# Recommender Systems

Tulasi 7/6/2017

## DATA643 Project 5

by Tulasi Ramarao

## Description

In this project, a recommender was implemented on a distributed system. The Performance of this recommender was then compared with the recommender that was created on Apache Spark using ALS.

#### Dataset

The dataset from MovieLens(Ref#:5) listed under Recommended for Education and Research was used. This dataset contains 100,000 ratings for 9,000 movies by 700 users. The ratings.csv and movies.csv were used to build the recommendation system.

#### **Installations:**

The first attempt was to build this recommender on databricks.com that runs on Amazon Web Services (AWS)[Fig.1]. The Spark connection was successfully established and the data was loaded successfully. The queries worked until the calls called sdf\_copy\_to() and sdf\_import() were made. The error messages that appeared were not easy to debug even with the help of google search.

So, this project on databricks was abandoned due to time restrictions and an alternative approach was chosen, which was to install the recommender on Apache Spark on a single node (Ref#:4).

Sparklyr was recommended by Ref# 1 because it is an easier environment to work with, when compared to SparkR. The data manipulation in Sparklyr uses the same verbage as dplyr, so the learning curve is said to be easier for R programmers. Also, Sparkly is faster than R and help documentation are easily available in R (?function\_name).

The returned spark connection(sc) below provides a remote dplyr data source to the Spark cluster

```
## * Using Spark: 2.1.0
```

