

Project Report-- ##Project 3: Data Analysis using Map/Reduce + Query Processing##
Spring 2023

Division of Labor:

Team 31

- [Sakshi](#)(sx3702@mavs.uta.edu, 1001993702) (installation of Hadoop, MapReduce function, query plan) ~20hours+
- [Spoorthi Seri](#)(sxs2684@mavs.uta.edu,1002032680) (documentation, UTM-virtual machine)
~10hours
- [Neha Thokala](#)(nxt0631@mavs.uta.edu, 1002030631) (trend analysis, oracle run)
~10hours

Overall Status: Successfully completed both the tasks and documented the results.

- **[Task 1]** Completed the pre-requisite steps to setup a development environment by installing Hadoop on an Ubuntu OS
 - Used UTM virtual box on MacBook.
- **[Task 1]** Implemented a Map/Reduce program to count number of movies with different genre combinations in JAVA by following the steps provided.
 - Initially ran the basic Wordcount program – created, compiled jar (javac *.java file)
 - Incrementally developed to generate the required results (part-r-00000 file).
 - After final changes, generated a CSV file/graph from the output and analyzed the counts to determine trends.
 - After completion, verified the counts matched with the oracle table as well.
- **[Task 2]** Completed the pre-requisite to access oracle database for task 2.
 - Followed the steps provided on how to access oracle on Omega
- **[Task 2]** Created an SQL query to get top 5 movies based on movie rating with the expected criteria (at least 100k votes within a given date range and genre combination)
 - Ran an explain statement of the query to understand how the query performed.

Validated the results by filing simple select query on the title_basics table on the imdb00 database on Omega with the correct where clauses for titletype, startyear and genres;

```
SELECT count(*) FROM imdb00.TITLE_BASICS tb WHERE tb.STARTYEAR >= '2000' AND tb.STARTYEAR <= '2006' AND tb.GENRES LIKE '%Action%' AND tb.GENRES LIKE '%Drama%' AND tb.titletype = 'movie'
```

Analysis Results

[Task 1] Analysis of genre combinations with given date ranges

Task1 (mapReduce)		Output File: part-r-00000	
Year Range	Genre	Number of Movies	
2000-2006	Action/Drama	1240	
2000-2006	Adventure/Sci-Fi	67	
2000-2006	Comedy/Romance	1620	
2007-2013	Action/Drama	1811	
2007-2013	Adventure/Sci-Fi	153	
2007-2013	Comedy/Romance	2402	
2014-2020	Action/Drama	2326	
2014-2020	Adventure/Sci-Fi	247	
2014-2020	Comedy/Romance	2784	

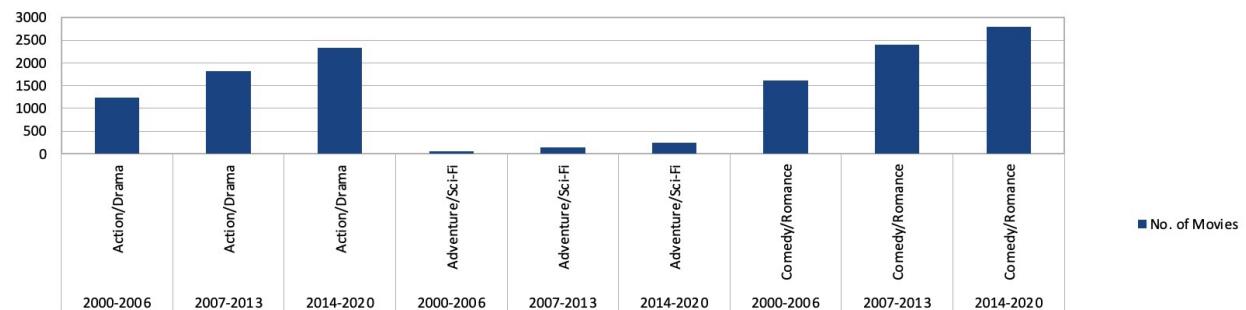
Scenario1:

Data Range	Genre Combinations	No. of Movies
2000-2006	Action/Drama	1240
2007-2013	Action/Drama	1811
2014-2020	Action/Drama	2326
2000-2006	Adventure/Sci-Fi	67
2007-2013	Adventure/Sci-Fi	153
2014-2020	Adventure/Sci-Fi	247
2000-2006	Comedy/Romance	1620
2007-2013	Comedy/Romance	2402
2014-2020	Comedy/Romance	2784

Inference of Scenario1:

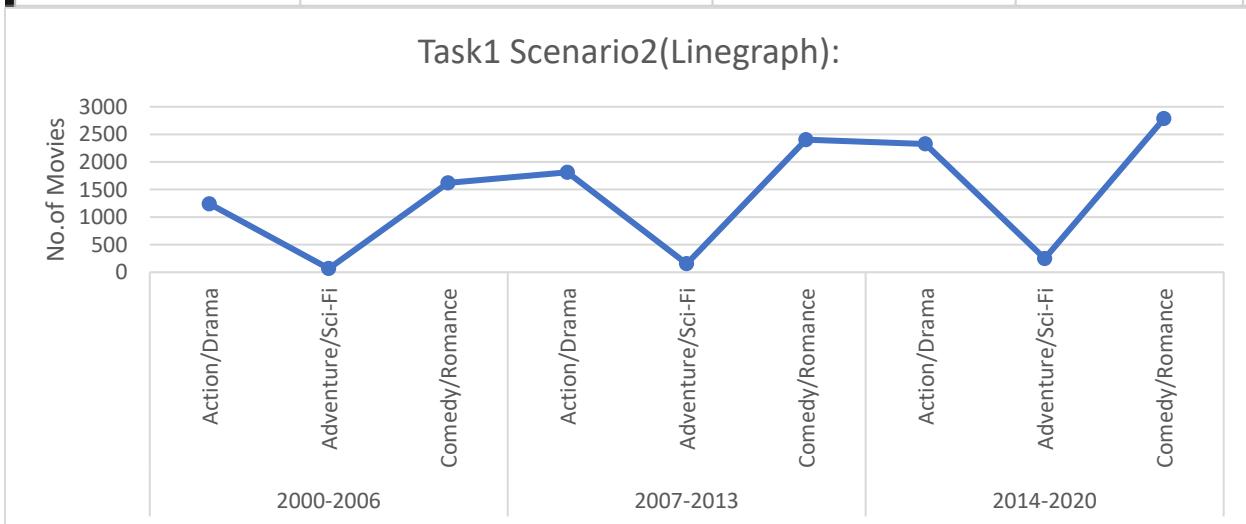
- All three genres (Action/Drama, Adventure/Sci-Fi, and Comedy/Romance) show an increasing trend in the number.
- Comedy/Romance is the most popular genre combination, with consistently high counts in all three time periods. It also shows the strongest increasing trend over time, with a total count of 2784 movies in the 2014-2020 period.
- Action/Drama and Adventure/Sci-Fi are less popular than Comedy/Romance, but still show an increasing trend in the number of movies produced over time.
- Adventure/Sci-Fi is the least popular genre combination, with the lowest counts in all three time periods. It shows a relatively flat trend over time, with a total count of only 467 movies across all three periods.
- Overall, there is an upward trend in the number of movies produced across all genres and time periods

