



IT-478 MILESTONE 4

Movie Ratings Analysis



To

Dr. Bryan Hosack

Madhuri Choudhary
Mayank Shrivansh
Mayank Agrawal
Yesh Polishetty

Table of Contents

Goal.....	3
Business Problem.....	3
What we have done	
• Implementation.....	3
Data flow of the Project.....	4
Team Project Planning.....	13
Lesson Learnt.....	14
Conclusion.....	14

Goal:

To Implement Big Data Architecture in an organization which makes educational movies, to create increase revenue by analyzing the relevant data sets and to generate meaningful data from different type of data sets that will help client to make better decision.

Business Problem:

The organization is not able to generate enough revenue by making educational movies and they want to increase their financial status without changing their main goal. They want to

What we have done:

Implementation:

Step 1: Installation of Cloudera CDH3 virtual platform with VM Player.



Step 2: Then we have done demo for our understanding of the system.

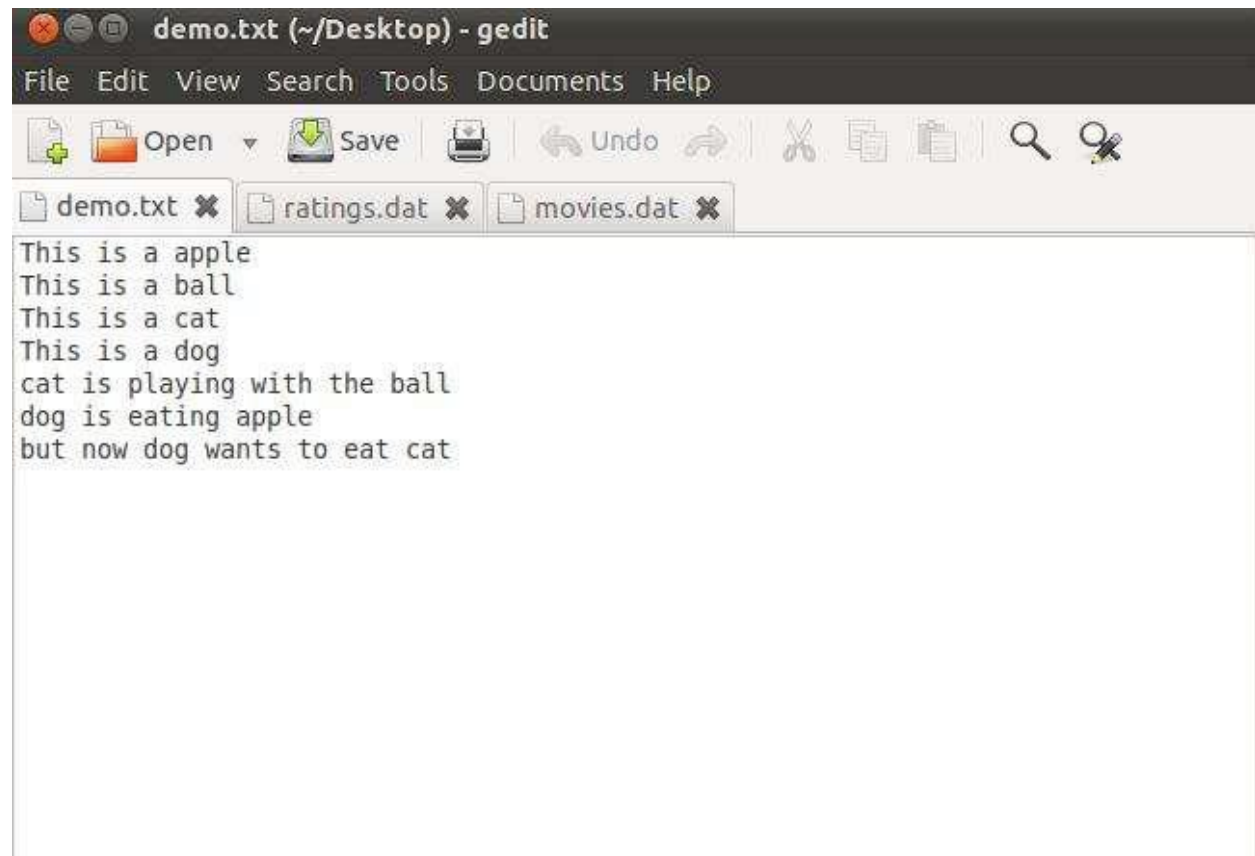
Here is the files that we inserted in to the HDFS:

```

cloudera@cloudera-vm:/usr/lib/hadoop$ hadoop dfs -ls
Found 4 items
-rw-r--r-- 1 cloudera cloudera 171308 2014-10-27 13:17 /user/cloudera/adb
-rw-r--r-- 1 cloudera cloudera 134368 2014-10-27 09:17 /user/cloudera/adb1
-rw-r--r-- 1 cloudera cloudera 24594131 2014-10-27 09:18 /user/cloudera/adb2
drwxr-xr-x - cloudera supergroup 0 2014-10-27 13:38 /user/cloudera/wordcount

