimp.jar merge dataset.MergeCSV /home/sagarshah95/Desktop/csv /BigDataProject/data/youtubeDataset  
3.csv  
4.csv  
1.csv  
2.csv  
0.csv  
sagarshah95@ubuntu:/usr/local/bin/hadoop-3.3.0/bin$ ./hadoop fs -ls /BigDataProject/data  
Found 1 items  
-rw-r--r-- 1 sagarshah95 supergroup 219815043 2021-04-20 15:48 /BigDataProject/data/youtubeDataset  
sagarshah95@ubuntu:/usr/local/bin/hadoop-3.3.0/bin$ ./hadoop fs -head /BigDataProject/data/youtubeDataset  
t  
PF_ZMlw4rHs,DisneyUnleashed,729,Film & Animation,288,96,3.5,2,0,1xbSFrHzFQ0,4VP4qSjDNQs,RJgGeYiJrj0,Mqy  
Sp7Nq5j0,MC--VwVTHAM,eVdiIbwT60M,Da80HD18tp0,mfp6Z8z1cI,tdRBH7VBrSY,xKvzxLeYoiQ,_I2EZYCdUXI,gompU_uhYq0  
,ClPPxBxPB0w,MeFi3SD1_n8,YRi20cWMyOM,v2UEfmW06z8,2t2Fe_ixWpI,_Wud5vSIQ01,bBml9opQxnc,1ZU_ytaZTxg  
c3XKOAxKc_w,TvBride,495,Entertainment,80,334,4,1,0,NVRkDihhHB0,wL1-yb-vb1o,7qeWjy9hPok,-30GaT2E200,gbFqa  
hWCyqQ,GZguki6X3nE,UcaCn0caIxQ,mkTWQrkEdTs,gloJygJc2aw,y_iew26B5JQ,KXq4j2aStGw,BL6hcqdRzs8,tQEeqTCxxdyk,2  
PFMouZoKGw,LS-xTWSgb1Y,yePQqn_YE9c,-pyPi5T6QNq,-B8UWZ6xb1o,utJi9Ha7yPc,U2hhxHny55w  
yr73064qrGE,jamesacisco3rd,491,Music,125,262,0,0,0,rD9zwdFmAwg,m6_GLY46Ee8,Pc8X4YLAFlc,NBLxyiuBwKE,-RMOM  
EZeuGc,bGimrcYWhkk,ZkNTZ4pUd8,lNCz2uSxiqw,cVwGA-4lE2U,rEDuBZNoCqo,HbRPbgtYASs,C5_8rgDsFHE,a2kiShr0r7I,t  
n_1-2UanZE,x1yJvXA5xg,VENsQsE4ZzY,HISio2WVxcI,wQoz4uaMSJs,EPQ_7CuDJBI,oreH2bC2Ixk  
atfNL0_KAcS,f0xmuld3r,454,Howto & DIY,102,47718,4.84,101,44,tt3W6X8971o,kTfYttriolI,pdyYe7sDlhA,WCzaeeah  
saqarshah95@ubuntu:/usr/local/bin/hadoop-3.3.0/bin$
```

Top Viewed Video: Returns the video with highest number of views received (youtubeanalysis)

This analysis comprises of :

1. Driver Class : Sets the configuration for Mapper and Reducer class
2. Mapper Class : Emits videoIDs of type Text and number of views of type FloatWritable
3. Reducer Class : Emits videoID of type Text and views of type FloatWritable

Driver Class

```
package top_viewed_video;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.FloatWritable;
import org.apache.hadoop.io.IntWritable;
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.input.TextInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;
import java.io.IOException;

public class DriverClass {

    public static void main(String[] args) throws IOException, InterruptedException, ClassNotFoundException{
        Configuration conf = new Configuration();
        Job job = new Job(conf, "Top Videos");
        job.setJarByClass(DriverClass.class);

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

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

        job.setMapOutputKeyClass(Text.class);
        job.setMapOutputValueClass(FloatWritable.class);

        job.setMapperClass(MapperClass.class);
        job.setReducerClass(ReducerClass.class);

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

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

Mapper Class

```
package top_viewed_video;

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

public class MapperClass extends Mapper<LongWritable, Text, Text, FloatWritable> {

    private Text video_name = new Text();
    private FloatWritable views = new FloatWritable();

    @Override
    public void map(LongWritable key, Text value, Context context)
        throws IOException, InterruptedException {
        String line = value.toString();
        String str[] = line.split(",");
        if (str.length >= 5) {
            video_name.set(str[0]);
            try {
                float temp = Float.parseFloat(str[5]); //typecasting string to Integer
                views.set(temp);
                context.write(video_name, views);
            }catch(Exception e){
                e.printStackTrace();
            }
            //views.set(temp);
        }
    }
}
```

Reducer Class

```
package top_viewed_video;

import java.io.IOException;

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

public class ReducerClass extends Reducer<Text, FloatWritable, Text, FloatWritable> {
    static float max = 0;
    static int sum = 0;
    static String finalKey = "";

    @Override
    public void reduce(Text key, Iterable<FloatWritable> values, Context context)
        throws IOException, InterruptedException {
        for (FloatWritable val : values) {
            sum += val.get();
        }
        if (sum > max) {
            max = sum;
            finalKey = key.toString();
        }
    }

    @Override
    protected void cleanup(Reducer<Text, FloatWritable, Text, FloatWritable>.Context context)
        throws IOException, InterruptedException {
        context.write(new Text(finalKey), new FloatWritable(max));
        // TODO Auto-generated method stub
    }
}
```

Execution

```

File Edit View Search Terminal Help
  File Output Format Counters
    Bytes Written=25
sagarshah95@ubuntu:/usr/local/bin/hadoop-3.3.0/bin$ ./hadoop jar /home/sagarshah95/Desktop/searchByVideo-0.0.1-SNAPSHOT.jar top_viewed_video.DriverClass /BigDataProject/data/youtubeDataset /BigDataProject/data/o
sgrep@ip-172-31-10-10:~$ ./hadoop jar /home/sagarshah95/Desktop/searchByVideo-0.0.1-SNAPSHOT.jar top_viewed_video.DriverClass /BigDataProject/data/youtubeDataset /BigDataProject/data/o
2021-04-28 20:03:41,441 INFO client.DefaultNoHDFSProxyProvider: Connecting to ResourceManager at /0.0.0.0:8032
2021-04-28 20:03:42,024 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
2021-04-28 20:03:42,062 INFO mapreduce.JobResourceUploader: Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/sagarshah95_.staging/job_1619663038748_0004
2021-04-28 20:03:42,338 INFO InputFormat: Total input files to process : 1
2021-04-28 20:03:42,340 INFO InputFormat: Input splits:
2021-04-28 20:03:42,183 INFO mapreduce.JobSubmitter: number of splits:1
2021-04-28 20:03:43,103 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1619663038748_0004
2021-04-28 20:03:43,103 INFO mapreduce.JobSubmitter: Executing with tokens: []
2021-04-28 20:03:43,438 INFO conf.Configuration: resource-types.xml not found
2021-04-28 20:03:43,439 INFO resource.ResourceUtil: Unable to find 'resource-types.xml'.
2021-04-28 20:03:43,629 INFO lpm.YarnClientImpl: Submitted application application_1619663038748_0004
2021-04-28 20:03:43,630 INFO lpm.YarnClientImpl: Application has been assigned job: http://172.31.10.10:8088/proxy/application_1619663038748_0004/
2021-04-28 20:03:43,715 INFO mapreduce.Job: Running job: job_1619663038748_0004
2021-04-28 20:03:58,995 INFO mapreduce.Job: Job job_1619663038748_0004 running in uber mode : false
2021-04-28 20:03:58,996 INFO mapreduce.Job: map 0% reduce 0%
2021-04-28 20:04:03,772 INFO mapreduce.Job: map 50% reduce 0%
2021-04-28 20:04:03,773 INFO mapreduce.Job: map 100% reduce 0%
2021-04-28 20:04:11,852 INFO mapreduce.Job: map 100% reduce 100%
2021-04-28 20:04:12,881 INFO mapreduce.Job: Job job_1619663038748_0004 completed successfully
2021-04-28 20:04:13,029 INFO mapreduce.Job: Counters: 55
  File System Counters
    FILE: Number of bytes read=13749937
    FILE: Number of bytes written=3239402
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=21919381
    HDFS: Number of bytes written=25
    HDFS: Number of read operations=11
    HDFS: Number of large read operations=0
    HDFS: Number of write operations=2
    HDFS: Number of bytes read erasure-coded=0
Job Counters
  Killed map tasks=1
  Launched map tasks=2
  Launched reduce tasks=1
  Data-local map tasks=2
  Total time spent by all map in occupied slots (ms)=20504
  Total time spent by all reduces in occupied slots (ms)=5583
  Total time spent by all map tasks (ms)=8594
  Total time spent by all reduce tasks (ms)=5583
  Total vcore-milliseconds taken by all map tasks=20504
  Total vcore-milliseconds taken by all reduce tasks=5583
  Total map-milliseconds taken by all map tasks=20996096
  Total map-bytes-milliseconds taken by all reduce tasks=3716992
Map-Reduce Framework
  Map Input records=769739

```

```
File Edit View Search Terminal Help
Job Counter
    Killed map tasks=1
    Launched map tasks=2
    Launched reduce tasks=1
    Data-local map tasks=2
    Total time spent by all maps in occupied slots (ms)=20504
    Total time spent by all reduces in occupied slots (ms)=5583
    Total time spent by all map tasks (ms)=20504
    Total time spent by all reduce tasks (ms)=5583
    Total vcore-milliseconds taken by all map tasks=20504
    Total vcore-milliseconds taken by all reduce tasks=5583
    Total megabyte-milliseconds taken by all map tasks=209966096
    Total megabyte-milliseconds taken by all reduce tasks=5716992
Map-Reduce Framework
    Map input records=769739
    Map output records=763885
    Map output bytes=12222161
    Map output materialized bytes=13749943
    Input splits bytes=242
    Combine output bytes=8
    Combine output records=0
    Reduce input groups=740427
    Reduce shuffle bytes=13749943
    Reduce input records=763885
    Reduce output records=1
    Spills local records=527770
    Shuffled Maps =2
    Failed Shuffles=0
    Merged Map outputs=2
    GC time elapsed (ms)=478
    CPU time spent (ms)=1990
    Physical memory (bytes) snapshot=1086612160
    Virtual memory (bytes) snapshot=7763783688
    Total committed heap usage (bytes)=959447840
    Peak Map Physical memory (bytes)=461369344
    Peak Map Virtual memory (bytes)=2586222592
    Peak Reduce Physical memory (bytes)=220241920
    Peak Reduce Virtual memory (bytes)=2592538624
Shuffle Errors
    BAD_ID=0
    CONNECTION=0
    IO_ERROR=0
    WRONG_LENGTH=0
    WRONG_TYPE=0
    WRONG_REDUCE=0
File Input Format Counters
    Bytes Read=219819139
File Output Format Counters
    Bytes Written=25
sagarshah95@ubuntu:/usr/local/bin/hadoop-3.3.0/bin$
```

```
sagarshah95@ubuntu:/usr/local/bin/hadoop-3.3.0/bin$ ./hadoop fs -cat /BigDataProject/data/outputFiles/topViewed_Video/part-r-00000
DDvq79XTiww      2.14748365E9
```

Average Rating: Returns the average number of ratings of videos along with comments

This analysis comprises of :

1. Driver Class : Sets up configuration for Mapper, Reducer, Combiner and Tuple class
2. Mapper Class : Emits VideoID of type Text and Tuple of custom type
Averagerating_CommentTuple
3. Reducer Class : Emits VideoID of type Text and Tuple of custom type
Averagerating_CommentTuple
4. Combiner Class : Emits Average rating along with comment count
5. Tuple Class : Pojo class consisting of Rating and Comment Count

Driver Class

```
package averagerating_youtube;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.conf.Configured;
import org.apache.hadoop.fs.Path;
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;
import org.apache.hadoop.util.Tool;
import org.apache.hadoop.util.ToolRunner;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.FloatWritable;
import org.apache.hadoop.io.IntWritable;
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.input.TextInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;
import java.io.IOException;

public class AverageRating_Youtube {

    /**
     * @param args the command line arguments
     */
    //@Override
    public static void main(String[] args) throws IOException, InterruptedException, ClassNotFoundException {
        Configuration conf = new Configuration();
        Job job = new Job(conf, "AverageRating_Youtube");
        //Job job = new Job(getConf());
        //job.setJobName("AverageRating_Youtube");

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

        job.setMapperClass(AvgRating_CommCountMapper.class);
        job.setMapOutputKeyClass(Text.class);
        job.setMapOutputValueClass(AverageRating_CommentCountTuple.class);
        job.setCombinerClass(AvgRating_CommCountCombiner.class);
        job.setReducerClass(AvgRating_CommCountReducer.class);
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(AverageRating_CommentCountTuple.class);
        System.exit(job.waitForCompletion(true)?0:1);
    }

}
```

Mapper Class

```
package averagerating_youtube;

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

public class AvgRating_CommCountMapper extends Mapper<Object, Text, Text, AverageRating_CommentCountTuple> {

    // Our output key and value Writables
    private Text video_name = new Text();
    private float v_rate;
    private AverageRating_CommentCountTuple outTuple = new AverageRating_CommentCountTuple();

    @Override
    protected void map(Object key, Text value, Context context) throws IOException, InterruptedException {
        String[] fields = value.toString().split(",");
        String videoId = (fields[0]);
        try {
            if(fields.length > 6)
                if (!fields[6].isEmpty())
                    this.v_rate = Float.parseFloat(fields[6]);
            else
                this.v_rate = 0;
            video_name.set(videoId);
            outTuple.setComment_count(1);
            outTuple.setVideo_rating(this.v_rate);
            context.write(video_name, outTuple);
        }catch (Exception e){
            e.printStackTrace();
        }
    }
}
```

Reducer Class

```
package averagerating_youtube;

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

public class AvgRating_CommCountReducer extends Reducer<Text, AverageRating_CommentCountTuple, Text, AverageRating_CommentCountTuple> {

    private AverageRating_CommentCountTuple result = new AverageRating_CommentCountTuple();

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

        float sum = 0;
        int count = 0;

        for (AverageRating_CommentCountTuple val : values) {
            sum += val.getComment_count() * val.getVideo_rating();
            count += val.getComment_count();
        }

        result.setVideo_rating(sum / count);
        //result.setComment_count(count);
        context.write(key, result);
    }
}
```

Combiner Class

```
package averagerating_youtube;

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

public class AvgRating_CommCountCombiner extends Reducer<Text, AverageRating_CommentCountTuple, Text, AverageRating_CommentCountTuple> {

    private AverageRating_CommentCountTuple result = new AverageRating_CommentCountTuple();

    protected void reduce(Text key, Iterable<AverageRating_CommentCountTuple> values, Reducer.Context context) throws IOException, InterruptedException {

        float sum = 0;
        int count = 0;

        for (AverageRating_CommentCountTuple val : values) {
            sum += val.getComment_count() * val.getVideo_rating();
            count += val.getComment_count();
        }

        result.setVideo_rating(sum / count);
        result.setComment_count(count);
        context.write(key, result);
    }
}
```

Tuple Class

```
package averagerating_youtube;

import java.io.DataInput;
import java.io.DataOutput;
import java.io.IOException;
import org.apache.hadoop.io.Writable;

public class AverageRating_CommentCountTuple implements Writable {

    private int comment_count = 0;
    private double video_rating = 0;

    public int getComment_count() {
        return comment_count;
    }

    public void setComment_count(int comment_count) {
        this.comment_count = comment_count;
    }

    public double getVideo_rating() {
        return video_rating;
    }

    public void setVideo_rating(double video_rating) {
        this.video_rating = video_rating;
    }

    public void write(DataOutput d) throws IOException {
        d.writeInt(comment_count);
        d.writeDouble(video_rating);
    }

    public void readFields(DataInput di) throws IOException {
        comment_count = di.readInt();
        video_rating = di.readDouble();
    }

    @Override
    public String toString() {
        return Integer.toString(comment_count) + " " + Double.toString(video_rating);
    }
}
```

Execution

```
sagarshah@ubuntu:~/usr/local/bin/hadoop-3.3.0/bin$ ./hadoop jar /home/sagarshah95/Desktop/AverageRating_Youtube-0.0.1-SNAPSHOT.jar averagerating_youtube.AverageRating_Youtube /BigDataProject/data/youtubeDataset /BigDataProject/data/outputFiles/averagerating
2021-04-22 21:32:10,641 INFO client.DefaultNoHARMFailoverProxyProvider: Connecting to ResourceManager at /0.0.0.0:8032
2021-04-22 21:32:11,638 WARN mapreduce.JobContextUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
2021-04-22 21:32:11,783 INFO mapreduce.JobResourceUploader: ResourceUploader: Erasing Configuration for path: /tmp/hadoop-yarn/staging/sagarshah95/.staging/job_1619148809882_0006
2021-04-22 21:32:12,124 INFO InputFormat: Total input files in process : 1
2021-04-22 21:32:12,686 INFO mapreduce.JobSubmitter: number of splits:2
2021-04-22 21:32:13,156 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1619148809882_0006
2021-04-22 21:32:13,157 INFO mapreduce.JobSubmitter: Executing with tokens:[]
2021-04-22 21:32:13,161 INFO mapreduce.JobSubmitter: Configuration resource-type.xml found
2021-04-22 21:32:13,774 INFO resource.ResourceUtil: Unable to find 'resource-types.xml'.
2021-04-22 21:32:13,994 INFO impl.YarnClientImpl: Submitted application application_1619148809882_0006
2021-04-22 21:32:14,198 INFO mapreduce.Job: The url to track the job: http://ubuntu:8088/proxy/application_1619148809882_0006/
2021-04-22 21:32:14,198 INFO mapreduce.Job: Job job_1619148809882_0006 running in uber mode : false
2021-04-22 21:32:27,892 INFO mapreduce.Job: map 0% reduce 0%
2021-04-22 21:32:52,926 INFO mapreduce.Job: map 33% reduce 0%
2021-04-22 21:32:55,063 INFO mapreduce.Job: map 50% reduce 0%
2021-04-22 21:32:56,071 INFO mapreduce.Job: map 100% reduce 0%
2021-04-22 21:33:05,198 INFO mapreduce.Job: 100% reduce 100%
2021-04-22 21:33:05,219 INFO mapreduce.Job: Job job_1619148809882_0006 completed successfully
2021-04-22 21:33:05,376 INFO mapreduce.Job: Counters: 55
File System Counters
    FILE: Number of bytes read=28013117
    FILE: Number of bytes written=40019991
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=19939381
    HDFS: Number of bytes written=18106621
    HDFS: Number of read operations=11
    HDFS: Number of large read operations=2
    HDFS: Number of write operations=2
    HDFS: Number of bytes read erasure-coded=0
Job Counters
    Killed map tasks=1
    Launched map tasks=2
    Launched reduce tasks=1
    Data locality miss=2
    Total time spent by all maps in occupied slots (ms)=49800
    Total time spent by all reduces in occupied slots (ms)=6293
    Total time spent by all map tasks (ms)=49800
    Total time spent by all reduce tasks (ms)=6293
    Total time spent by all reduce tasks (ms)=6293
    Total vcore-milliseconds taken by all map tasks=49800
```

```

Job Counters
  Killed map tasks=1
  Launched map tasks=2
  Launched reduce tasks=1
  Data-local map tasks=2
  Total time spent by all maps in occupied slots (ms)=49800
  Total time spent by all reduces in occupied slots (ms)=6293
  Total time spent by all map tasks (ms)=49800
  Total time spent by all reduce tasks (ms)=6293
  Total vcore-milliseconds taken by all map tasks=49800
  Total vcore-milliseconds taken by all reduce tasks=6293
  Total megabyte-milliseconds taken by all map tasks=50995200
  Total megabyte-milliseconds taken by all reduce tasks=6444032

Map-Reduce Framework
  Map input records=769739
  Map output records=769735
  Map output bytes=18473641
  Map output materialized bytes=20013123
  Input split bytes=242
  Combine input records=769735
  Combine output records=769735
  Reduce input groups=746192
  Reduce shuffle bytes=20013123
  Reduce input records=769735
  Reduce output records=746192
  Spilled Records=1539470
  Shuffled Maps=2
  Failed Shuffles=0
  Merged Map outputs=2
  GC time elapsed (ms)=1306
  CPU time spent (ms)=25090
  Physical memory (bytes) snapshot=1093885952
  Virtual memory (bytes) snapshot=7762518016
  Total committed heap usage (bytes)=962592768
  Peak Map Physical memory (bytes)=461099008
  Peak Map Virtual memory (bytes)=2586202112
  Peak Reduce Physical memory (bytes)=220381184
  Peak Reduce Virtual memory (bytes)=2595991552

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=219819139
File Output Format Counters
  Bytes Written=18106621
sagarshah95@ubuntu:/usr/local/bin/hadoop-3.3.0/bin$ ■

sagarshah95@ubuntu:/usr/local/bin/hadoop-3.3.0/bin$ ./hadoop fs -cat /BigDataProject/data/outputFiles/averagerating/part-r-00000 | head
PF_2Nlw4rhS 1 3.5
---mkyh90bc 1 5.0
---nh-hn_3E 1 4.880000114440918
---x4K1JvQ0 1 3.799999952316284
-09WIapU0c 1 1.0
-0R69A3CVU 1 0.0
-0VHTCNYzs 1 4.550000190734863
-0eZhhav08 1 1.0
-0ts5lqos 1 4.380000114440918
-1K0JeTg2I 1 4.440000057220459
cat: Unable to write to output stream.
sagarshah95@ubuntu:/usr/local/bin/hadoop-3.3.0/bin$ ■

```

Youtuber based on number of videos uploaded: Returns the number of video uploaded by Youtuber (youtubeuploader)

This analysis comprises of :

1. Driver Class : Sets up configuration for Mapper, Reducer class
2. Mapper Class : Emits uploader of type Text and increments of 1 every time its occurs
3. Reducer Class : Emits uploader as a key of type Text and its total occurrence of type IntWritable

Mapper Class

```
public class Youtubetopuploader {

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

        private Text uploader = new Text();
        private final static IntWritable occurrence = new IntWritable(1);

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

            String record = value.toString();
            String str[] = record.split(",");
            if (str.length >= 7) {
                uploader.set(str[1]);
            }
            context.write(uploader, occurrence);
        }
    }
}
```

Reducer Class

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

    @Override
    public void reduce(Text key, Iterable<IntWritable> values,
                       Context context) throws IOException, InterruptedException {
        int totaloccurrence = 0;

        for (IntWritable value : values) {
            totaloccurrence += value.get();
        }
        context.write(key, new IntWritable(totaloccurrence));
    }
}
```

Driver Class

```
public static void main(String[] args) throws IOException, ClassNotFoundException, InterruptedException {  
  
    Configuration conf1 = new Configuration();  
  
    @SuppressWarnings("deprecation")  
    Job job = new Job(conf1, "myyoutube");  
  
    job.setJarByClass(Youtubetopuploader.class);  
    job.setMapperClass(Map.class);  
    job.setReducerClass(Reduce.class);  
  
    job.setOutputKeyClass(Text.class);  
    job.setOutputValueClass(IntWritable.class);  
  
    job.setInputFormatClass(TextInputFormat.class);  
    job.setOutputFormatClass(TextOutputFormat.class);  
  
    FileInputFormat.addInputPath(job, new Path(args[0]));  
    FileOutputFormat.setOutputPath(job, new Path(args[1]));  
  
    System.exit(job.waitForCompletion(true) ? 0 : 1);  
}
```

Execution



Total Views on a Video : Returns the total number views on a video (Youtube_VIEWS)

This analysis comprises of :

1. Driver Class : Sets up configuration for Mapper, Reducer class
2. Mapper Class : Emits videoID of type Text and views of type FloatWritable
3. Reducer Class : Emits videoID as key of type Text and total occurrence of type FloatWritable

Mapper

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

    private Text video_name = new Text();
    private FloatWritable views = new FloatWritable();

    @Override
    public void map(LongWritable key, Text value, Context context)
        throws IOException, InterruptedException {
        String line = value.toString();
        String str[] = line.split(",");
        try {
            if (str.length >= 5) {
                video_name.set(str[0]);
                float temp = Float.parseFloat(str[5]); //typecasting string to Integer
                views.set(temp);
            }

            context.write(video_name, views);
        } catch (Exception e) {
            e.printStackTrace();
        }
    }
}
```