Task1:Scenario1(Histogram)



Scenario2:		
Data Range	Genre Combinations	No. of Movies
2000-2006	Action/Drama	1240
	Adventure/Sci-Fi	67
	Comedy/Romance	1620
2007-2013	Action/Drama	1811
	Adventure/Sci-Fi	153
	Comedy/Romance	2402
2014-2020	Action/Drama	2326
	Adventure/Sci-Fi	247
	Comedy/Romance	2784

Inference of Scenario2:
-In all three time periods, Comedy/Romance is the most popular genre combination with the highest count of movies produced. It is followed by Action/Drama, which has the second-highest count in all three time periods.
-Adventure/Sci-Fi is the least popular genre combination in all three time periods, with the lowest count of movies produced.
-Over time, the count of movies produced in all three genre combinations has increased. However, the rate of increase is not the same for all three genres. Comedy/Romance shows the highest rate of increase, followed by Action/Drama and then Adventure/Sci-Fi.
-The largest increase in the number of movies produced was seen in the Comedy/Romance genre combination, which saw a 71% increase in the number of movies produced between the 2000-2006 and 2014-2020 time periods.
-Adventure/Sci-Fi saw the smallest increase in the number of movies produced over time, with only a 268% increase between 2000-2006 and 2014-2020. This is likely due to the relatively low popularity of this genre combination compared to the other two genres.



[Task 2] Analysis of query explanation

```
4331-5331_Proj3Spring23_team_31.sql ×

Users > sakhisrivastava > Desktop > uta > Sem3Spring > DBMS > project3 > 433
1 -- Query 1: Task2 Get top five movie from 2000<=startyear<=2006
2 set echo on
3 spool 4331-5331_Proj3Spring23_team_31.txt
4 SELECT
5   tb.PRIMARYTITLE,
6   tr.AVERAGEGATING,
7   nb.PRIMARYNAME
8   FROM
9   imdb00.TITLE_BASICS tb
10  JOIN imdb00.TITLE_RATINGS tr ON tb.TCONST = tr.TCONST
11  JOIN imdb00.TITLE_PRINCIPALS tp ON tb.TCONST = tp.TCONST
12  AND tp.ORDERING = 1
13  JOIN imdb00.NAME_BASICS nb ON tp.NCONST = nb.NCONST
14 WHERE
15   tb.STARTYEAR >= 2000
16   AND tb.STARTYEAR <= 2006
17   AND tb.GENRES LIKE '%Action%'
18   AND tb.GENRES LIKE '%Drama%'
19   AND tr.NUMVOTES >= 100000
20   AND tb.titletype = 'movie'
21 ORDER BY
22   tr.AVERAGEGATING DESC
23 FETCH FIRST
24   5 ROWS ONLY;
```

```
4331-5331_Proj3Spring23_team_31.sql X
Users > sakshisrivastava > Desktop > uta > Sem3Spring > DBMS > project3 > 4
25 | EXPLAIN PLAN
26 | SET
27 |   statement_id = 'Query_1' FOR
28 | SELECT
29 |   tb.PRIMARYTITLE,
30 |   tr.AVERAGERATING,
31 |   nb.PRIMARYNAME
32 | FROM
33 |   imdb00.TITLE_BASICS tb
34 |   JOIN imdb00.TITLE_RATINGS tr ON tb.TCONST = tr.TCONST
35 |   JOIN imdb00.TITLE_PRINCIPALS tp ON tb.TCONST = tp.TCONST
36 |   AND tp.ORDERING = 1
37 |   JOIN imdb00.NAME_BASICS nb ON tp.NCONST = nb.NCONST
38 | WHERE
39 |   tb.STARTYEAR >= 2000
40 |   AND tb.STARTYEAR <= 2006
41 |   AND tb.GENRES LIKE '%Action%'
42 |   AND tb.GENRES LIKE '%Drama%'
43 |   AND tr.NUMVOTES >= 100000
44 |   AND tb.titletype = 'movie'
45 | ORDER BY
46 |   tr.AVERAGERATING DESC
47 |   FETCH FIRST
48 |   5 ROWS ONLY;
49 |
50 |
51 | SELECT
52 |   PLAN_TABLE_OUTPUT
53 | FROM
54 |   TABLE(DBMS_XPLAN.DISPLAY(NULL, 'Query_1', 'BASIC'));
55 |
56 | SPOOL OFF
57 | set echo off
```

