

# DATA-643 Assignment - 03

*Mohamed Elmoudni*

*Shazia Khan*

*Senthil Dhanapal*

## Contents

1 Requirements: . . . . .	2
2.Introduction . . . . .	2
<b>Toy Dataset</b>	<b>2</b>
Testing . . . . .	4
<b>MovieLense dataset</b>	<b>5</b>
2 Data Load . . . . .	5
<b>3 Data Exploration</b>	<b>5</b>
3.2 Data visualization . . . . .	7
<b>4 Model Evaluation</b>	<b>7</b>

# 1 Requirements:

## Project 03

The goal of this assignment is give you practice working with Singular Value Decomposition. Your task is implement a matrix factorization method-such as singular value decomposition (SVD) or Alternating Least Squares (ALS)-in the context of a recommender system. You may approach this in a large number of ways. You are welcome to start with an existing recommender system written by yourself or someone else (always citing your sources, so that you can be graded on what you added, not what you found). Here is one example. Suppose you start with (or create) a collaborative filtering system against (a subset of) the MovieLens database or our toy dataset. You could create a content-based system, where you populate your item profiles by pulling text information for specific movies from a source like imdb, applying text processing techniques (like TF-IDF), then using SVD and topic modeling to create a set of features derived from the text. An extra intermediate step could be to take text that was pre-classified, e.g. “fighting” or “singing” and build out two “explainable” features. SVD builds features that may or may not map neatly to movie genres or news topics. You may work in a small group (2 or 3 people) on this assignment.

## 2.Introduction

### Recommendation Using SVD

The goal of CF-based recommendation algorithms is to suggest new products or to predict the utility of a product for a customer, based on the customer's previous behavior and other similar customers 'opinions. However, these systems have some problems like sparsity, scalability, and synonymy. The weakness of CF algorithms for large, sparse databases led the researchers to alternative ways. In order to remove noise data from a large and sparse database, some dimensionality reduction techniques are proposed. Latent Semantic Indexing (LSI), which is a dimensionality reduction technique that used in information retrieval (IR), is a widely used technique to reduce the dimensionality of user-item ratings matrix. LSI, which uses singular value decomposition (SVD) as its underlying dimension reduction algorithm, maps nicely into the collaborative filtering recommender algorithm challenge

## Toy Dataset

First let's show an example of dimension reduction, consider the rating matrix A

```
User1 <- c(1, 5, 0, 5, 4)
User2<- c(5, 4, 4, 3, 2)
User3<- c(0, 4, 0, 0, 5)
User4<- c(4, 4, 1, 4, 0)
User5<- c(0, 4, 3, 5, 0)
User6<- c(2, 4, 3, 5, 3 )

A<- rbind(User1,User2, User3, User4, User5, User6)
colnames(A)<- c("Movie1","Movie2", "Movie3", "Movie4", "Movie5")
kable(A, caption = 'Original Matrix')
```

Table 1: Original Matrix

	Movie1	Movie2	Movie3	Movie4	Movie5
User1	1	5	0	5	4
User2	5	4	4	3	2
User3	0	4	0	0	5
User4	4	4	1	4	0
User5	0	4	3	5	0
User6	2	4	3	5	3

Applying SVD to matrix A:

```
A.svd <- svd(A)
U<- matrix(A.svd$u,6,6)
kable(U, caption = 'SVD, U Matrix')
```

Table 2: SVD, U Matrix

-0.4599561	0.3964320	0.3000847	-0.4318831	0.3225994	-0.4599561
-0.4608062	-0.3069318	-0.6474524	0.2836982	0.0197989	-0.4608062
-0.2487966	0.7546451	-0.2790170	0.1562937	-0.4645854	-0.2487966
-0.3834392	-0.3454649	-0.1318702	-0.6833112	-0.3218441	-0.3834392
-0.3762169	-0.2443748	0.6212002	0.3802641	-0.5020945	-0.3762169
-0.4750090	-0.0089592	0.0981056	0.3115279	0.5692235	-0.4750090

```
S<- diag(A.svd$d, 6,5)
kable(S, caption = 'Diagonal Matrix, S')
```

