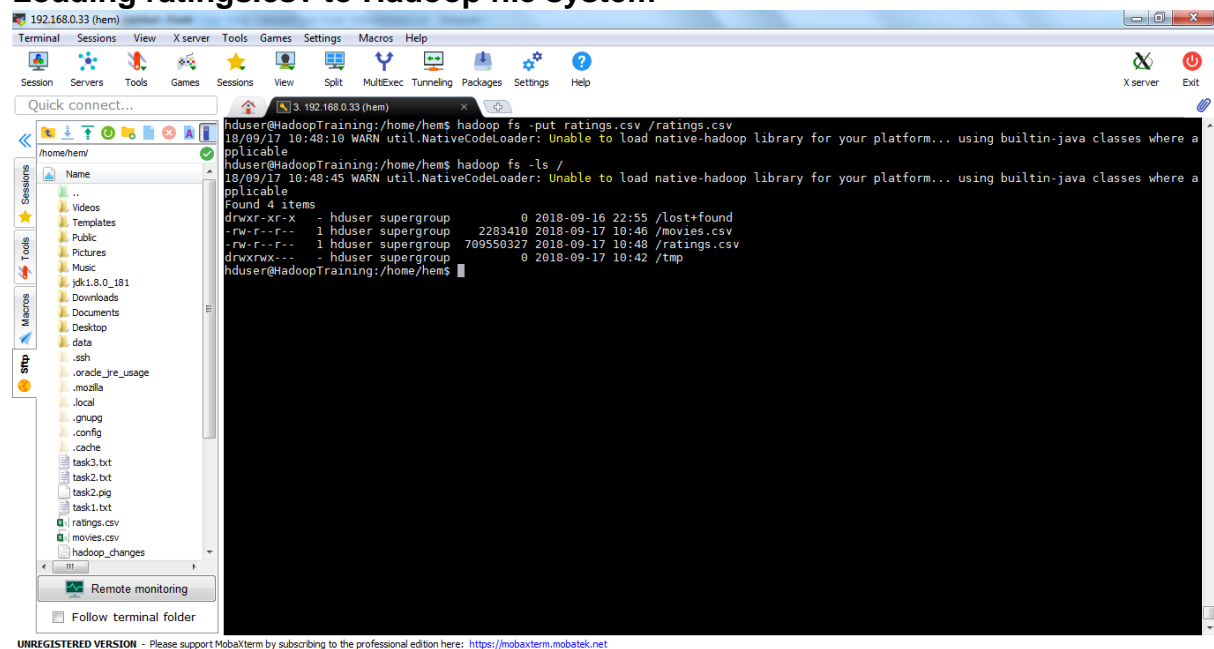## Problem Statement:

## Load Rating.csv and Movies.csv

## Loading movies.csv to Hadoop file system



```
hduser@HadoopTraining:/home/hem$ hadoop fs -put movies.csv /movies.csv
18/09/17 10:46:27 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where a
pplicable
hduser@HadoopTraining:/home/hem$
```

## Loading ratings.csv to Hadoop file system



```
hduser@HadoopTraining:/home/hem$ hadoop fs -put ratings.csv /ratings.csv
18/09/17 10:48:10 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where a
pplicable
hduser@HadoopTraining:/home/hem$ hadoop fs -ls /
18/09/17 10:48:45 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where a
pplicable
Found 4 items
drwxr-xr-x   - hduser supergroup          0 2018-09-16 22:55 /lost+found
-rw-r--r--   1 hduser supergroup    2283410 2018-09-17 10:46 /movies.csv
-rw-r--r--   1 hduser supergroup  709550327 2018-09-17 10:48 /ratings.csv
drwxrwx---   - hduser supergroup          0 2018-09-17 10:42 /tmp
hduser@HadoopTraining:/home/hem$
```