[Task 2] Oracle run screenshot.

```
Enter user-name: sx3702
Enter password:

Connected to:
Oracle Database 19c Enterprise Edition Release 19.0.0.0.0 - Production

SQL> start execute.sql
SQL> spool 4331-5331_Proj3Spring23_team_31.txt
SQL> SELECT
  2      tb.PRIMARYTITLE,
  3      tr.AVERAGERATING,
  4      nb.PRIMARYNAME
  5  FROM
  6      imdb00.TITLE_BASICS tb
  7      JOIN imdb00.TITLE_RATINGS tr ON tb.TCONST = tr.TCONST
  8      JOIN imdb00.TITLE_PRINCIPALS tp ON tb.TCONST = tp.TCONST
  9      AND tp.ORDERING = 1
10      JOIN imdb00.NAME_BASICS nb ON tp.NCONST = nb.NCONST
11 WHERE
12      tb.STARTYEAR >= 2000
13      AND tb.STARTYEAR <= 2006
14      AND tb.GENRES LIKE '%Action%'
15      AND tb.GENRES LIKE '%Drama%'
16      AND tr.NUMVOTES >= 100000
17      AND tb.titletype = 'movie'
18 ORDER BY
19      tr.AVERAGERATING DESC
20 FETCH FIRST
21      5 ROWS ONLY;
```

```
PRIMARYTITLE
-----
AVERAGERATING
-----
PRIMARYNAME
-----
The Lord of the Rings: The Return of the King
      9
Elijah Wood

The Lord of the Rings: The Fellowship of the Ring
      8.8
Elijah Wood

PRIMARYTITLE
-----
AVERAGERATING
-----
PRIMARYNAME
-----
The Lord of the Rings: The Two Towers
      8.8
Elijah Wood

Gladiator
      8.5
```

```

PRIMARYTITLE
-----
AVERAGERATING
-----
PRIMARYNAME
-----
Russell Crowe

Oldboy          8.4
Choi Min-sik

SQL> EXPLAIN PLAN
 2  SET
 3      statement_id = 'Query_1' FOR
 4  SELECT
 5      tb.PRIMARYTITLE,
 6      tr.AVERAGERATING,
 7      nb.PRIMARYNAME
 8  FROM
 9      imdb00.TITLE_BASICS tb
10     JOIN imdb00.TITLE_RATINGS tr ON tb.TCONST = tr.TCONST
11     JOIN imdb00.TITLE_PRINCIPALS tp ON tb.TCONST = tp.TCONST
12     AND tp.ORDERING = 1
13     JOIN imdb00.NAME_BASICS nb ON tp.NCONST = nb.NCONST
14 WHERE
15     tb.STARTYEAR >= 2000
16     AND tb.STARTYEAR <= 2006
17     AND tb.GENRES LIKE '%Action%'
18     AND tb.GENRES LIKE '%Drama%'
19     AND tr.NUMVOTES >= 100000
20     AND tb.titletype = 'movie'
21 ORDER BY
22     tr.AVERAGERATING DESC
23 FETCH FIRST
24     5 ROWS ONLY;

Explained.

SQL>
SQL> SELECT
 2      PLAN_TABLE_OUTPUT
 3  FROM
 4      TABLE(DBMS_XPLAN.DISPLAY(NULL, 'Query_1', 'BASIC'));

```

```

PLAN_TABLE_OUTPUT
-----
Plan hash value: 2239367442

| Id  | Operation          | Name   |
|---|---|---|
| 0  | SELECT STATEMENT  |        |
| 1  |   VIEW             |        |
| 2  |   WINDOW SORT PUSHED RANK |
| 3  |   NESTED LOOPS    |        |
| 4  |     NESTED LOOPS  |        |
| 5  |       HASH JOIN    |        |

PLAN_TABLE_OUTPUT
-----
| 6  | NESTED LOOPS      |        |
| 7  |   NESTED LOOPS    |        |
| 8  |     TABLE ACCESS FULL | TITLE_RATINGS |
| 9  |     INDEX UNIQUE SCAN | SYS_C00547784 |
| 10 |     TABLE ACCESS BY INDEX ROWID | TITLE_BASICS |
| 11 |     TABLE ACCESS FULL | TITLE_PRINCIPALS |
| 12 |     INDEX UNIQUE SCAN | SYS_C00547785 |
| 13 |     TABLE ACCESS BY INDEX ROWID | NAME_BASICS |

-----
```

20 rows selected.

```

SQL>
SQL> SPOOL OFF
SQL> set echo off
```

Task2: SQL Queries Get top five movie from 2000<=startyear<=2006 GENRES='Action,Drama'

```

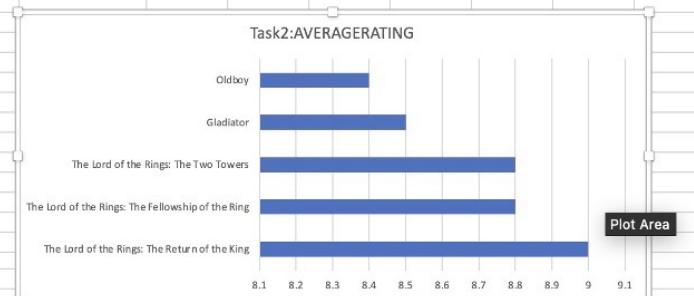
SELECT
tb.PRIMARYTITLE,
tr.AVERAGERATING,
nb.PRIMARYNAME
FROM
imdb00.TITLE_BASICS tb
JOIN imdb00.TITLE_RATINGS tr ON tb.TCONST = tr.TCONST
JOIN imdb00.TITLE_PRINCIPALS tp ON tb.TCONST = tp.TCONST
AND tp.ORDERING = 1
JOIN imdb00.NAME_BASICS nb ON tp.NCONST = nb.NCONST
WHERE
tb.STARTYEAR >= 2000
AND tb.STARTYEAR <= 2006
AND tb.GENRES LIKE '%Action%'
AND tb.GENRES LIKE '%Drama%'
AND tr.NUMVOTES >= 100000
AND tb.titletype = 'movie'
ORDER BY
tr.AVERAGERATING DESC
FETCH FIRST
5 ROWS ONLY;
```

OutputReceived:

PRIMARYTITLE	AVERAGERATING	PRIMARYNAME
The Lord of the Rings: The Return of the King	9	Elijah Wood
The Lord of the Rings: The Fellowship of the Ring	8.8	Elijah Wood
The Lord of the Rings: The Two Towers	8.8	Elijah Wood
Gladiator	8.5	Russell Crowe
Oldboy	8.4	Choi Min-sik

Inference of Task2: SQL Queries:

- 1) The top 3 movies are all part of "The Lord of the Rings" trilogy, with "The Return of the King" being the highest-rated movie with an average rating of 9.0.
- 2) The other two movies in the top 5 are "Gladiator" with an average rating of 8.5 and "Oldboy" with an average rating of 8.4.
- 3) Elijah Wood is the first-billed cast member in all three "The Lord of the Rings" movies, and Russell Crowe and Choi Min-sik are the first-billed cast members in "Gladiator" and "Oldboy", respectively.