```

Sample Text File demo.txt contains the set of statements as below:



Then we load this file into HDFS:

```

cloudera@cloudera-vm:~$ hadoop dfs -copyFromLocal /home/cloudera/Desktop/demo.tx
t /usr/training/demo.txt

```

After that we apply Map task on the file to get the word count of a particular word. The source of Map program is

http://www.cloudera.com/content/cloudera/en/documentation/HadoopTutorial/CDH4/Hadoop-Tutorial/ht_wordcount1_source.html

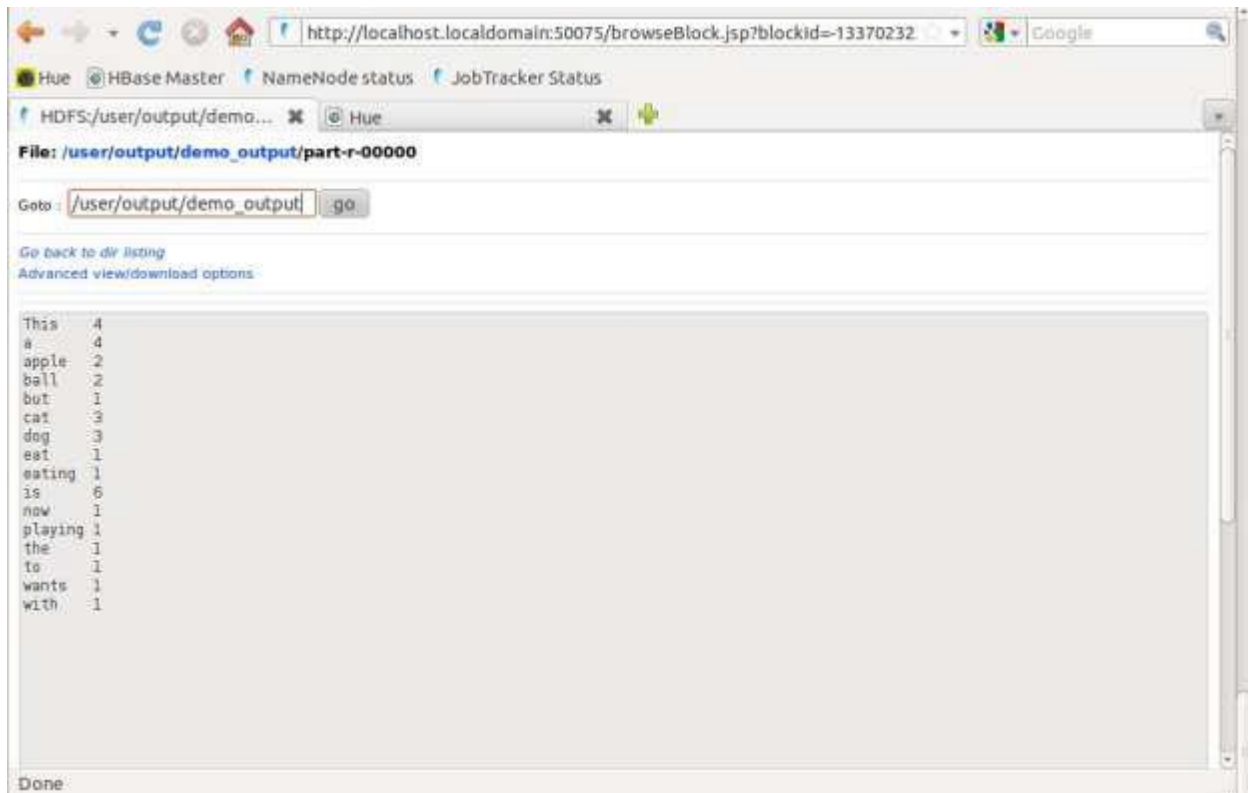
The detailed execution of Map task is shown below:

```

cloudera@cloudera-vm:~/usr/lib/hadoop$ hadoop jar hadoop-examples.jar wordcount /
usr/training/demo.txt /user/output/demo_output
14/10/27 15:02:24 INFO input.FileInputFormat: Total input paths to process : 1
14/10/27 15:02:25 INFO mapred.JobClient: Running job: job_201410210013_0001
14/10/27 15:02:26 INFO mapred.JobClient: map 0% reduce 0%
14/10/27 15:02:36 INFO mapred.JobClient: map 100% reduce 0%
14/10/27 15:02:45 INFO mapred.JobClient: map 100% reduce 100%
14/10/27 15:02:46 INFO mapred.JobClient: Job complete: job_201410210013_0001
14/10/27 15:02:46 INFO mapred.JobClient: Counters: 22
14/10/27 15:02:46 INFO mapred.JobClient: Job Counters
14/10/27 15:02:46 INFO mapred.JobClient:   Launched reduce tasks=1
14/10/27 15:02:46 INFO mapred.JobClient:   SLOTS_MILLIS_MAPS=7287
14/10/27 15:02:46 INFO mapred.JobClient:   Total time spent by all reduces waiting after reserving slots (ms)=0
14/10/27 15:02:46 INFO mapred.JobClient:   Total time spent by all maps waiting after reserving slots (ms)=0
14/10/27 15:02:46 INFO mapred.JobClient:   Launched map tasks=1
14/10/27 15:02:46 INFO mapred.JobClient:   Data-local map tasks=1
14/10/27 15:02:46 INFO mapred.JobClient:   SLOTS_MILLIS_REDUCES=9570
14/10/27 15:02:46 INFO mapred.JobClient: FileSystemCounters
14/10/27 15:02:46 INFO mapred.JobClient:   FILE_BYTES_READ=176
14/10/27 15:02:46 INFO mapred.JobClient:   HDFS_BYTES_READ=240
14/10/27 15:02:46 INFO mapred.JobClient:   FILE_BYTES_WRITTEN=106676
14/10/27 15:02:46 INFO mapred.JobClient:   HDFS_BYTES_WRITTEN=106
14/10/27 15:02:46 INFO mapred.JobClient: Map-Reduce Framework
14/10/27 15:02:46 INFO mapred.JobClient:   Reduce input groups=16
14/10/27 15:02:46 INFO mapred.JobClient:   Combine output records=16
14/10/27 15:02:46 INFO mapred.JobClient:   Map input records=7
14/10/27 15:02:46 INFO mapred.JobClient:   Reduce shuffle bytes=176
14/10/27 15:02:46 INFO mapred.JobClient:   Reduce output records=16
14/10/27 15:02:46 INFO mapred.JobClient:   Spilled Records=32
14/10/27 15:02:46 INFO mapred.JobClient:   Map output bytes=269
14/10/27 15:02:46 INFO mapred.JobClient:   Combine input records=33
14/10/27 15:02:46 INFO mapred.JobClient:   Map output records=33
14/10/27 15:02:46 INFO mapred.JobClient:   SPLIT_RAW_BYTES=103
14/10/27 15:02:46 INFO mapred.JobClient:   Reduce input records=16