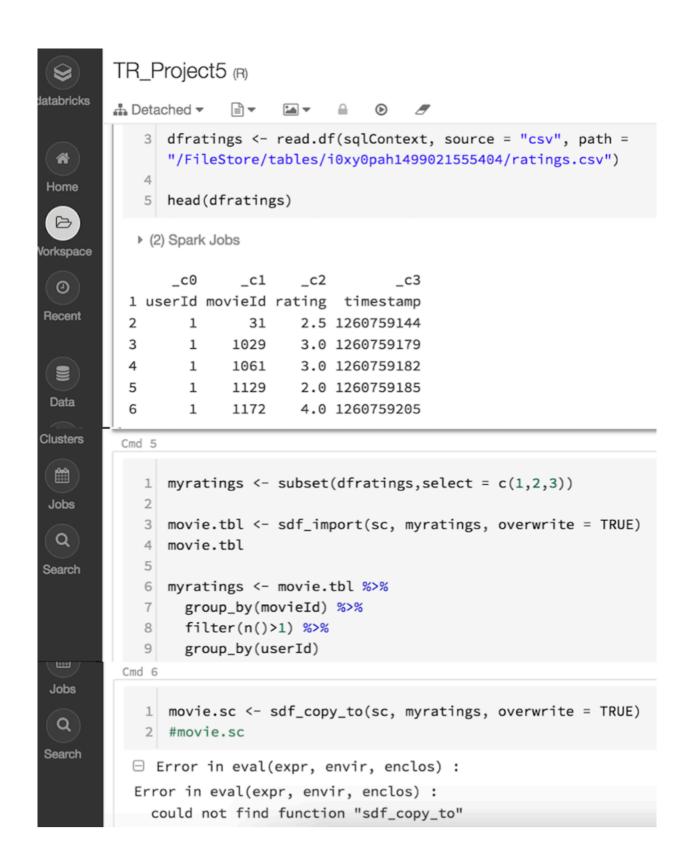


Fig. 1: Screenshot from a workspace on Databricks

```
set.seed(3445) # to keep #s from the results the same
# Set the working directory
setwd("/Users/tulasiramarao/Documents/Tulasi/CUNYProjects/DATA643/RPrograms")
tic()
# load data from local drive
dfratings <- read.csv("MovieRatingsData/ml-latest-small/ratings10M.csv", header = TRUE, sep =",",
                     stringsAsFactors = FALSE)
colnames(dfratings)
## [1] "X1..122..5..838985046"
exectime <- toc()</pre>
## 61.646 sec elapsed
exectimeALS <- exectime$toc - exectime$tic</pre>
# Time the 10 million dataset load
tic()
# load data from local drive
dfmovies <- read.csv("MovieRatingsData/ml-latest-small/movies10M.csv", header = TRUE, sep =",",
                     stringsAsFactors = FALSE)
## Warning in scan(file = file, what = what, sep = sep, quote = quote, dec =
## dec, : EOF within quoted string
exectime <- toc()</pre>
## 0.09 sec elapsed
exectimeALS <- exectime$toc - exectime$tic</pre>
Ten million rows for each csv files loaded in less than a minute. With R, the program hung when loading
these huge datasets.
Going back to 1M for simplicty
dfratings <- read.csv("MovieRatingsData/ml-latest-small/ratings.csv", header = TRUE, sep =",",
                     stringsAsFactors = FALSE)
colnames(dfratings)
## [1] "userId"
                    "movieId"
                                "rating"
                                             "timestamp"
dfmovies <- read.csv("MovieRatingsData/ml-latest-small/movies.csv", header = TRUE, sep =",",
                     stringsAsFactors = FALSE)
colnames(dfmovies)
## [1] "movieId" "title"
                            "genres"
Now, the two datasets are merged - to get a movie name for each movieId.
combinedData <- merge(dfratings,dfmovies, by=c("movieId"))</pre>
colnames(combinedData)
                                             "timestamp" "title"
## [1] "movieId"
                    "userId"
                                "rating"
                                                                       "genres"
#Select the relevant columns - skip timestamp and genres, movieID
dfratings <- subset(combinedData, select = c(1,2,3,5))
colnames(dfratings)
```

```
## [1] "movieId" "userId" "title"
myratings <- as.data.frame(dfratings)
kable(summary(dfratings))</pre>
```

movieId	userId	rating	title
Min. : 1	Min. : 1	Min. :0.500	Length:100004
1st Qu.: 1028	1st Qu.:182	1st Qu.:3.000	Class :character
Median: 2406	Median $:367$	Median $:4.000$	Mode :character
Mean: 12549	Mean $:347$	Mean $:3.544$	NA
3rd Qu.: 5418	3rd Qu.:520	3rd Qu.:4.000	NA
Max. :163949	Max. :671	Max. $:5.000$	NA

The ratings run from 0.5 to 5.

Copy myratings into Spark and return an R object wrapping the copied object (a spark dataframe) colnames (myratings)

```
## [1] "movieId" "userId" "rating" "title"

#which_train <- sample(x=c(TRUE,FALSE), size = nrow(myratings), #replace=TRUE,prob=c(0.8,0.2))
##trainData <- myratings[which_train,]
#testData <- myratings[!which_train,]
#head(trainData)

## seelct 3 columns here
# copy the table to Spark
movie.tbl <- sdf_copy_to(sc, myratings, overwrite = TRUE)
#movie.tbl
#movie.tbl <- sdf_copy_to(sc, trainData, overwrite = TRUE)

mv <- movie.tbl %>% filter(userId == 428)
```

## kable(head(mv,2))

movieId	userId	rating	title
1	428	5	Toy Story (1995)
2	428	3	Jumanji (1995)
SVD was at	tempted t	o impleme	nt, but due to errors, switched to ALS instead.

This SVD command gave a sparkException for the Java.lang.OutofMemoryError in heap space. "Error: org.apache.spark.SparkException: Job aborted due to stage failure: Task 0 in stage 18.0 failed 1 times, most recent failure: Lost task 0.0 in stage 18.0 (TID 34, localhost, executor driver): java.lang.OutOfMemoryError: Java heap space at java.nio.HeapByteBuffer.(HeapByteBuffer.java:57)"

So chose Alternating Least Squares(ALS) to perform matrix factorization on a Spark Dataframe.

Create an ALS model:

```
##
                    Length Class
                                       Mode
## item.factors
                           data.frame list
                    11
## user.factors
                    11
                            data.frame list
## data
                     2
                            spark_jobj environment
## ml.options
                     6
                           ml options list
## model.parameters 2
                            -none-
                                       list
## .call
                     6
                            -none-
                                       call
                            spark_jobj environment
## .model
                     2
```

It took  $\sim 2$  seconds in Spark. In R, UBCF needed  $\sim 9$  seconds and IBCF needed  $\sim 89$  seconds to evaluate the models.

