

# Machine Learning Approach for Unstructured Data Using Hive

Neha Mangla([apj.neha@gmail.com](mailto:apj.neha@gmail.com)), Shanthi Mahesh([shanthi\\_md@yahoo.com](mailto:shanthi_md@yahoo.com))

(ChhayaM([chhaya.11ms@gmail.com](mailto:chhaya.11ms@gmail.com)), Vidyashree G([victorius.vidya@gmail.com](mailto:victorius.vidya@gmail.com)),  
Vikas([vikasdstar@gmail.com](mailto:vikasdstar@gmail.com)), Atria Institute of Technology, Bangalore.

**Abstract**—Voluminous amount of structured, semi-structured and unstructured data sets that have the potential to learn the relationship among data in the area of business is being collected rapidly; termed as big data. The storage of large chunks of data is difficult as even terabytes and petabytes of traditional data warehousing solutions is insufficient and exorbitant. [1][2]

It is viable to store and process these ransom amount of data on Hadoop; which is a low cost, reliable, scalable and fault tolerant Java-based programming framework that supports the processing of large data sets in a distributed computing environment. Hadoop implements MapReduce programing model for storing and processing large data sets with a parallel, distributed algorithm on commodity hardware. Nevertheless, the programming model expects the developers to write bespoke programs that are less flexible, time consuming, hard to code; maintain and reuse. This challenging task of writing complex MapReduce codes was rationalized by making use of HiveQL.

Hive is the platform required to run HiveQL. Hive is built on top of Hadoop to query Big Data. Internally the Hive queries are converted into the corresponding MapReduce task. [3][4]

In this paper, by making use of machine learning algorithm a movie rating prediction system is built based on MovieLens dataset.

**keywords** - Big Data, HDFS, Hadoop, Hive, MapReduce, linear regression

## I. INTRODUCTION

The prediction system is built using machine learning algorithm. This system employs sentiment analysis that identifies and extracts subjective information based on selected training sets from MovieLens dataset.

At present, cinema has the greatest potential to be the most effective and entertaining mass media instrument. Various software applications and websites such as Bookmyshow, Moviefone, etc., which are the biggest online movie brand don't just assist for ticket booking but also aim at reaching users' satisfaction by providing them with the facility to rate the movies that they have watched and also access feedback on yet to watch movies based on others' ratings. Hence prediction on movie ratings is remarkable.

The system predicts and provides the users with suggestions based on their previous ratings recorded and other users' ratings. These predictions provide an opportunity for movie makers to have better understanding about the viewer's expectations which in turn is beneficial for marketing. This is done by determining the relationship between viewers' and their ratings. Further by making use of effective BigData analysis tools as in this paper Hadoop and Hive are made use of, larger datasets can be analyzed which provides statistically accurate outcomes.[5] These findings provide better understanding about viewers' expectations and hence movie choice.

In this paper, we use MovieLens dataset which is an open dataset collected by GroupLensresearch; University of Minnesota. This dataset is made available on the website for the users to rate movies. MovieLens Dataset comprises of 100K, 1M, 10M datasets having 100 thousand ratings on 1,700 movies from 1,000 users, 1 million ratings on 4,000 movies from 6,000 users and 100 thousand ratings on 10,000 movies from 72,000 users respectively.[6][7][8]

Further HiveQL is used to analyze the dataset which is elaborated in section 2.

## II. DATASET PREPROCESSING USING HIVE

### 2.1 MovieLens Dataset schema

MovieLens Dataset is collected and stored into HDFS (Hadoop Distributed File System) from the website

<http://grouplens.org/datasets/movielens>. For the ease of analysis 100K data set has been chosen. The movies.dat, ratings.dat and the users.dat files have [movieID, tile, gednre], [userID, movieID, rating, timestamp] and [userID, gender, age, occupation, zipCode] fields respectively;[9] with each field delimited from the other by # symbol.

## 2.2 Creating tables and loading data

Tables with same schema as that of the data is created for each of the three files. Hive query to create ratings table and result for the same is as shown in Fig 1.0

Similarly, tables have been created for movies and users files based on their attributes respectively.