```

The result of the Map task is as follows:

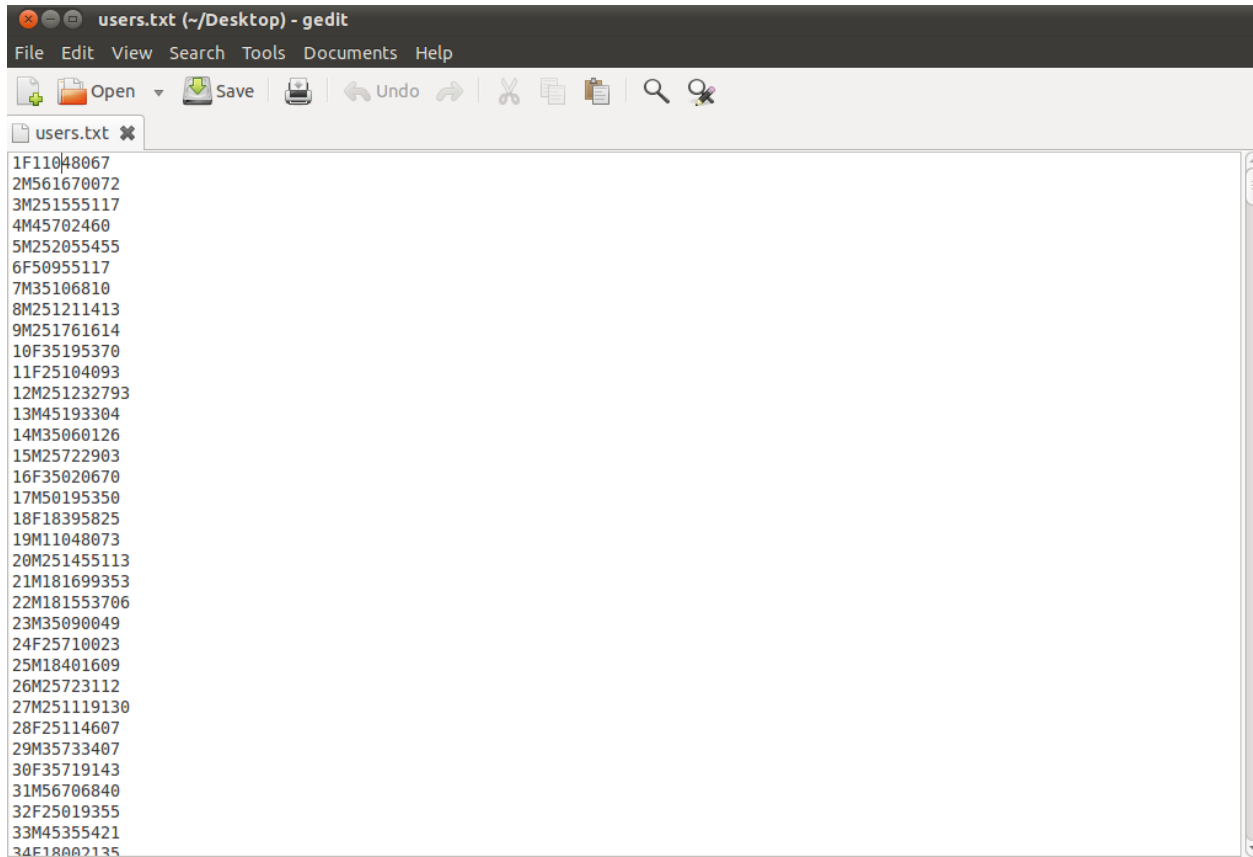


Step 3: Now we started using the system according to our requirements.

At first, we transfer our datasets from our local machine to Cloudera's VM through FileZilla FTP client.

And then we stored this unstructured data files in to the HDFS.