File descriptions:

[Task 1] Map/Reduce mapper step which receives an input movie row and map's it into a key with date range as the prefix and genre combination as the suffix and a value of 1.

```
J IMDbData.java X
Users > sakshisrivastava > Desktop > uta > Sem3Spring > DBMS > project3 > J IMDbData.java
 4  import java.util.List;
 5  import org.apache.hadoop.conf.Configuration;
 6  import org.apache.hadoop.fs.Path;
 7  import org.apache.hadoop.io.IntWritable;
 8  import org.apache.hadoop.io.Text;
 9  import org.apache.hadoop.mapreduce.Job;
10  import org.apache.hadoop.mapreduce.Mapper;
11  import org.apache.hadoop.mapreduce.Reducer;
12  import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
13  import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
14  import org.apache.hadoop.mapreduce.lib.input.FileSplit;
15
16 public class IMDbData {
17     public static class TokenizerMapper extends Mapper<Object, Text, Text, IntWritable> {
18         // Define the date ranges and genre categories
19         int[][] dateRanges = { { 2000, 2006 }, { 2007, 2013 }, { 2014, 2020 } };
20         String[][] genreCategories = { { "Comedy", "Romance" }, { "Action", "Drama" }, { "Adventure", "Sci-Fi" } };
21         private final static IntWritable one = new IntWritable(1);
22         private Text word = new Text();
23         // mapper is transformer
24         // object key associated with this mapper execution is unique,value=sentence
25         // context knows which file it is and which split it is
26         // mapper emits all the occurrence of the word inside one sentence
27         public void map(Object key, Text value, Context context) throws IOException, InterruptedException {
28             FileSplit fileSplit = (FileSplit) context.getInputSplit();
29             // we can use log that will tell which file name is running inside this map
30             String filename = fileSplit.getPath().getName();
31             // StringTokenizer itr = new StringTokenizer(value.toString());
32             // Split the sentence into individual parts
33             String[] parts = value.toString().split(";");
34             // Errro solution: Null /N is coming
35             if (parts.length > 3 && parts[3].matches("-?\\d+(\\.\\d+)?"))
36                 && parts[1].equals("movie") ) {
37                 // Extract the year value from the sentence
38                 int year = Integer.parseInt(parts[3]);
39                 // Extract the genres from the sentence
40                 List<String> genres = Arrays.asList(parts[4].split(","));

```

```
J IMDbData.java X
Users > sakshisrivastava > Desktop > uta > Sem3Spring > DBMS > project3 > J IMDbData.java
39      // Extract the genres from the sentence
40      List<String> genres = Arrays.asList(parts[4].split(","));
41      // Find the date range that the year falls in
42      String dateRangeLabel = "Unknown";
43      for (int i = 0; i < dateRanges.length; i++) {
44          if (year >= dateRanges[i][0] && year <= dateRanges[i][1]) {
45              dateRangeLabel = dateRanges[i][0] + "-" + dateRanges[i][1];
46              break;
47          }
48      }
49      // Find the genre category that matches the genres in the sentence
50      String genreCategoryLabel = "Unknown";
51      for (int i = 0; i < genreCategories.length; i++) {
52          List<String> categoryGenres = Arrays.asList(genreCategories[i]);
53          if (genres.containsAll(categoryGenres)) {
54              genreCategoryLabel = String.join("/", categoryGenres);
55              break;
56          }
57      }
58      // Create the key with the specified format
59      if (!dateRangeLabel.equals("Unknown") && !genreCategoryLabel.equals("Unknown")) {
60          String key1 = dateRangeLabel + " | " + genreCategoryLabel;
61          System.out.println(key1); // Output: 2007-2013 | Drama/Action
62          // word is text type , key1 is in string format
63          word.set(key1);
64          context.write(word, one);
65      } else {
66          System.out.println("Unable to create key");
67      }
68  }
69 }
70 }
71 }
```

[Task 1] Map/Reduce reducer step which aggregates the occurrences of the keys generated in the mapper step.

```

J IMDbData.java ×
Users > sakshisrivastava > Desktop > uta > Sem3Spring > DBMS > project3 > J IMDbData.java
71 }
72 // reducer accepts the mapper output and aggregates them to a final result
73 public static class IntSumReducer
74     extends Reducer<Text, IntWritable, Text, IntWritable> {
75     private IntWritable result = new IntWritable();
76     // key=word
77     // values=list of occurrences reading all the occurancance and summing
78     public void reduce(Text key, Iterable<IntWritable> values,
79         Context context) throws IOException, InterruptedException {
80         int sum = 0;
81         for (IntWritable val : values) {
82             sum += val.get();
83         }
84         result.set(sum);
85         // key=word of text type
86         // result=Iterable
87         context.write(key, result);
88     }
89 }
90
91 public static void main(String[] args) throws Exception {
92     Configuration conf = new Configuration();
93     Job job = Job.getInstance(conf, "word count");
94     job.setJarByClass(IMDbData.class);
95     job.setMapperClass(TokenizerMapper.class);
96     job.setCombinerClass(IntSumReducer.class);
97     job.setReducerClass(IntSumReducer.class);
98     job.setOutputKeyClass(Text.class);
99     job.setOutputValueClass(IntWritable.class);
100    FileInputFormat.addInputPath(job, new Path(args[0]));
101    FileOutputFormat.setOutputPath(job, new Path(args[1]));
102    System.exit(job.waitForCompletion(true) ? 0 : 1);
103 }
104 }

```

[Task 1] Map/Reduce program summary

- IMDDBData.java is a program for processing IMDb movie data with the Apache Hadoop MapReduce framework.
- The program consists of a main class IMDDBData and two nested classes TokenizerMapper and IntSumReducer.
- The main function configures and runs the MapReduce job, serving as the entry point of the program.
- The TokenizerMapper class extends the Mapper class provided by Hadoop and overrides the map function to generate intermediate key-value pairs from input text data where the prefix is a date range and suffix is a combination of genres
 - **2000-2006/Action/Drama**
- Each line of text data is split into individual parts, the year and genres of the movie are extracted, and the movie is categorized based on these factors, with each category represented as a key-value pair.