Table 3: Diagonal Matrix, S

16.46644	0.000000	0.000000	0.000000	0.000000
0.00000	6.210013	0.000000	0.000000	0.000000
0.00000	0.000000	4.399085	0.000000	0.000000
0.00000	0.000000	0.000000	2.903364	0.000000
0.00000	0.000000	0.000000	0.000000	1.584456
0.00000	0.000000	0.000000	0.000000	0.000000

```
V_t<- t(A.svd$v)
#V_t
kable(V_t, caption = 'SVD, V transpose Matrix')
```

Table 4: SVD, V transpose Matrix

-0.3186943	-0.6119629	-0.2903081	-0.5752356	-0.3297887
-0.4086956	0.2218704	-0.3757148	-0.2555815	0.7597749
-0.7429835	0.0327995	-0.1281545	0.5971956	-0.2717235

-0.3869928	-0.1258970	0.8703205	-0.2006657	0.1914731
0.1720872	-0.7478965	-0.0260415	0.4548327	0.4510941

---

Matrix U (6x6), matrix S (6x5), and matrix V (5x5) are calculated. Now, we will collapse this matrix from a (6x5) space into a 2-Dimensional one. To do this, we simply take the first two columns of U, S and V. The end result

```
a<- U[, c(1,2)]
b<- S[, c(1,2)]
c<- V_t[, c(2,2)]
```

For a recommender system based on SVD, here is one very simple strategy: find the most similar user using the 2-Dimensional matrixes above with one of the similarity calculation algorithms and compare his/her items against that of the new user; take the items that the similar user has rated and the new user has not and return them for the new user. Similar to this, for a new item, find the most similar item using the 2-Dimensional matrixes above with one of the similarity calculation algorithms and compare the users rated similar item against the new item; take the users that rate similar item but not the new item and return the ratings for the new item.

## Testing

### Selecting k = 2

```
svd <- svd(A,nu=2,nv=2)
S <- diag(svd$d[1:2])
kable(svd$u %*% S %*% t(svd$v), caption = 'SVD for k =2')
```

Table 5: SVD for k =2

1.4075937	5.181121	1.2737945	3.727541	4.3682179
3.1971956	4.220580	2.9189424	4.851947	1.0542144
-0.6096677	3.546850	-0.5714007	1.158876	4.9116525
2.8889890	3.387872	2.6390068	4.180278	0.4522691
2.5945209	3.454378	2.3686189	3.951421	0.8900136
2.5154726	4.774251	2.2916088	4.513545	2.5372394

New user with rating: 0,0,2,0,4

```
b <-matrix(c(0,0,2,0,4,1),byrow=T, ncol=6)
S1 <- diag(1/svd$d[1:2])
b %*% svd$u %*% S1
```

```
##           [,1]      [,2]
## [1,] -0.1504557 0.08419174
```

Selecting  $k = 3$

```
svd <- svd(A,nu=3,nv=3)
S <- diag(svd$d[1:3])
kable(svd$u %*% S %*% t(svd$v), caption = 'SVD for k =3')
```

Table 6: SVD for  $k = 3$

0.4267829	5.224419	1.1046180	4.5158972	4.0095163
5.3133595	4.127161	3.2839519	3.1510159	1.8281366
0.3022846	3.506591	-0.4141014	0.4258669	5.2451711
3.3199999	3.368845	2.7133503	3.8338400	0.6098981
0.5641610	3.544009	2.0184095	5.5833848	0.1474716
2.1948197	4.788406	2.2363005	4.7712798	2.4199704

New user with rating: 0,0,2,0,4

```
b <-matrix(c(0,0,2,0,4,1),byrow=T, ncol=6)
S1 <- diag(1/svd$d[1:3])
b %*% svd$u %*% S1
```

```
##           [,1]      [,2]      [,3]
## [1,] -0.1504557 0.08419174 0.460294
```

## MovieLens dataset

Now that we have a basic understanding on how to apply SVD in our toy data, let's apply it in the to bigger set such the Movie Lens dataset.

We will be using the MovieLens dataset that comes with the recommenderlab package.