Reducer

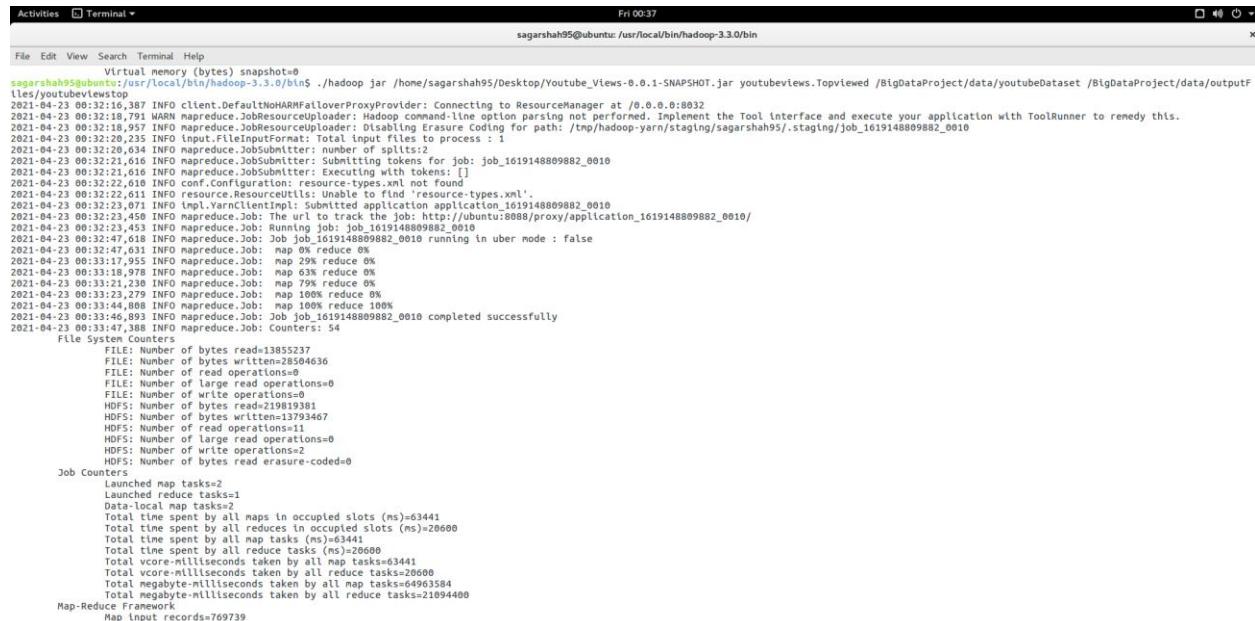
```
public static class Reduce| extends Reducer<Text, FloatWritable, Text, FloatWritable> {

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

Driver

```
@SuppressWarnings("deprecation")
public static void main(String[] args) throws Exception {
    Configuration conf = new Configuration();
    Job job = new Job(conf, "Top Videos");
    job.setJarByClass(Topviewed.class);
    job.setMapOutputKeyClass(Text.class);
    job.setMapOutputValueClass(FloatWritable.class);
    job.setOutputKeyClass(Text.class);
    job.setOutputValueClass(FloatWritable.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);
}
```

Execution



The screenshot shows a terminal window titled 'Terminal' with the command 'Virtual memory (bytes) snapshot@' entered. The output is as follows:

```
sagarshah95@ubuntu:/usr/local/bin/hadoop-3.3.0/bin$ ./hadoop jar /home/sagarshah95/Desktop/Youtube_Views-0.0.1-SNAPSHOT.jar youtubeviews.Topviewed /BigDataProject/data/youtubeDataset /BigDataProject/data/outputFiles/youtubebevewstop
2021-04-23 00:32:16,367 WARN mapreduce.JobResourceUploader: Connecting to ResourceManager at /0.0.0.0:8032
2021-04-23 00:32:18,791 INFO mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool Interface and execute your application with ToolRunner to remedy this.
2021-04-23 00:32:18,957 INFO mapreduce.JobResourceUploader: Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/sagarshah95/.staging/job_1619148889882_0010
2021-04-23 00:32:20,235 INFO input.FileInputFormat: Total input files to process : 1
2021-04-23 00:32:20,634 INFO mapreduce.JobSubmitter: number of splits:2
2021-04-23 00:32:21,000 INFO mapreduce.JobSubmitter: Submitting token for job: job_1619148889882_0010
2021-04-23 00:32:21,618 INFO mapreduce.JobSubmitter: Number of的心 tokens: []
2021-04-23 00:32:22,618 INFO conf.Configuration: resource-types.xml not found
2021-04-23 00:32:22,618 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.
2021-04-23 00:32:23,071 INFO impl.YarnClientImpl: Submitted application application_1619148889882_0010
2021-04-23 00:32:23,071 INFO mapreduce.Job: The url to track the job: http://ubuntu:8088/proxy/application_1619148889882_0010/
2021-04-23 00:32:23,451 INFO mapreduce.Job: Running in uber mode : false
2021-04-23 00:32:47,618 INFO mapreduce.Job: Job job_1619148889882_0010 running in uber mode : false
2021-04-23 00:32:47,631 INFO mapreduce.Job: map 0% reduce 0%
2021-04-23 00:33:17,955 INFO mapreduce.Job: map 29% reduce 0%
2021-04-23 00:33:18,238 INFO mapreduce.Job: map 79% reduce 0%
2021-04-23 00:33:23,238 INFO mapreduce.Job: map 79% reduce 0%
2021-04-23 00:33:23,279 INFO mapreduce.Job: map 100% reduce 0%
2021-04-23 00:33:44,808 INFO mapreduce.Job: Job job_1619148889882_0010 completed successfully
2021-04-23 00:33:46,893 INFO mapreduce.Job: Job job_1619148889882_0010 completed successfully
2021-04-23 00:33:47,140 INFO mapreduce.Job: Counters: 54
  File System Counters:
    FILE: Number of bytes read=13855237
    FILE: Number of bytes written=28504636
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=21919381
    HDFS: Number of bytes written=13793467
    HDFS: Number of read operations=1
    HDFS: Number of large read operations=0
    HDFS: Number of write operations=2
    HDFS: Number of bytes read erasure-coded=0
  Job Counters:
    Launched map tasks=2
    Launched reduce tasks=1
    Data-local map tasks=2
    Total time spent by all maps in occupied slots (ms)=63441
    Total time spent by all reduces in occupied slots (ms)=20600
    Total time spent by all map tasks (ms)=26800
    Total time spent by all reduce tasks (ms)=26800
    Total vcore-milliseconds taken by all map tasks=63441
    Total vcore-milliseconds taken by all reduce tasks=20600
    Total negabyte-milliseconds taken by all map tasks=64963584
    Total negabyte-milliseconds taken by all reduce tasks=21094400
  Map-Reduce Framework
    Map Input records=769739
```

```

Launched map tasks=2
Launched reduce tasks=1
Data-local map tasks=2
Total time spent by all maps in occupied slots (ms)=63441
Total time spent by all reduces in occupied slots (ms)=20600
Total time spent by all map tasks (ms)=63441
Total time spent by all reduce tasks (ms)=20600
Total vcore-milliseconds taken by all map tasks=63441
Total vcore-milliseconds taken by all reduce tasks=20600
Total megabyte-milliseconds taken by all map tasks=64963584
Total megabyte-milliseconds taken by all reduce tasks=21094400
Map-Reduce Framework
Map Input records=769739
Map Output records=769735
Map Input bytes=1539470
Map output materialized bytes=13855243
Input split bytes=242
Combine input records=0
Combine output records=0
Reduce input groups=740427
Reduce shuffle bytes=13855243
Reduce output bytes=13855243
Reduce output records=740427
Spilled Records=1539470
Shuffled Maps =2
Failed Shuffles=0
Merged Map outputs=2
GC time elapsed (ms)=1047
CPU time spent (ms)=30620
Physical memory (bytes) snapshot=1129230336
Virtual memory (bytes) snapshot=7768645632
Total committed heap usage (bytes)=1007681536
Peak Map Physical memory (bytes)=468340736
Peak Map Virtual memory (bytes)=2587058176
Peak Reduce Physical memory (bytes)=249479168
Peak Reduce Virtual memory (bytes)=2595913728
Shuffle Errors
  BAD_ID=0
  CONNECTION=0
  IO_ERROR=0
  WRONG_LENGTH=0
  WRONG_MAP=0
  WRONGReducer=0
File Input Format Counters
  Bytes Read=219819139
File Output Format Counters
  Bytes Written=13793467
sagarshah95@ubuntu:/usr/local/bin/hadoop-3.3.0/bin$ ■

```

```

sagarshah95@ubuntu:/usr/local/bin/hadoop-3.3.0/bin$ ./hadoop fs -cat /BigDataProject/data/outputFiles/youtubeviewstop/part-r-00000 | head
PF_ZMlw4rHs    96.0
---mykyh90bc   1616.0
---nH-hN_3E    1715.0
---x4KIjVQ0    7594.0
--09WIapUUC   154.0
--OR69A3CVU   202.0
--0VhtCnzyz   2357.0
--0eZhhAv08   892.0
--0tsSllqos   1923.0
--1K0JeTg2I    2010.0
cat: Unable to write to output stream.

```

Total category occurance : Returns the total number of occurrence of each category (Youtube Categories)

This analysis comprises of :

1. Driver Class : Sets up configuration for Mapper, Reducer class
2. Mapper Class : Emits category of type Text and 1 on every occurrence of type IntWritable
3. Reducer Class : Emits category as a key of type Text and total occurrence of type IntWritable

Mapper

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

    private Text category = new Text();
    private final static IntWritable occurrence = new IntWritable(1);

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

        String record = value.toString();
        String str[] = record.split(",");
        if (str.length > 5) {
            category.set(str[3]);
        }
        context.write(category, occurrence);
    }
}
```

Reducer

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

    @Override
    public void reduce(Text key, Iterable<IntWritable> values,
                       Context context) throws IOException, InterruptedException {
        int totaloccurrence = 0;

        for (IntWritable value : values) {
            totaloccurrence += value.get();
        }
        context.write(key, new IntWritable(totaloccurrence));
    }
}
```

Driver

```
@SuppressWarnings("deprecation")
public static void main(String[] args) throws IOException, ClassNotFoundException, InterruptedException {

    Configuration conf = new Configuration();
    Job job = new Job(conf, "myyoutube");
    job.setJarByClass(Youtube_DataAnalysis.class);
    job.setMapperClass(Map.class);
    job.setReducerClass(Reduce.class);

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

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

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

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

Execution

```
...  
sagarshah95@ubuntu:~$ /usr/local/bin/hadoop-3.3.0/bin$ ./hadoop jar /home/sagarshah95/Desktop/Youtube_Categories-0.0.1-SNAPSHOT.jar youtube_dataanalysis.Youtube_DataAnalysis /BigDataProject/data/youtubeDataset /Big  
DataProject/data/outputFiles/youtubeanalysis  
2021-04-23 12:37:03.646 INFO client.DefaultHttpAclFollowerProxyProvider: Connecting to ResourceManager at /0.0.0.0:8032  
2021-04-23 12:37:04.089 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.  
2021-04-23 12:37:04.112 INFO mapreduce.JobResourceUploader: Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/sagarshah95/.staging/job_1619206134623_0003  
2021-04-23 12:37:04.421 INFO Input.FileInputFormat: Total input files to process : 1  
2021-04-23 12:37:05.000 INFO mapreduce.Job: Job job_1619206134623_0003 running in uber mode : false  
2021-04-23 12:37:05.315 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1619206134623_0003  
2021-04-23 12:37:05.315 INFO mapreduce.JobSubmitter: Executing with tokens: []  
2021-04-23 12:37:05.621 INFO conf.Configuration: resource-types.xml not found  
2021-04-23 12:37:05.621 INFO resource.ResourceCells: Unable to find 'resource-types.xml'.  
2021-04-23 12:37:05.621 INFO resource.ResourceCells: Using default application_1619206134623_0003  
2021-04-23 12:37:06.113 INFO mapreduce.Job: Submitter updated application_1619206134623_0003  
2021-04-23 12:37:06.191 INFO mapreduce.Job: The url to track the job: http://ubuntu:8088/proxy/application_1619206134623_0003/  
2021-04-23 12:37:06.192 INFO mapreduce.Job: Running job: job_1619206134623_0003  
2021-04-23 12:37:15.403 INFO mapreduce.Job: map 10% reduce 0%  
2021-04-23 12:37:15.403 INFO mapreduce.Job: map 100% reduce 0%  
2021-04-23 12:37:27.687 INFO mapreduce.Job: map 100% reduce 0%  
2021-04-23 12:37:34.692 INFO mapreduce.Job: map 100% reduce 100%  
2021-04-23 12:37:35.717 INFO mapreduce.Job: Job job_1619206134623_0003 completed successfully  
2021-04-23 12:37:35.717 INFO mapreduce.Job: Counters:  
  File System Counters  
    FILE: Number of bytes read=13268663  
    FILE: Number of bytes written=27330468  
    FILE: Number of read operations=0  
    FILE: Number of large read operations=0  
    FILE: Number of write operations=0  
    HDFS: Number of bytes read=19819381  
    HDFS: Number of bytes written=266  
    HDFS: Number of read operations=11  
    HDFS: Number of large read operations=0  
    HDFS: Number of write operations=2  
    HDFS: Number of bytes read erasure-coded=0  
  Job Counters  
    Killed map tasks=1  
    Launched map tasks=2  
    Launched reduce tasks=1  
    Data-local map tasks=2  
    Total time spent by all map tasks in occupied slots (ms)=18110  
    Total time spent by all reducers in occupied slots (ms)=5479  
    Total time spent by all map tasks (ms)=18110  
    Total time spent by all reduce tasks (ms)=5479  
    Total vcore-milliseconds taken by all map tasks=18110  
    Total vcore-milliseconds taken by all reduce tasks=5479  
    Total megabyte-milliseconds taken by all map tasks=18544640  
    Total megabyte-milliseconds taken by all reduce tasks=5610496  
  Map-Reduce Framework  
    Map input records=769739  
  ...  
sagarshah95@ubuntu:~$
```

```
Job Counters  
  Killed map tasks=1  
  Launched map tasks=2  
  Launched reduce tasks=1  
  Data-local map tasks=2  
  Total time spent by all maps in occupied slots (ms)=18110  
  Total time spent by all reduces in occupied slots (ms)=5479  
  Total time spent by all map tasks (ms)=18110  
  Total time spent by all reduce tasks (ms)=5479  
  Total vcore-milliseconds taken by all map tasks=18110  
  Total vcore-milliseconds taken by all reduce tasks=5479  
  Total megabyte-milliseconds taken by all map tasks=18544640  
  Total megabyte-milliseconds taken by all reduce tasks=5610496  
Map-Reduce Framework  
  Map input records=769739  
  Map output records=769739  
  Map output bytes=11729179  
  Map output materialized bytes=13268669  
  Input split bytes=242  
  Combine input records=0  
  Combiner input records=0  
  Reduce input groups=15  
  Reduce shuffle bytes=13268669  
  Reduce input records=769739  
  Redundant input records=15  
  Spilled Records=1539478  
  Shuffled Maps =2  
  Failed Shuffles=0  
  Merged Map outputs=2  
  GC (bytes)=1086967808  
  CPU time spent (ms)=9150  
  Physical memory (bytes) snapshot=1086967808  
  Virtual memory (bytes) snapshot=7761457152  
  Total committed heap usage (bytes)=945614416  
  Peak committed physical memory (bytes)=103338632  
  Peak Map Virtual memory (bytes)=2586505216  
  Peak Reduce Physical memory (bytes)=191303680  
  Peak Reduce Virtual memory (bytes)=2589937664  
Shuffle Bytes  
  BAD ID=0  
  CONNECTION=0  
  IO_ERROR=0  
  WRONG_LENGTH=0  
  WRONG_TYPE=0  
  WRONG_REDUCE=0  
File Input Format Counters  
  Bytes Read=219819139  
File Output Format Counters  
  Bytes Written=266  
sagarshah95@ubuntu:~$
```

```
sagarshah95@ubuntu:/usr/local/bin/hadoop-3.3.0/bin$ ./hadoop fs -cat /BigDataProject/data/outputFiles/youtubeanalysis/part-r-00000 | tail  
Entertainment 132087  
Film & Animation 76010  
Gadgets & Games 61419  
Howto & DIY 18887  
Music 184957  
News & Politics 37739  
People & Blogs 51245  
Pets & Animals 10782  
Sports 69113  
Travel & Places 15093  
sagarshah95@ubuntu:/usr/local/bin/hadoop-3.3.0/bin$
```

Top Youtube Videos based on Ratings : Returns top Rated videos in the sorted manner descendingly (Top_Youtuber)

This analysis comprises of :

1. Driver Class : Sets up configuration for Mapper, Reducer class
2. Mapper Class : Emits videoID and rating
3. Reducer Class : Emits videoID as a key of type Text and ratings of type Floatwritable sorted in the descending order

Mapper

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

    private FloatWritable video_rating = new FloatWritable();
    private Text video_id = new Text();

    public void map(Object key, Text value, Mapper.Context context
    ) throws IOException, InterruptedException {
        String[] fields = value.toString().split(",");
        video_id = new Text(fields[0]);
        try {
            if(fields.length > 6) {
                //if (!fields[6].isEmpty()) {
                    video_rating = new FloatWritable(Float.parseFloat(fields[7]));
                //}
            }
            context.write(video_id, video_rating);
        }catch (Exception e) {
            e.printStackTrace();
        }
    }
}
```

Reducer

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

    private Map<Text, FloatWritable> countMap = new HashMap<>();

    @Override
    public void reduce(Text key, Iterable<FloatWritable> values, Context context) throws IOException, InterruptedException {
        // computes the number of occurrences of a single word
        float sum = 0.0f;
        int count = 0;

        for (FloatWritable val : values) {
            sum += val.get();
            count++;
        }

        countMap.put(new Text(key), new FloatWritable(sum / count));
    }
}
```

Driver

```
public static void main(String[] args) throws Exception {
    Configuration conf = new Configuration();
    Job job = Job.getInstance(conf, "Top50");
    job.setJarByClass(Top_Youtuber.class);
    job.setMapperClass(TopNMapper.class);
    job.setReducerClass(TopNReducer.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);
}
```

Execution

```
Player || Activities Terminal Fri 16:19 sagarshah95@ubuntu: /usr/local/bin/hadoop-3.3.0/bin
File Edit View Search Terminal Help
    at java.lang.reflect.Method.invoke(Method.java:498)
    at org.apache.hadoop.util.RunJar.run(RunJar.java:233)
    at org.apache.hadoop.util.RunJar.main(RunJar.java:236)
sagarshah95@ubuntu:~/local/bin/hadoop-3.3.0/bin$ ./hadoop jar /home/sagarshah95/Desktop/Top_Youtuber-0.0.1-SNAPSHOT.jar com.neu.bigdata.Top_Youtuber /BigDataProject/data/youtubeDataset /BigDataProject/data/outputFiles/topyoutubeur
:putFiles/topyoutubeur
021-04-23 16:13:59,632 INFO Client: DeafultHDFSFallbackProxyProvider: Connecting to ResourceManager at /0.0.0.0:8032
021-04-23 16:14:00,467 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool Interface and execute your application with ToolRunner to remedy this.
021-04-23 16:14:00,536 INFO mapreduce.JobResourceUploader: Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/sagarshah95.staging/job_1619206134623_0007
021-04-23 16:14:00,540 INFO InputFormatFactory: Total number of splits: 1
021-04-23 16:14:01,298 INFO mapreduce.JobSubmitter: number of splits:2
021-04-23 16:14:02,088 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1619206134623_0007
021-04-23 16:14:02,081 INFO mapreduce.JobSubmitter: Executing with tokens: []
021-04-23 16:14:02,528 INFO ConfigurationResourceManager: Found configuration resource types.xml
021-04-23 16:14:02,654 INFO YarnClientImpl: Submitted application application_1619206134623_0007
021-04-23 16:14:02,793 INFO mapreduce.Job: The url to track the job: http://ubuntu:8088/proxy/application_1619206134623_0007/
021-04-23 16:14:02,793 INFO mapreduce.Job: Running job: job_1619206134623_0007
021-04-23 16:14:02,793 INFO mapreduce.Job: 100% complete
021-04-23 16:14:35,188 INFO mapreduce.Job: map 50% reduce 0%
021-04-23 16:14:35,776 INFO mapreduce.Job: map 50% reduce 0%
021-04-23 16:14:36,787 INFO mapreduce.Job: map 100% reduce 0%
021-04-23 16:14:48,925 INFO mapreduce.Job: map 100% reduce 100%
021-04-23 16:14:50,169 INFO mapreduce.Job: Job: job_1619206134623_0007 completed successfully
021-04-23 16:14:50,169 INFO mapreduce.Job: Counters: 55
    File System Counters
        FILE: Number of bytes read=13855237
        FILE: Number of bytes written=2850269
        FILE: Number of read operations=0
        FILE: Number of large read operations=0
        FILE: Number of write operations=0
        HDFS: Number of bytes read=219819381
        HDFS: Number of bytes written=800
        HDFS: Number of read operations=1
        HDFS: Number of large read operations=0
        HDFS: Number of write operations=2
        HDFS: Number of bytes read erasure-coded=0
    Job Counters
        Launched map tasks=1
        Launched map tasks=2
        Launched reduce tasks=1
        Data-local map tasks=2
        Total time spent by all map tasks in occupied slots (ms)=38281
        Total time spent by all map tasks in occupied slots (ms)=10186
        Total time spent by all map tasks (ms)=58281
        Total time spent by all reduce tasks (ms)=10186
        Total vcore-milliseconds taken by all map tasks=38281
        Total vcore-milliseconds taken by all reduce tasks=10186
        Total megabytes-milliseconds taken by all map tasks=39199744
        Total megabytes-milliseconds taken by all reduce tasks=10430464
Map-Reduce Framework
```

```

Killed map tasks=1
Launched map tasks=2
Launched reduce tasks=1
Data-local map tasks=2
Total time spent by all maps in occupied slots (ms)=38281
Total time spent by all map tasks (ms)=38281
Total time spent by all reduce tasks (ms)=10186
Total time spent by all reduce tasks (ns)=10186
Total vcore-milliseconds taken by all map tasks=38281
Total vcore-milliseconds taken by all reduce tasks=10186
Total vcore-milliseconds taken by all map tasks=39199744
Total megabyte-milliseconds taken by all reduce tasks=10430464
Map-Reduce Framework
  Map input records=769739
  Map output records=769735
  Map output bytes=13855243
  Map output materialized bytes=13855243
  Input split bytes=242
  Combine input records=0
  Combine output records=0
  Reduce input records=1539470
  Reduce output records=50
  Spilled Records=1539470
  Shuffled Maps=0
  Failed Shuffles=0
  Merged Map outputs=2
  GC time elapsed (ms)=1459
  CPU time spent (ms)=1836
  Physical memory snapshot=1173762948
  Virtual memory (bytes) snapshot=7771832320
  Total committed heap usage (bytes)=1024983040
  Peak Map Physical memory (bytes)=451727360
  Peak Map Virtual memory (bytes)=2586181632
  Peak Reduce Physical memory (bytes)=302112768
  Peak Reduce Virtual memory (bytes)=2600153088
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=219819139
File Output Format Counters
  Bytes Written=888
sagarshah95@ubuntu:/usr/local/bin/hadoop-3.3.0/bin$ 

```

```

sagarshah95@ubuntu:/usr/local/bin/hadoop-3.3.0/bin$ ./hadoop fs -cat /BigDataProject/data/outputFiles/Top_Youtuber/part-r-00000 | head
QjA5faZf1A8    120506.0
dMHObHeIRNg   87520.0
0XI-hvPRRA    80710.0
R0049_tDAUB   70972.0
noJWJ6-XmeQ    62265.0
Jahdn0Q9XCA   59008.0
VcQIwbvGRKU   46472.0
sdUUXsFdySS   42417.0
pv5zWaTEVkI   42386.0
D2kJZOfq7zk   42162.0
sagarshah95@ubuntu:/usr/local/bin/hadoop-3.3.0/bin$ 