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

**Pig Script-:**

REGISTER '/home/acadgild/piggybank.jar';
A = LOAD '/movies.csv' USING
org.apache.pig.piggybank.storage.CSVExcelStorage(',', 'NO_MULTILINE', 'UNIX',
'SKIP_INPUT_HEADER');
B = foreach A generate (int)$0 as movieid, (chararray)$1 as title;
C = filter B by title is not null;
D = order C by movieid asc;
E = LOAD '/ratings.csv' USING
org.apache.pig.piggybank.storage.CSVExcelStorage(',', 'NO_MULTILINE', 'UNIX',
'SKIP_INPUT_HEADER');
F = foreach E generate (int)$1 as movieid,(float)$2 as rating;
G = filter F by (rating > 0 );
H = order G by movieid asc;
I = COGROUP D by movieid  , H by movieid;
J = order I by $0 asc;
K = distinct J;
L = join C by $0 , K by $0 ;
M = foreach L generate $1;
dump M;

**Script Ilustration-:**
**In Line 1**: We are registering the *piggybank* jar in order to use the CSVExcelStorage class.
In relation **A**, we are loading the dataset using CSVExcelStorage because of its effective technique to handle double quotes and headers.
In relation **B**, we are generating the columns that are required for processing and explicitly typecasting each of them.
In relation **C**, we are filtering the null values from the "title" column.
In relation **D**, we are ordering result by "movieid" column
In relation **E**, we are loading the dataset using CSVExcelStorage because of its effective technique to handle double quotes and headers
In relation **F** we are generating the columns that are required for processing and explicitly typecasting each of them.
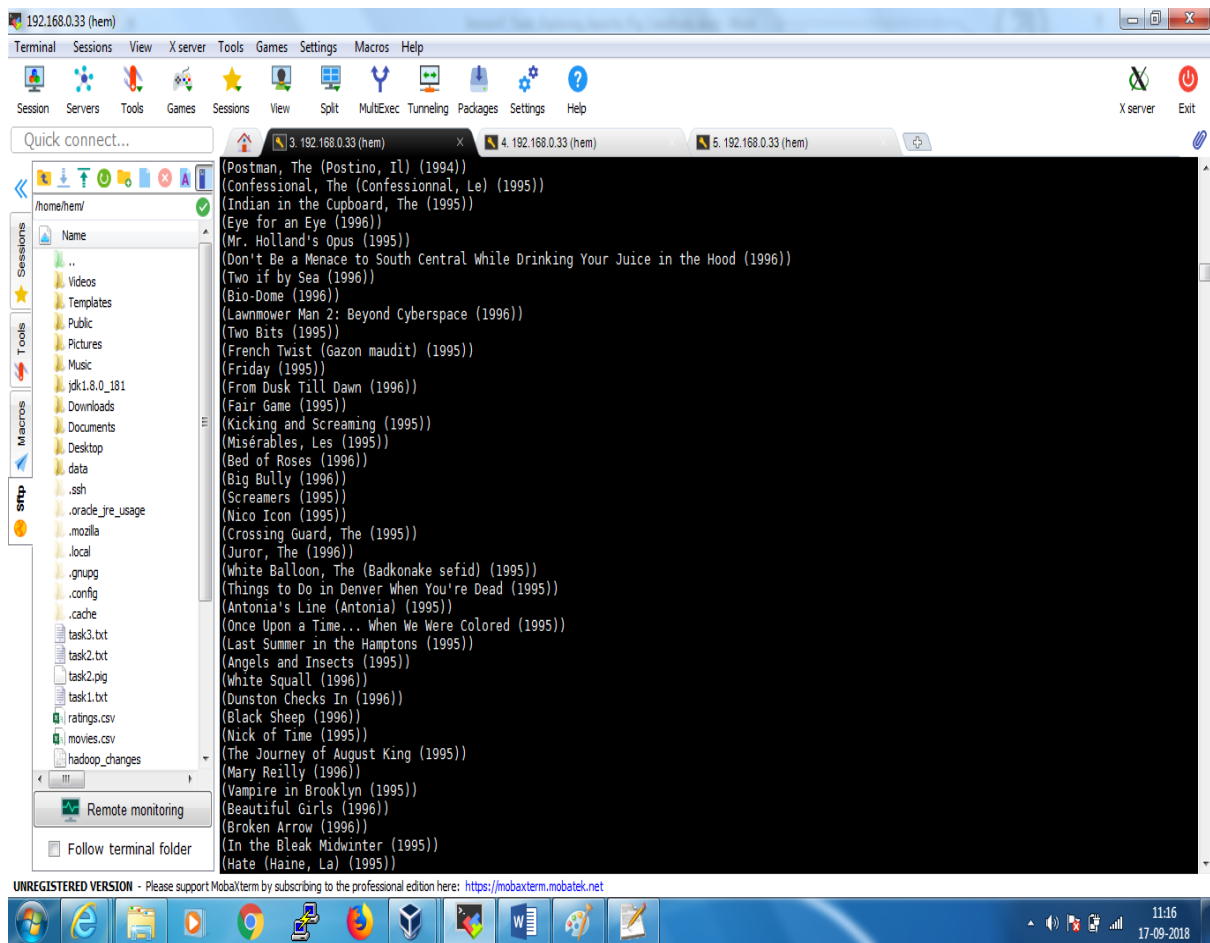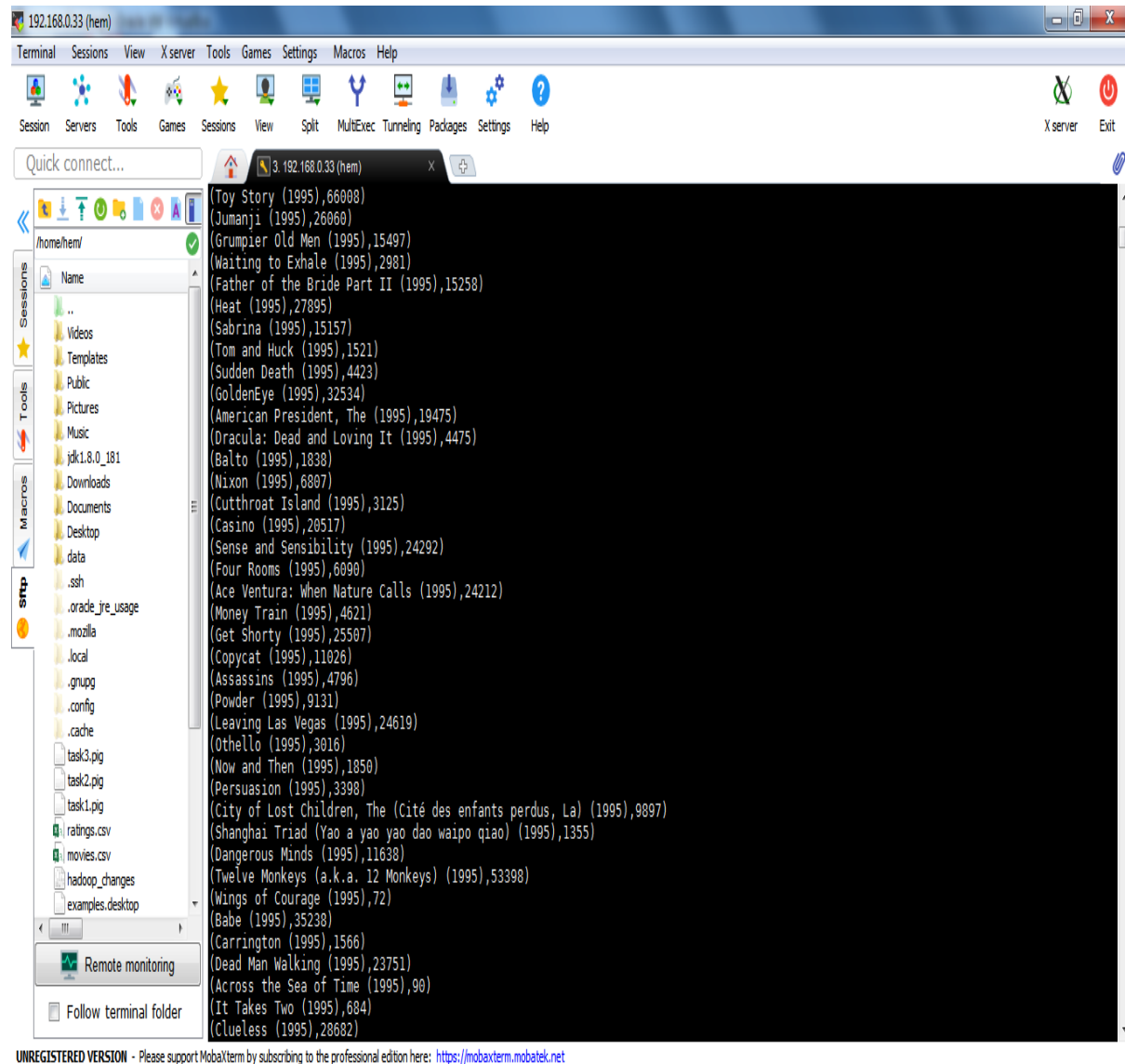Relation **G** and **H** is used to filter rating is greater than 0 and order by movie id in asc order and co grouping D by movieid and H by movieid.
Relation **J** is ordering group by movie id in asc.
Relation **K** is removing duplicates.
In final steps **L**, **M** we are joining both group and **C** by movie id and generate the required column