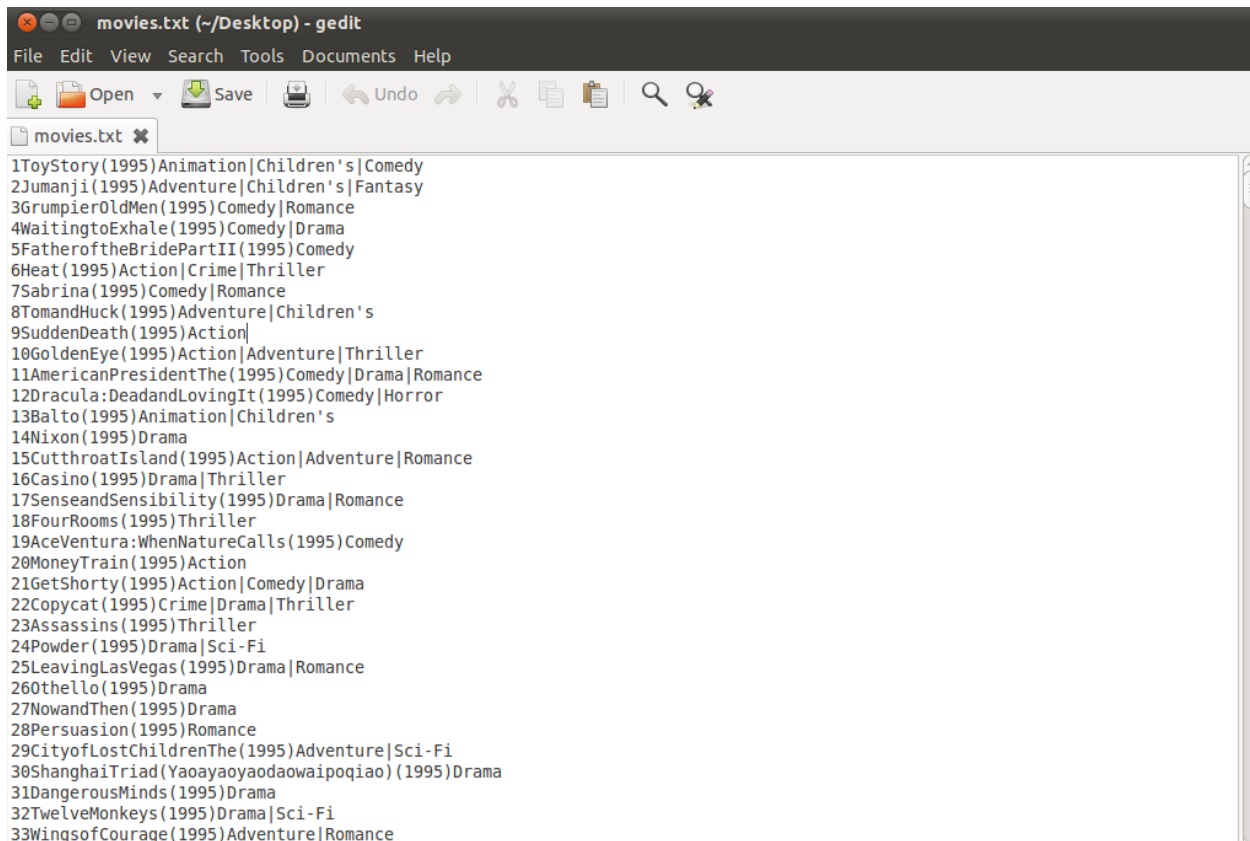
1. User.txt



The image shows a screenshot of a gedit text editor window. The title bar reads "users.txt (~/Desktop) - gedit". The menu bar includes "File", "Edit", "View", "Search", "Tools", "Documents", and "Help". The toolbar contains icons for "Open", "Save", "Print", "Undo", "Redo", "Cut", "Copy", "Paste", "Find", and "Replace". The text area displays a list of 34 user IDs, each consisting of a number followed by a letter and a six-digit hexadecimal string. The list is as follows:

```
1F11048067
2M561670072
3M251555117
4M45702460
5M252055455
6F50955117
7M35106810
8M251211413
9M251761614
10F35195370
11F25104093
12M251232793
13M45193304
14M35060126
15M25722903
16F35020670
17M50195350
18F18395825
19M11048073
20M251455113
21M181699353
22M181553706
23M35090049
24F25710023
25M18401609
26M25723112
27M251119130
28F25114607
29M35733407
30F35719143
31M56706840
32F25019355
33M45355421
34F18002135
```