```

Rating Summarization: Provides a summarization of ratings, rate and comments

This analysis comprises of :

1. Driver Class : Sets up configuration for Mapper, Reducer, Tuple class
2. Mapper Class : Emits videoID of type Text and Tuple of custom type MinMaxCountTuple
3. Reducer Class : Emits videoID as key of type Text and Tuple of custom type MinMaxCountTuple
4. Tuple Class : Pojo class for Tuple ratings, rate and comments

Mapper

```
class MapperClass extends Mapper<Object, Text, Text, MinMaxCountTuple> {

    private Text video_ID = new Text();
    private MinMaxCountTuple outTuple = new MinMaxCountTuple();

    protected void map(Object key, Text value, Context context) throws IOException, InterruptedException {
        String[] input = value.toString().split(",");
        video_ID.set(input[0]);
        if(input.length > 8) {
            try {
                outTuple.setTotalRating(Float.valueOf(input[7]));
                outTuple.setAverageRating(Float.valueOf(input[6]));
                outTuple.setTotalComment(Float.valueOf(input[8]));

                context.write(video_ID, outTuple);
            }catch(Exception e) {
                e.printStackTrace();
            }
        }
    }
}
```

Reducer

```
class ReducerClass extends Reducer<Text, MinMaxCountTuple, Text, MinMaxCountTuple> {

    private MinMaxCountTuple result = new MinMaxCountTuple();

    @Override
    protected void reduce(Text key, Iterable<MinMaxCountTuple> values, Context context) throws IOException, InterruptedException {
        // Initialize our result
        result.setAverageRating(0);
        result.setTotalRating(0);
        result.setTotalComment(0);
        int sum = 0;

        for (MinMaxCountTuple val : values) {
            //max
            if (result.getAverageRating() == 0 || val.getAverageRating() < result.getAverageRating()) {
                result.setAverageRating(val.getAverageRating());
            }
            //min
            if (result.getTotalRating() == 0
                || val.getTotalRating() > (result.getTotalRating())) {
                result.setTotalRating(val.getTotalRating());
            }
            //sum
            sum += val.getTotalComment();
        }
        result.setTotalComment(sum);
        context.write(key, result);
    }
}
```

Driver

```
'public static void main(String[] args) throws IOException {
    try {
        Configuration conf = new Configuration();
        Job job = Job.getInstance(conf, "VideoMinMaxRating");
        job.setJarByClass(Rating_Summarization.class);
        job.setMapperClass(MapperClass.class);
        job.setCombinerClass(ReducerClass.class);
        job.setReducerClass(ReducerClass.class);
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(MinMaxCountTuple.class);
        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));
        System.exit(job.waitForCompletion(true) ? 0 : 1);
    } catch (InterruptedException | ClassNotFoundException ex) {
        Logger.getLogger(Rating_Summarization.class.getName()).log(Level.SEVERE, null, ex);
    }
}
```

Execution

```
sagarshah95@ubuntu:/usr/local/bin/hadoop-3.3.0/bin$ ./hadoop jar /home/sagarshah95/Desktop/Rating_Summarization-0.0.1-SNAPSHOT.jar com.neu.bigdata.Rating_Summarization /BigDataProject/data/youtubeDataset /BigDataProject/data/outputfiles/ratingsummari
1921-04-23 18:20:30,422 INFO Client: DatanodeManager: Connecting to ResourceManager at 0.0.0.0:8082
1921-04-23 18:20:30,422 WARN mapreduce.JobResourceUploader: Hadoop Command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
1921-04-23 18:20:30,432 INFO mapreduce.JobResourceUploader: Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/sagarshah95/.staging/job_1619206134623_0010
1921-04-23 18:20:31,023 INFO FileInputFormat: Total Input files to process : 1
1921-04-23 18:20:31,265 INFO mapreduce.JobSubmitter: number of splits:2
1921-04-23 18:20:31,265 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1619206134623_0010
1921-04-23 18:20:31,576 INFO mapreduce.JobSubmitter: Executing with tokens: []
1921-04-23 18:20:32,213 INFO conf.Configuration: resource-types.xml not found
1921-04-23 18:20:32,214 INFO conf.ResourceUtils: Unable to find 'resource-types.xml'.
1921-04-23 18:20:32,429 INFO impl.YarnClientImpl: Submitted application application_1619206134623_0010
1921-04-23 18:20:32,430 INFO impl.YarnClientImpl: Application report from cluster: application_1619206134623_0010
1921-04-23 18:20:32,563 INFO mapreduce.Job: Running job: job_1619206134623_0010
1921-04-23 18:20:46,114 INFO mapreduce.Job: Job job_1619206134623_0010 running in uber mode : False
1921-04-23 18:21:05,692 INFO mapreduce.Job: map 0% reduce 0%
1921-04-23 18:21:05,692 INFO mapreduce.Job: map 50% reduce 0%
1921-04-23 18:21:05,692 INFO mapreduce.Job: map 100% reduce 0%
1921-04-23 18:21:20,971 INFO mapreduce.Job: map 100% reduce 100%
1921-04-23 18:21:21,992 INFO mapreduce.Job: Job job_1619206134623_0010 completed successfully
1921-04-23 18:21:22,178 INFO mapreduce.Job: Counters
File System Counters
  File system bytes read=19221109
  FILE: Number of bytes written=39294541
  FILE: Number of read operations=0
  FILE: Number of large read operations=0
  FILE: Number of write operations=0
  HDFS: Number of bytes read=18628531
  HDFS: Number of bytes written=18628537
  HDFS: Number of read operations=11
  HDFS: Number of large read operations=2
  HDFS: Number of write operations=2
  HDFS: Number of bytes read erasure-coded=0
Job Counters
  Killed map tasks=1
  Launched map tasks=2
  Launched reduce tasks=1
  Data-local map tasks=2
  Total time spent by all maps in occupied slots (ms)=36350
  Total time spent by all reduces in occupied slots (ms)=10990
  Total time spent by all map tasks (ms)=36350
```

```
Job Counters
  Killed map tasks=1
  Launched map tasks=2
  Launched reduce tasks=1
  Data-local map tasks=2
  Total time spent by all maps in occupied slots (ms)=36350
  Total time spent by all reduces in occupied slots (ms)=10990
  Total time spent by all map tasks (ms)=36350
  Total time spent by all reduce tasks (ms)=10990
  Total vcore-milliseconds taken by all map tasks=36350
  Total vcore-milliseconds taken by all reduce tasks=10990
  Total megabyte-milliseconds taken by all map tasks=37222460
  Total megabyte-milliseconds taken by all reduce tasks=11253760

Map-Reduce Framework
  Map input records=769739
  Map output records=763885
  Map output bytes=18333241
  Map output materialized bytes=19251115
  Input splits=1
  Combine input records=763885
  Combine output records=740427
  Reduce input groups=740427
  Reduce shuffle bytes=19251115
  Reduce input records=740427
  Reduce output records=740427
  Spilled Records=1480854
  Shuffled Maps =2
  Failed Maps =0
  Merged Map outputs=2
  GC time elapsed (ms)=626
  CPU time spent (ns)=20110
  Physical memory (bytes) snapshot=1072840704
  Virtual memory (bytes) snapshot=7770677248
  Total committed heap usage (bytes)=10728416
  Peak Virtual memory (bytes)=1072843220
  Peak Map Virtual memory (bytes)=2588856320
  Peak Reduce Physical memory (bytes)=229715968
  Peak Reduce Virtual memory (bytes)=2594189312

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=219819139
File Output Format Counters
  Bytes Written=18628537
agarshah95@ubuntu:/usr/local/bin/hadoop-3.3.0/bin$ █
```

```
agarshah@5gubuntu:~/usr/local/bin/hadoop-3.3.0/bin$ ./hadoop fs -cat /BigDataProject/data/outputFiles/ratingsummarization/part-r-00000 | head  
PF_ZhW4rHs 3.5 2.0 3.0  
--mkyh90b 5.0 2.0 1.0  
--HN_HN_3E 4.88 8.0 3.0  
--091190 3.0 16.0 3.0  
--09W1apUUC 1.0 1.0 1.0  
--0R69A3CVU 0.0 0.0 0.0  
--0VhtCNYzs 4.55 51.0 6.0  
--0eZhhvaV08 1.0 1.0 3.0  
--0ts5Llqos 4.38 8.0 5.0  
--1K0Jefg2I 4.44 9.0 9.0  
cat: Unable to write to output stream.  
agarshah@5gubuntu:~/usr/local/bin/hadoop-3.3.0/bin$
```

Binning by Categories : Performed binning based on categories to output multiple files per bin

This analysis comprises of :

1. Driver Class : Sets up configuration for Mapper class
 2. Mapper Class : Uses MultipleOutput class to write to multiple output files based on the categories. Number of bins determine how many output files will be created in Binning

Mapper

```
public static class YouTubeBinMapper extends Mapper<Object, Text, Text, NullWritable> {  
    private MultipleOutputs<Text, NullWritable> mos = null;  
  
    @Override  
    protected void setup(Mapper.Context context) throws IOException, InterruptedException {  
        mos = new MultipleOutputs<Text, NullWritable>(context);  
    }  
  
    @Override  
    protected void map(Object key, Text value, Mapper.Context context)  
        throws IOException, InterruptedException {  
        String[] input = value.toString().split(",");  
        if (input.length > 2) {  
            Text Name = new Text(input[3]);  
            String line = Name.toString();  
            if (line.contains("UNA ")) {  
                mos.write("bins", value, NullWritable.get(), "UNA");  
            } else if (line.contains("Autos & Vehicles")) {  
                mos.write("bins", value, NullWritable.get(), "Autos & Vehicles");  
            } else if (line.contains("Comedy")) {  
                mos.write("bins", value, NullWritable.get(), "Comedy");  
            } else if (line.contains("Entertainment")) {  
                mos.write("bins", value, NullWritable.get(), "Entertainment");  
            } else if (line.contains("Film & Animation")) {  
                mos.write("bins", value, NullWritable.get(), "Film & Animation");  
            } else if (line.contains("Gadgets & Games")) {  
                mos.write("bins", value, NullWritable.get(), "Gadgets & Games");  
            } else if (line.contains("Howto & DIY")) {  
                mos.write("bins", value, NullWritable.get(), "Howto & DIY");  
            } else if (line.contains("Music")) {  
                mos.write("bins", value, NullWritable.get(), "Music");  
            } else if (line.contains("News & Politics")) {  
                mos.write("bins", value, NullWritable.get(), "News & Politics");  
            } else if (line.contains("People & Blogs")) {  
                mos.write("bins", value, NullWritable.get(), "People & Blogs");  
            } else if (line.contains("Pets & Animals")) {  
                mos.write("bins", value, NullWritable.get(), "Pets & Animals");  
            } else if (line.contains("Sports")) {  
                mos.write("bins", value, NullWritable.get(), "Sports");  
            } else if (line.contains("Travel & Places")) {  
                mos.write("bins", value, NullWritable.get(), "Travel & Places");  
            } else {  
                mos.write("bins", value, NullWritable.get(), "UnCatogized");  
            }  
        }  
    }  
    @Override  
    protected void cleanup(Mapper.Context context)  
        throws IOException, InterruptedException {  
        mos.close();  
    }  
}
```

Driver

```
public static void main(String[] args) throws Exception {
    Configuration conf = new Configuration();
    Job job = new Job(conf, "Binning");
    job.setJarByClass(BinningByCategories.class);
    job.setMapperClass(YouTubeBinMapper.class);
    job.setNumReduceTasks(0);

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

    MultipleOutputs.addNamedOutput(job, "bins", TextOutputFormat.class,
        Text.class, NullWritable.class);

    MultipleOutputs.setCountersEnabled(job, true);

    System.exit(job.waitForCompletion(true) ? 0 : 2);
}
```

Execution

```
File Edit View Search Terminal Help
Virtual memory (bytes) snapshot=0
agarshah95@ubuntu:/usr/local/bin/hadoop-3.3.0/bin$ ./hadoop jar /home/sagarshah95/Desktop/BinningByCategories-0.0.1-SNAPSHOT.jar com.neu.bigdata.BinningByCategories /BigDataProject/data/youtubeDataset /BigDataProject/dataset/BinnedDataset
'021-04-23 22:46:19,181 INFO client.DefaultHARMFailoverProxyProvider: Connecting to ResourceManager at /0.0.0.0:8032
'021-04-23 22:46:19,181 INFO client.DefaultHARMFailoverProxyProvider: Connecting to ResourceManager at /0.0.0.0:8032
'021-04-23 22:46:20,148 INFO mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
'021-04-23 22:46:20,148 INFO mapreduce.JobResourceUploader: Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/sagarshah95/.staging/job_1619206134623_0012
'021-04-23 22:46:21,243 INFO input.FileInputFormat: Total Input files to process : 1
'021-04-23 22:46:21,723 INFO mapreduce.JobSubmitter: number of attempts:2
'021-04-23 22:46:22,125 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1619206134623_0012
'021-04-23 22:46:22,436 INFO mapreduce.JobSubmitter: Executing with tokens: []
'021-04-23 22:46:22,436 INFO mapreduce.JobSubmitter: Submitting token for job: job_1619206134623_0012
'021-04-23 22:46:23,485 INFO conf.Configuration: resource-types.xml not found
'021-04-23 22:46:23,487 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.
'021-04-23 22:46:23,779 INFO yarn.Client: Submitted application application_1619206134623_0012
'021-04-23 22:46:24,045 INFO mapreduce.Job: User code provided no map or reduce functions; falling back to identity functions for these stages.
'021-04-23 22:46:24,018 INFO mapreduce.Job: Running job: job_1619206134623_0012
'021-04-23 22:46:46,852 INFO mapreduce.Job: Job job_1619206134623_0012 running in uber mode : False
'021-04-23 22:46:46,856 INFO mapreduce.Job: map 0% reduce 0%
'021-04-23 22:46:46,922 INFO mapreduce.Job: map 100% reduce 0%
'021-04-23 22:47:23,022 INFO mapreduce.Job: map 41% reduce 0%
'021-04-23 22:47:31,594 INFO mapreduce.Job: map 52% reduce 0%
'021-04-23 22:47:32,642 INFO mapreduce.Job: map 69% reduce 0%
'021-04-23 22:47:38,158 INFO mapreduce.Job: map 81% reduce 0%
'021-04-23 22:47:45,152 INFO mapreduce.Job: map 100% reduce 0%
'021-04-23 22:47:45,474 INFO mapreduce.Job: map 100% reduce 0%
'021-04-23 22:47:51,474 INFO mapreduce.Job: job_1619206134623_0012 completed successfully
'021-04-23 22:47:51,862 INFO mapreduce.Job: Counters: 47
File System Counters
  File read=0 of bytes read=0
  FILE: Number of bytes written=528326
  FILE: Number of read operations=0
  FILE: Number of large read operations=0
  FILE: Number of write operations=0
  HDFS: Number of bytes read=2222230381
  HDFS: Number of bytes written=219061999
  HDFS: Number of read operations=42
  HDFS: Number of large read operations=0
  HDFS: Number of write operations=0
  HDFS: Number of bytes read erasure-coded=0
Job Counters
  Launched map tasks=2
  Data-local map tasks=2
  Total time spent by all maps in occupied slots (ms)=114789
  Total time spent by all map tasks in occupied slots (ms)=0
  Total time spent by all map tasks (ms)=114789
  Total vcore-milliseconds taken by all map tasks=114789
  Total megabyte-milliseconds taken by all map tasks=117543936
Map-Reduce Framework
  Input records=769739
  Map output records=8
  Input split bytes=242
```

```

HDFS: Number of read operations=42
HDFS: Number of large read operations=0
HDFS: Number of write operations=60
HDFS: Number of bytes read erasure-coded=0
Job Counters
Launched map tasks=2
Data-local map tasks=2
Total time spent by all maps in occupied slots (ms)=114789
Total time spent by all reducers in occupied slots (ms)=0
Total time spent by all map tasks (ms)=114789
Total vcore-milliseconds taken by all map tasks=114789
Total negabyte-milliseconds taken by all map tasks=117543936
Map-Reduce Framework
  Map input records=769739
  Map output records=0
  Input splits=242
  Spilled Records=0
  Failed Shuffles=0
  Merged Map outputs=0
  GC time elapsed (ms)=1593
  CPU time spent (ms)=38690
  Physical memory (bytes) snapshot=5253589128
  Virtual memory (bytes) snapshot=620756992
  Total committed heap usage (bytes)=620756992
  Peak Map Physical memory (bytes)=374652928
  Peak Map Virtual memory (bytes)=2628857856
File Input Format Counters
  Bytes Read=19819139
File Output Format Counters
  Bytes Written=0
org.apache.hadoop.mapreduce.lib.output.MultipleOutputs
  Autos & Vehicles=14537
  Comedy=89911
  Entertainment=131052
  Film & Animation=75487
  Gadgets & Games=1068
  Howto & DIY=18774
  Music=183411
  News & Politics=37469
  People & Blogs=50454
  Pets & Animals=10736
  Sports=8710
  Travel & Places=15012
  UNA=6264
  Uncategorized=4
agarshah95@ubuntu:/usr/local/bin/hadoop-3.3.0/bin$ 

```

```

agarshah95@ubuntu:/usr/local/bin/hadoop-3.3.0/bin$ ./hadoop fs -ls /BigDataProject/data/outputFiles/binning
Found 31 items
-rw-r--r-- 1 agarshah95 supergroup 2474833 2021-04-23 22:47 /BigDataProject/data/outputFiles/binning/Autos & Vehicles-m-00000
-rw-r--r-- 1 agarshah95 supergroup 1779437 2021-04-23 22:47 /BigDataProject/data/outputFiles/binning/Auto & Vehicles-m-00001
-rw-r--r-- 1 agarshah95 supergroup 1855272 2021-04-23 22:47 /BigDataProject/data/outputFiles/binning/Comedy-m-00000
-rw-r--r-- 1 agarshah95 supergroup 2936268 2021-04-23 22:47 /BigDataProject/data/outputFiles/binning/Comedy-m-00001
-rw-r--r-- 1 agarshah95 supergroup 24837226 2021-04-23 22:47 /BigDataProject/data/outputFiles/binning/Entertainment-m-00000
-rw-r--r-- 1 agarshah95 supergroup 13482200 2021-04-23 22:47 /BigDataProject/data/outputFiles/binning/Entertainment-m-00001
-rw-r--r-- 1 agarshah95 supergroup 11918688 2021-04-23 22:47 /BigDataProject/data/outputFiles/binning/Film & Animation-m-00000
-rw-r--r-- 1 agarshah95 supergroup 19456978 2021-04-23 22:47 /BigDataProject/data/outputFiles/binning/Film & Animation-m-00001
-rw-r--r-- 1 agarshah95 supergroup 3586113 2021-04-23 22:47 /BigDataProject/data/outputFiles/binning/Gadgets & Games-m-00000
-rw-r--r-- 1 agarshah95 supergroup 9696564 2021-04-23 22:47 /BigDataProject/data/outputFiles/binning/Gadgets & Games-m-00001
-rw-r--r-- 1 agarshah95 supergroup 4859479 2021-04-23 22:47 /BigDataProject/data/outputFiles/binning/Howto & DIY-m-00000
-rw-r--r-- 1 agarshah95 supergroup 2167533 2021-04-23 22:47 /BigDataProject/data/outputFiles/binning/Howto & DIY-m-00001
-rw-r--r-- 1 agarshah95 supergroup 3572334 2021-04-23 22:47 /BigDataProject/data/outputFiles/binning/Howto & DIY-m-00000
-rw-r--r-- 1 agarshah95 supergroup 1810449 2021-04-23 22:47 /BigDataProject/data/outputFiles/binning/Howto & DIY-m-00001
-rw-r--r-- 1 agarshah95 supergroup 32624226 2021-04-23 22:47 /BigDataProject/data/outputFiles/binning/Music-m-00000
-rw-r--r-- 1 agarshah95 supergroup 19491779 2021-04-23 22:47 /BigDataProject/data/outputFiles/binning/Music-m-00001
-rw-r--r-- 1 agarshah95 supergroup 15001113 2021-04-23 22:47 /BigDataProject/data/outputFiles/binning/News & Politics-m-00000
-rw-r--r-- 1 agarshah95 supergroup 3586013 2021-04-23 22:47 /BigDataProject/data/outputFiles/binning/News & Politics-m-00001
-rw-r--r-- 1 agarshah95 supergroup 9816583 2021-04-23 22:47 /BigDataProject/data/outputFiles/binning/People & Blogs-m-00000
-rw-r--r-- 1 agarshah95 supergroup 4859479 2021-04-23 22:47 /BigDataProject/data/outputFiles/binning/People & Blogs-m-00001
-rw-r--r-- 1 agarshah95 supergroup 2167533 2021-04-23 22:47 /BigDataProject/data/outputFiles/binning/Pets & Animals-m-00000
-rw-r--r-- 1 agarshah95 supergroup 937932 2021-04-23 22:47 /BigDataProject/data/outputFiles/binning/Pets & Animals-m-00001
-rw-r--r-- 1 agarshah95 supergroup 9816583 2021-04-23 22:47 /BigDataProject/data/outputFiles/binning/Sports-m-00000
-rw-r--r-- 1 agarshah95 supergroup 30381209 2021-04-23 22:47 /BigDataProject/data/outputFiles/binning/Sports-m-00001
-rw-r--r-- 1 agarshah95 supergroup 1518986 2021-04-23 22:47 /BigDataProject/data/outputFiles/binning/Travel & Places-m-00000
-rw-r--r-- 1 agarshah95 supergroup 247595 2021-04-23 22:47 /BigDataProject/data/outputFiles/binning/UNA-m-00000
-rw-r--r-- 1 agarshah95 supergroup 167141 2021-04-23 22:47 /BigDataProject/data/outputFiles/binning/UNA-m-00001
-rw-r--r-- 1 agarshah95 supergroup 765 2021-04-23 22:47 /BigDataProject/data/outputFiles/binning/UNCategorized-m-00000
-rw-r--r-- 1 agarshah95 supergroup 0 2021-04-23 22:47 /BigDataProject/data/outputFiles/binning/SUCCESS
-rw-r--r-- 1 agarshah95 supergroup 0 2021-04-23 22:47 /BigDataProject/data/outputFiles/binning/part-m-00000
-rw-r--r-- 1 agarshah95 supergroup 0 2021-04-23 22:47 /BigDataProject/data/outputFiles/binning/part-m-00001

```

Top 25 Categories : Returns the top 25 based on rating

2 MR jobs, 1st MR job calculates avg rating, 2nd MR job gets the top 25 records with the help of CustomKeyComparator (Top10_Categories)

This analysis comprises of :

1. Driver Class : Sets up configuration for Mapper1, Reducer1, Mapper2, Reducer2 class
2. Mapper 1 Class : Emits videoID and rating
3. Reducer 1 Class : Returns average rating
4. Mapper 2 Class : Emits rating and videoID
5. Reducer 2 Class : Emits top 25 values in descending order, videoID as a key of type Text and rating of type FloatWritable
6. CustomKeyComparator : Used for sorting by implementing the comparable method

Mapper 1

```
public static class Map1 extends Mapper<Object, Text, Text, FloatWritable> {  
    private FloatWritable video_rating = new FloatWritable();  
    private Text video_id = new Text();  
  
    public void map(Object key, Text value, Mapper.Context context  
    ) throws IOException, InterruptedException {  
  
        String[] fields = value.toString().split(",");  
        video_id = new Text(fields[0]);  
        try {  
            if (fields.length > 6) {  
                //if (!fields[6].isEmpty()) {  
                    video_rating = new FloatWritable(Float.parseFloat(fields[6]));  
                }  
  
                context.write(video_id, video_rating);  
            } catch (Exception e) {  
                e.printStackTrace();  
            }  
    }  
}
```

Reducer 1

```
public static class Reducel extends Reducer<Text, FloatWritable, Text, FloatWritable> {  
  
    private FloatWritable result = new FloatWritable();  
  
    @Override  
    protected void reduce(Text key, Iterable<FloatWritable> values, Context context)  
    throws IOException, InterruptedException {  
  
        int count = 0;  
        float sum = 0, avg = 0;  
  
        for (FloatWritable val : values) {  
            sum += val.get();  
            ++count;  
        }  
  
        avg = sum / count;  
        result.set(avg);  
        context.write(key, result);  
    }  
}
```