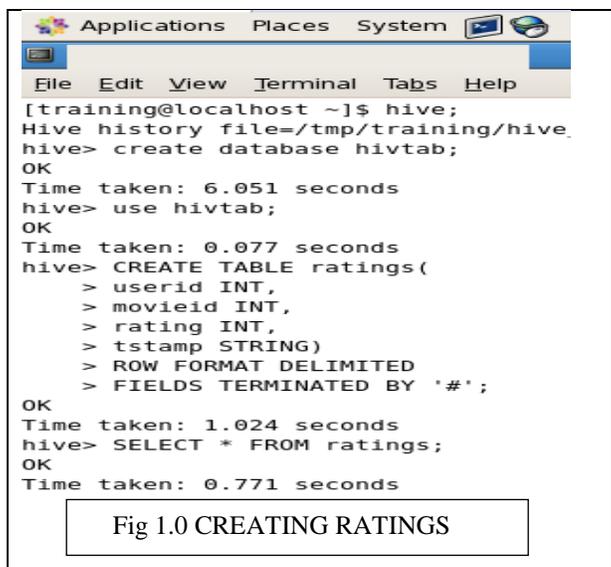


Fig 1.0 CREATING RATINGS

The next step is to load the data into the tables after creating all the three tables. Hive provides us with the utilities to load datasets from flat files stored on HDFS using the LOAD DATA command. The following is the command signature:

```
LOAD DATA LOCAL INPATH <'path_to_flat_file'>
OVERWRITE INTO TABLE <table_name>;
```

The result is as shown below:



Fig 1.1 LOADING DATA INTO RATINGS

The same process of loading the data is carried out for all the three tables.

The verification of data getting loaded into the table can be carried out by displaying the table contents or in the following way making use of SELECT COUNT.

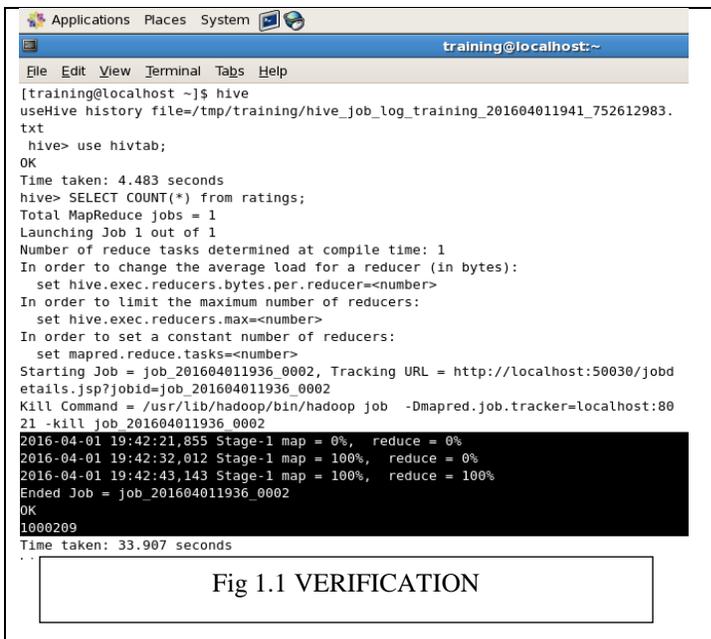


Fig 1.1 VERIFICATION

The highlighted region conveys 2 important points. Firstly, the number of rows present in the table which is as expected approximately equal to 10K. secondly, how hive is internally converted into map-reduce tasks and only after the completion of map phase there reduce phase begins.

### III. APPLYING HIVE QUERIES ON DATASETS

Now that the tables have been created and loaded with the respective datasets successfully, they can be queried using hiveQL which is depicted in the following sections.

#### 3.1 Differential rating based on gender

The subsequent hive query determines the number of people who have rated 5 for the movies based on gender.

```
hive> select users.gender, count(*)
from ratings join users on(users.userid=ratings.userid)
where rating=5 group by users.gender;
```

The result is as shown below:

```
2016-02-29 23:31:22,096 Stage-2 map = 0%, reduce = 0%
2016-02-29 23:31:27,164 Stage-2 map = 100%, reduce = 0%
2016-02-29 23:31:39,913 Stage-2 map = 100%, reduce = 100%
Ended Job = job_201602292316_0003
OK
F      58546
M      167764
```

Fig 1.2: Gender based ratings

From the result obtained we can infer that more number of males rate a movie 5 than female.