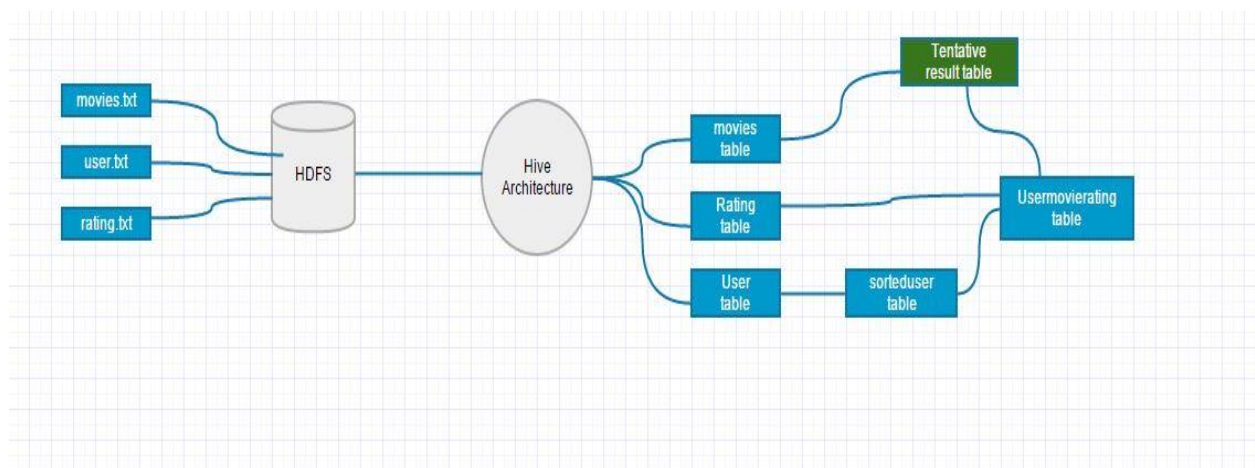
2. Movies.txt



3. Rating.txt

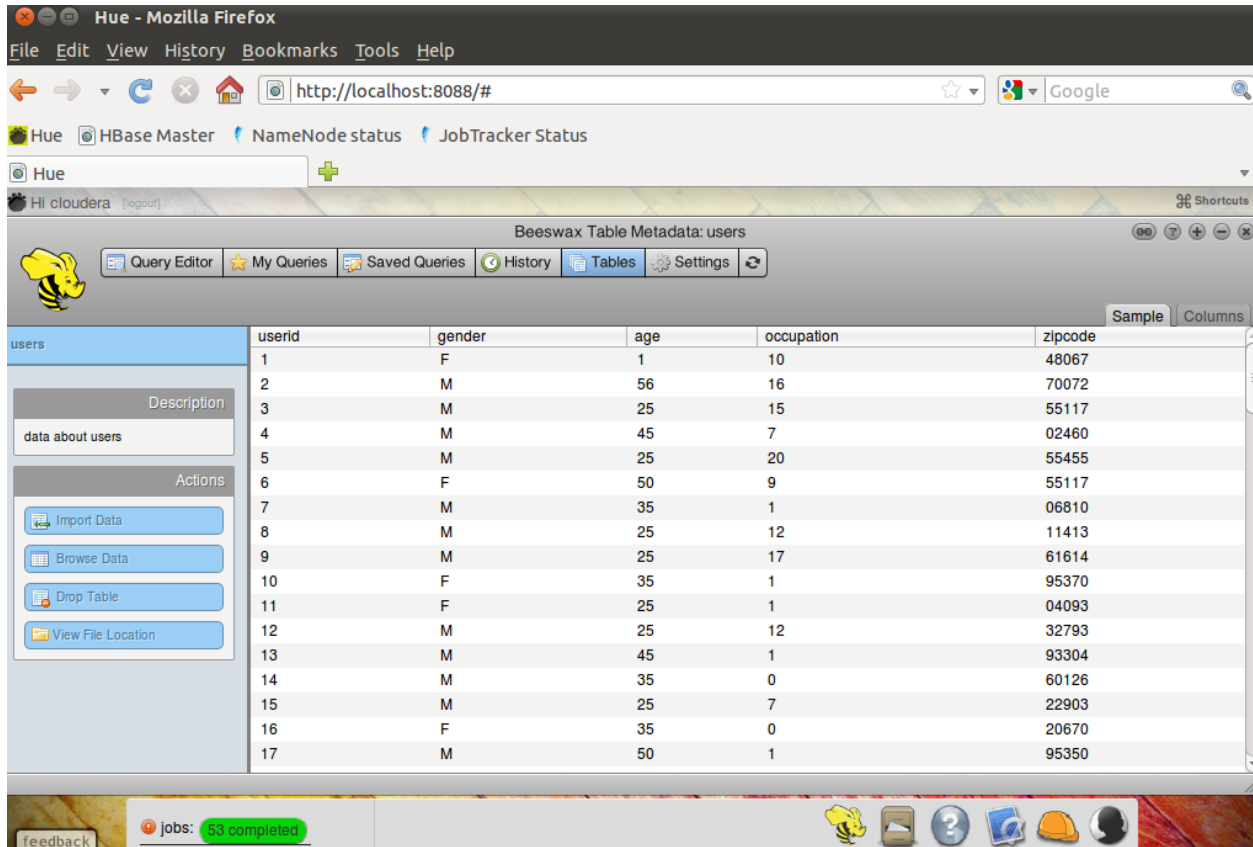
```
ratings.txt (~/Desktop) - gedit
File Edit View Search Documents Help
Open Save Undo
ratings.txt
229163978299809
234685978298542
212104978298151
217923978299941
216873978300174
212132978298458
235785978298958
228813978300002
230304978298434
212173978298151
231054978298673
24342978300174
221263978300123
231072978300002
231083978299712
230354978298625
212533978299120
216105978299809
22923978300123
222365978299220
230714978299120
29022978298905
23684978300002
212595978298841
231475978298652
215444978300174
212935978298261
211884978299620
232554978299321
232562978299839
232573978300073
21105978298625
222783978299889
```

Here is the Data flow of the Overall Project:



According to the data flow we proceed and convert the unstructured datasets in to a structured format like tables.

1. User Table.



users

userid	gender	age	occupation	zipcode
1	F	1	10	48067
2	M	56	16	70072
3	M	25	15	55117
4	M	45	7	02460
5	M	25	20	55455
6	F	50	9	55117
7	M	35	1	06810
8	M	25	12	11413
9	M	25	17	61614
10	F	35	1	95370
11	F	25	1	04093
12	M	25	12	32793
13	M	45	1	93304
14	M	35	0	60126
15	M	25	7	22903
16	F	35	0	20670
17	M	50	1	95350

feedback jobs: 53 completed

2. Movies Table.

Hue - Mozilla Firefox

File Edit View History Bookmarks Tools Help

http://localhost:8088/#

Hue HBase Master NameNode status JobTracker Status

Hi cloudera [logout]