Mapper 2

```
public static class Map2 extends Mapper<Object, Text, FloatWritable, Text> {  
  
    @Override  
    protected void map(Object key, Text value, Mapper.Context context) throws IOException, InterruptedException {  
  
        String row[] = value.toString().split("\t");  
        Text video_id = new Text(row[0]);  
        String rating = row[1];  
  
        try {  
            FloatWritable ratingg = new FloatWritable(Float.parseFloat(rating));  
            context.write(ratingg, video_id);  
        } catch (Exception e) {  
            e.printStackTrace();  
        }  
    }  
}
```

Reducer 2

```
public static class Reduce2 extends Reducer<FloatWritable, Text, Text, FloatWritable> {  
    private static int count = 25;  
  
    @Override  
    protected void reduce(FloatWritable key, Iterable<Text> values, Context context) throws IOException, InterruptedException {  
        for (Text val : values) {  
            if (count > 0) {  
                context.write(val, key);  
                --count;  
            } else {  
                break;  
            }  
        }  
    }  
}
```

Driver

```
public static void main(String[] args) throws IOException, InterruptedException, ClassNotFoundException {  
    Configuration conf1 = new Configuration();  
    Configuration conf = new Configuration();  
  
    Job job1 = Job.getInstance(conf1, "Chaining");  
    job1.setJarByClass(Top10_Categories.class);  
  
    job1.setMapperClass(Map1.class);  
    job1.setMapOutputKeyClass(Text.class);  
    job1.setMapOutputValueClass(FloatWritable.class);  
  
    job1.setReducerClass(Reduce1.class);  
    job1.setOutputKeyClass(Text.class);  
    job1.setOutputValueClass(DoubleWritable.class);  
    job1.setCombinerClass(Reduce1.class);  
  
    FileInputFormat.addInputPath(job1, new Path(args[0]));  
    FileOutputFormat.setOutputPath(job1, new Path(args[1]));  
  
    boolean complete = job1.waitForCompletion(true);  
  
    Configuration conf2 = new Configuration();  
    Job job2 = Job.getInstance(conf2, "Chaining");  
  
    if (complete) {  
        job2.setJarByClass(Top10_Categories.class);  
        job2.setMapperClass(Map2.class);  
        job2.setMapOutputKeyClass(FloatWritable.class);  
        job2.setMapOutputValueClass(Text.class);  
  
        job2.setReducerClass(Reduce2.class);  
        job2.setOutputKeyClass(Text.class);  
        job2.setOutputValueClass(FloatWritable.class);  
  
        job2.setSortComparatorClass(SortKeyComparator.class);  
  
        job2.setNumReduceTasks(1);  
  
        FileInputFormat.addInputPath(job2, new Path(args[2]));  
        FileOutputFormat.setOutputPath(job2, new Path(args[3]));  
  
        System.exit(job2.waitForCompletion(true) ? 0 : 1);  
    }  
}
```

Execution

```
sagarshah95@ubuntu: /usr/local/bin/hadoop-3.3.0/bin$ ./hadoop jar /home/sagarshah95/Desktop/Top10_Categories-0.8.1-SNAPSHOT.jar com.neu.bigdata.Top10_Categories /BigDataProject/data/youtubeDataset /BigDataProject
File Edit View Search Terminal Help
sagarshah95@ubuntu: /usr/local/bin/hadoop-3.3.0/bin
separshah95@ubuntu: /usr/local/bin/hadoop-3.3.0/bin$ ./hadoop jar /home/sagarshah95/Desktop/Top10_Categories-0.8.1-SNAPSHOT.jar com.neu.bigdata.Top10_Categories /BigDataProject/data/outputfiles/tempOutput /BigDataProject/data/outputfiles/tempOutput/part-r-00000 /BigDataProject/data/outputfiles/top10Output
2021-04-28 11:52:41,370 INFO client.DefaultNoHARMFailoverProxyProvider: Connecting to ResourceManager at /0.0.0.0:8032
2021-04-28 11:52:42,144 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
2021-04-28 11:52:42,189 INFO mapreduce.JobResourceUploader: Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/sagarshah95/.staging/job_1619549114880_0006
2021-04-28 11:52:42,190 INFO mapreduce.JobResourceUploader: Number of files to process : 1
2021-04-28 11:52:42,941 INFO mapreduce.JobSubmitter: number of splits:2
2021-04-28 11:52:43,312 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1619549114880_0006
2021-04-28 11:52:43,313 INFO mapreduce.JobSubmitter: Executing with tokens: []
2021-04-28 11:52:43,728 INFO conf.Configuration: resource-types.xml not found
2021-04-28 11:52:43,730 INFO conf.Configuration: Resource types: []
2021-04-28 11:52:43,911 INFO impl.YarnClientImpl: Submitted application application_1619549114880_0006
2021-04-28 11:52:44,010 INFO mapreduce.Job: The url to track the job: http://ubuntu:8088/proxy/application_1619549114880_0006
2021-04-28 11:52:44,012 INFO mapreduce.Job: Running job: job_1619549114880_0006
2021-04-28 11:52:45,253 INFO mapreduce.Job: Job job_1619549114880_0006 running in uber mode : false
2021-04-28 11:52:45,547 INFO mapreduce.Job: map 0% reduce 0%
2021-04-28 11:53:14,677 INFO mapreduce.Job: map 100% reduce 0%
2021-04-28 11:53:24,795 INFO mapreduce.Job: map 100% reduce 100%
2021-04-28 11:53:25,250 INFO mapreduce.Job: Job job_1619549114880_0006 completed successfully
2021-04-28 11:53:25,251 INFO mapreduce.Job: Counters: 55
  File System Counters
    FILE: Number of bytes read=13431463
    FILE: Number of bytes written=2756446
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=19819381
    HDFS: Number of bytes written=12246195
    HDFS: Number of read operations=1
    HDFS: Number of large read operations=0
    HDFS: Number of write operations=2
    HDFS: Number of bytes read erasure-coded=0
  Job Counters
    Killed map tasks=1
    Launched map tasks=2
    Launched reduce tasks=1
    Data-local map tasks=2
    Total time spent by all map tasks in occupied slots (ms)=31834
    Total time spent by all reduce tasks in occupied slots (ms)=8573
    Total time spent by all map tasks (ms)=31834
    Total time spent by all reduce tasks (ms)=8573
    Total vcore-milliseconds taken by all map tasks=31834
    Total vcore-milliseconds taken by all reduce tasks=8573
    Total megabyte-milliseconds taken by all map tasks=32598016
    Total megabyte-milliseconds taken by all reduce tasks=8778752
  Map-Reduce Framework
    Map input records=769739
    Map output Records=769735
    Map output bytes=12315761
    Map output materialized bytes=13431469
    Input split bytes=242
    Combine output records=769735
    Combine output bytes=746192
    Reduce input groups=746192
    Reduce shuffle bytes=13431469
    Reduce input records=746192
    Reduces input records=746192
    Spilled Records=1493284
    Shuffled Maps =2
    Failed Shuffles=0
    Merged Map outputs=2
    GC time (ms)=1817
    CPU time spent (ms)=13850
    Physical memory (bytes) snapshot=1062842368
    Virtual memory (bytes) snapshot=7764004864
    Total committed memory (bytes)=1380996
    Peak Map Physical memory (bytes)=1581532
    Peak Map Virtual memory (bytes)=2587725824
    Peak Reduce Physical memory (bytes)=185225216
    Peak Reduce Virtual memory (bytes)=2590793728
  Shuffle Errors
    IO=0
    CONNECTION=0
    IO_ERROR=0
    WRONG_LENGTH=0
    WRONG_MAGIC=0
    WRONG_REDUCE=0
  File Input Format Counters
    Bytes Read=219819139
  File Output Format Counters
    Bytes Written=12246195
2021-04-28 11:53:26,181 INFO client.DefaultNoHARMFailoverProxyProvider: Connecting to ResourceManager at /0.0.0.0:8032
2021-04-28 11:53:26,216 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
2021-04-28 11:53:26,238 INFO mapreduce.JobResourceUploader: Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/sagarshah95/.staging/job_1619549114880_0007
2021-04-28 11:53:26,445 INFO mapreduce.JobSubmitter: Number of splits:1
2021-04-28 11:53:26,465 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1619549114880_0007
2021-04-28 11:53:26,579 INFO mapreduce.JobSubmitter: Executing with tokens: []
2021-04-28 11:53:26,677 INFO impl.YarnClientImpl: Submitted application application_1619549114880_0007
```

```
sagarshah95@ubuntu: /usr/local/bin/hadoop-3.3.0/bin
File Edit View Search Terminal Help
sagarshah95@ubuntu: /usr/local/bin/hadoop-3.3.0/bin
Total time spent by all map tasks (ms)=31834
Total time spent by all reduce tasks (ms)=8573
Total vcore-milliseconds taken by all map tasks=31834
Total vcore-milliseconds taken by all reduce tasks=8573
Total megabyte-milliseconds taken by all map tasks=32598016
Total megabyte-milliseconds taken by all reduce tasks=8778752
Map-Reduce Framework
  Map input records=769739
  Map output records=769735
  Map output bytes=12315761
  Map output materialized bytes=13431469
  Input split bytes=242
  Combine output records=769735
  Combine output bytes=746192
  Reduce input groups=746192
  Reduce shuffle bytes=13431469
  Reduce input records=746192
  Reduces input records=746192
  Spilled Records=1493284
  Shuffled Maps =2
  Failed Shuffles=0
  Merged Map outputs=2
  GC time (ms)=1817
  CPU time spent (ms)=13850
  Physical memory (bytes) snapshot=1062842368
  Virtual memory (bytes) snapshot=7764004864
  Total committed memory (bytes)=1380996
  Peak Map Physical memory (bytes)=1581532
  Peak Map Virtual memory (bytes)=2587725824
  Peak Reduce Physical memory (bytes)=185225216
  Peak Reduce Virtual memory (bytes)=2590793728
  Shuffle Errors
    IO=0
    CONNECTION=0
    IO_ERROR=0
    WRONG_LENGTH=0
    WRONG_MAGIC=0
    WRONG_REDUCE=0
  File Input Format Counters
    Bytes Read=219819139
  File Output Format Counters
    Bytes Written=12246195
2021-04-28 11:53:26,181 INFO client.DefaultNoHARMFailoverProxyProvider: Connecting to ResourceManager at /0.0.0.0:8032
2021-04-28 11:53:26,216 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
2021-04-28 11:53:26,238 INFO mapreduce.JobResourceUploader: Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/sagarshah95/.staging/job_1619549114880_0007
2021-04-28 11:53:26,445 INFO mapreduce.JobSubmitter: Number of splits:1
2021-04-28 11:53:26,465 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1619549114880_0007
2021-04-28 11:53:26,579 INFO mapreduce.JobSubmitter: Executing with tokens: []
2021-04-28 11:53:26,677 INFO impl.YarnClientImpl: Submitted application application_1619549114880_0007
```

```

Activities Terminal ▾ sagarshah95@ubuntu: /usr/local/bin/hadoop-3.3.0/bin
Wed 11:59

File Edit View Search Terminal Help
sagarshah95@ubuntu: /usr/local/bin/hadoop-3.3.0/bin

2021-04-28 11:53:26,698 INFO mapreduce.Job: The url to track the job: http://ubuntu:8088/proxy/application_1619549114880_0007/
2021-04-28 11:53:26,699 INFO mapreduce.Job: Running job: job_1619549114880_0007
2021-04-28 11:53:27,000 INFO mapreduce.Job: Job job_1619549114880_0007 running in uber mode : false
2021-04-28 11:53:43,289 INFO mapreduce.Job: map 0% reduce 0%
2021-04-28 11:53:52,402 INFO mapreduce.Job: map 100% reduce 0%
2021-04-28 11:54:01,587 INFO mapreduce.Job: map 100% reduce 100%
2021-04-28 11:54:03,665 INFO mapreduce.Job: Job job_1619549114880_0007 completed successfully
2021-04-28 11:54:03,666 INFO mapreduce.Job: Counters: 54
  File System Counters
    FILE: Number of bytes read=13431463
    FILE: Number of bytes written=27391969
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=22246338
    HDFS: Number of bytes written=400
    HDFS: Number of read operations=8
    HDFS: Number of large read operations=0
    HDFS: Number of write operations=2
    HDFS: Number of bytes read erasure-coded=0
  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)=6883
    Total time spent by all reduces in occupied slots (ms)=5582
    Total time spent by all map tasks (ms)=6883
    Total time spent by all reduce tasks (ms)=5582
    Total vcore-milliseconds taken by all map tasks=6883
    Total vcore-milliseconds taken by all reduce tasks=7048192
    Total megabyte-milliseconds taken by all map tasks=7048192
    Total megabyte-milliseconds taken by all reduce tasks=5715968
  Map-Reduce Framework
    Map input records=746192
    Map output records=746192
    Map output bytes=11939873
    Map output materialized bytes=13431463
    Input split bytes=143
    Combine input records=0
    Reduce input groups=400
    Reduce shuffle bytes=13431463
    Reduce input records=746192
    Reduce output records=25
    Spilled Records=1492384
    Shuffled Maps =1
    Failed Shuffles=0
    Merged Map outputs=1
    GC time elapsed (ms)=183
    CPU time spent (ms)=5590

```

```

Activities Terminal ▾ sagarshah95@ubuntu: /usr/local/bin/hadoop-3.3.0/bin
File Edit View Search Terminal Help
sagarshah95@ubuntu: /usr/local/bin/hadoop-3.3.0/bin

File Edit View Search Terminal Help
sagarshah95@ubuntu: /usr/local/bin/hadoop-3.3.0/bin

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)=6883
  Total time spent by all reduces in occupied slots (ms)=5582
  Total time spent by all map tasks (ms)=6883
  Total time spent by all reduce tasks (ms)=5582
  Total vcore-milliseconds taken by all map tasks=6883
  Total vcore-milliseconds taken by all reduce tasks=5582
  Total megabyte-milliseconds taken by all map tasks=7048192
  Total megabyte-milliseconds taken by all reduce tasks=5715968
Map-Reduce Framework
  Map input records=746192
  Map output records=746192
  Map output bytes=11939873
  Map output materialized bytes=13431463
  Input split bytes=143
  Combine input records=0
  Reduce input groups=400
  Reduce shuffle bytes=13431463
  Reduce input records=746192
  Reduce output records=25
  Spilled Records=1492384
  Shuffled Maps =1
  Failed Shuffles=0
  Merged Map outputs=1
  GC time elapsed (ms)=183
  CPU time spent (ms)=5590
  Physical memory (bytes) snapshot=609832960
  Virtual memory (bytes) snapshot=5176967168
  Total committed heap usage (bytes)=486014976
  Peak Map Physical memory (bytes)=426663936
  Peak Map Virtual memory (bytes)=258514048
  Peak Reduce Physical memory (bytes)=183169834
  Peak Reduce Virtual memory (bytes)=2591817728
  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=1246195
  File Output Format Counters
sagarshah95@ubuntu: /usr/local/bin/hadoop-3.3.0/bin$ 

```

```

sagarshah95@ubuntu: /usr/local/bin/hadoop-3.3.0/bin$ ./hadoop fs -cat /BigDataProject/data/outputFiles/top10output/part-r-00000
Gby0saVfq4U 5.0
HBH2hNvQTAc 5.0
GbzTV7VnEA 5.0
HB0CqkLJ4L 5.0
Cc-PBkQ9nc 5.0
H8A6GCUvxt8Y 5.0
H99eSeVsV4E 5.0
-7wC2LSHSjk 5.0
H85MKRKTh74 5.0
H82xR51vE7Y 5.0
H815hzX0uQg 5.0
H7yVRQLx3yo 5.0
H7y51NUuuhi 5.0
Gc6E0LEuDk 5.0
H7xf1fUgchvo 5.0
Gc7EM7JhdNw 5.0
H7vneSg9qc 5.0
Gc7Wdxt1ao 5.0
HT75mf14vFA 5.0
H75rM13Us 5.0
H75lk1YtME 5.0
zb4k3P9s-U 5.0
Gc9yjwv34I 5.0
GcAP0bx9ovg 5.0
GcAV-9us4_8 5.0
sagarshah95@ubuntu: /usr/local/bin/hadoop-3.3.0/bin$ 

```

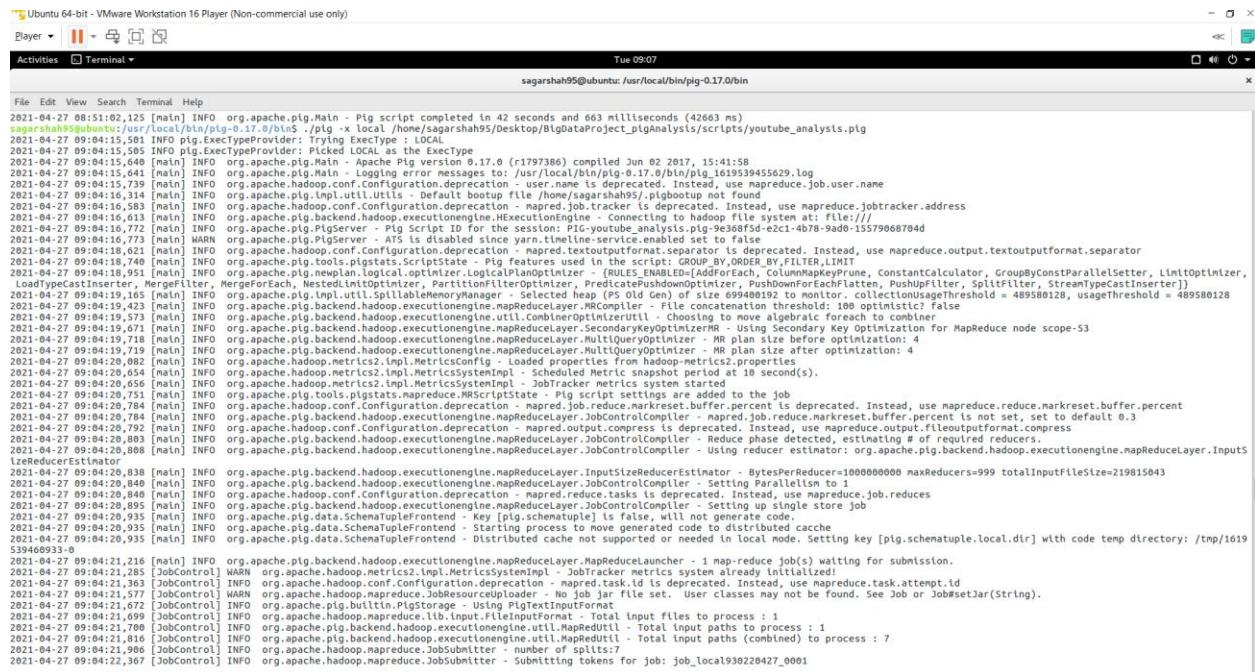
Pig Analysis

1)Top 5 Categories : Returns the Top 5 Categories of Youtube Videos

Script

```
infiles = load '/home/sagarshah95/Desktop/BigDataProject_pigAnalysis/data/youtubeDataset.csv' using PigStorage(',') as (videoid:chararray,uploader:chararray,age:int,category:chararray,length:int,views:int,rate:int,rating:int,comments:int,related_id:chararray);
files = FILTER infiles BY category is not null;
grpns_for_categories = group files by category;
cnt_for_categories = foreach grpns_for_categories generate group, COUNT(files.videoid) as counting;
sorted_for_categories_desc = order cnt_for_categories by counting desc;
top5_for_categories = limit sorted_for_categories_desc 5;
STORE top5_for_categories INTO 'Top5Categories.txt' using PigStorage('|');
```