#### 3.2 Differential rating based on occupation

The subsequent hive query determines the number of people who have rated 5 for the movies based on occupation:

```
hive> select occupations.occupation,count(*)
from users join occupations on(occupation.id=users.occupation)
join ratings on(ratings.userid=users.userid)
where ratings=5
group by occupation.occupation;
```

```
2016-03-01 03:06:00,582 Stage-3 map = 0%, reduce = 0%
2016-03-01 03:06:06,649 Stage-3 map = 100%, reduce = 0%
2016-03-01 03:06:20,400 Stage-3 map = 100%, reduce = 100%
Ended Job = job_201603010226_0008
OK
K-12 student      5822
academic/educator 18603
artist 11702
clerical/admin 7825
college/grad student 30272
customer service 4655
doctor/health care 9269
executive/managerial 23044
farmer 489
homemaker 2555
lawyer 5069
other/not specified 28178
programmer 13670
retired 3839
sales/marketing 11315
scientist 5654
self-employed 9902
technician/engineer 16209
tradesman/craftsman 2315
unemployed 3179
writer 12744
Time taken: 124.077 seconds
hive>
```

Fig 1.3: Occupation based ratings

#### 3.3 Differential rating based on age

The following hive query determines the number of people who have rated 5 for the movies based on age:

```
hive> select users.age, count(*)
from ratings join users on(ratings.userid=users.userid)
where rating=5
GROUP BY users.age;
```

The outcome of the query is as shown below:

```
2016-03-01 03:43:30,936 Stage-2 map = 0%, reduce = 0%
2016-03-01 03:43:37,002 Stage-2 map = 100%, reduce = 0%
2016-03-01 03:43:51,131 Stage-2 map = 100%, reduce = 100%
Ended Job = job_201603010226_0028
OK
1      6802
18     40558
25     85730
35     44710
45     19142
50     18600
56     10768
Time taken: 75.08 seconds
```

Fig 1.4: Age based ratings

From the result obtained we can conclude that viewers around the age group 25years rate movies the highest (rate movies 5).

#### 3.4 Differential rating based on occupation and gender

The following query determines the number of people who have rated 5 for the movies based on occupation and gender.

```
hive> SELECT occupations.occupation, count(*)
from users join ratings on(ratings.userid=users.userid)
join occupations on(users.occupation=occupations.id)
where ratings=5
GROUP BY occupations.occupation, gender
```

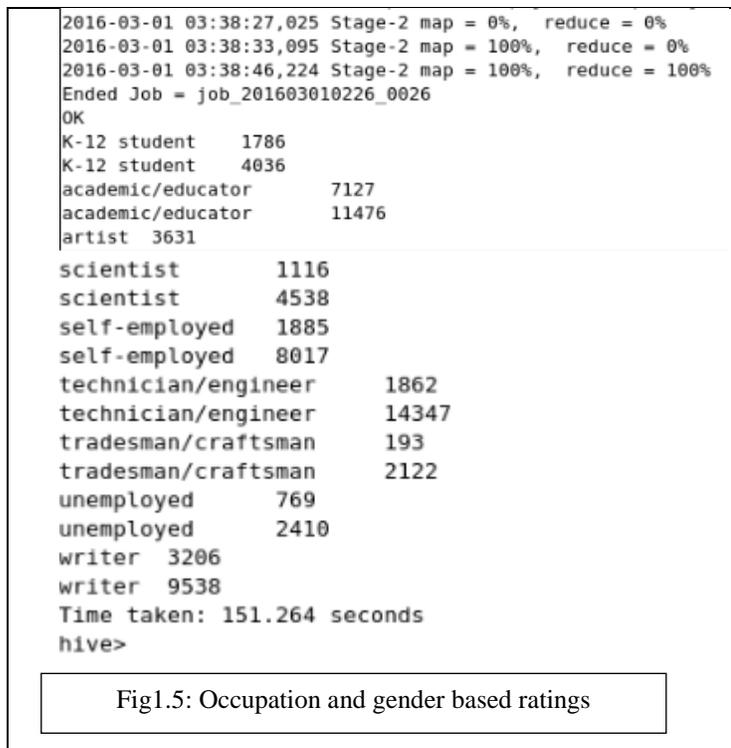


Fig1.5: Occupation and gender based ratings

From the above shown outcome we can conclude that each occupation's rating is mentioned twice with respect to the gender females rating followed by the males rating.