- The IntSumReducer class extends the Reducer class provided by Hadoop and overrides the reduce function to aggregate intermediate key-value pairs and produce final output key-value pairs.
- The main function sets up the Hadoop job by configuring input and output paths, map and reduce classes, and output key and value classes.
- The job is executed using the job.waitForCompletion() method, which submits the job to the Hadoop cluster and waits for it to complete. -The program returns 0 if the job completes successfully and 1 if not.
- The user provides input and output paths as command-line arguments. -The Hadoop cluster runs the MapReduce job and produces the output.

Logical errors and how we handled them:

- Encountered new lines (\N) which caused array out of bound errors. **Solution:** Added below if condition to prevent processing any data with length less than 3 in IMDbData.java file.

```

39
40          // If Null is coming
41          if (parts.length > 3 && parts[3].matches("-?\\d+(\\.\\d+)?") .
42              && parts[1].equals("movie")) {|

```

- Initially the Map/Reduce program generated very high counts, we resolved the logical error by adding 'movie' in the JAVA filtering logic. After this change the counts matched with the SQL query counts.
- Giving wrong argument for the input data led to empty data problem even though Hadoop ran successfully without any errors (0%map and 100% reducer)
- Resolving JAR file build and compilation errors. We faced dependency and type mismatch errors when using java.util.* instead of the exact imports for java.util.Arrays and java.util.List

Installation errors and how we handled them:

- Resolved JAVA class path errors when setting up HADOOP by installing usr/local/bin directory and adding exports in ~/.bashrc file
- While installing Ubuntu faced issues as MacBook required java-11-openjdk-arm64 unlike all other systems which used amd64. Downloaded all dependency file with arm64.

```

sx3702@sakshiubuntu:~$ source ~/.bashrc
bash: export: `/usr/lib/jvm/java-11-openjdk-arm64': not a valid identifier
bash: export: `/usr/local/hadoop/': not a valid identifier
sx3702@ubuntu:~/Downloads$ hadoop-3.3.5/bin/hdfs namenode -format
ERROR: JAVA_HOME /usr/lib/jvm/java-8-openjdk-amd64 does not exist.

```

```

117 fi
118 export JAVA_HOME=/usr/lib/jvm/java-11-openjdk-arm64
119 export HADOOP_HOME=/usr/local/hadoop
120 export PATH=$PATH:$HADOOP_HOME/bin
121 export PATH=$PATH:$HADOOP_HOME/sbin
122 export HADOOP_MAPRED_HOME=$HADOOP_HOME
123 export HADOOP_COMMON_HOME=$HADOOP_HOME
124 ----- HADOOP_HOME $HADOOP_HOME

```

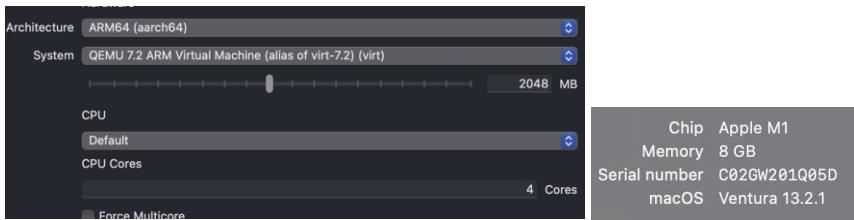
- Got permission denied in Ubuntu while doing pdsh, solution re-configure ubuntu with SSH external.

```

sx3702@ubuntu:/usr/local$ export PDSh_RCMD_TYPE=ssh
sx3702@ubuntu:/usr/local$ start-all.sh
WARNING: Attempting to start all Apache Hadoop daemons as sx3702 in 10 seconds.
WARNING: This is not a recommended production deployment configuration.
WARNING: Use CTRL-C to abort.
Starting namenodes on [localhost]
localhost: sx3702@localhost: Permission denied (publickey,password).
pdsh@ubuntu: localhost: ssh exited with exit code 255
Starting datanodes
localhost: sx3702@localhost: Permission denied (publickey,password).
pdsh@ubuntu: localhost: ssh exited with exit code 255
Starting secondary namenodes [ubuntu]
ubuntu: sx3702@ubuntu: Permission denied (publickey,password).
pdsh@ubuntu: ubuntu: ssh exited with exit code 255
2023-04-17 21:51:51,144 WARN util.NativeCodeLoader: Unable to load native-hadoop
library for your platform... using builtin-java classes where applicable
Starting resourcemanager
resourcemanager is running as process 7531. Stop it first and ensure /tmp/hadoop-
p-sx3702-resourcemanager.pid file is empty before retry.
Starting nodemanagers
localhost: sx3702@localhost: Permission denied (publickey,password).
pdsh@ubuntu: localhost: ssh exited with exit code 255

```

- Mac Book keeps on getting hang due to storage shortage. iCloud Drive is always full. After multiple attempts found balanced memory resource on UTM with 2048 MB and 4 core



Future Improvements

Here are some possible future improvements

1. Use of Combiner Class: Currently, the program uses the same reducer class for combiner as well. Using a separate combiner class can help in reducing the amount of data sent from mapper to reducer, thus improving the overall performance.
2. Better Error Handling: The current program only checks for a few specific conditions before emitting the output. However, there can be other cases that might result in an error. Therefore, adding more error handling code can improve the robustness of the program.