Execution



```
Ubuntu 64-bit - VMware Workstation 16 Player (Non-commercial use only)
Player | ||| | Activities Terminal Tue 09:07
sagarshah95@ubuntu: /usr/local/bin/pig-0.17.0/bin

File Edit View Search Terminal Help
2021-04-27 08:51:02,125 [main] INFO org.apache.pig.Main - Pig script completed in 42 seconds and 663 milliseconds (42663 ms)
sagarshah95@ubuntu:/usr/local/bin/pig-0.17.0/bin$ ./pig -x local /home/sagarshah95/Desktop/BigDataProject_pigAnalysis/scripts/youtube_analysis.pig
2021-04-27 09:04:15,153 [main] INFO pig.ExectypeProvider - Trying ExecType : LOCAL
2021-04-27 09:04:15,153 [main] INFO pig.ExectypeProvider - Selected ExecType : LOCAL as Local
2021-04-27 09:04:15,683 [main] INFO org.apache.pig.Main - Apache Pig version 0.17.0 (r1797306) compiled Jun 02 2017, 15:41:58
2021-04-27 09:04:15,683 [main] INFO org.apache.pig.Main - Logging error messages to: /usr/local/bin/pig-0.17.0/bin/pig_1619539455629.log
2021-04-27 09:04:15,739 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - user.name is deprecated. Instead, use mapreduce.job.user.name
2021-04-27 09:04:15,739 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-04-27 09:04:15,739 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-04-27 09:04:16,383 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-04-27 09:04:16,613 [main] INFO org.apache.pig.backend.hadoop.executionengine.HExecutionEngine - Connecting to hadoop file system at: file:///
2021-04-27 09:04:16,772 [main] INFO org.apache.pig.PigServer - Pig Script ID for the session: PIG-youtube.analysis.pig-9e368f5d-e2c1-4b78-9a0e-15579068704d
2021-04-27 09:04:16,773 [main] WARN org.apache.pig.PigServer - ATS is disabled since yarn.timeline-service.enabled set to false
2021-04-27 09:04:16,773 [main] INFO org.apache.pig.tools.pigscript.parser.PigScriptParser - Using Secondary Optimization - see https://pig.apache.org/docs/r0.12.0/secondary.html
2021-04-27 09:04:16,740 [main] INFO org.apache.pig.tools.pigscript.parser.PigScriptParser - Using Secondary Optimization - see https://pig.apache.org/docs/r0.12.0/secondary.html
2021-04-27 09:04:18,951 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer - RULES_ENABLED=[AddForEach, ColumnMapKeyPrune, ConstantCalculator, GroupByConstParallelSetter, LimitOptimizer, LoadTypeCastInserter, MergeFilter, MergeForEach, NeedLimitOptimizer, PartitionFilterOptimizer, PredicatePushdownOptimizer, PushdownForEachFlatten, PushupFilter, SplitFilter, StreamTypeCastInserter]
2021-04-27 09:04:19,105 [main] INFO org.apache.pig.impl.util.SplittableMemoryManager - Selected heap (PS OR Gen) of size 0994000192 to monitor, collectionUsageThreshold = 489580128, usageThreshold = 489580128
2021-04-27 09:04:19,105 [main] INFO org.apache.pig.impl.util.SplittableMemoryManager - Collection usage threshold is 0.0000000000000001, usage threshold is 0.0000000000000001, false
2021-04-27 09:04:19,573 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.CombineOptimizerUtil - Choosing to move algebraic forward to combine
2021-04-27 09:04:19,671 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.CombineOptimizerUtil - Using Secondary Key Optimization for MapReduce node scope-53
2021-04-27 09:04:19,718 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultitQueryOptimizer - MR plan size before optimization: 4
2021-04-27 09:04:19,719 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultitQueryOptimizer - MR plan size after optimization: 4
2021-04-27 09:04:19,719 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultitQueryOptimizer - Local file size: 0 properties
2021-04-27 09:04:20,654 [main] INFO org.apache.hadoop.metrics2.impl.MetricsSystemImpl - Scheduled Metric snapshot period at 10 second(s).
2021-04-27 09:04:20,656 [main] INFO org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system started
2021-04-27 09:04:20,751 [main] INFO org.apache.pig.tools.pigscript.MRScriptState - Pig script settings are added to the job
2021-04-27 09:04:20,784 [main] INFO org.apache.pig.tools.pigscript.MRScriptState - mapreduce.buffer.percent is deprecated. Instead, use mapreduce.reduce.mapred.buffer.percent. mapreduce.buffer.percent is deprecated. Instead, use mapreduce.reduce.mapred.buffer.percent. mapred.buffer.percent is not set, set to default 0.3
2021-04-27 09:04:20,784 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Setting Parallelism to 1
2021-04-27 09:04:20,792 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.output.compress is deprecated. Instead, use mapreduce.output.fileoutputformat.compress
2021-04-27 09:04:20,803 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Reduce phase detected, estimating # of required reducers.
2021-04-27 09:04:20,808 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Using reduce estimator: org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.InputsizerReducerEstimator
2021-04-27 09:04:20,838 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.InputsizerReducerEstimator - BytesPerReducer=1000000000 maxReducers=999 totalInputFileSize=219815043
2021-04-27 09:04:20,849 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Setting Parallelism to 1
2021-04-27 09:04:20,849 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.reduce.tasks is deprecated. Instead, use mapreduce.job.reduces
2021-04-27 09:04:20,849 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Setting up single store job
2021-04-27 09:04:20,850 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - KILL [pig.schematuple.local.dir] with code temp
2021-04-27 09:04:20,935 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Status changed to no generated code to distributed cache
2021-04-27 09:04:20,935 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Distributed cache not supported or needed in local mode. Setting key [pig.schematuple.local.dir] with code temp directory: /tmp/1619546993-0
2021-04-27 09:04:21,216 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 1 map-reduce job(s) waiting for submission.
2021-04-27 09:04:21,285 [JobControl] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:04:21,363 [JobControl] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.task.id is deprecated. Instead, use mapreduce.task.attempt.id
2021-04-27 09:04:21,577 [JobControl] WARN org.apache.hadoop.mapreduce.JobResourceUploader - No job jar file set. User classes may not be found. See Job or Job#setJar(String).
2021-04-27 09:04:21,672 [JobControl] INFO org.apache.pig.builtin.PigStorage - Using PigTextInputFormat
2021-04-27 09:04:21,700 [JobControl] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input files to process : 1
2021-04-27 09:04:21,700 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceUtil - Total input paths to process : 1
2021-04-27 09:04:21,816 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceUtil - Total input paths (combined) to process : 7
2021-04-27 09:04:21,906 [JobControl] INFO org.apache.hadoop.mapreduce.JobSubmitter - number of splits:7
2021-04-27 09:04:22,367 [JobControl] INFO org.apache.hadoop.mapreduce.JobSubmitter - Submitting tokens for job: job_local93028427_0001
```

```

Activities Terminal Tue 09:07
sagarshah95@ubuntu: /usr/local/bin/pig-0.17.0/bin

File Edit View Search Terminal Help
2021-04-27 09:04:28,264 [LocalJobRunner Map Task Executor #0] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.PigCombiner$Combine - Aliases being processed per job phase (AliasName[line,offset ])
]: M: infiles[1,10],infles[-1,-1],files[3,8],cnt_for_catagories[5,21],grpn_for_catagories[4,22] C: cnt_for_catagories[5,21],R: cnt_for_catagories[5,21]
2021-04-27 09:04:28,854 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.MapTask - Finished spill 0
2021-04-27 09:04:28,858 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.Task - Task attempt local930220427_0001_m_000000 is done. And is in the process of committing
2021-04-27 09:04:28,985 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.Task - Task attempt local930220427_0001_m_000000 map
2021-04-27 09:04:28,988 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.Task - Task attempt local930220427_0001_m_000000_0 done.
2021-04-27 09:04:29,063 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.Task - Final Counters for attempt_local930220427_0001_m_000000: Counters: 20
  File System Counters
    FILE: Number of bytes read=3383176
    FILE: Number of bytes written=621758
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
  Map-Reduce Framework
    Map Input records=118289
    Map output records=117223
    Map output bytes=2129248
    Map output materialized bytes=315
    Input split bytes=117223
    Combine input records=117223
    Combine output records=13
    Spilled Records=13
    Failed Shuffles=0
    Merged Map outputs=0
    GC time elapsed (ms)=57
    Total committed heap usage (bytes)=381681664
  File Input Format Counters
    Bytes Read=0
    BYTES_READ_TO_HEAP=0
    org.apache.hadoop.mapreduce.ACCESsing_NON_EXISTENT_FIELD=1295
    FIELD_DISCARDED_TYPE_CONVERSION_FAILED=1020
2021-04-27 09:04:29,063 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.LocalJobRunner - Finishing task: attempt_local930220427_0001_m_000000_0
2021-04-27 09:04:29,063 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.LocalJobRunner - Starting task: attempt_local930220427_0001_m_000001_0
2021-04-27 09:04:29,077 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter - File Output Committer Algorithm Version ls 2
2021-04-27 09:04:29,077 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter - FileOutputCommitter skip cleanup _temporary folders under output directory=false,
ignore cleanup failures: false
2021-04-27 09:04:29,080 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.Task - Using ResourceCalculatorProcessTree : []
2021-04-27 09:04:29,080 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.MapTask - (EQUATOR) e kvl 26214396(104857584)
Total Length = 33554432
Input split[0]:
  Length = 33554432
  ClassName: org.apache.hadoop.mapreduce.lib.input.FileSplit
  Locations:

-----
2021-04-27 09:04:29,117 [LocalJobRunner Map Task Executor #0] INFO org.apache.pig.builtin.PigStorage - Using PigTextInputFormat
2021-04-27 09:04:29,118 [LocalJobRunner Map Task Executor #0] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.PigRecordReader - Current split being processed file:/home/sagarshah95/Desktop/BIG
DataProject_piganalysis/data/youtubeDataset.csv:33554432+33554432
2021-04-27 09:04:29,134 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.MapTask - (EQUATOR) e kvl 26214396(104857584)

```

```

Activities Terminal Tue 09:08
sagarshah95@ubuntu: /usr/local/bin/pig-0.17.0/bin

File Edit View Search Terminal Help
1021-04-27 09:04:29,438 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 12% complete
1021-04-27 09:04:29,444 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - MapReduceLauncher - Running jobs are [job_local930220427_0001]
1021-04-27 09:04:33,099 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.LocalJobRunner -
1021-04-27 09:04:33,100 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.MapTask - Starting flush of map output
1021-04-27 09:04:33,101 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.MapTask - Spilling map output
1021-04-27 09:04:33,102 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.MapTask - Spilled 0.00 bytes in 0.000 seconds(s). binned = 2140775; bufvalid = 104857600
1021-04-27 09:04:33,103 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.MapTask - kvstart = 26214396(104857584); kvend = 25747668(102990672); length = 466729/6553600
1021-04-27 09:04:33,346 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.MapTask - Finished spill 0
1021-04-27 09:04:33,348 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.Task - Task attempt_local930220427_0001_m_000001 is done. And is in the process of committing
1021-04-27 09:04:33,349 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.Task - Task attempt_local930220427_0001_m_000001 map
1021-04-27 09:04:33,347 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.LocalJobRunner - map
1021-04-27 09:04:33,348 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.Task - Task attempt_local930220427_0001_m_000001_0 done.
1021-04-27 09:04:33,348 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.Task - Final Counters for attempt_local930220427_0001_m_000001: Counters: 20
  File System Counters
    FILE: Number of bytes read=67122999
    FILE: Number of bytes written=622139
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
  Map-Reduce Framework
    Map Input records=117261
    Map output records=116683
    Map output bytes=2140775
    Map output materialized bytes=349
    Input split bytes=466729
    Combine output records=15
    Spilled Records=15
    Failed Shuffles=0
    Merged Map outputs=0
    GC time elapsed (ms)=318
    Total committed heap usage (bytes)=554696704
  File Input Format Counters
    Bytes Read=0
    org.apache.pig.PigLearning
    org.apache.pig.PigLearning
    ACCESsing_NON_EXISTENT_FIELD=5496
    FIELD_DISCARDED_TYPE_CONVERSION_FAILED=788
1021-04-27 09:04:33,349 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.LocalJobRunner - Finishing task: attempt_local930220427_0001_m_000001_0
1021-04-27 09:04:33,349 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.LocalJobRunner - Starting task: attempt_local930220427_0001_m_000002_0
1021-04-27 09:04:33,350 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter - File Output Committer Algorithm Version ls 2
1021-04-27 09:04:33,355 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter - FileOutputCommitter skip cleanup _temporary folders under output directory=false,
.ignore cleanup failures: false
1021-04-27 09:04:33,386 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.Task - Using ResourceCalculatorProcessTree : []
1021-04-27 09:04:33,389 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.MapTask - Processing split: Number of splits : 1
Total Length = 33554432
Input split[0]:
  Length = 33554432
  ClassName: org.apache.hadoop.mapreduce.lib.input.FileSplit
  Locations:

```

```

Player - Terminal
Activities Terminal
Tue 09:09
sagarshah95@ubuntu: /usr/local/bin/pig-0.17.0/bin

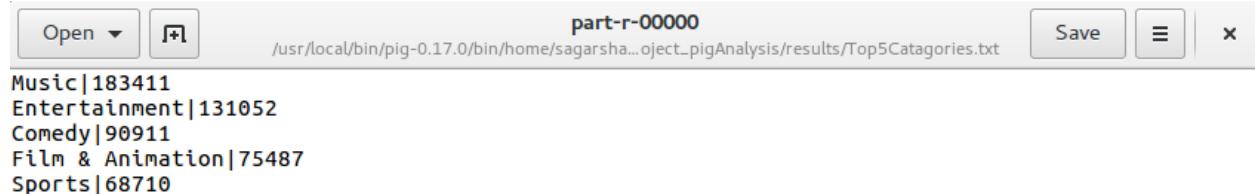
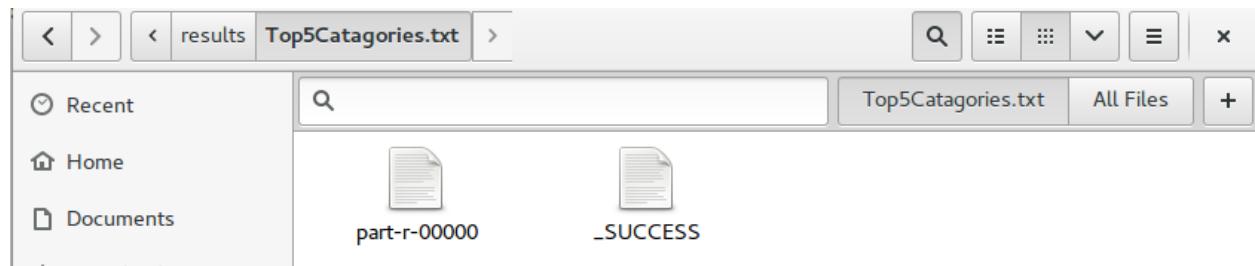
File Edit View Search Terminal Help
021-04-27 09:04:48,292 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 100% complete
021-04-27 09:04:48,316 [main] INFO org.apache.pig.tools.pigstats.mapreduce.SimplePigStats - Script Statistics:
ladopVersion PlgVersion UserId StartedAt FinishedAt Features
.3.0 0.17.0 sagarshah95 2021-04-27 09:04:20 2021-04-27 09:04:48 GROUP_BY,ORDER_BY,FILTER,LIMIT
uccess!
ob Stats (time in seconds):
obId Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime MedianReduceTime Alias Feature Outputs
ob_local1176657791_0002 1 1 n/a n/a n/a n/a n/a n/a sorted_for_catagories_desc SAMPLER
ob_local1688923654_0004 1 1 n/a n/a n/a n/a n/a n/a sorted_for_catagories_desc file:///usr/local/bin/pig-0.17.0/bin/Top5Categories.txt,
ob_local1850735436_0003 1 1 n/a n/a n/a n/a n/a n/a sorted_for_catagories_desc ORDER_BY,COMBINER
ob_local19393220427_0001 7 1 n/a n/a n/a n/a n/a n/a cnt_for_catagories,files,grpn_for_catagories,infiles GROUP_BY,COMBINER
nput(s);
uccessfully read 769739 records from: "/home/sagarshah95/Desktop/BigDataProject_pigAnalysis/data/youtubeDataset.csv"
utput(s);
uccessfully stored 5 records in: "file:///usr/local/bin/pig-0.17.0/bin/Top5Categories.txt"

counters:
otal records written : 5
otal bytes written : 0
pillable Memory Manager spill count : 0
otal bags proactively spilled: 0
otal records proactively spilled: 0

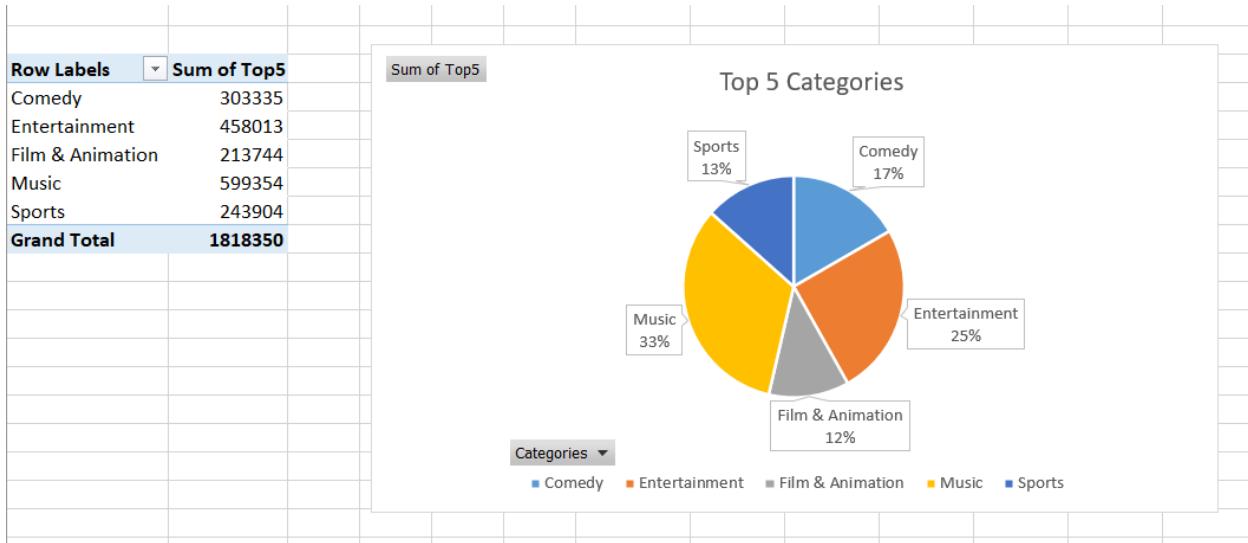
ob DAG:
ob_local1930220427_0001 -> job_local1176657791_0002,
ob_local1176657791_0002 -> job_local1850735436_0003,
ob_local1688923654_0004,
ob_local1850735436_0003 -> job_local1688923654_0004,
ob_local19393220427_0001 -> job_local19393220427_0001

021-04-27 09:04:48,323 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
021-04-27 09:04:48,346 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
021-04-27 09:04:48,339 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
021-04-27 09:04:48,353 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
021-04-27 09:04:48,367 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
021-04-27 09:04:48,385 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
021-04-27 09:04:48,397 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
021-04-27 09:04:48,398 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
021-04-27 09:04:48,404 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
021-04-27 09:04:48,410 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
021-04-27 09:04:48,416 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
021-04-27 09:04:48,424 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
021-04-27 09:04:48,427 [main] WARN org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Encountered Warning FIELD_DISCARDED_TYPE_CONVERSION_FAILED 1826 time(s).
021-04-27 09:04:48,427 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Encountered Warning ACCESSING_NON_EXISTENT_FIELD 17591 time(s).
021-04-27 09:04:48,427 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!

```



Pivot Table



2) **Top10 Rated :** Returns the Top 10 most rated Youtube videos

Script

```
infiles = load '/home/sagarshah95/Desktop/BigDataProject_pigAnalysis/data/youtubeDataset.csv' using PigStorage(',') as  
(videoid:chararray,uploader:chararray,age:int,category:chararray,length:int,views:int,rate:int,rating:int,comments:int,related_id:chararray);  
files = FILTER infiles BY category is not null;  
order_rated_video = order files by rating desc;  
top10_rated_video = limit order_rated_video 10;  
final_top10_rated_video = foreach top10_rated_video generate $0,$3,$7;  
STORE final_top10_rated_video INTO 'Top10Rated.txt' using PigStorage('|');
```

Execution

```

File Edit View Search Terminal Help
021-04-27 09:14:25,480 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.lib.SplittableMemoryManager - Selected heap (P5 Old Gen) of size 699400192 to monitor. collectionUsageThreshold = 489580128
usageThreshold = 489580128
021-04-27 09:14:25,491 [LocalJobRunner Map Task Executor #0] WARN org.apache.pig.data.SchemaTupleBackend - SchemaTupleBackend has already been initialized
021-04-27 09:14:25,502 [LocalJobRunner Map Task Executor #0] INFO org.apache.pig.backend.hadoop.executionengine.mapreduceLayer.PigGenericMapReduce$Map - Aliases being processed per job phase (AliasName[line,of set]): M: order_rated_video[4,28] C: R: final_top10_rated_video[6,26]
021-04-27 09:14:25,505 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.LocalJobRunner -
021-04-27 09:14:25,511 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.MapTask - Starting flush of map output
021-04-27 09:14:25,512 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.MapTask - Spilling map output
021-04-27 09:14:25,513 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.MapTask - Map output spilled to 1 file(s) total length = 340; bufUsed = 104857600
021-04-27 09:14:25,516 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.MapTask - kvstart = 0; kvend = 1044396(104857584); kvend = 26214360(104857440); length = 37/6553600
021-04-27 09:14:25,516 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.MapTask - Finished spill 0
021-04-27 09:14:25,537 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.Task - Taskattempt_local428338920_0004_m_000000_0 ls done. And is in the process of committing
021-04-27 09:14:25,548 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.LocalJobRunner - map
021-04-27 09:14:25,548 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.Task - Task 'attempt_local428338920_0004_m_000000_0' done.
021-04-27 09:14:25,549 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.Task - Final Counters for attempt_local428338920_0004_m_000000_0: Counters: 17
  File System Counters
    FILE: Number of bytes read=270739464
    FILE: Number of bytes written=27852141
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
Map-Reduce Framework
  Map Input records=10
  Map output records=10
  Map output bytes=0
  Map output materialized bytes=366
  Input split bytes=377
  Combine input records=0
  Spilled Records=10
  Failed Shuffles=0
  Merged Map outputs=0
  GC time elapsed (ms)=0
  Total committed heap usage (bytes)=629145600
  File Input Format Counters
    Bytes read=0
    Bytes written=0
021-04-27 09:14:25,542 [LocalJobRunner Map Task Executor #0] INFO org.apache.hadoop.mapred.LocalJobRunner - Finishing task: attempt_local428338920_0004_m_000000_0
021-04-27 09:14:25,542 [Thread-23] INFO org.apache.hadoop.mapred.LocalJobRunner - map task executor complete.
021-04-27 09:14:25,545 [Thread-23] INFO org.apache.hadoop.mapred.LocalJobRunner - Waiting for reduce tasks
021-04-27 09:14:25,554 [pool-14-thread-1] INFO org.apache.hadoop.mapred.LocalJobRunner - Starting task attempt_local428338920_0004_r_000000_0
021-04-27 09:14:25,598 [pool-14-thread-1] INFO org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter - File Output Committer Algorithm version is 2
021-04-27 09:14:25,598 [pool-14-thread-1] INFO org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter - FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cleanup fail res: false
021-04-27 09:14:25,608 [pool-14-thread-1] INFO org.apache.hadoop.mapred.Task - Using ResourceCalculatorProcessTree : []
021-04-27 09:14:25,608 [pool-14-thread-1] INFO org.apache.hadoop.mapred.Task.ReduceTask - Using ShuffleConsumerPlugin: org.apache.hadoop.mapreduce.task.reduce.Shuffle$Sba0db2
021-04-27 09:14:25,608 [pool-14-thread-1] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
021-04-27 09:14:25,608 [pool-14-thread-1] INFO org.apache.hadoop.mapreduce.task.ReduceManagerImpl - ReduceManager: memoryLimit:652528832, maxStringShuffleLimit:163132208, mergeThreshold:430669056, toSortactor:10, nonTempMergeOutputSizeThreshold:10
021-04-27 09:14:25,637 [EventFetcher for fetching Map Completion Events] INFO org.apache.hadoop.mapreduce.task.reduce.EventFetcher - attempt_local428338920_0004_r_000000_0 Thread started: EventFetcher for fetch Map Completion Events
021-04-27 09:14:25,643 [localfetcher#3] INFO org.apache.hadoop.mapreduce.task.reduce.LocalFetcher - localfetcher#3 about to shuffle output of map attempt_local428338920_0004_m_000000_0 decomp: 362 len: 366 to

```

Tue 09:22

sagarshah95@ubuntu: /usr/local/bin/pig-0.17.0/bin

```

Activities Terminal Tue 09:22
File Edit View Search Terminal Help
WRONG_REDUCE=0
Bytes Written=0
2021-04-27 09:14:25,757 [Thread-14-thread-1] INFO org.apache.hadoop.mapred.LocalJobRunner - Finishing task: attempt_local428338920_0004_r_000000_0
2021-04-27 09:14:25,761 [Thread-23] INFO org.apache.hadoop.mapred.LocalJobRunner - reduce task executor complete.
2021-04-27 09:14:25,893 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:14:25,908 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:14:25,915 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:14:25,945 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapreduceLayer.MapReduceLauncher - 100% complete
2021-04-27 09:14:25,984 [main] INFO org.apache.pig.tools.pigstats.Napreduce.SimplePigStats - Script Statistics:
HadoopVersion PigVersion UserId StartedAt FinishedAt Features
3.3.0 0.17.0 sagarshah95 2021-04-27 09:13:53 2021-04-27 09:14:25 ORDER_BY,FILTER,LIMIT
Success!
Job Stats (time in seconds):
JobID Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime OutputSize Alias Feature Outputs
job_local1884688449_0002 7 1 n/a n/a n/a n/a n/a n/a 0 final_top10_rated_video,sampler
job_local1884688449_0002 1 1 n/a n/a n/a n/a n/a n/a 0 order_rated_video SAMPLER
job_local428338920_0004 1 n/a n/a n/a n/a n/a n/a n/a 0 final_top10_rated_video,order_rated_video file:///usr/local/bin/pig-0.17.0/bin/Top10Rated.txt,ORDER_BY,COMBINER
job_local428338926_0003 1 n/a n/a n/a n/a n/a n/a n/a 0 order_rated_video

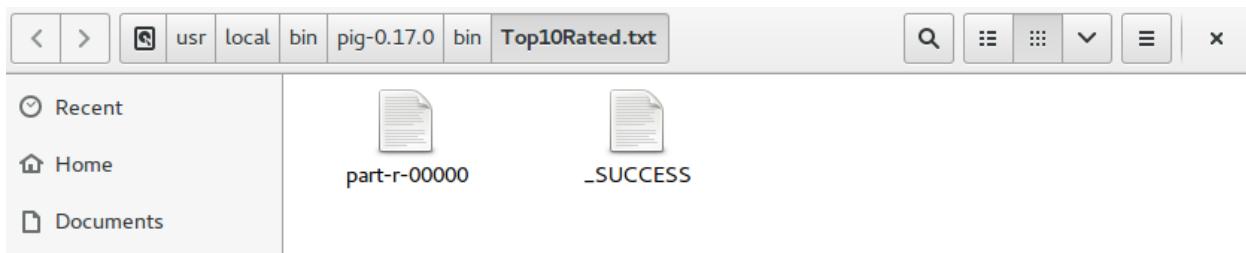
Input(s):
Successfully read 769739 records from: "/home/sagarshah95/Desktop/BigDataProject_pigAnalysis/data/youtubeDataset.csv"

Output(s):
Successfully stored 10 records in: "file:///usr/local/bin/pig-0.17.0/bin/Top10Rated.txt"

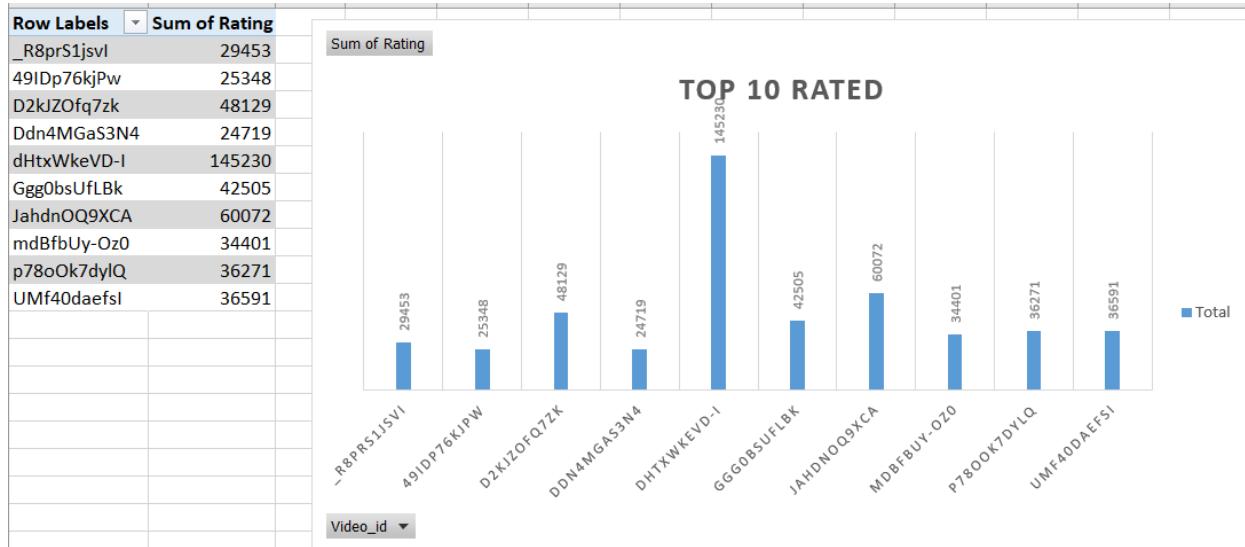
Counters:
Total records written : 10
Total bytes written : 8
Spilled local memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job Data:
job_local1421158876_0001 -> job_local1884688449_0002,
job_local1884688449_0002 -> job_local95538266_0003,
job_local195538266_0003 -> job_local428338920_0004,
job_local428338926_0004

2021-04-27 09:14:25,997 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:14:26,009 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:14:26,012 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:14:26,015 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:14:26,035 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:14:26,038 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:14:26,079 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:14:26,079 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
```