#### IV. ALGORITHMS

##### 4.1 Introduction

Machine Learning is the field of study that gives computers the ability to learn without being explicitly programmed. It grew out of Artificial Intelligence.

In literature we have many learning algorithms which comes under either supervised or unsupervised learning.

1) *Supervised Learning*: is a type of machine learning algorithm that uses a known training dataset to make predictions.

2) *Unsupervised learning*: is a type of machine learning algorithm used to draw inferences from datasets consisting of input data without labelled responses. [10][11]

Through supervised learning we can learn what makes the rating a certain value from the selected training dataset. In our

paper we will focus on linear regression which is a type of supervised learning.

##### 4.1.1 Linear Regression

Linear Regression is an approach for modelling the relationship between a scalar dependent variable  $y$  and one or more explanatory variables denoted by  $x$ . In subsequent section we have described the mathematical description of linear regression for our problem statement.

##### 4.1.2 Variables description

Let,  
 $m$ =number of training examples  
 $x$ =input variables  
 $y$ =output/target variables  
 $i$ =an index to training set

$(x^i, y^i)$  implies  $i^{\text{th}}$  training example

In the following equation,

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

where,

$h_{\theta}(x)$  is the numerically calculated values based on chosen parameters also termed as hypothesis function

- $\theta_i$  are parameters
- $\theta_0$  is zero condition
- $\theta_1$  is gradient

##### 4.1.3 Cost Function

Cost function lets us figure out how to fit the best straight line to our data by choosing values for  $\theta_1$ .

Based on the training set values for the parameters have to be generated so as to fit in the best possible straight line.

Values for the parameters are chosen such that  $h_{\theta}(x)$  is close to  $y$  for the training example. Basically, uses  $x$ s in training set with  $h_{\theta}(x)$  to give output which is as close as possible to the actual  $y$  value.  $h_{\theta}(x)$  can be considered as a "y imitator" - it tries to convert the  $x$  into  $y$ , and considering we already have  $y$  we can evaluate how well  $h_{\theta}(x)$  converges with  $y$ .

The cost function is given by:

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

where,

$J(\theta_0, \theta_1)$  is the cost function

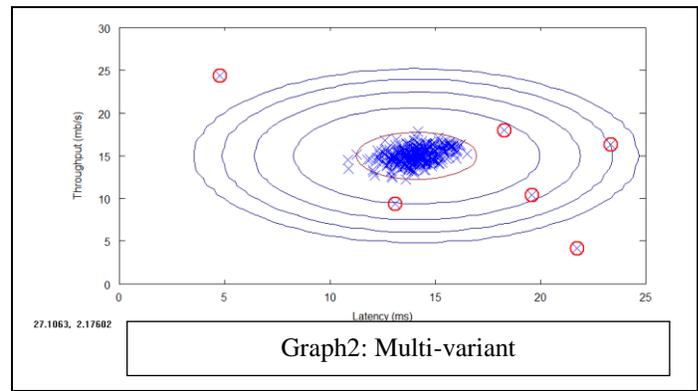
$y$  is the linear function of  $x$

$i$  varies from  $i=1$  to  $m$

$(h_{\theta}(x^{(1)}) - y^{(i)})^2$  implying trying to minimize squared difference between predicted ratings and actual ratings called in general as minimization problem.

- $1/2m$ 
  - $1/m$  – average determination
  - $1/2m$  the 2 doesn't change the constant value negligibly.
- Minimizing  $\theta_0, \theta_1$  means finding the values of  $\theta_0$  and  $\theta_1$ , which find on average the minimal deviation of  $x$  from  $y$  when the parameters are used in hypothesis function.

Above  $\theta_0, \theta_1$  is taken only for one input and one output. But our problem statement is having multiple attributes so we have more theta values for experiment.



The above graph is with respect to multiple variables as implemented in this paper.