**Output-:**

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

**Pig Script-:**

REGISTER '/usr/local/pig/lib/piggybank.jar';

A = LOAD '/movies.csv' USING
org.apache.pig.piggybank.storage.CSVExcelStorage(',', 'NO_MULTILINE', 'UNIX',
'SKIP_INPUT_HEADER');

B = foreach A generate (int)$0 as movieid, (chararray)$1 as title;

C = filter B by title is not null;

D = order C by movieid asc;

E = LOAD '/ratings.csv' USING
org.apache.pig.piggybank.storage.CSVExcelStorage(',', 'NO_MULTILINE', 'UNIX',
'SKIP_INPUT_HEADER');

F = foreach E generate (int)$1 as movieid,(float)$2 as rating;

G = filter F by (rating > 0 );

H = order G by movieid asc;

I = COGROUP D by movieid  , H by movieid;

J = order I by $0 asc;

K = foreach J generate group, COUNT(H.$1) as cnt;

L = join C by $0, K by $0;

M = foreach L generate $1 , $3;

dump M;

**Script Ilustration-:**

**In Line 1**: We are registering the *piggybank* jar in order to use the CSVExcelStorage class.
In relation **A**, we are loading the dataset using CSVExcelStorage because of its effective
technique to handle double quotes and headers.
In relation **B**, we are generating the columns that are required for processing and
explicitly typecasting each of them.
In relation **C**, we are filtering the null values from the "title" column.
In relation **D**, we are ordering result by "movieid" column
In relation **E**, we are loading the dataset using CSVExcelStorage because of its effective
technique to handle double quotes and headers
In relation **F** we are generating the columns that are required for processing and
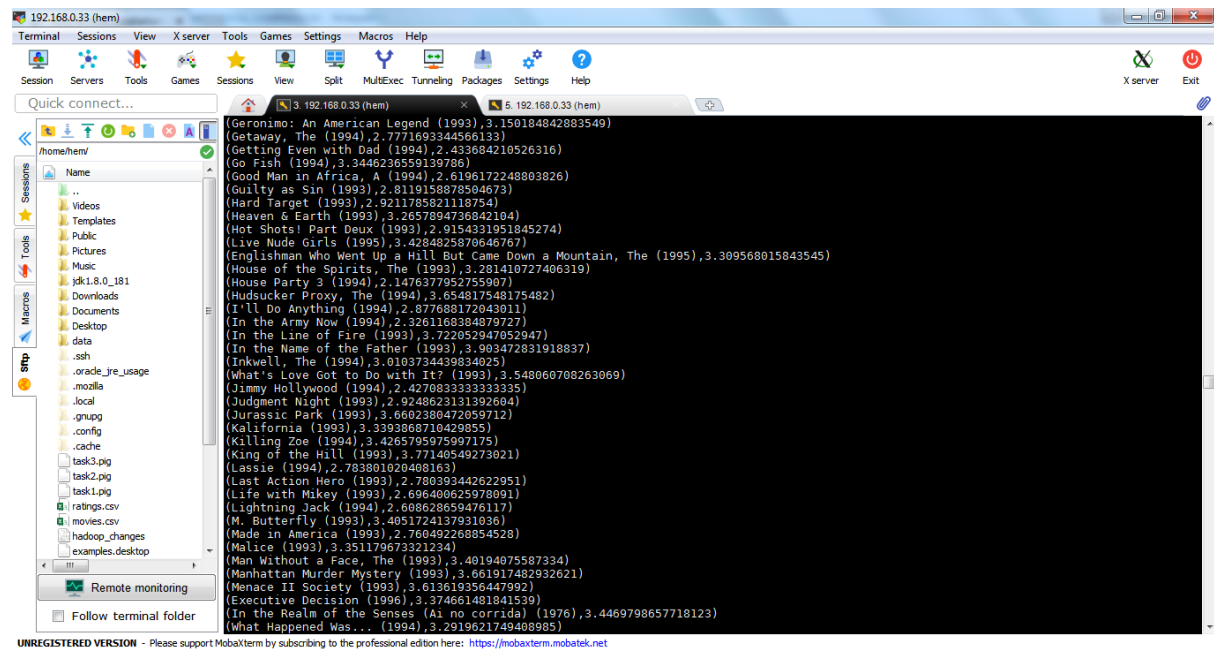explicitly typecasting each of them.
Relation **G** and **H** is used to filter rating is greater than 0 and order by movie id in asc

order and co grouping D by movieid and H by movieid.

Relation **J** is ordering group by movie id in asc.

Relation **K** is used to generate group and count.

In final steps **L**, **M** we are joining both group and **C** by movie id and generate the required columns