Pivot Table



3)Top10 Viewed : Return Top 10 most viewed Youtube videos

Script

```
infiles = load '/home/sagarshah95/Desktop/BigDataProject_pigAnalysis/data/youtubeDataset.csv' using PigStorage(',') as
(videoId:chararray,uploader:chararray,age:int,category:chararray,length:int,views:int,rate:int,rating:int,comments:int,related_id:chararray);
files = FILTER infiles BY category IS NOT null;
order_viewed_video = order files BY views desc;
top10_viewed_video = limit order_viewed_video 10;
final_top10_viewed_video = foreach top10_viewed_video generate $0,$3,$5;
STORE final_top10_viewed_video INTO 'Top10Viewed.txt' using PigStorage('|');
```

Execution

```
2021-04-27 09:14:26,175 [main] WARN org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Encountered Warning FIELD_DISCARDED_TYPE_CONVERSION_FAILED 4 time(s).
2021-04-27 09:14:26,176 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-04-27 09:14:26,176 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Job completed in 42 seconds and 43 milliseconds (42843 ms)
sagarshah@shahs-MacBook-Pro:~/Desktop$ pig -x local /home/sagarshah95/Desktop/BigDataProject_pigAnalysis/scripts/top10Viewed.pig
2021-04-27 09:27:53,352 INFO pig.ExectypeProvider: Trying ExecType : LOCAL
2021-04-27 09:27:53,353 INFO pig.ExectypeProvider: Picked Local as the ExecType
2021-04-27 09:27:53,525 [main] INFO org.apache.pig.Main - Apache Pig version 0.17.0 (r1797386) compiled Jun 02 2017, 15:41:58
2021-04-27 09:27:53,526 [main] INFO org.apache.pig.Main - Copyright 2009 The Apache Software Foundation
2021-04-27 09:27:53,759 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - user.name is deprecated. Instead, use mapreduce.job.user.name
2021-04-27 09:27:54,755 [main] INFO org.apache.pig.impl.util.Util - Default bootstrap file /home/sagarshah95/.pigbootstrap not found
2021-04-27 09:27:55,164 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.job.tracker.address
2021-04-27 09:27:55,165 [main] INFO org.apache.hadoop.mapred.JobClient - Connecting to hadoop file system at: file:///tmp
2021-04-27 09:27:55,413 [main] INFO org.apache.pig.PigServer - Pig is running in local mode
2021-04-27 09:27:55,414 [main] WARN org.apache.pig.PigServer - ATIS is disabled since yarn.timeline-service.enabled set to false
2021-04-27 09:27:57,938 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.textoutputformat.separator is deprecated. Instead, use mapreduce.output.textoutputformat.separator
2021-04-27 09:27:58,006 [main] INFO org.apache.pig.tools.scripts.PigScriptState - Pig features used in the script: ORDER_BY,FILTER,LIMIT
2021-04-27 09:27:58,301 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer - [RULES_ENABLED:[AllRules], ColumnMapKeyPrune, ConstantCalculator, GroupByConstParallelSetter, LimitOptimizer, LocalOptimizer, MapJoinOptimizer, MetricsOptimizer, NullOptimizer, ParallelizationOptimizer, PartitionOptimizer, PushUpFilter, SplitFilter, StreamTypeCastInserter]]
2021-04-27 09:27:58,401 [main] INFO org.apache.pig.newplan.logical.rules.ColumnPruneVisitor - Columns pruned for infiles: $1, $2, $4, $6, $7, $8, $9
2021-04-27 09:27:58,568 [main] INFO org.apache.pig.impl.util.SplittableMemoryManager - (PS Old Gen) 6994001920 to monitor. collectionUsageThreshold = 489580128, usageThreshold = 489580128
2021-04-27 09:27:58,569 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRCCompiler - File concatenation threshold: 100 optimistic? false
2021-04-27 09:27:58,571 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRCCompiler - Using Key Optimization for MapReduce node scope=28
2021-04-27 09:27:58,124 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size before optimization: 4
2021-04-27 09:27:59,124 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size after optimization: 4
2021-04-27 09:27:59,585 [main] INFO org.apache.hadoop.metrics2.lib.MetricsConfig - Loaded properties from hadoop-metrics2.properties
2021-04-27 09:28:00,289 [main] INFO org.apache.hadoop.metrics2.lib.MetricSnapshotImpl - Scheduling metric snapshot period at 10 second(s).
2021-04-27 09:28:00,306 [main] INFO org.apache.hadoop.metrics2.lib.MetricSnapshotImpl - Metrics snapshot started
2021-04-27 09:28:00,389 [main] INFO org.apache.pig.tools.pigstats.mapreduce.RSScriptState - Pig script settings are added to the job
2021-04-27 09:28:00,439 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.reduce.narkreset.buffer.percent is deprecated. Instead, use mapreduce.reduce.narkreset.buffer.percent
2021-04-27 09:28:00,439 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - mapred.job.reduce.narkreset.buffer.percent is not set, set to default 0.3
2021-04-27 09:28:00,439 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.output.compression is deprecated. Instead, use mapreduce.output.fileoutputformat.compress
2021-04-27 09:28:00,571 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.RSScriptState - Setting up signal store job
2021-04-27 09:28:00,641 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Key [pig.schematuple] is false, will not generate code
2021-04-27 09:28:00,642 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Starting process to move generated code to distributed cache
2021-04-27 09:28:00,642 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Distributed cache not supported or needed in local mode. Setting key [pig.schematuple.local.dir] with code temp directory: /tmp/161954800447
2021-04-27 09:28:00,846 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 1 map-reduce job(s) waiting for submission.
2021-04-27 09:28:00,904 [JobControl] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:28:01,353 [JobControl] INFO org.apache.hadoop.mapreduce.JobResourceUploader - No job jar file set. User classes may not be found. See Job or Job#setJar(String).
2021-04-27 09:28:01,353 [JobControl] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input files to process : 1
2021-04-27 09:28:01,508 [JobControl] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2021-04-27 09:28:01,576 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths (combined) to process : 1
2021-04-27 09:28:01,576 [JobControl] INFO org.apache.hadoop.mapreduce.JobResourceUploader - Number of inputs?
2021-04-27 09:28:01,576 [JobControl] INFO org.apache.hadoop.mapreduce.JobResourceUploader - Number of outputs?
2021-04-27 09:28:02,482 [JobControl] INFO org.apache.hadoop.mapreduce.Job - Executing with tokens: []
2021-04-27 09:28:03,152 [JobControl] INFO org.apache.hadoop.mapreduce.Job - The url to track the job: http://localhost:8080/
2021-04-27 09:28:03,154 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - HadoopJobId: job_local100907917_0001
2021-04-27 09:28:03,154 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Processing aliases files,infiles
```

```
sagarshah95@ubuntu: /usr/local/bin/pig-0.17.0/bin
File Edit View Search Terminal Help
2021-04-27 09:28:37,105 [pool-14-thread-1] INFO org.apache.hadoop.mapreduce.task.reduce.MergeManagerImpl - Merging 0 segments, 0 bytes from memory into reduce
2021-04-27 09:28:37,105 [pool-14-thread-1] INFO org.apache.hadoop.mapred.Merger - Merging 1 sorted segments
2021-04-27 09:28:37,108 [pool-14-thread-1] INFO org.apache.hadoop.mapred.Merger - Down to the last merge-pass, with 1 segments left of total size: 361 bytes
2021-04-27 09:28:37,108 [pool-14-thread-1] INFO org.apache.hadoop.mapred.LocalJobRunner - 1 / 1 copied.
2021-04-27 09:28:37,112 [pool-14-thread-1] INFO org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter - File Output Committer Algorithm version is 2
2021-04-27 09:28:37,119 [pool-14-thread-1] INFO org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter - FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cleanup fall
ures: false
2021-04-27 09:28:37,134 [pool-14-thread-1] INFO org.apache.pig.impl.util.SpillableMemoryManager - Selected heap (P9 Old Gen) of size 699400192 to monitor. collectionUsageThreshold = 489580128
2021-04-27 09:28:37,141 [pool-14-thread-1] WARN org.apache.pig.data.SchemaTupleBackend - SchemaTupleBackend has already been initialized
2021-04-27 09:28:37,162 [pool-14-thread-1] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.PigMapReduceReduce - Aliases being processed per job phase (AliasName[line,offset]): M: order_viewed
video[4,21] C: R: final_top10 viewed video[6,27]
2021-04-27 09:28:37,177 [pool-14-thread-1] INFO org.apache.hadoop.mapred.Task - Task attempt_local093968640_0004_r_000000_0 is done. And is in the process of committing
2021-04-27 09:28:37,183 [pool-14-thread-1] INFO org.apache.hadoop.mapred.LocalJobRunner - 1 / 1 copied.
2021-04-27 09:28:37,203 [pool-14-thread-1] INFO org.apache.hadoop.mapreduce.Task - Task attempt_local093968640_0004_r_000000_0 is allowed to commit now
2021-04-27 09:28:37,203 [pool-14-thread-1] INFO org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter - Saved output of task 'attempt_local093968640_0004_r_000000_0' to file:/usr/local/bin/pig-0.17.0/bin/T
op10Viewed.txt
2021-04-27 09:28:37,217 [pool-14-thread-1] INFO org.apache.hadoop.mapred.LocalJobRunner - Reduce > reduce
2021-04-27 09:28:37,228 [pool-14-thread-1] INFO org.apache.hadoop.mapred.Task - Task attempt_local093968640_0004_r_000000_0 done.
2021-04-27 09:28:37,230 [Thread-27] INFO org.apache.hadoop.mapred.Task - Final Counters for attempt_local093968640_0004_r_000000_0: Counters: 24
File System Counters
FILE: Number of bytes read=72605823
FILE: Number of bytes written=28787544
FILE: Number of small file operations=0
FILE: Number of large read operations=0
FILE: Number of write operations=0
Map-Reduce Framework
Combine input records=0
Combine output records=0
Reduce input groups=18
Reduce shuffle bytes=373
Reduce input records=10
Reduce output records=10
Split input records=16
Shuffled Maps =1
Failed Shuffles=0
Merged Map outputs=1
GC Time elapsed (ms)=0
Total committed heap usage (bytes)=634912768
Shuffle Errors
BAD_ID=0
CONNECTION=0
IO_ERROR=0
WRONG_LENGTH=0
WRONG_MAP=0
WRONG_REDUCE=0
File Output Format Counters
Bytes written=0
2021-04-27 09:28:37,234 [pool-14-thread-1] INFO org.apache.hadoop.mapred.LocalJobRunner - Finishing task: attempt_local093968640_0004_r_000000_0
2021-04-27 09:28:37,235 [Thread-23] INFO org.apache.hadoop.mapred.LocalJobRunner - reduce task executor complete.
```

```
Activities Terminal sagarshah95@ubuntu: /usr/local/bin/pig-0.17.0/bin
File Edit View Search Terminal Help
2021-04-27 09:28:37,452 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 100% complete
2021-04-27 09:28:37,476 [main] INFO org.apache.pig.tools.pigstats.Script Statistics:
HadoopVersion PigVersion UserId StartedAt FinishedAt Features
3.3.0 0.17.0 sagarshah95 2021-04-27 09:28:00 2021-04-27 09:28:37 ORDER_BY,FILTER,LIMIT
Success!
Job Stats (time in seconds):
JobId Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime MedianReducetime Alias Feature Outputs
job_local1090907917_0001 0 n/a n/a n/a 0 0 0 0 files,inf files MAP_ONLY
job_local1389838244_0002 1 1 n/a n/a n/a n/a n/a n/a order_viewed_video SAMPLER
job_local1693968640_0004 1 1 n/a n/a n/a n/a n/a n/a final_top10_viewed_video,order_viewed_videofile:///usr/local/bin/pig-0.17.0/bin/Top10Viewed.txt,
job_local1710553747_0003 1 n/a n/a n/a n/a n/a n/a n/a order_viewed_video ORDER_BY,COMBINER

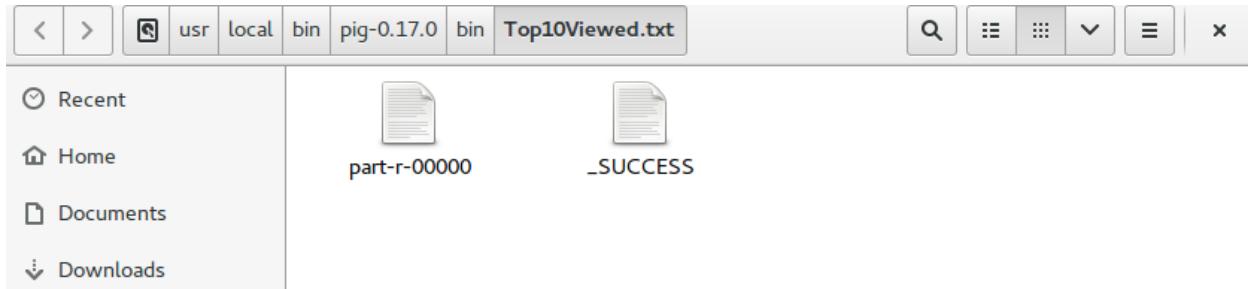
Input(s):
Successfully read 769739 records from: "/home/sagarshah95/Desktop/BigDataProject_pigAnalysis/data/youtubeDataset.csv"

Output(s):
Successfully stored 10 records in: "file:///usr/local/bin/pig-0.17.0/bin/Top10Viewed.txt"

Counters:
Total records written : 10
Total bytes written : 0
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

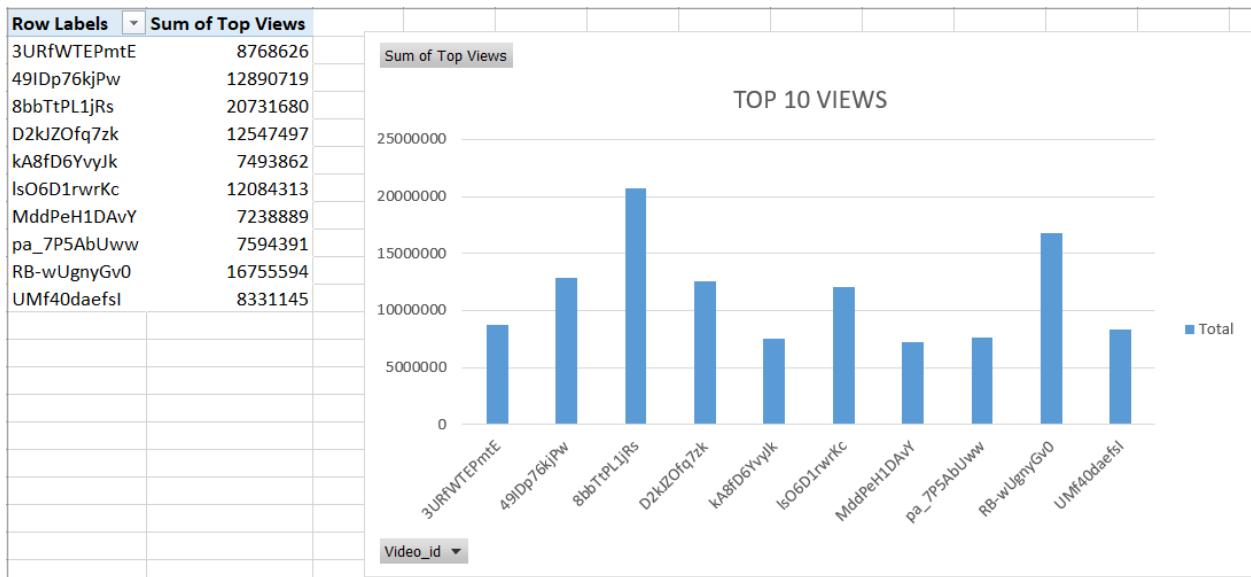
Job DAG:
job_local1090907917_0001 -> job_local1389838244_0002,
job_local1389838244_0002 -> job_local1710553747_0003,
job_local1710553747_0003 -> job_local1693968640_0004,
job_local1693968640_0004

2021-04-27 09:28:37,493 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:28:37,501 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:28:37,509 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:28:37,548 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:28:37,551 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:28:37,554 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:28:37,557 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:28:37,586 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:28:37,590 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:28:37,594 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:28:37,628 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:28:37,629 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:28:37,629 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:28:37,640 [main] WARN org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Encountered Warning FIELD_DISCARDED_TYPE_CONVERSION_FAILED 4 time(s).
2021-04-27 09:28:37,640 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
```



part-r-00000	
/usr/local/bin/pig-0.17.0/bin/Top10Viewed.txt	
dMH0bHeiRNg	Comedy 42513417
0XXI-hvPRRA	Comedy 20282464
1dmVU08zVpA	Entertainment 16087899
RB-wUgnyGv0	Entertainment 15712924
QJASfaZF1A8	Music 15256922
_CSo1gOd48	People & Blogs 13199833
49IDp76kjPw	Comedy 11970018
tYnn51C3X_w	Music 11823701
pv5zWaTEVkI	Music 11672017
D2kJZOfq7zk	People & Blogs 11184051

Pivot Table



4)Top10 rated by categories : Returns top 10 rated Youtube videos based on categories

Script

```

infiles = load '/home/sagarshah95/Desktop/BigDataProject_pigAnalysis/data/youtubeDataset.csv' using PigStorage(',') as
(videoid:chararray,uploader:chararray,age:int,category:chararray,length:int,views:int,rate:int,rating:int,comments:int,related_id:chararray);
files = FILTER infiles BY category is not null;
grpns_for_categories = group files by category;
top10_rated_categories = foreach grpns_for_categories{
    sorted = order files by rating desc;
    top10 = limit sorted 10;
    generate flatten(top10);
};
top10_rated_by_categories = foreach top10_rated_categories generate $0,$3,$7;
STORE top10_rated_by_categories INTO 'Top10RatedByCategories.pig' using PigStorage('|');

```

Execution

```
sagarshah95@ubuntu: /usr/local/bin/pig-0.17.0/bin
File Edit View Search Terminal Help
2021-04-27 09:28:37,649 [main] WARN org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Encountered Warning ACCESSING_NON_EXISTENT_FIELD 3026 time(s).
2021-04-27 09:28:37,649 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-04-27 09:28:37,733 [main] INFO org.apache.pig.Main - Pig script completed in 45 seconds and 876 milliseconds (45876 ms)
sagarshah95@ubuntu: /usr/local/bin/pig-0.17.0/bin> pig -x local /home/sagarshah95/Desktop/BIGDATA/project_pigAnalysts/scripts/top10RatedByCategories.pig
2021-04-27 09:43:23,832 [main] INFO org.apache.pig.ExcuteTypeProvider - PigExcuteTypeProvider - LOCAL as the ExcuteType
2021-04-27 09:43:23,848 [main] INFO org.apache.pig.ExcuteTypeProvider - PigExcuteTypeProvider - LOCAL as the ExcuteType
2021-04-27 09:43:24,026 [main] INFO org.apache.pig.Main - Apache Pig version 0.17.0 (r1797386) compiled Jun 02 2017 15:41:58
2021-04-27 09:43:24,026 [main] INFO org.apache.pig.Main - Logging error messages to: /usr/local/bin/pig-0.17.0/bin/1619541803999.log
2021-04-27 09:43:24,026 [main] INFO org.apache.pig.Main - To set log4j appenders in your application, use log4j.appender.job.user.name
2021-04-27 09:43:24,026 [main] INFO org.apache.pig.impl.util.Unix - Default bootstrap file /home/sagarshah95/pigbootstrap not found
2021-04-27 09:43:25,284 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-04-27 09:43:25,284 [main] INFO org.apache.pig.backend.hadoop.executionengine.HExecutionEngine - Connecting to hadoop file system at: file:///
2021-04-27 09:43:25,394 [main] INFO org.apache.pig.PigServer - Pig Script ID for the session: PIG-top10RatedByCategories.pig-b926dbf-01d7-4578-8125-ae559cc7e94f
2021-04-27 09:43:25,394 [main] INFO org.apache.pig.PigServer - YARN Session ID for the session: 1619541803999
2021-04-27 09:43:27,879 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used in the script: GROUP_BY_FILTER
2021-04-27 09:43:27,998 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used in the script: GROUP_BY_FILTER
2021-04-27 09:43:27,998 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer - [RULES_ENABLED=[AddForEach, ColumnWisePrune, ConstantCalculator, GroupByConstParallelSetter, LimitOptimizer, LoadTypeCastInserter, MergeFilter, MergeForJoin, NeededWithOptimizer, PartitionOptimizer, PredictivePushdownOptimizer, PushdownEachLimiter, PushUpFilter, SplitFilter, StreamTypeCastInserter]]
2021-04-27 09:43:28,018 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer - [RULES_DISABLE=[AddForEach, ColumnWisePrune, ConstantCalculator, GroupByConstParallelSetter, LimitOptimizer, LoadTypeCastInserter, MergeFilter, MergeForJoin, NeededWithOptimizer, PartitionOptimizer, PredictivePushdownOptimizer, PushdownEachLimiter, PushUpFilter, SplitFilter, StreamTypeCastInserter]]
2021-04-27 09:43:28,018 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer - [RULES_ENABLE=[AddForEach, ColumnWisePrune, ConstantCalculator, GroupByConstParallelSetter, LimitOptimizer, LoadTypeCastInserter, MergeFilter, MergeForJoin, NeededWithOptimizer, PartitionOptimizer, PredictivePushdownOptimizer, PushdownEachLimiter, PushUpFilter, SplitFilter, StreamTypeCastInserter]]
2021-04-27 09:43:28,839 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRCCompiler - File concatenation threshold= 100 optimistic? false
2021-04-27 09:43:29,019 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRCCompiler - M plan size before optimization: 1
2021-04-27 09:43:29,021 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MRCCompiler - M plan size after optimization: 1
2021-04-27 09:43:29,507 [main] INFO org.apache.hadoop.metrics2.impl.MetricsConfig - Loaded properties from hadoop-metrics2.properties
2021-04-27 09:43:29,507 [main] INFO org.apache.hadoop.metrics2.impl.MetricsConfig - Loaded properties from hadoop-metrics2.properties
2021-04-27 09:43:30,204 [main] INFO org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system started
2021-04-27 09:43:30,316 [main] INFO org.apache.hadoop.tools.pigstats.MRScriptState - Pig script settings are added to the job
2021-04-27 09:43:30,335 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.reduce.markreset.buffer.percent is deprecated. Instead, use mapreduce.reduce.markreset.buffer.percent
2021-04-27 09:43:30,335 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.reduce.markreset.buffer.percent is not set, so default 0.3
2021-04-27 09:43:30,344 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.output.compression.codec is deprecated. Instead, use mapreduce.output.compress
2021-04-27 09:43:30,351 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Reduce phase detected, estimating # of required reducers
2021-04-27 09:43:30,355 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Using reducer estimator: org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.InputsLzeReducerEstimator
2021-04-27 09:43:30,360 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.InputsPerReducer - maxReducers=1000000000 maxReducers=999 totalInputFileSize=219815043
2021-04-27 09:43:30,369 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Setting Parallelism to 1
2021-04-27 09:43:30,369 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.reduce.tasks is deprecated. Instead, use mapreduce.job.reduces
2021-04-27 09:43:30,412 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Setting single store job
2021-04-27 09:43:30,509 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Key [pig.schematuple] is false, will not generate code.
2021-04-27 09:43:30,509 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Starting process to move generated code to distributed cache
2021-04-27 09:43:30,509 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Distributed cache not supported or needed in local mode. Setting key [pig.schematuple.local.dir] with code temp directory: /tmp/16195808-6
2021-04-27 09:43:30,954 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 1 map-reduce job(s) waiting for submission.
2021-04-27 09:43:31,019 [JobControl] INFO org.apache.hadoop.mapred.LocalJobController - JobTracker metrics system already initialized
2021-04-27 09:43:31,019 [JobControl] INFO org.apache.hadoop.mapred.LocalJobController - Configuration deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.job.tracker
2021-04-27 09:43:31,548 [JobControl] WARN org.apache.hadoop.mapreduce.JobResourceUploader - No Job jar file set. User classes may not be found. See Job or JobSetJar(String).
2021-04-27 09:43:31,609 [JobControl] INFO org.apache.pig.PigStorage - Using PLGTextInputFormat
2021-04-27 09:43:31,651 [JobControl] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total Input files to process : 1
2021-04-27 09:43:31,652 [JobControl] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCompiler - Total input paths to process : 1
2021-04-27 09:43:31,652 [JobControl] INFO org.apache.hadoop.mapreduce.Job - Total Input paths (combined) to process : 7
2021-04-27 09:43:31,984 [JobControl] INFO org.apache.hadoop.mapreduce.JobSubmitter - number of splits?
2021-04-27 09:43:32,625 [JobControl] INFO org.apache.hadoop.mapreduce.JobSubmitter - submitting tokens for job: job_local2046767052_0001
sagarshah95@ubuntu: /usr/local/bin/pig-0.17.0/bin
```