3. Use of Custom Data Types: The program uses Text and IntWritable as input and output types respectively. However, in some cases, using custom data types can provide better abstraction and make the code more readable and maintainable.
4. Use of Advanced Techniques: This program uses basic MapReduce techniques to analyze the data. However, there are other advanced techniques like Apache Spark, Apache Flink, or Hadoop Streaming, which can help in processing large datasets more efficiently.
5. More Complex Analysis: The current program only finds the count of movies based on their genre and release date. However, there are many other interesting insights that can be obtained from the IMDb dataset. For example, analyzing the ratings of movies based on their genre and release date can help in identifying popular genres or trends over time.

Screenshots

Mounted mac file directory to UTM /media/host folder

```
Remote disk share fromm mac to UTM -->
sudo apt install spice-vdagent spice-vdagent spice-webdavd -y
(make sure project3 folder is shared before starting UTM)
sudo mkdir /media/host
sudo mount -t 9p -o trans=virtio share /media/host -oversion=9p2000.L
```

Hadoop restart and connect with server.

```
hadoop restart----->
hadoop namenode -format
start-yarn.sh
start-dfs.sh
jps
export HADOOP_CLASSPATH=$(hadoop classpath)
hadoop fs -mkdir /Count
hadoop fs -mkdir /Count/Word
hadoop fs -mkdir /Count/Imdb
hadoop fs -mkdir /Count/Input
hadoop fs -put Desktop/Lab/Input/Spring2023-Project3-IMDbData.txt /Count/Input
copy imdb.java from hostshared to download folder

cd Desktop/Lab
cp ../../Downloads/IMDbData.java .
javac -classpath $HADOOP_CLASSPATH -d imdbcoundbms IMDbData.java
jar -cvf Imdb.jar -C imdbcoundbms .
hadoop jar Imdb.jar IMDbData /Count/Input /Count/Imdb/Output
hadoop dfs -cat /Count/Imdb/Output/*
```

The terminal window shows the following commands and their outputs:

```
sx3702@ubuntu:~$ ls
configuration-steps  Documents  Music      Public  Templates
Desktop             Downloads  Pictures   snap    Videos
sx3702@ubuntu:~$ sudo apt install spice-vdagent spice-vdagent spice-webdavd -y
[sudo] password for sx3702:
Reading package lists... Done
Building dependency tree... Done
Reading state information... Done
spice-vdagent is already the newest version (0.22.1-1).
spice-webdavd is already the newest version (2.5-1).
0 upgraded, 0 newly installed, 0 to remove and 4 not upgraded.
sx3702@ubuntu:~$ sudo mkdir /media/host
mkdir: cannot create directory '/media/host': File exists
sx3702@ubuntu:~$ sudo mount -t 9p -o trans=virtio share /media/host -oversion=9p
2000.L
```

Below the terminal, a file browser window is open. The left sidebar shows "Recent", "Starred", "Home", "Documents", and "Downloads". The main area, titled "On This Computer", shows three entries:

HostShared		26.8 GB / 45.1 GB available	/
Computer		9.9 GB / 245.1 GB available	/media/host
host			