When comparing SVD, the performance of the ALS model performed 5 times faster than R( 2.8 seconds versus 11.076 seconds)

Using documentation from

https://spark.apache.org/docs/latest/ml-collaborative-filtering.html, the following predictions were created.

```
predictions <- MLSmodel$.model %>%
  invoke("transform", spark_dataframe(movie.tbl)) %>%
  collect()

#predictions[predictions$userId == 8,]
#dim(predictions)
```

```
tic()
#Feed the matrix form of data
svd.form <- svd(user_item,nu=3,nv=3)
exectime <- toc()
## 11.076 sec elapsed</pre>
```

Fig 2: Execution time for SVD calculated in Project 3

Figure 2:

```
pred_RMSE <- sqrt(mean(with(predictions, prediction-rating)^2))
pred_RMSE</pre>
```

## [1] 0.6364164

The RMSE for the ALS model is calculated and it is: 0.6093105

Comparing the ALS RMSE to the RMSE calculated from R irlba/svd in project 3, it can be seen that there is a great improvement in RMSE for ALS (0.6093105 vs. 3.43567) [Fig.2]

The executime time is very fast when running the recommender in Spark. R recommender took several minutes to run.

```
# exclude missing values NA from analysis with na.rm = True
(RMSE.svd <- sqrt(mean((predict.svd - user_item)^2, na.rm=T)))
## [1] 3.43567</pre>
```

Fig. 2: RMSE for SVD from Project 3

Figure 3:

ALS is a method where the entire loss function is minimized at once, changing half the parameters at a time[Ref#:6]. So, half the parameters are fixed and the other half is recomputed and the process is repeated. ALS uncovers latent features.

Next, the matrix for the predictions from ALS is calculated. First two data frames are created - one for userid and one for movie ID  $\rm w/title$ .

```
usernames <- myratings %>%
  distinct(userId) %>%
  arrange(userId)

userratings <- myratings %>%
  distinct(userId,rating) %>%
  arrange(userId,rating)
```

```
movienames <- myratings %>%
    group_by(title) %>%
    distinct(title)

u.df <- MLSmodel$user.factors[,-1]
m.df <- MLSmodel$item.factors[,-1]
u.matrix <- as.matrix(u.df)
m.matrix <- as.matrix(m.df)

# now predict
predict.ALS <- u.matrix %*% t(m.matrix)

rownames(predict.ALS) = usernames$userId
colnames(predict.ALS) = movienames$title</pre>
```

	Toy Story (1995)	Jumanji (1995)	Heat (1995)	GoldenEye (1995)	Leaving Las Vegas (1995)
4	4.611692	3.918420	4.250472	4.081249	4.050529
8	3.870543	3.187401	3.896780	3.353761	3.933417
15	2.831116	2.071788	4.040609	2.408394	2.906512
17	3.196486	2.395453	3.989754	2.751173	4.027206
19	3.608474	3.089160	3.551563	3.118113	3.430007

# Performance testing:

```
movieStats <- myratings %>%
  group_by(userId) %>%
  summarise(
    meanreview = mean(rating)
)

kable(head(movieStats,3))
```

userId	meanreview
4	4.714286
8	4.214286
15	3.263514

```
predict.ALS.df <- as.data.frame(predict.ALS)
predict.ALS.df[5,1]

## [1] 3.608474
kable(predict.ALS.df[1:6, 1:6])</pre>
```

	Toy Story (1995)	Jumanji (1995)	Heat (1995)	GoldenEye (1995)	Leaving Las Vegas (1995)	Twelve Monkeys
4	4.611692	3.918420	4.250472	4.081249	4.050529	
8	3.870543	3.187401	3.896780	3.353761	3.933417	
15	2.831116	2.071788	4.040609	2.408394	2.906512	

	Toy Story (1995)	Jumanji (1995)	Heat (1995)	GoldenEye (1995)	Leaving Las Vegas (1995)	Twelve Monkeys
17	3.196486	2.395453	3.989754	2.751173	4.027206	
19	3.608474	3.089160	3.551563	3.118113	3.430007	
21	3.832396	2.952647	3.314175	2.826197	3.357965	

```
# remove dates and numbers from the movie name ( easy to query without dates and numbers)
colnames(predict.ALS.df) <- sub("\\([0-9][0-9][0-9][0-9]\\)", "", colnames(predict.ALS.df) )
predict.ALS.df["21", "Jumanji "]  # Keep - an example
## [1] 2.952647
predict.ALS.df["15", "Toy Story "]</pre>
```