**Output-:**

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

**Pig Script-:**

REGISTER '/usr/local/pig/lib/piggybank.jar';

A = LOAD '/movies.csv' USING org.apache.pig.piggybank.storage.CSVExcelStorage(',', 'NO_MULTILINE', 'UNIX', 'SKIP_INPUT_HEADER');

B = foreach A generate (int)$0 as movieid, (chararray)$1 as title;

C = filter B by title is not null;

D = order C by movieid asc;

E = LOAD '/ratings.csv' USING org.apache.pig.piggybank.storage.CSVExcelStorage(',', 'NO_MULTILINE', 'UNIX', 'SKIP_INPUT_HEADER');

F = foreach E generate (int)$1 as movieid,(float)$2 as rating;

G = filter F by (rating > 0 );

H = order G by movieid asc;

I = COGROUP D by movieid , H by movieid;

J = order I by $0 asc;

K = foreach J generate group, AVG(H.$1) as cnt;

L = join C by $0, K by $0;

M = foreach L generate $1 , $3;

dump M;

**Script Ilustration-:**

**In Line 1**: We are registering the *piggybank* jar in order to use the CSVExcelStorage class.
In relation **A**, we are loading the dataset using CSVExcelStorage because of its effective technique to handle double quotes and headers.
In relation **B**, we are generating the columns that are required for processing and explicitly typecasting each of them.
In relation **C**, we are filtering the null values from the "title" column.
In relation **D**, we are ordering result by "movieid" column
In relation **E**, we are loading the dataset using CSVExcelStorage because of its effective technique to handle double quotes and headers
In relation **F** we are generating the columns that are required for processing and explicitly typecasting each of them.
Relation **G** and **H** is used to filter rating is greater than 0 and order by movie id in asc order and co grouping D by movieid and H by movieid.
Relation **J** is ordering group by movie id in asc.
Relation **K** is used to generate group and avg

In final steps **L**, **M** we are joining both group and **C** by movie id and generate the required columns

**Output-:**