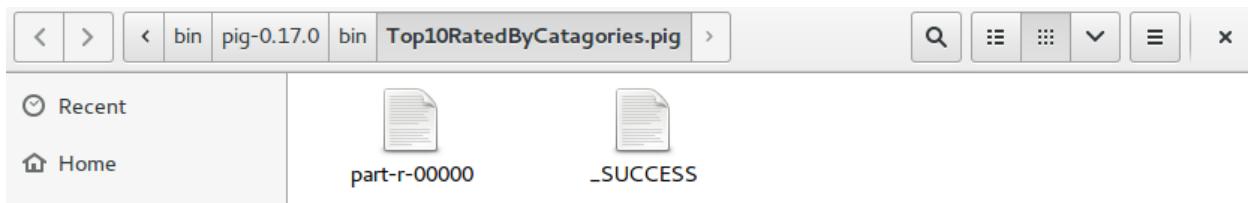
```
sagarshah95@ubuntu: /usr/local/bin/pig-0.17.0/bin
File Edit View Search Terminal Help
2021-04-27 09:44:03,137 [pool-4-thread-1] INFO org.apache.hadoop.mapred.local.JobRunner - 7 / 7 copied.
2021-04-27 09:44:03,171 [pool-4-thread-1] INFO org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter - File Output Committer Algorithm version is 2
2021-04-27 09:44:03,171 [pool-4-thread-1] INFO org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter - FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cleanup fallus: false
2021-04-27 09:44:03,181 [pool-4-thread-1] INFO org.apache.pig.impl.SplittableMemoryManager - Selected heap (PS Old Gen) of size 699400192 to monitor. collectionUsageThreshold = 44958018
2021-04-27 09:44:03,181 [pool-4-thread-1] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.PigMapReduceSReduce - mapred.skip.on is deprecated. Instead, use mapreduce.job.skipprecords
2021-04-27 09:44:03,214 [pool-4-thread-1] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.skip.on is deprecated. Instead, use mapreduce.job.skipprecords
2021-04-27 09:44:03,214 [pool-4-thread-1] INFO org.apache.pig.impl.SplittableMemoryManager - Selected heap (PS Old Gen) of size 699400192 to monitor. collectionUsageThreshold = 44958018
2021-04-27 09:44:03,236 [pool-4-thread-1] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.PigMapReduceSReduce - Aliases being processed per job phase (AliasName[line,offset]): M: [infiles[1,10], infiles[-1,-1], files[3,8], grp[0,0], categories[4,22], C: R: top10_rated_categories[5,25], sorted[0,30], top10_rated_by_categories[10,28]
2021-04-27 09:44:03,236 [pool-4-thread-1] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.PigMapReduceSReduce - Reducer 0: [infiles[1,10], files[3,8], grp[0,0], categories[4,22], C: R: top10_rated_categories[5,25], sorted[0,30], top10_rated_by_categories[10,28]
2021-04-27 09:44:03,236 [pool-4-thread-1] INFO org.apache.hadoop.mapred.LocalJobRunner - 7 / 7 copied
2021-04-27 09:44:03,236 [pool-4-thread-1] INFO org.apache.hadoop.mapred.Task - Task attempt_local2046767052_0001_r_000000_0 is done. And is in the process of committing
2021-04-27 09:44:03,236 [pool-4-thread-1] INFO org.apache.hadoop.mapred.Task - Task attempt_local2046767052_0001_r_000000_0 is allowed to commit now
2021-04-27 09:44:12,306 [pool-4-thread-1] INFO org.apache.hadoop.mapred.Task - Task attempt_local2046767052_0001_r_000000_0 is saved
2021-04-27 09:44:12,306 [pool-4-thread-1] INFO org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter - Saved output of task 'attempt_local2046767052_0001_r_000000_0' to file:/usr/local/bin/pig-0.17.0/bin/T_ip100000ByCategories.pig
2021-04-27 09:44:12,315 [pool-4-thread-1] INFO org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter - Saved output of task 'attempt_local2046767052_0001_r_000000_0' to file:/usr/local/bin/pig-0.17.0/bin/T_ip100000ByCategories.pig
2021-04-27 09:44:12,335 [pool-4-thread-1] INFO org.apache.hadoop.mapred.Task - Task attempt_local2046767052_0001_r_000000_0 done.
2021-04-27 09:44:12,336 [pool-4-thread-1] INFO org.apache.hadoop.mapred.Task - Final Counters for attempt_local2046767052_0001_r_000000_0: Counters: 24
File System Counters
FILE: Number of bytes read=31279387
FILE: Number of bytes written=12840620
FILE: Number of read operations=0
FILE: Number of large read operations=0
FILE: Number of write operations=0
Map-Reduce Framework
  Combine input records=0
  Combine output records=0
  Reduce input records=15
  Reduce input groups=15
  Reduce shuffle bytes=5712909
  Reduce shuffle records=5712909
  Reduce output records=144
  Spilled Records=763889
  Shuffled Maps = 7
  Failed Shuffles=0
  Merged Map outputs=7
  GC time elapsed (ms)=532
  Total committed heap usage (bytes)=822083584
Shuffle Errors
  Bytes Written=0
  CONNECTION=0
  IO_ERROR=0
  WRONG_LENGTH=0
  WRONG_MAP=0
  WRONG_REDUCE=0
  File Output Format Counters
    Bytes Written=0
2021-04-27 09:44:12,339 [pool-4-thread-1] INFO org.apache.hadoop.mapred.LocalJobRunner - Finishing task: attempt_local2046767052_0001_r_000000_0
2021-04-27 09:44:12,346 [Thread-6] INFO org.apache.hadoop.mapred.LocalJobRunner - reduce task executor complete.
2021-04-27 09:44:12,481 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:44:12,537 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
```

```

sagarshah95@ubuntu: /usr/local/bin/pig-0.17.0/bin
File Edit View Search Terminal Help
IO_ERROR=0
WRONG_LENGTH=0
WRONG_MAP=0
WRONG_REDUCE=0
File Output Format Counters
Bytes Written:0
2021-04-27 09:44:12,339 [pool-4-thread-1] INFO org.apache.hadoop.mapred.LocalJobRunner - Finishing task: attempt_local2046767052_0001_r_000000_0
2021-04-27 09:44:12,346 [Thread-6] INFO org.apache.hadoop.mapred.LocalJobRunner - reduce task executor complete.
2021-04-27 09:44:12,346 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:44:12,357 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - Metrics system already initialized!
2021-04-27 09:44:12,542 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.map.tasks is deprecated. Instead, use mapreduce.job.maps
2021-04-27 09:44:12,546 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:44:12,752 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - MapReduceLauncher - 100% complete
2021-04-27 09:44:12,767 [main] INFO org.apache.pig.tools.pigstats.mapreduce.SimplePigStats - Script Statistics:
HadoopVersion PigVersion Userid Startedat Finishedat Features
3.3.0 0.17.0 sagarshah95 2021-04-27 09:43:30 2021-04-27 09:44:12 GROUP_BY,FILTER
Success!
Job Stats (time in seconds):
JobId Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime MedianReduceTime Alias Feature Outputs
job_local2046767052_0001 7 1 n/a n/a n/a n/a n/a n/a n/a files,grpn_for_catagories,infiles,sorted,top10_rated_by_catagories,top10_rated_catagories GRO
UP_BY file:///usr/local/bin/pig-0.17.0/bin/Top10RatedByCatagories.pig,
Input(s):
Successfully read 769739 records from: "/home/sagarshah95/Desktop/BIGDataProject_pigAnalysis/data/youtubeDataset.csv"
Output(s):
Successfully stored 144 records in: "file:///usr/local/bin/pig-0.17.0/bin/Top10RatedByCatagories.pig"
Counters:
Total records written : 144
Total bytes written : 0
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0
Job DAG:
job_local2046767052_0001

2021-04-27 09:44:12,770 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:44:12,774 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:44:12,780 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:44:12,798 [main] WARN org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Encountered Warning FIELD_DISCARDED_TYPE_CONVERSION_FAILED 1826 time(s).
2021-04-27 09:44:12,798 [main] WARN org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Encountered Warning ACCESSING_NON_EXISTENT_FIELD 17591 time(s).
2021-04-27 09:44:12,800 [main] INFO org.apache.pig.Main - Pig script completed in 50 seconds and 122 milliseconds (50122 ms)
2021-04-27 09:44:12,910 [main] INFO org.apache.pig.Main - Pig script completed in 50 seconds and 122 milliseconds (50122 ms)

sagarshah95@ubuntu:/usr/local/bin/pig-0.17.0/bin$
```



Open		part-r-00000	Save		
		/usr/local/bin/pig-0.17.0/bin/Top10RatedByCatagories.pig			
R0049_TDAU8 UNA 70972					
aRNzWyD7C9o UNA 36815					
3jLRNik6oWY UNA 3128					
jFoyIcuhqos UNA 2728					
jFoyIcuhqos UNA 2728					
LIhbap3FlGc UNA 2704					
LIhbap3FlGc UNA 2704					
c8sW2Qzhh38 UNA 2568					
X8w_-9M5_SA UNA 2234					
QjA5faZF1A8 Music 120506					
pv5zWaTEVkI Music 42386					
UMf40daefsI Music 31886					
tYnn51C3X_w Music 29479					
FLn45-7Pn2Y Music 21249					
FLn45-7Pn2Y Music 21249					
-xEzGIuY7kw Music 20828					
HSoVKUVOnfQ Music 19803					
HSoVKUVOnfQ Music 19803					
ARHyRI9_NB4 Music 19243					
dMH0bHeiRNg Comedy 87520					
0XxI-hvPRRA Comedy 80710					
nojWJ6-XmeQ Comedy 62265					
CQ03K8BcyGM Comedy 36460					
5P6UU6m3cqk Comedy 34972					
_R8prS1jsvI Comedy 29453					
49IDp76kjPw Comedy 22579					
N6j475XI1Xg Comedy 20593					
PKnloiM-0Ns Comedy 19728					
6D9p-wmtIJc Comedy 19043					
sMEnn0xJ0l0 Sports 11478					
Ugrlz7fySE Sports 8541					
8X2_zsnPkq8 Sports 7359					
qjWQNwv-GJ4 Sports 6377					
KtdYsd_QRmc Sports 5819					

Open	Save	≡	X
part-r-00000			
/usr/local/bin/pig-0.17.0/bin/Top10RatedByCatagories.pig			
0DolbU5njEM News & Politics 9722			
xDh_pvv1tUM News & Politics 9543			
QCVxQ_3Ejkg News & Politics 8120			
qdSSlkeN8_8 News & Politics 6798			
bNF_P281Uu4 Travel & Places 29152			
RIH1I1doUI4 Travel & Places 4077			
_5QUdvUhCZc Travel & Places 4007			
WXL-CTMku1o Travel & Places 3260			
lqbt6X4ZgEI Travel & Places 2966			
8L7SxcBiDOY Travel & Places 2615			
tIzycq8252Q Travel & Places 2223			
bGkZSiENKIA Travel & Places 2159			
m9A_vxIOB-I Travel & Places 1966			
dVRUBIyRAYk Travel & Places 1906			
RjreQaG5jPM Autos & Vehicles 5034			
46LQd9dXF瑞 Autos & Vehicles 3852			
46LQd9dXF瑞 Autos & Vehicles 3852			
cv157ZIInUk Autos & Vehicles 2850			
8c1GGgXLepY Autos & Vehicles 2487			
8c1GGgXLepY Autos & Vehicles 2487			
aCamHfJwSGU Autos & Vehicles 2108			
aCamHfJwSGU Autos & Vehicles 2108			
3Xo5GY4kDXg Autos & Vehicles 1770			
SLTvSUCCqPo Autos & Vehicles 1730			
sduUx5FdySs Film & Animation 42417			
6B26asyGKDo Film & Animation 31792			
JzqumbhfxRo Film & Animation 27545			
h7svw0m-w00 Film & Animation 27372			
h7svw0m-w00 Film & Animation 27372			
AJzU3NjDikY Film & Animation 22973			
PnCVZozHTG8 Film & Animation 14309			
hquCf-sS2sU Film & Animation 13377			
o9698TqtY4A Film & Animation 12457			
o9698TqtY4A Film & Animation 12457			

5) Top10 Viewed by categories : Returns top 10 viewed Youtube videos based on categories

Script

```

infiles = load '/home/sagarshah95/Desktop/BigDataProject_pigAnalysis/data/youtubeDataset.csv' using PigStorage(',') as
(videoid:chararray,uploader:chararray,age:int,category:chararray,length:int,views:int,rate:int,rating:int,comments:int,related_id:chararray);
files = FILTER infiles BY category IS NOT NULL;
grpn_for_categories = group files by category;
top10_viewed_categories = foreach grpn_for_categories{
    sorted = order files by views desc;
    top10 = limit sorted 10;
    generate flatten(top10);
};
top10_viewed_by_categories = foreach top10_viewed_categories generate $0,$3,$5;
STORE top10_viewed_by_categories INTO 'Top10ViewedByCategories.txt' using PigStorage(' ');

```

Execution

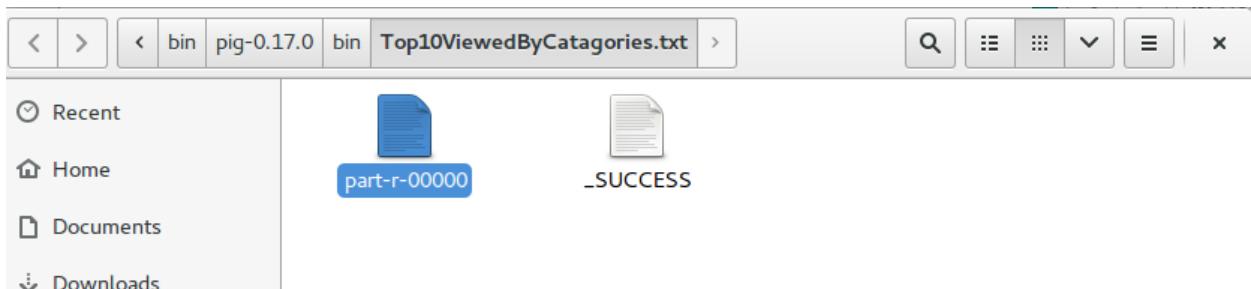
```

sagarshah95@ubuntu: /usr/local/bin/pig-0.17.0/bin
File Edit View Search Terminal Help
2021-04-27 09:44:12,798 [main] WARN org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Encountered Warning ACCESSING_NON_EXISTENT_FIELD 17591 time(s).
2021-04-27 09:44:12,798 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-04-27 09:44:12,910 [main] INFO org.apache.pig.Main - Pig script completed in 50 seconds and 122 milliseconds (50122 ms)
2021-04-27 09:51:22,514 INFO pig.Executor - ExecutorProvider: TryLocal ExecutorType: LOCAL
2021-04-27 09:51:22,517 INFO pig.ExecutorTypeProvider: Picked LOCAL as the ExecType
2021-04-27 09:51:22,690 [main] INFO org.apache.pig.Main - Apache Pig version 0.17.0 (r1797608) compiled Jun 02 2017, 15:41:58
2021-04-27 09:51:22,690 [main] INFO org.apache.pig.Main - Logging error messages to: /usr/local/bin/pig-0.17.0/bin/pig_161954228679.log
2021-04-27 09:51:22,793 [main] INFO org.apache.pig.Main - Logging error messages to: /usr/local/bin/pig-0.17.0/bin/pig_161954228679.log
2021-04-27 09:51:23,011 [main] INFO org.apache.pig.impl.util.Util - Default bootstrap file: /home/sagarshah95/Desktop/BigDataProject_piganalysis/scripts/top10ViewedByCategories.pig
2021-04-27 09:51:23,619 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2021-04-27 09:51:23,626 [main] INFO org.apache.pig.backend.hadoop.executionengine.HExecutionEngine - Connecting to hadoop file system at: file:///
2021-04-27 09:51:23,785 [main] INFO org.apache.pig.PigServer - Pig Script ID for the session: pig_top10ViewedByCategories.pig-14e78c82-dc4d-4448-a229-a133f29cd9c9
2021-04-27 09:51:23,785 [main] INFO org.apache.pig.PigServer - ATS is disabled since yarn.resourcemanager.service.address is not provided
2021-04-27 09:51:23,785 [main] INFO org.apache.pig.PigServer - InputFormat is deprecated. Instead, use mapreduce.input.textinputformat.separator
2021-04-27 09:51:23,785 [main] INFO org.apache.pig.PigServer - OutputFormat is deprecated. Instead, use mapreduce.output.textoutputformat.separator
2021-04-27 09:51:23,785 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used in the script! : GROUP_BY_FILTER
2021-04-27 09:51:23,927 [main] INFO org.apache.pig.newplan.logical.optimizer.LogicalPlanOptimizer - [RULES_ENABLEDForAdd, ColumnByKeyPrune, ConstantCalculator, GroupByConstParallelSetter, LimitOptimizer, LoadTypeCastInserter, MergeFilter, MergeForEach, NestedLmToOptimizer, PartitionFilterOptimizer, PredicatePushdownOptimizer, PushdownForEachFlatlen, PushupFilter, SplitFilter, StreamTypeCastInserter]
2021-04-27 09:51:26,128 [main] INFO org.apache.pig.impl.util.SplittableRecordReader - Reader heap size: 6992000192 to monitor. collectionUsageThreshold = 489580128, usageThreshold = 489580128
2021-04-27 09:51:26,130 [main] INFO org.apache.pig.impl.io.NullableRecordReader - Reader heap size: 6992000192 to monitor. collectionUsageThreshold = 489580128, usageThreshold = 489580128
2021-04-27 09:51:26,513 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size before optimization: 1
2021-04-27 09:51:26,513 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MultiQueryOptimizer - MR plan size after optimization: 1
2021-04-27 09:51:26,925 [main] INFO org.apache.hadoop.metrics2.sink.inpl.MetricsConfig - Loaded properties from hadoop-metrics2.properties
2021-04-27 09:51:27,475 [main] INFO org.apache.hadoop.metrics2.sink.inpl.MetricsSystemImpl - Scheduled Metrics snapshot period at 10 second(s).
2021-04-27 09:51:27,475 [main] INFO org.apache.hadoop.metrics2.sink.inpl.MetricsSystemImpl - JobContext metrics system started
2021-04-27 09:51:27,595 [main] INFO org.apache.pig.tools.pigstats.MetricsScriptState - MetricsScriptState class is added to the job
2021-04-27 09:51:27,617 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.reduce.markRset.buffer.percent is deprecated. Instead, use mapreduce.reduce.markrset.buffer.percent
2021-04-27 09:51:27,617 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.reduce.markRset.buffer.percent is not set, set to default 0.3
2021-04-27 09:51:27,626 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.output.compress is deprecated. Instead, use mapreduce.output.fileoutputformat.compress
2021-04-27 09:51:27,637 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.output.compress is deprecated. Instead, use mapreduce.output.fileoutputformat.compress
2021-04-27 09:51:27,641 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCommitter - Reduc phase detected, estimating # of required reducers.
2021-04-27 09:51:27,641 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCommitter - Using reducer estimator: org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.InputsListReducerEstimator
2021-04-27 09:51:27,655 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.InputSizeReducerEstimator - BytesPerReducer=1000000000 maxReducers=999 totalInputFileSize=219815843
2021-04-27 09:51:27,655 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.JobControlCommitter - Setting Parallelism to 1
2021-04-27 09:51:27,655 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.reduce.tasks is deprecated. Instead, use mapreduce.job.reduces
2021-04-27 09:51:27,751 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Key [pig.schematuple] is false, will not generate code.
2021-04-27 09:51:27,752 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Starting process to move generated code to distributed casche
2021-04-27 09:51:27,753 [main] INFO org.apache.pig.data.SchemaTupleFrontend - Distributed cache not supported or needed in local mode. Setting key [pig.schematuple.local.dir] with code temp directory: /tmp/161954228679
542287751
2021-04-27 09:51:28,039 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 1 map-reduce job(s) waiting for submission.
2021-04-27 09:51:28,118 [JobControl] WARN org.apache.hadoop.metrics2.sink.inpl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:51:28,233 [JobControl] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.task.id is deprecated. Instead, use mapreduce.task.attempt.id
2021-04-27 09:51:28,495 [JobControl] WARN org.apache.hadoop.mapreduce.JobResourceUploader - No jar file set. User classes may not be found. See Job or Job#setJar(String).
2021-04-27 09:51:28,532 [JobControl] INFO org.apache.pig.backend.pigStorage - Using PigTextInputFormat
2021-04-27 09:51:28,565 [JobControl] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Input files to process : 1
2021-04-27 09:51:28,565 [JobControl] INFO org.apache.hadoop.mapreduce.lib.map.ReduceUtil - Total input paths to process : 1
2021-04-27 09:51:28,656 [JobControl] INFO org.apache.hadoop.mapreduce.lib.map.ReduceUtil - Total input paths (combined) to process : 7
2021-04-27 09:51:28,885 [JobControl] INFO org.apache.hadoop.mapreduce.JobSubmitter - number of splits?: 7
2021-04-27 09:51:29,260 [JobControl] INFO org.apache.hadoop.mapreduce.JobSubmitter - Submitting tokens for job: job_local423179355_0001