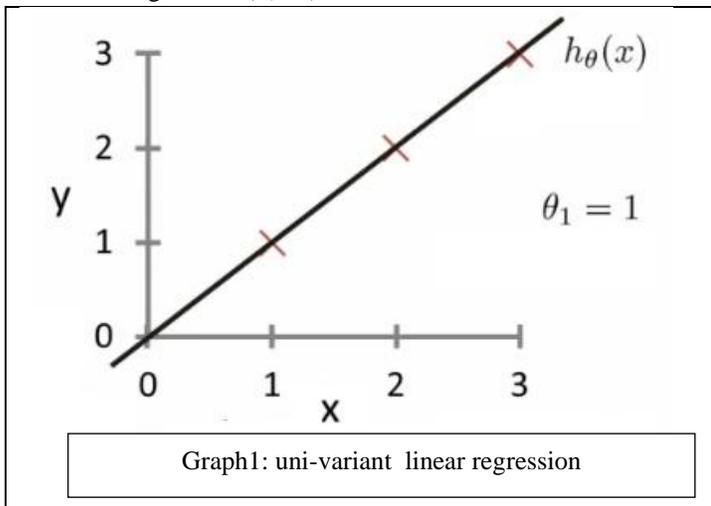
#### 4.2 Gradient descent

Gradient descent is an optimization method for minimizing an objective function that is written as a sum of differentiable functions. Used in machine learning for minimization of cost function.

Gradient Descent is all about:

We have  $J(\theta_0, \theta_1)$

We want to get  $\min J(\theta_0, \theta_1)$



As shown in the above graph  $y$  directly depends on the value of  $x$ . Repeated computation of the hypothesis function  $h_{\theta}(x)$  and applying minimization to the hence obtained cost function, most accurate graph for the prediction system can be determined.

#### 4.3 Implementation

In this paper, we make use of GNU Octave which is a high-level interpreted language intended for numerical computations. It provides capabilities for the numerical solution of linear and non-linear problems and for performing other numerical experiments. It also provides extensive graphic capabilities for data visualization and manipulation.

The Octave language is quite similar to Matlab so that programs are easily portable. [12]

Consider the following rating table for the movie dataset:

MOVIE	Ana	Joe	Mia	Matt	
RIO	5	5	0	0	★ ★ ★ ★ ★
CROODS	5	?	?	0	★ ★ ★ ★ ★
SMURFS	?	4	0	?	★ ★ ★ ★ ★
BRAVE	0	0	5	4	★ ★ ★ ★ ★
CARS	0	0	5	?	★ ★ ★ ★ ★

Fig 2.1: ratings table

Applying linear regression on the previously obtained and stored ratings we can predict the possible unknown ratings.

##### 4.3.1 Implementation steps

- Loading movie dataset: We will start by loading the movie ratings dataset to understand the structure of the data using `load('movies.m')`

```
ans =
<
[1,1] = Toy Story (1995)
[2,1] = GoldenEye (1995)
[3,1] = Four Rooms (1995)
[4,1] = Get Shorty (1995)
[5,1] = Copycat (1995)
[6,1] = Shanghai Triad (Yao a yao dao waipo qiao) (1995)
[7,1] = Twelve Monkeys (1995)
[8,1] = Babe (1995)
[9,1] = Dead Man Walking (1995)
[10,1] = Richard III (1995)
```

Fig 2.2: A part of loading process result

here Y is considered a matrix,containing rating (1-5) and R is a matrix,where R(i,j)=1 if and only if user j gives rating to movie i.From these matrices we can calculate statics like average rating using mean(Y(1,R(1,:)))

```
Loading movie ratings dataset.
Average rating for movie 1 (Toy Story): 3.878319 / 5
```

Fig 2.3: Computation of mean

We can visualize the ratings matrix by plotting it with imagec function as:

```
imagec(Y);
ylabel('Movies')
xlabel('Users')
```

- Collaborative Filtering: now we implement the collaborative filtering for cost function. The cost function is evaluated using:

J = cofiCostFunc([X(:) ; Theta(:)], Y, R, num\_users, num\_movies,num\_features, 0);

```
Cost at loaded parameters: 22.224604
(this value should be about 22.22)
```