Hadoop execution steps

```
sx3702@ubuntu:~$ hadoop namenode -format
WARNING: Use of this script to execute namenode is deprecated.
WARNING: Attempting to execute replacement "hdfs namenode" instead.

2023-04-25 21:58:15,879 INFO namenode.NameNode: STARTUP_MSG:
/*****
STARTUP_MSG: Starting NameNode
STARTUP_MSG: host = ubuntu/127.0.1.1
STARTUP_MSG: args = [-format]
STARTUP_MSG: version = 3.3.5
STARTUP_MSG: classpath = /usr/local/hadoop-3.3.5/etc/hadoop:/usr/local/hadoop-3.3.5/share/hadoop/common/lib/netty-codec-smtp-4.1.77.Final.jar:/usr/local/hadoop-3.3.5/share/hadoop/common/lib/netty-common-4.1.77.Final.jar:/usr/local/hadoop-3.3.5/share/hadoop/common/lib/netty-handler-4.1.77.Final.jar:/usr/local/hadoop-3.3.5/share/hadoop/common/lib/netty-handler-proxy-4.1.77.Final.jar:/usr/local/hadoop-3.3.5/share/hadoop/common/lib/netty-resolver-4.1.77.Final.jar:/usr/local/hadoop-3.3.5/share/hadoop/common/lib/netty-transport-4.1.77.Final.jar:/usr/local/hadoop-3.3.5/share/hadoop/common/lib/netty-transport-native-epoll-4.1.77.Final.jar:/usr/local/hadoop-3.3.5/share/hadoop/common/lib/netty-transport-native-unix-4.1.77.Final.jar:/usr/local/hadoop-3.3.5/share/hadoop/common/lib/zookeeper-3.5.14
```

```
sx3702@ubuntu:~
```

```
*****  
SHUTDOWN_MSG: Shutting down NameNode at ubuntu/127.0.1.1  
*****  
sx3702@ubuntu:~$ start-dfs.sh —  
Starting namenodes on [localhost]  
Starting datanodes  
Starting secondary namenodes [ubuntu]  
2023-04-25 21:59:22,669 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable  
sx3702@ubuntu:~$ start-yarn.sh —  
Starting resourcemanager  
Starting nodemanagers  
sx3702@ubuntu:~$ jps —  
3920 ResourceManager  
3345 NameNode  
3731 SecondaryNameNode  
4422 Jps  
4206 NodeManager  
3518 DataNode  
sx3702@ubuntu:~$ export HADOOP_CLASSPATH=$(hadoop classpath) —  
sx3702@ubuntu:~$ echo HADOOP_CLASSPATH  
HADOOP_CLASSPATH  
sx3702@ubuntu:~$ hadoop fs -mkdir /Count —  
2023-04-25 22:00:32,911 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable  
sx3702@ubuntu:~$ hadoop fs -mkdir /Count/Word —  
2023-04-25 22:00:51,930 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable  
sx3702@ubuntu:~$ hadoop fs -mkdir /Count/Imdb —  
2023-04-25 22:01:18,535 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable  
sx3702@ubuntu:~$ hadoop fs -mkdir /Count/Input —  
2023-04-25 22:01:25,549 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
```

```
sx3702@ubuntu:~$ hadoop fs -put Desktop/Lab/Input/Spring2023-Project3-IMDbData.txt /Count/Input/ —  
2023-04-25 22:13:25,233 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform...  
re applicable  
sx3702@ubuntu:~/Desktop/Lab$ ls —  
imdbcoundbms Input wordcountdbms WordCount.java  
sx3702@ubuntu:~/Desktop/Lab$ cp ../../Downloads/IMDbData.java . —  
sx3702@ubuntu:~/Desktop/Lab$ javac -classpath $HADOOP_CLASSPATH -d imdbcoundbms IMDbData.java —  
sx3702@ubuntu:~/Desktop/Lab$ ls —  
imdbcoundbms IMDbData.java Input wordcountdbms WordCount.java  
sx3702@ubuntu:~/Desktop/Lab$ jar -cvf Imdb.jar -C imdbcoundbms . —  
added manifest  
adding: IMDbData$IntSumReducer.class(in = 1736) (out= 741)(deflated 57%)  
adding: IMDbData$TokenizerMapper.class(in = 3515) (out= 1756)(deflated 50%)  
adding: IMDbData.class(in = 1487) (out= 815)(deflated 45%)  
sx3702@ubuntu:~/Desktop/Lab$ ls —  
imdbcoundbms IMDbData.java Imdb.jar Input wordcountdbms WordCount.java  
sx3702@ubuntu:~/Desktop/Lab$ hadoop jar Imdb.jar IMDbData /Count/Input /Count/Imdb/Output —
```

```

1682459971055_0003
2023-04-25 22:14:29,006 INFO input.FileInputFormat: Total input files to process : 1
2023-04-25 22:14:29,518 INFO mapreduce.JobSubmitter: number of splits:4
2023-04-25 22:14:29,663 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1682459971055_0003
2023-04-25 22:14:29,663 INFO mapreduce.JobSubmitter: Executing with tokens: []
2023-04-25 22:14:29,791 INFO conf.Configuration: resource-types.xml not found
2023-04-25 22:14:29,791 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.
2023-04-25 22:14:30,176 INFO impl.YarnClientImpl: Submitted application application_1682459971055_0003
2023-04-25 22:14:30,226 INFO mapreduce.Job: The url to track the job: http://ubuntu:8088/proxy/application_1682459971055_0003
2023-04-25 22:14:30,227 INFO mapreduce.Job: Running job: job_1682459971055_0003
2023-04-25 22:14:35,427 INFO mapreduce.Job: Job job_1682459971055_0003 running in uber mode : false
2023-04-25 22:14:35,429 INFO mapreduce.Job: map 0% reduce 0%
2023-04-25 22:14:55,045 INFO mapreduce.Job: map 52% reduce 0%
2023-04-25 22:14:57,067 INFO mapreduce.Job: map 100% reduce 0%
2023-04-25 22:15:00,094 INFO mapreduce.Job: map 100% reduce 100%
2023-04-25 22:15:01,143 INFO mapreduce.Job: Job job_1682459971055_0003 completed successfully
2023-04-25 22:15:01,636 INFO mapreduce.Job: Counters: 55
    File System Counters
        FILE: Number of bytes read=1124
        FILE: Number of bytes written=1384076
        FILE: Number of read operations=0
        FILE: Number of large read operations=0
        FILE: Number of write operations=0
        HDFS: Number of bytes read=497433541
        HDFS: Number of bytes written=284
        HDFS: Number of read operations=17
        HDFS: Number of large read operations=0
        HDFS: Number of write operations=2
        HDFS: Number of bytes read erasure-coded=0
    Job Counters

```

Hadoop final output

```

sx3702@ubuntu:~/Desktop/Lab$ hadoop dfs -cat /Count/Imdb/Output2/*
WARNING: Use of this script to execute dfs is deprecated.
WARNING: Attempting to execute replacement "hdfs dfs" instead.

```

2023-04-25 22:23:32,154 WARN util.NativeCodeLoader: Unable to load native library, see stacktrace below.		
2000-2006	Action/Drama	1240
2000-2006	Adventure/Sci-Fi	67
2000-2006	Comedy/Romance	1620
2007-2013	Action/Drama	1811
2007-2013	Adventure/Sci-Fi	153
2007-2013	Comedy/Romance	2402
2014-2020	Action/Drama	2326
2014-2020	Adventure/Sci-Fi	247
2014-2020	Comedy/Romance	2784

Hadoop file system

Browsing HDFS New Tab

localhost:9870/explorer.html#/Count/Imdb/Output

Hadoop Overview Datanodes Datanode Volume Failures Snapshot Startup Progress Utilities

Browse Directory

/Count/Imdb/Output

Show 25 entries

	Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
<input type="checkbox"/>	-rW-r--r--	sx3702	supergroup	0 B	Apr 25 22:06	1	128 MB	_SUCCESS
<input type="checkbox"/>	-rW-r--r--	sx3702	supergroup	0 B	Apr 25 22:06	1	128 MB	part-r-00000

Showing 1 to 2 of 2 entries

Search:

Previous 1 Next

Hadoop 2023.