## [1] 2.831117

A confusion matrix was created to visualize the performance of the algorithm. Choosing a threshold value of 3, the predictions and the ratings are partitioned.

```
# choose threshold value of 3
confusionMat <- predictions %>%
  mutate(actual = if_else(rating >= 3, 1, 0),
  predicted = if_else(prediction >= 3, 1, 0))

confusionMat <- subset(confusionMat,select = c(6,7))
#colnames(confusionMat)

confusionTable<- table(confusionMat)
confusionTable</pre>
```

```
## predicted
## actual 0 1
## 0 874 446
## 1 638 9257
```

The diagonal represents the cases where the ratings are correctly predicted. Now the precision, recall and Fscores are calculated.

Precision is the ability of the classifier to not label as positive when its actually negative. In the formula Precision : tp/(tp + fp) to is the number of true positives and fp is the number of true negatives. Recall is the ability of the classifier to find all the positive samples.

FScore is the weighted harmonic mean of the precision and recall, where the best score is at 1 and worse at 0.

```
# Precision: tp/(tp+fp):
(precision <- confusionTable[1,1]/sum(confusionTable[1,1:2]))

## [1] 0.6621212

# Recall: tp/(tp + fn):
(recall <- confusionTable[1,1]/sum(confusionTable[1:2,1]))

## [1] 0.5780423

# F-Score: 2 * precision * recall / (precision + recall):
(F_Score <- 2 * precision * recall / (precision + recall))</pre>
```

## [1] 0.6172316

To improve the FScore, a different threshold value (2) is chosen.

```
confusionMatLow <- predictions %>%
  mutate(actual = if_else(rating >= 2, 1, 0),
  predicted = if_else(prediction >= 2, 1, 0))
confusionMatLow <- subset(confusionMatLow,select = c(6,7))</pre>
colnames(confusionMatLow)
## [1] "actual"
                    "predicted"
confusionTableLow<- table(confusionMatLow)</pre>
confusionTableLow
         predicted
##
## actual
              0
##
        0
              97
                   300
        1
              28 10790
# Precision: tp/(tp+fp):
(precision L \gets confusion Table Low[1,1]/sum(confusion Table Low[1,1:2]))\\
## [1] 0.2443325
# Recall: tp/(tp + fn):
(recallL <- confusionTableLow[1,1]/sum(confusionTableLow[1:2,1]))</pre>
## [1] 0.776
# F-Score: 2 * precision * recall /(precision + recall):
(F_ScoreL <- 2 * precisionL * recallL / (precisionL + recallL))</pre>
## [1] 0.3716475
The F-Score got worse, so the accuracy of this recommender is not too great.
Now disconnect from spark gracefully.
# disconnect from Spark
spark_disconnect(sc)
The results are tabulated and it can be seen that Spark is much faster, eventhough the RMSE didn't improve
much.
ubcf \leftarrow c('0.938', '9')
ibcf <-c('1.067', '85')
svd \leftarrow c('3.43','9')
als <-c('3.12','2')
myresults.df <- data.frame(ubcf,ibcf,svd,als)</pre>
str(myresults.df)
                     2 obs. of 4 variables:
## 'data.frame':
## $ ubcf: Factor w/ 2 levels "0.938", "9": 1 2
## $ ibcf: Factor w/ 2 levels "1.067", "85": 1 2
## $ svd : Factor w/ 2 levels "3.43", "9": 1 2
## $ als : Factor w/ 2 levels "2", "3.12": 2 1
colnames(myresults.df) <- c("UBCF", "IBCF", "SVD", "ALS")</pre>
rownames(myresults.df) <- c("RMSE-Distance", "Executime Time(secs)")</pre>
kable(myresults.df, type = "html",caption="Results - R vs Spark")
```