Fig 2.4: J value

- Collaborative Filtering Gradient: Once our cost function matches with expected value as shown in Fig 2.4,Collaborative Filtering Gradient function should be implemented where in we check Gradients by running checkNNGradients and checkCostFunction. (without using regularization)

```
Checking Gradients (without regularization) ...
1.077681 1.077681
3.076620 3.076620
-0.063313 -0.063313
-0.016618 -0.016618
0.118171 0.118171
3.489001 3.489001
-0.071853 -0.071853
-2.784864 -2.784864
-0.119073 -0.119073
2.620050 2.620050
-0.010861 -0.010861
-1.454231 -1.454231
0.134090 0.134090
0.904398 0.904398
0.056093 0.056093
1.162020 1.162020
-0.848104 -0.848104
1.892171 1.892171
1.181381 1.181381
-0.360925 -0.360925
-0.499161 -0.499161
-3.816202 -3.816202
```

Fig 2.5: Gradients without Regularization

- Collaborative Filtering Gradient with Regularization: now we implement regularization for cost function for collaborative filtering this is done by adding the cost of regularization to the original cost computation. It is evaluated as follows:

J = cofiCostFunc([X(:) ; Theta(:)], Y, R, num\_users, num\_movies, num\_features, 1.5);

```
Checking Gradients (with regularization) ...
-1.419257 -1.419257
1.415663 1.415663
3.054467 3.054467
1.371528 1.371528
2.039074 2.039074
-0.812317 -0.812317
-10.125630 -10.125630
2.343506 2.343506
-0.155398 -0.155398
1.256389 1.256389
-15.459471 -15.459471
-4.177575 -4.177575
-2.067246 -2.067246
-3.298391 -3.298391
-1.248892 -1.248892
0.498188 0.498188
2.411015 2.411015
-1.081186 -1.081186
9.713643 9.713643
-1.268081 -1.268081
-4.106953 -4.106953
-7.291396 -7.291396
```

Fig: Gradients with Regularization

Collaborative Filtering Gradient Regularization:

As the cost matches as shown in Fig 2.7 we proceed to implement regularization for the gradient.

We check the gradient by running:

checkNNGradientscheckCostFunction(1.5);

- Enter ratings for a new user: We would train the collaborative filtering model first by adding ratings that correspond to new users, by using:  
movieList=loadMovieList();  
and we initialize the ratings for the new movies:

```

New user ratings:
Rated 4 for Toy Story (1995)
Rated 3 for Twelve Monkeys (1995)
Rated 5 for Usual Suspects, The (1995)
Rated 4 for Outbreak (1995)
Rated 5 for Shawshank Redemption, The (1994)
Rated 3 for While You Were Sleeping (1995)
Rated 5 for Forrest Gump (1994)
Rated 2 for Silence of the Lambs, The (1991)
Rated 4 for Alien (1979)
Rated 5 for Die Hard 2 (1990)
Rated 5 for Sphere (1998)
Program paused. Press enter to continue.

Training collaborative filtering...
Iteration 100 ! Cost: 7.205218e+004
Recommender system learning completed.
  
```

my\_ratings  
(u)=v;  
where u  
represents  
the movie  
ID and v  
represents  
rating(1-5).

Opportunities”, DASFAA Workshops 2013, LNCS 7827, pp. 1–15, 2013  
[ii] SurajitMohanty, KedarNathRout, ShekhareshBarik, Sameer Kumar Das, ”A Study on Evolution of Data in Traditional RDBMS to Big Data Analytics”, International Journal of Advanced Research in Computer and Communication Engineering ISSN 2278-1021; Vol 4, pp.230-232, October 2015  
[iii] Barkha Jain, Manish Km. Kakhani, “Query Optimization in Hadoop Distributed File System”, IJCSIS 2015, 9915; Volume 2,

Fig 2.8 training collaborative filtering

- Learning movie ratings and recommendations:

Now the collaborative filtering model similarly and complete the recommender system leaning as shown in Fig 2.8.