Additional Work

[Task 1] Wrote python script to automate manual work of downloading part-r-0000 map reduce output from Hadoop localhost that has code to convert csv format to json and directly display analysis graph without human interweaving.

```
➊ DownloadAndPlotHistogram.py ×

Users > sakshisrivastava > Desktop > uta > Sem3Spring > DBMS > project3 > ➋ DownloadAndPlotHistogram.py

1  #!pip3 install matplotlib
2  import matplotlib.pyplot as plt
3  data = {}
4  with open('/home/sx3702/Downloads/part-r-00000', 'r') as f:
5      for line in f:
6          year_genre, count = line.strip().split('\t')
7          year, genre = year_genre.split(' | ')
8          count = int(count)
9          if year not in data:
10              data[year] = {}
11              data[year][genre] = count
12  print(data)
13  # data = {
14  #     '2000-2006': {'Action/Drama': 167371, 'Adventure/Sci-Fi': 11891, 'Comedy/Romance': 100000},
15  #     '2007-2013': {'Action/Drama': 330421, 'Adventure/Sci-Fi': 27804, 'Comedy/Romance': 100000},
16  #     '2014-2020': {'Action/Drama': 474792, 'Adventure/Sci-Fi': 49363, 'Comedy/Romance': 100000}
17  # }
18  def plot_histogram(start_date, end_date, genres):
19      #genres = ['Action/Drama', 'Adventure/Sci-Fi', 'Comedy/Romance']
20      data_range = {}
21      for year in data:
22          if start_date <= year <= end_date:
23              for genre in genres:
24                  if genre in data[year]:
25                      if year in data_range:
26                          data_range[year][genre] = data[year][genre]
27                      else:
28                          data_range[year] = {genre: data[year][genre]}
29      fig, ax = plt.subplots()
30      for genre in genres:
31          genre_data = [data_range[year][genre] for year in data_range]
32          ax.bar(data_range.keys(), genre_data, label=genre)
33      ax.set_xlabel('Year Range')
34      ax.set_ylabel('Number of Movies')
35      ax.set_title('Movie Genres by Year Range')
36      ax.legend()
37      ax.legend()
38      plt.show()
39  plot_histogram('2000-2006', '2014-2020', ['Action/Drama'])
40  plot_histogram('2000-2006', '2014-2020', ['Adventure/Sci-Fi'])
41  plot_histogram('2000-2006', '2014-2020', ['Comedy/Romance'])
```

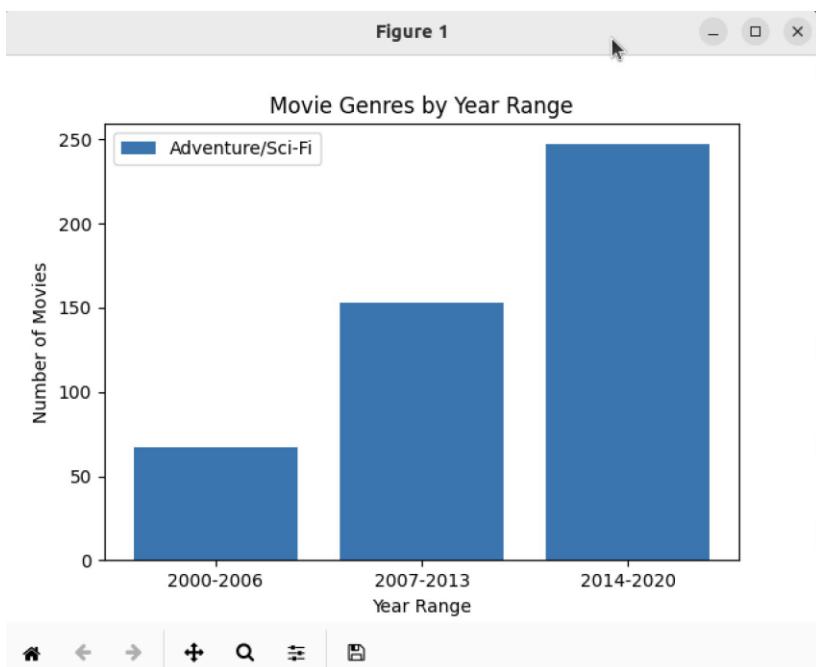
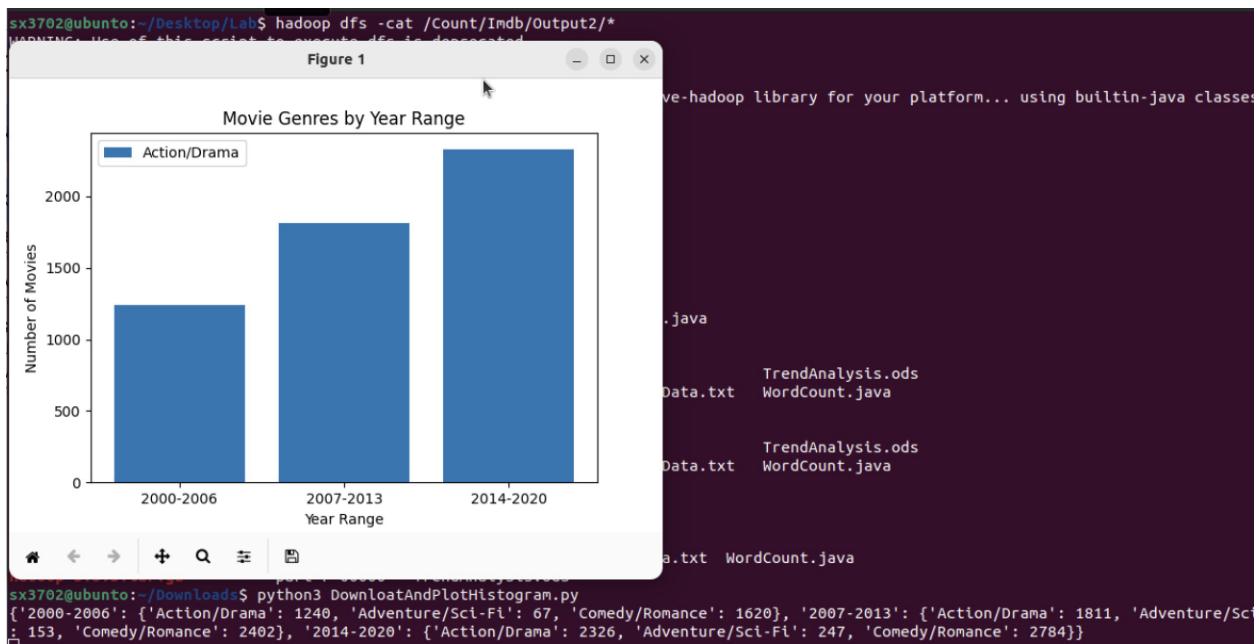


Figure 1