Beeswax Table Metadata: movies

Query Editor My Queries Saved Queries History Tables Settings

movies

Description

data about movies

Actions

Import Data

Browse Data

Drop Table

View File Location

Sample	Columns
null1	movieid null2 genre null3
1	Animation Children's Comedy
2	Adventure Children's Fantasy
3	Comedy Romance
4	Comedy Drama
5	Comedy
6	Action Crime Thriller
7	Comedy Romance
8	Adventure Children's
9	Action
10	Action Adventure Thriller
11	Comedy Drama Romance
12	Comedy Horror
13	Animation Children's
14	Drama
15	Action Adventure Romance
16	Drama Thriller
17	Drama Romance

feedback

Jobs: 53 completed

3. Ratings Table.

Beeswax Table Metadata: ratings

userid	movieid	rating	timestamp
1	1193	5	978300760
1	661	3	978302109
1	914	3	978301968
1	3408	4	978300275
1	2355	5	978824291
1	1197	3	978302268
1	1287	5	978302039
1	2804	5	978300719
1	594	4	978302268
1	919	4	978301368
1	595	5	978824268
1	938	4	978301752
1	2398	4	978302281
1	2918	4	978302124
1	1035	5	978301753
1	2791	4	978302188
1	2687	3	978824268

After that we apply some queries to the tables and came with the certain sets of Tables, by using certain delimiters and Serializable and Deserializable interfaces.

Some queries that we used are:-

- To sort the user with age group 18 (18 to 24yrs) and generate **sorteduser** tables.

Select * from **usertable** where age ='18';

- To create a new table called **usermovierating** from **sorteduser** and **rating** table.

Select userid, gender, age, movieid, rating from **sortedusers** join **rating** on (users.userid=movies.userid) where ratings.rating='5';

- To generate the **tentativeresult** table from **movies** and **usermovierating** table.

Select movieid, userid, gender, age, rating, genre from **movies** join **usermovierating** on (movies.movieid=usermovierating.movieid) where movies.movieid=usermovierating.movieid;

- To calculate the number of user under the age group 18 (18 to 24 yrs). i.e. **1098**

Select distinct userid, gender, age, rating, genre from tentativeresult;

- To calculate the number of females from (18 to 24 yrs). i.e. **298**

Select * from tentativeresult where gender='F';

- To calculate the number of males from (18 to 24 yrs) i.e. **800**

Select * from tentativeresult where gender='M';

Some of the Screenshots of our result tables are:

1. Tentative Result.

movieid	userid	gender	age	rating	genre
1	1664	F	18	5	Animation Children's Comedy
1	1676	M	18	5	Animation Children's Comedy
1	3163	M	18	5	Animation Children's Comedy
1	1765	M	18	5	Animation Children's Comedy
1	3174	F	18	5	Animation Children's Comedy
1	3185	M	18	5	Animation Children's Comedy
1	1778	M	18	5	Animation Children's Comedy
1	301	M	18	5	Animation Children's Comedy
1	3006	M	18	5	Animation Children's Comedy
1	2996	M	18	5	Animation Children's Comedy
1	1117	M	18	5	Animation Children's Comedy
1	2968	F	18	5	Animation Children's Comedy
1	2931	M	18	5	Animation Children's Comedy
1	2924	F	18	5	Animation Children's Comedy
1	1125	F	18	5	Animation Children's Comedy
1	2885	F	18	5	Animation Children's Comedy
1	2873	F	18	5	Animation Children's Comedy

2. Male Table.

Hue - Mozilla Firefox

File Edit View History Bookmarks Tools Help

http://localhost:8088/#

Hue HBase Master NameNode status JobTracker Status

Hue

Hi cloudera [logout]

Beeswax Table Metadata: tmalesresult

Query Editor My Queries Saved Queries History Tables Settings

Sample Columns

movielid	userid	gender	age	rating	genre
2406	3163	M	18	5	Action Adventure Comedy Romance
2406	2186	M	18	5	Action Adventure Comedy Romance
2406	3029	M	18	5	Action Adventure Comedy Romance
2406	1889	M	18	5	Action Adventure Comedy Romance
2406	272	M	18	5	Action Adventure Comedy Romance
2406	1420	M	18	5	Action Adventure Comedy Romance
2406	3391	M	18	5	Action Adventure Comedy Romance
2406	3195	M	18	5	Action Adventure Comedy Romance
2406	817	M	18	5	Action Adventure Comedy Romance
2406	714	M	18	5	Action Adventure Comedy Romance
2406	4635	M	18	5	Action Adventure Comedy Romance
2406	3653	M	18	5	Action Adventure Comedy Romance
2406	4041	M	18	5	Action Adventure Comedy Romance
2406	4374	M	18	5	Action Adventure Comedy Romance
2406	4535	M	18	5	Action Adventure Comedy Romance
2407	5706	M	18	5	Comedy Sci-Fi
2407	764	M	18	5	Comedy Sci-Fi