Table 6: Results - R vs Spark

	UBCF	IBCF	SVD	ALS
RMSE-Distance	0.938	1.067	3.43	3.12
Executime Time(secs)	9	85	9	2

```
x <- c("R is a very familiar language", "has negligible learning curve", "Most datasets fit in memory", "
y <- c("Spark is unfamiliar to programmers", "has a teep learning curve", "is scalable to handle big data

tradeoffs.df <- data.frame(x,y)
str(tradeoffs.df)

## 'data.frame': 7 obs. of 2 variables:
## $ x: Factor w/ 7 levels "Abundance of help in google search",..: 5 3 4 1 7 2 6

## $ y: Factor w/ 7 levels "Accuracy is not too good",..: 7 2 3 4 5 1 6

rownames(tradeoffs.df) <- c("1","2","3","4","5","6","7")
colnames(tradeoffs.df) <- c("R","Scale")

kable(tradeoffs.df, type = "html",caption="Comparison: R vs. Spark")</pre>
```

Table 7: Comparison: R vs. Spark

R	Scale
R is a very familiar language	Spark is u
has negligible learning curve	has a teep
Most datasets fit in memory	is scalable
Abundance of help in google search	Not much
Takes a lot of time to process large samples in the dataset. Hence performance of the entire process is affected	runs on m
Accuracy in prediction is good	Accuracy i
Slow executime time for large datasets	Spark has

#### Conclusion:

R is single threaded and comparing its performance with Spark is not a fair one. R does better with smaller samples of dataset. With its enormous statistical computation and visual libraries, R has a lot to offer. And its easier to test and verify results. But in the past few projects, the biggest frustration was with the time it took for the R program to run for 1M+ datasets. It was several hours each time to run some of the commands. However, it can be argued that using a small sample of dataset is not representative of the whole population in real world scenarios. So R may compromise accuracy. Massive datasets are ideal for Spark mllib and with R tools in Spark, it can be used to explore datasets on distributed systems. In the end, increased productivity is a huge deal for businesses with massive datasets and utilizing Spark will help in that area. So, if the dataset is below 1M rows, then R is a better choice with its abundance of tools.

This project was mostly comprised of installation of distributed systems on a single node, working on Databricks on top of Amazon AWS and researching several errors when making the recommender to work on Spark. The mllib in Spark supports collaborative filtering, where users and the movies are described by a small set of latent factors (used to predict the missing entries). Spark's mllib uses ALS with its parameter, lambda, scaled in solving each least squares problem and (Ref#.8) so lambda is less dependent on the dataset scale. So a smilar performance can be expected in a large dataset. It is worth the learning time

to get comfortable with Spark in a distrubuted environment, so that the data science skills fit the business requirements.

#### References:

Ref#1: http://www.lyzander.com/r/spark/2016/11/26/spark\_and\_r

Ref#2: https://blog.rstudio.org/author/javierrstudiocom/

Ref#3: http://spark.rstudio.com/h2o.html

Ref#4: https://github.com/rstudio/sparklyr/blob/master/README.md

Ref#5: http://grouplens.org/datasets/movielens Ref#6: https://www.quora.com/What-is-the-Alternating-Least-Squares-n

Ref#7: http://spark.rstudio.com/h2o.html

 $Ref\#8: \ https://spark.apache.org/docs/latest/mllib-collaborative-filtering.html \ Ref\#9: \ https://github.com/rstudio/sparklyr/blob/master/man/ml_als_factorization.Rd \ Ref\#10: \ https://rdrr.io/cran/sparklyr/latest/mllib-collaborative-filtering.html \ Ref\#9: \ https://github.com/rstudio/sparklyr/blob/master/man/ml_als_factorization.Rd \ Ref\#10: \ https://rdrr.io/cran/sparklyr/latest/mllib-collaborative-filtering.html \ Ref\#9: \ https://github.com/rstudio/sparklyr/blob/master/man/ml_als_factorization.Rd \ Ref\#10: \ https://rdrr.io/cran/sparklyr/latest/mllib-collaborative-filtering.html \ Ref\#9: \ https://rdrr.io/cran/sparklyr/latest/mllib-collaborative-filtering.html \ Ref\#9: \ https://rdrr.io/cran/sparklyr/latest/mllib-collaborative-filtering.html \ Ref\#9: \ https://rdrr.io/cran/sparklyr/latest/mllib-collaborative-filtering.html \ Ref\#9: \ https://rdrr.io/cran/sparklyr/latest/mllib-collaborative-filtering.html \ Ref#10: \ https://rdrr.io/cran/sparklyr/latest/mllib-collaborative-filtering.html \ Ref#9: \ https://rdrr.io/cran/sparkl$ 

man/ml als factorization.html