```

```
sagarshah95@ubuntu: /usr/local/bin/pig-0.17.0/bin
File Edit View Search Terminal Help
2021-04-27 09:51:48,391 [localfetcher#1] INFO org.apache.hadoop.mapreduce.task.reduce.LocalFetcher - localfetcher#1 about to shuffle output of map attempt_local423179355_0001_m_000002_0 decomps: 8499037 len: 849
9041 _MEMORY
2021-04-27 09:51:48,412 [localfetcher#1] INFO org.apache.hadoop.mapreduce.task.reduce.InMemoryMapOutput - Read 8499037 bytes from map-output for attempt_local423179355_0001_m_000002_0
2021-04-27 09:51:48,412 [localfetcher#1] INFO org.apache.hadoop.mapreduce.task.reduce.MergeManagerImpl - closeInMemoryfile -> map-output of size: 8499037, inMemoryMapOutputs.size() -> 3, commitMemory -> 1702576
6, usedMemory -> 25524803
2021-04-27 09:51:48,430 [localfetcher#1] INFO org.apache.hadoop.mapreduce.task.reduce.LocalFetcher - localfetcher#1 about to shuffle output of map attempt_local423179355_0001_m_000005_0 decomps: 8461601 len: 846
1605 _MEMORY
2021-04-27 09:51:48,453 [localfetcher#1] INFO org.apache.hadoop.mapreduce.task.reduce.InMemoryMapOutput - Read 8461601 bytes from map-output for attempt_local423179355_0001_m_000005_0 decomps: 8461601 len: 846
2021-04-27 09:51:48,453 [localfetcher#1] INFO org.apache.hadoop.mapreduce.task.reduce.MergeManagerImpl - closeInMemoryfile -> map-output of size: 8461601, inMemoryMapOutputs.size() -> 4, commitMemory -> 2552480
3, usedMemory -> 31986404
2021-04-27 09:51:48,462 [localfetcher#1] INFO org.apache.hadoop.mapreduce.task.reduce.LocalFetcher - localfetcher#1 about to shuffle output of map attempt_local423179355_0001_m_000006_0 decomps: 4691089 len: 469
1093 _MEMORY
2021-04-27 09:51:48,470 [localfetcher#1] INFO org.apache.hadoop.mapreduce.task.reduce.InMemoryMapOutput - Read 4691089 bytes from map-output for attempt_local423179355_0001_m_000006_0
2021-04-27 09:51:48,470 [localfetcher#1] INFO org.apache.hadoop.mapreduce.task.reduce.MergeManagerImpl - closeInMemoryfile -> map-output of size: 4691089, inMemoryMapOutputs.size() -> 5, commitMemory -> 3398640
4, usedMemory -> 38677493
2021-04-27 09:51:48,488 [localfetcher#1] INFO org.apache.hadoop.mapreduce.task.reduce.LocalFetcher - localfetcher#1 about to shuffle output of map attempt_local423179355_0001_m_000008_0 decomps: 8480182 len: 848
0186 _MEMORY
2021-04-27 09:51:48,499 [localfetcher#1] INFO org.apache.hadoop.mapreduce.task.reduce.InMemoryMapOutput - Read 8480182 bytes from map-output for attempt_local423179355_0001_m_000008_0 decomps: 8480182 len: 848
2021-04-27 09:51:48,500 [localfetcher#1] INFO org.apache.hadoop.mapreduce.task.reduce.MergeManagerImpl - closeInMemoryfile -> map-output of size: 8480182, inMemoryMapOutputs.size() -> 6, commitMemory -> 3867749
3, usedMemory -> 47157675
2021-04-27 09:51:48,526 [localfetcher#1] INFO org.apache.hadoop.mapreduce.task.reduce.LocalFetcher - localfetcher#1 about to shuffle output of map attempt_local423179355_0001_m_000009_0 decomps: 8555206 len: 855
5210 _MEMORY
2021-04-27 09:51:48,555 [localfetcher#1] INFO org.apache.hadoop.mapreduce.task.reduce.InMemoryMapOutput - Read 8555206 bytes from map-output for attempt_local423179355_0001_m_000009_0 decomps: 8555206 len: 855
2021-04-27 09:51:48,556 [localfetcher#1] INFO org.apache.hadoop.mapreduce.task.reduce.MergeManagerImpl - closeInMemoryfile -> map-output of size: 8555206, inMemoryMapOutputs.size() -> 7, commitMemory -> 4715767
5, usedMemory -> 55712881
2021-04-27 09:51:48,568 [EventFetcher for fetching Map Completion Events] INFO org.apache.hadoop.mapreduce.task.reduce.EventFetcher - EventFetcher is interrupted.. Returning
2021-04-27 09:51:48,562 [pool-4-thread-1] INFO org.apache.hadoop.mapreduce.TaskAttemptLocalJobRunner - 7 copied
2021-04-27 09:51:48,562 [pool-4-thread-1] INFO org.apache.hadoop.mapreduce.TaskAttemptLocalJobRunner - finalMerge called with 7 in-memory map-outputs and 0 on-disk map-outputs
2021-04-27 09:51:48,571 [pool-4-thread-1] INFO org.apache.hadoop.mapred.Merger - Herging 7 sorted segments
2021-04-27 09:51:48,571 [pool-4-thread-1] INFO org.apache.hadoop.mapred.Merger - Down to the last merge-pass, with 7 segments left of total size: 55712811 bytes
2021-04-27 09:51:49,444 [pool-4-thread-1] INFO org.apache.hadoop.mapreduce.task.reduce.MergeManagerImpl - Merged 7 segments, 55712881 bytes to disk to satisfy reduce memory limit
2021-04-27 09:51:49,445 [pool-4-thread-1] INFO org.apache.hadoop.mapreduce.task.reduce.MergeManagerImpl - Merging 1 files, 55712873 bytes from disk
2021-04-27 09:51:49,446 [pool-4-thread-1] INFO org.apache.hadoop.mapreduce.TaskAttemptLocalJobRunner - Reducers, MergeManagerImpl - Merging 0 segments, 0 bytes from memory into reduce
2021-04-27 09:51:49,451 [pool-4-thread-1] INFO org.apache.hadoop.mapred.Merger - Herging 1 sorted segments
2021-04-27 09:51:49,451 [pool-4-thread-1] INFO org.apache.hadoop.mapred.Merger - Down to the last merge-pass, with 1 segments left of total size: 55712859 bytes
2021-04-27 09:51:49,452 [pool-4-thread-1] INFO org.apache.hadoop.mapred.LocalJobRunner - 7 / 7 copied.
2021-04-27 09:51:49,496 [pool-4-thread-1] INFO org.apache.hadoop.mapred.Merger - lib.output.FileOutputCommitter - File Output Committer Algorithm version is 2
2021-04-27 09:51:49,496 [pool-4-thread-1] INFO org.apache.hadoop.mapred.lib.output.FileOutputCommitter - lib.output.FileOutputCommitter skipp cleanup _temporary folders under output directory:false, ignore cleanup fail
res: False
2021-04-27 09:51:49,539 [pool-4-thread-1] INFO org.apache.hadoop.conf.Configuration.deprecation - napred.skip.on is deprecated. Instead, use mapreduce.job.skipped
2021-04-27 09:51:49,539 [pool-4-thread-1] INFO org.apache.ptg.impl.util.SplittableMemoryManager - Selected heap (PS Old Gen) of size 699400192 to monitor. collectionUsageThreshold = 489580128, usageThreshold = 4
89580128
2021-04-27 09:51:49,544 [pool-4-thread-1] WARN org.apache.ptg.data.SchemaTupleBackend - SchemaTupleBackend has already been initialized
2021-04-27 09:51:49,575 [pool-4-thread-1] INFO org.apache.ptg.backend.hadoop.executionengine.mapReduceLayer.PigMapReduceSReduce - Aliases being processed per job phase (AliasName[line,offset]): M: infiles[1,10],inffiles[-1,-1],ffiles[1,8],grpn_for_catagories[4,22] C: R: top10_viewed_catagories[5,26],sorted[6,36],top10_viewed_by_catagories[10,29]
2021-04-27 09:51:49,575 [pool-4-thread-1] INFO org.apache.hadoop.mapred.Task - Taskattempt_local423179355_0001_r_000000_0 is done. And it is in the process of committing
2021-04-27 09:51:49,577 [pool-4-thread-1] INFO org.apache.hadoop.mapred.LocalJobRunner - 7 / 7 copied.
2021-04-27 09:51:49,586 [pool-4-thread-1] INFO org.apache.hadoop.mapred.Task - task attempt_local423179355_0001_r_000000_0 is allowed to commit now
2021-04-27 09:51:49,586 [pool-4-thread-1] INFO org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter - Saved output of task 'attempt_local423179355_0001_r_000000_0' to file:/usr/local/bin/pig-0.17.0/bln/To
piViewedByCatagories.txt
```

```
sagarshah95@ubuntu: /usr/local/bin/pig-0.17.0/bin
File Edit View Search Terminal Help
IO_ERROR=0
WRONG_LENGTH=0
WRONG_MAP=0
WRONG_REDUCE=0
File Output Format Counters
Bytes Written=0
2021-04-27 09:51:57,208 [pool-4-thread-1] INFO org.apache.hadoop.mapred.LocalJobRunner - Finishing task: attempt_local423179355_0001_r_000000_0
2021-04-27 09:51:57,212 [Thread-0] INFO org.apache.hadoop.mapred.LocalJobRunner - reduce task executor complete.
2021-04-27 09:51:57,212 [Thread-0] INFO org.apache.hadoop.mapred.TaskAttemptLocalJobRunner - JobTracker metrics system already initialized!
2021-04-27 09:51:57,462 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:51:57,470 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.maps task is deprecated. Instead, use mapreduce.job.maps
2021-04-27 09:51:57,470 [main] INFO org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:51:57,674 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 100% complete
2021-04-27 09:51:57,685 [main] INFO org.apache.pig.tools.pigstats.mapreduce.SimplePigStats - Script Statistics:
HadoopVersion PigVersion UserId StartedAt FinishedAt Features
3.3.0 0.17.0 sagarshah95 2021-04-27 09:51:27 2021-04-27 09:51:57 GROUP_BY,FILTER
Success!
Job Stats (time in seconds):
JobId Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime MedianReduceTime Alias Feature Outputs
job_local423179355_0001 7 1 n/a n/a n/a n/a n/a n/a n/a files,grpn_for_catagories,infiles,sorted,top10_viewed_by_catagories,top10_viewed_catagories GROUP_BY
file:///usr/local/bin/pig-0.17.0/bin/Top10ViewedByCatagories.txt,
Inputs:
Successfully read 709739 records from: "/home/sagarshah95/Desktop/BigDataProject_pigAnalysis/data/youtubeDataset.csv"
Outputs:
Successfully stored 144 records in: "file:///usr/local/bin/pig-0.17.0/bin/Top10ViewedByCatagories.txt"
Counters:
Total records written : 144
Total bytes written : 0
Splittable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0
Job DAG:
job_local423179355_0001
2021-04-27 09:51:57,696 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:51:57,701 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:51:57,705 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2021-04-27 09:51:57,740 [main] WARN org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Encountered Warning ACCESSING_NON_EXISTENT_FIELD 17591 time(s).
2021-04-27 09:51:57,746 [main] WARN org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Encountered Warning FIELD_DISCARDED_TYPE_CONVERSION_FAILED 1826 time(s).
2021-04-27 09:51:57,746 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2021-04-27 09:51:57,848 [main] INFO org.apache.pig.Math - Pig script completed in 36 Seconds and 511 Milliseconds (30311 ms)
sagarshah95@ubuntu: /usr/local/bin/pig-0.17.0/bin
```

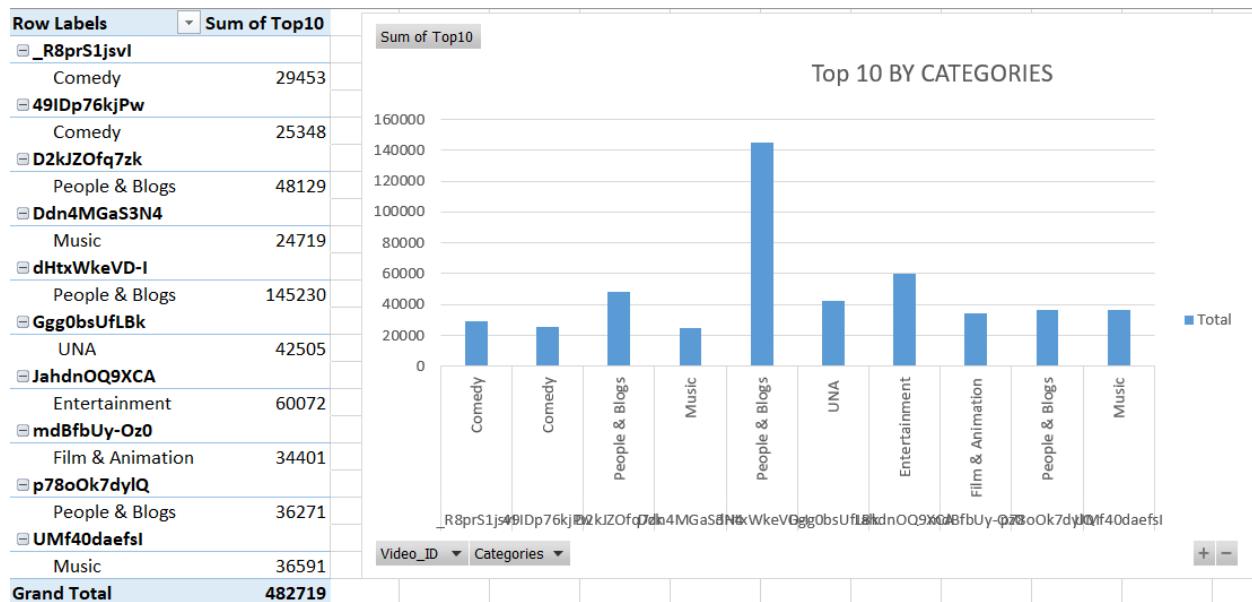


part-r-00000		
/usr/local/bin/pig-0.17.0/bin/Top10ViewedByCategories.txt		
aRNzWyD7C9o	UNA	8825788
PxNNR4symUE	UNA	4033376
LIhbap3FlGc	UNA	2849832
LIhbap3FlGc	UNA	2849832
lCSTULqmmYE	UNA	2179562
y6oXEWowirI	UNA	1666084
ByM_K-xkcag	UNA	1366452
783ynnbiFrI	UNA	1251129
R0049_tDAU8	UNA	1204982
R0049_tDAU8	UNA	1204982
QjA5faZF1A8	Music	15256922
tYnn51C3X_w	Music	11823701
pV5zWaTEVkI	Music	11672017
8bbTtPL1jRs	Music	9579911
UMf40daefsI	Music	7533070
-xEzGIuY7kw	Music	6946033
d6C0bNDqf3Y	Music	6935578
HSoVKUVOnfQ	Music	6193057
HSoVKUVOnfQ	Music	6193057
3URfWTEPmtE	Music	5581171
dMH0bHeiRNg	Comedy	42513417
0XXI-hvPRRA	Comedy	20282464
49IDp76kjPw	Comedy	11970018
SP6UU6m3cqk	Comedy	10107491
_BuRwH59oAo	Comedy	9566609
MNxwAU_xAMk	Comedy	7066676
pYak2F1hUYA	Comedy	6322117
h0zAlXr1UOs	Comedy	5826923
h0zAlXr1UOs	Comedy	5826923
C8rjr4jmWd0	Comedy	5587299
Ugrlzm7fySE	Sports	2867888
q8t7iSGAKik	Sports	2735003
q8t7iSGAKik	Sports	2735003
7vl19n8vl54	Sports	2527713

Open +/- part-r-00000 /usr/local/bin/pig-0.17.0/bin/Top10ViewedByCategories.txt Save ☰ X

xDh_pvV1tUM	News & Politics	2335060
p_YMigZmUuk	News & Politics	2326680
QCVxQ_3EjkG	News & Politics	2318782
a9WB_PXjTB0	News & Politics	2310583
bNF_P281Uu4	Travel & Places	5231539
s5ipz_0uC_U	Travel & Places	1198840
6jjW7aSNCzU	Travel & Places	1143287
dVRUBiYRAYk	Travel & Places	1000309
lqbt6X4ZgEI	Travel & Places	921593
RIH1I1doUI4	Travel & Places	879577
ALPql7IUT6M	Travel & Places	845180
ALPql7IUT6M	Travel & Places	845180
_5QUdvUhCzC	Travel & Places	819974
m9A_vxIOB-I	Travel & Places	677876
RjrEQaG5jPM	Autos & Vehicles	2803140
cv157ZIInUk	Autos & Vehicles	2773979
Gyg9U1YaVk8	Autos & Vehicles	1832224
6GNB7xT3rNE	Autos & Vehicles	1412497
tth9krDtxII	Autos & Vehicles	1347317
46LQd9dXFru	Autos & Vehicles	1262173
46LQd9dXFru	Autos & Vehicles	1262173
pdiuDXwgrjQ	Autos & Vehicles	1013697
kY_cDpENQLE	Autos & Vehicles	956665
YtxfbxGz1u4	Autos & Vehicles	942604
sdUUx5FdySs	Film & Animation	5840839
6B26asyGKDo	Film & Animation	5147533
H20dhY01Xjk	Film & Animation	3772116
55YYaJIRmzo	Film & Animation	3356163
55YYaJIRmzo	Film & Animation	3356163
JzqumbhfxRo	Film & Animation	3230774
eAhfZUZiwSE	Film & Animation	3114215
eAhfZUZiwSE	Film & Animation	3114215
h7svw0m-w00	Film & Animation	2866490
h7svw0m-w00	Film & Animation	2866490

Pivot Table



Hive Analysis

Create Schema/Table

```
hive> create table YouTube_data_table (vedioID STRING,uploader STRING, age INT, category STRING, length INT, noofviews INT, no_of_comments INT, IDs INT)
   > ROW FORMAT DELIMITED
   > FIELDS TERMINATED BY ','
   > STORED AS TEXTFILE;
OK
Time taken: 1.87 seconds
hive> show tables
   > ;
OK
youtube_data_table
Time taken: 0.241 seconds, Fetched: 1 row(s)
hive> set hive.cli.print.header=true;
```

Load Data into Schema

```
hive> LOAD DATA LOCAL INPATH '/home/sagarshah95/Desktop/BigDataProject_pigAnalysis/data/youtubedata.csv' OVERWRITE INTO TABLE YouTube_data_table;
Loading data to table default.youtube_data_table
OK
Time taken: 2.121 seconds
```

Execution

```

File Edit View Search Terminal Help
... 43 more
Job Submission failed with exception 'org.apache.hadoop.security.AccessControlException(Permission denied: user=root, access=EXECUTE, inode="/tmp/hadoop-yarn":sagarshah95:supergroup:drwx-----'
at org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.check(FSPermissionChecker.java:496)
at org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.check(FSPermissionChecker.java:412)
at org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.checkPermission(FSPermissionChecker.java:323)
at org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.checkPermissionWithContext(FSPermissionChecker.java:360)
at org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.checkPermission(FSPermissionChecker.java:239)
at org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.checkTraverse(FSPermissionChecker.java:703)
at org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.checkTraverse(FSPermissionChecker.java:104)
at org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.checkTraverse(FSPermissionChecker.java:1876)
at org.apache.hadoop.hdfs.server.namenode.FSDirectory.checkTraverse(FSDirectory.java:718)
at org.apache.hadoop.hdfs.server.namenode.FSDirectory.resolvePath(FSDirectory.java:718)
at org.apache.hadoop.hdfs.server.namenode.FSDirectory.listingOp.getFileInfo(FSDirStatAndListingOp.java:112)
at org.apache.hadoop.hdfs.server.namenode.FSNamesystem.getFileInfo(FSNamesystem.java:332)
at org.apache.hadoop.hdfs.server.namenode.FSNamesystem.listingOp.getFileInfo(FSNamesystem.java:1210)
at org.apache.hadoop.hdfs.protocol.ClientNamemodeProtocolServerSideTranslatorPB.callBlockingMethod(ClientNamemodeProtocolServerSideTranslatorPB.java:1041)
at org.apache.hadoop.hdfs.protocol.ClientNamemodeProtocolProtos$ClientNamemodeProtocol$2.callBlockingMethod(ClientNamemodeProtocolProtos.java)
at org.apache.hadoop.ipc.ProtobufRpcEngine$Server$ProtoBufRpcInvoker.call(ProtobufRpcEngine.java:532)
at org.apache.hadoop.ipc.RPC$Server.call(RPC.java:1070)
at org.apache.hadoop.ipc.Server$RPC$Call.run(Server.java:1020)
at org.apache.hadoop.ipc.Server$RPC$Call.run(Server.java:948)
at java.security.AccessController.doAsPrivileged(Native Method)
at javax.security.auth.Subject.doAs(Subject.java:422)
at org.apache.hadoop.security.UserGroupInformation.doAs(UserGroupInformation.java:1845)
at org.apache.hadoop.ipc.Server$Handler.run(Server.java:2952)
)'

FAILED: Execution Error, return code 1 from org.apache.hadoop.hive.ql.exec.mr.MappRedTask. Permission denied: user=root, access=EXECUTE, inode="/tmp/hadoop-yarn":sagarshah95:supergroup:drwx----- at org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.check(FSPermissionChecker.java:496)
at org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.checkTraverse(FSPermissionChecker.java:412)
at org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.check(FSPermissionChecker.java:412)
at org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.checkPermission(FSPermissionChecker.java:323)
at org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.checkPermissionWithContext(FSPermissionChecker.java:360)
at org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.checkPermission(FSPermissionChecker.java:239)
at org.apache.hadoop.hdfs.server.namenode.FSPermissionChecker.checkTraverse(FSPermissionChecker.java:703)
at org.apache.hadoop.hdfs.server.namenode.FSDirectory.checkTraverse(FSDirectory.java:1858)
at org.apache.hadoop.hdfs.server.namenode.FSDirectory.listingOp.getFileInfo(FSDirStatAndListingOp.java:112)
at org.apache.hadoop.hdfs.server.namenode.FSNamesystem.getFileInfo(FSNamesystem.java:332)
at org.apache.hadoop.hdfs.server.namenode.FSNamesystem.listingOp.getFileInfo(FSNamesystem.java:1210)
at org.apache.hadoop.hdfs.protocol.ClientNamemodeProtocolServerSideTranslatorPB.callBlockingMethod(ClientNamemodeProtocolServerSideTranslatorPB.java:1041)
at org.apache.hadoop.hdfs.protocol.ClientNamemodeProtocolProtos$ClientNamemodeProtocol$2.callBlockingMethod(ClientNamemodeProtocolProtos.java)
at org.apache.hadoop.ipc.ProtobufRpcEngine$Server$ProtoBufRpcInvoker.call(ProtobufRpcEngine.java:532)
at org.apache.hadoop.ipc.RPC$Server.call(RPC.java:1070)
at org.apache.hadoop.ipc.Server$RPC$Call.run(Server.java:1020)
at org.apache.hadoop.ipc.Server$RPC$Call.run(Server.java:948)
at java.security.AccessController.doAsPrivileged(Native Method)
at javax.security.auth.Subject.doAs(Subject.java:422)
at org.apache.hadoop.security.UserGroupInformation.doAs(UserGroupInformation.java:1845)
at org.apache.hadoop.ipc.Server$Handler.run(Server.java:2952)

hive> ■

```

```

:sagarshah95@ubuntu:/usr/local/bin/apache-hive-3.1.2-bin/bin$ sudo ./hive -f ~/Desktop/hive_analysis.hql
hive Session ID = 45689e3e-ebd7-4633-ac41-cf5d531d90f8

Logging initialized using configuration in jar:file:/usr/local/bin/apache-hive-3.1.2-bin/lib/hive-common-3.1.2.jar!/hive-log4j2.properties Async: true
hive Session ID = 1b03a501-0aac-498d-916c-6d3c5421f57f
)K
Time taken: 4.561 seconds
)K
Time taken: 1.233 seconds
>Loading data to table default.youtube_datatable
)K
Time taken: 1.18 seconds
:sagarshah95@ubuntu:/usr/local/bin/apache-hive-3.1.2-bin/bin$ ■

```

Queries

Calculate top 10 channels with maximum number of likes

```

select vedioid, uploader, no_of_comments FROM
YouTube_DataTable ORDER BY no_of_comments DESC LIMIT 10;

```

Calculate top 5 categories with maximum number of comments

```

select category, max(no_of_comments) as max_no_of_comments from YouTube_DataTable e GROUP
ORDER BY max_no_of_comments DESC LIMIT 5;

```

Appendix Section

The code snippets are provided below each analysis along with description of each analysis and its execution