Jobs: 53 completed

3. Female Table.

Beeswax Table Metadata: tfmaleresult

ID	Age	Gender	Genre
2406	2944	F	Action Adventure Comedy Romance
2406	2885	F	Action Adventure Comedy Romance
2406	3951	F	Action Adventure Comedy Romance
2406	4298	F	Action Adventure Comedy Romance
2407	3308	F	Comedy Sci-Fi
2409	101	F	Action Drama
241	3308	F	Children's Drama
2410	5530	F	Action Drama
2411	101	F	Action Drama
2411	4660	F	Action Drama
2412	101	F	Action Drama
2413	3196	F	Comedy Mystery
2413	1755	F	Comedy Mystery
2413	3308	F	Comedy Mystery
2413	3259	F	Comedy Mystery
2413	1645	F	Comedy Mystery
2413	92	F	Comedy Mystery
2413	4918	F	Comedy Mystery

Jobs: 53 completed

The above result tables are derived from the unstructured datasets that we have at the starting and with the help of this result tables the analyst can easily analyze the data according to the requirements which would be very hard to do with the unorganized data sets that we have earlier.

Some of the results that we have found are as follows:

- There are 1098 users in the datasets that are in the age range 18-24yrs.
- Among them 800 were males and most of the time they would like to watch Action/Adventure movies.
- And there are 298 females and most of the time they would like to watch Comedy/Drama type movies.

Team Project Planning:

Our job was to make meaningful analysis report from the available dataset that will help the organization to make better decision in their movie making choice. For that, our plan was to first identify the datasets based on the requirements of the organization. From those datasets we performed the required testing and proceed with the planned application. Firstly, data was inserted

in the HDFS and then we wrote map and reduce programs for it. Then we planned to apply multiple queries to the result data using map and reduce programs and to filter that data and then to analyze the filtered data and generate the final report for the end client.

We faced many problems in implementing this plan, like we had problems while creating map-reduce programs then as team we decided to switch over a tool called Beeswax Hive.

Hive can be defined as an infrastructure that runs on the top of Hadoop and it is used for the analysis of the structured as well as unstructured data. It works with HiveQL. HiveQL is a query based language similar to SQL with extra features like multi-table insert and create table as select. And Beeswax Hive is a tool which is similar to oracle developer, which is a UI to interact with Hive infrastructure (Data Warehouse).

Hive converts this HiveQL queries to the map reduce program and run accordingly. By using this tool we came up with certain sets of tables which are more organized and easy to understand.

After analyzing the sets of table we came up with the final report which contains the statistics related to movie viewers combined with movies reviews which will help the movie making company to take better decision in their movie making choice. We came up with the statistics related to the people, who are in the age group between 18years to 24years, like their favorite genre of movie based on the movie ratings and reviews. Now the movie making company may decide what movie genre they can opt to make the movie for the young population. We also categorized those tables and data according to the gender (Male and Female).

For the project planning and implementation as a team we collaborated very well and coordinated to work on this project. We faced some issues with time management; we took more time to do the tasks than we decided. But at the end we all did well according to our planning and completed the project on time. When we had issues while creating map-reduce programs, we decided as a team to switch over Beeswax Hive. During this project we also learnt to work in a team and skills to make our points or thoughts in a group.

Lesson Learnt:

- We have learned to work and collaborate in the team.
- Learned to present our ideas in a team.
- Learned the concepts of technologies relevant to the data technologies like Big Data, Hadoop, and the tools that works on top of Hadoop like Hive.
- Learned the project development cycle and the processes in each phase of the project.
- Learned time management and problem solving skills.

Conclusion:

Previously with low volumes of data, intuitive decisions would suffice for an organization. As the data size has grown human ability to make intuitive decisions has been completely reduced. This project is designed to retrieve the necessary information from the available data sets, analyze the information and provide to client in the form of reports.

Overall process and decisions are based on the quantitative methods which are cyclic process like

-

- Problem solving definition and identification

- Design and build Hadoop framework

- Data sets management

- Analysis to produce models

- Execution and testing

- Creating final reports for decision making