After obtaining the trained model, now recommendations can be computed using prediction matrix. And hence the following result for recommendation based on original is obtained.

```

Top recommendations for you:
Predicting rating 8.5 for movie Titanic (1997)
Predicting rating 8.5 for movie Star Wars (1977)
Predicting rating 8.3 for movie Shawshank Redemption, The (1994)
Predicting rating 8.3 for movie Schindler's List (1993)
Predicting rating 8.2 for movie Raiders of the Lost Ark (1981)
Predicting rating 8.2 for movie Good Will Hunting (1997)
Predicting rating 8.1 for movie Usual Suspects, The (1995)
Predicting rating 8.1 for movie Goodfellas, The (1992)
Predicting rating 8.0 for movie Braveheart (1995)
Predicting rating 8.0 for movie Empire Strikes Back, The (1980)

Original ratings provided:
Rated 4 for Toy Story (1995)
Rated 3 for Twelve Monkeys (1995)
Rated 5 for Usual Suspects, The (1995)
Rated 4 for Outbreak (1995)
Rated 5 for Shawshank Redemption, The (1994)
Rated 3 for While You Were Sleeping (1995)
Rated 5 for Forrest Gump (1994)
Rated 2 for Silence of the Lambs, The (1991)
Rated 4 for Alien (1979)
Rated 5 for Die Hard 2 (1990)
Rated 5 for Sphere (1998)
  
```

Fig 2.9 recommendation

[iv][https://hadoop.apache.org/docs/r1.2.1/mapred\\_tutorial.html](https://hadoop.apache.org/docs/r1.2.1/mapred_tutorial.html)  
[v]<https://cwiki.apache.org/confluence/display/Hive/Home;jsessionid=986ED81CB5EB4E18AB21A6BDDE4610EA>  
[vi]SurekhaSharadMuzumdar, JharnaMajumdar, “Big Data Analytics Framework using Machine Learning on Multiple Datasets” IJSR ISSN:2319-7064 Volume 4 Issue 8, August 2015 pp.414-418  
[vii] <http://grouplens.org/datasets/movielens/>  
[viii] <http://www.recsyswiki.com/wiki/MovieLens>  
[ix]<http://files.grouplens.org/datasets/movielens/ml-latest-README.html>  
[x]<http://machinelearningmastery.com/a-tour-of-machine-learning-algorithms/>  
[xi] [https://en.wikipedia.org/wiki/Machine\\_learning](https://en.wikipedia.org/wiki/Machine_learning)  
[xii] <https://www.gnu.org/software/octave/>  
[xiii] [http://www.iosrjen.org/Papers/vol2\\_issue8%20\(part-1\)/K0287882.pdf](http://www.iosrjen.org/Papers/vol2_issue8%20(part-1)/K0287882.pdf)  
[xiv][https://www.researchgate.net/publication/251935098\\_Optimization\\_of\\_IP\\_routing\\_with\\_content\\_delivery\\_network](https://www.researchgate.net/publication/251935098_Optimization_of_IP_routing_with_content_delivery_network)[xv][http://shodhganga.inflibnet.ac.in/bitstream/10603/10203/16/16\\_publications.pdf](http://shodhganga.inflibnet.ac.in/bitstream/10603/10203/16/16_publications.pdf) .

## V. CONCLUSION

MACHINE LEARNING IS A METHOD OF DATA ANALYSIS THAT AUTOMATES ANALYTICAL MODEL BUILDING. USING ALGORITHMS THAT ITERATIVELY LEARN FROM DATA, MACHINE LEARNING ALLOWS COMPUTERS TO FIND HIDDEN INSIGHTS WITHOUT BEING EXPLICITLY PROGRAMMED WHERE TO LOOK. MACHINE LEARNING SOLVES PROBLEM THAT CANNOT BE SOLVED BY OTHER NUMERICAL MEANS. WE HAVE SEEN, IN THIS PAPER BY APPLYING LINEAR REGRESSION, WE CAN PREDICT THE RATINGS FOR FUTURE MOVIES. THIS WAY MACHINE LEARNING HELPS IN IMPROVING PERFORMANCE FOR ANY SUCH APPLICATIONS. HERE WE HAVE MADE USE OF HUGE DATASET ON HADOOP PLATFORM WHICH HELP US PROCESS THESE DATASETS AT A FASTER RATE USING MAPREDUCE PROCESSING WHICH IS NOT POSSIBLE BY OTHER TRADITIONAL PROCESSING SYSTEM. HENCE HADOOP AND MACHINE LEARNING TOGETHER CAN BE USED TO SOLVE A VARIETY OF LEARNING PROBLEMS MORE EFFICIENTLY.

## REFERENCES

- [i] DunrenChe, MejdI Safran, and Zhiyong Peng, “From Big Data to Big Data Mining: Challenges, Issues, and