## 2 Data Load

## 3 Data Exploration

In this phase we will perform data analysis and summary of the main characteristics of a dataset. We will be describing the data by means of statistical and visualization techniques

Table 7: Data Summary

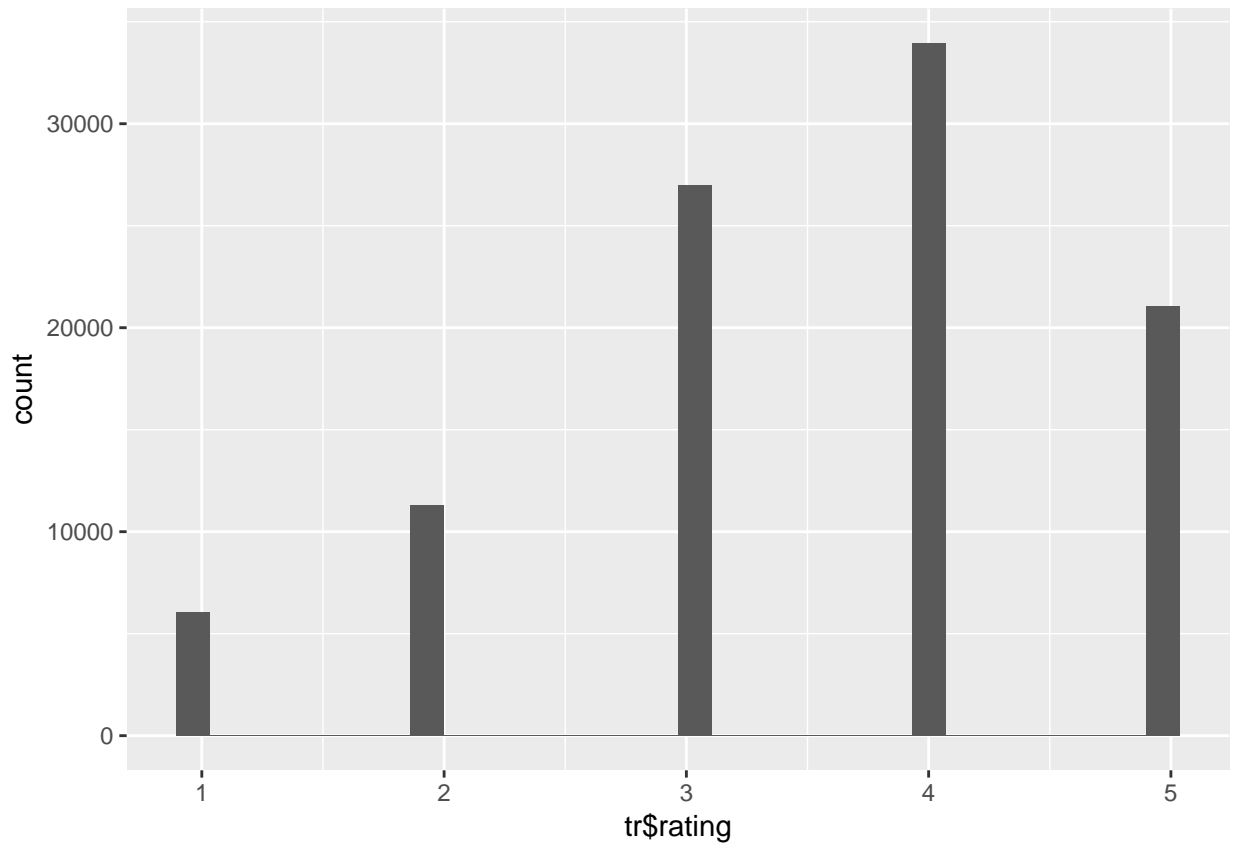
user_id	item_id	rating
1	101 Dalmatians (1996)	2
1	12 Angry Men (1957)	5
1	20,000 Leagues Under the Sea (1954)	3

user_id	item_id	rating
1	2001: A Space Odyssey (1968)	4
1	Abyss, The (1989)	3
1	Ace Ventura: Pet Detective (1994)	3

Table 8: Ratings Frequency Table

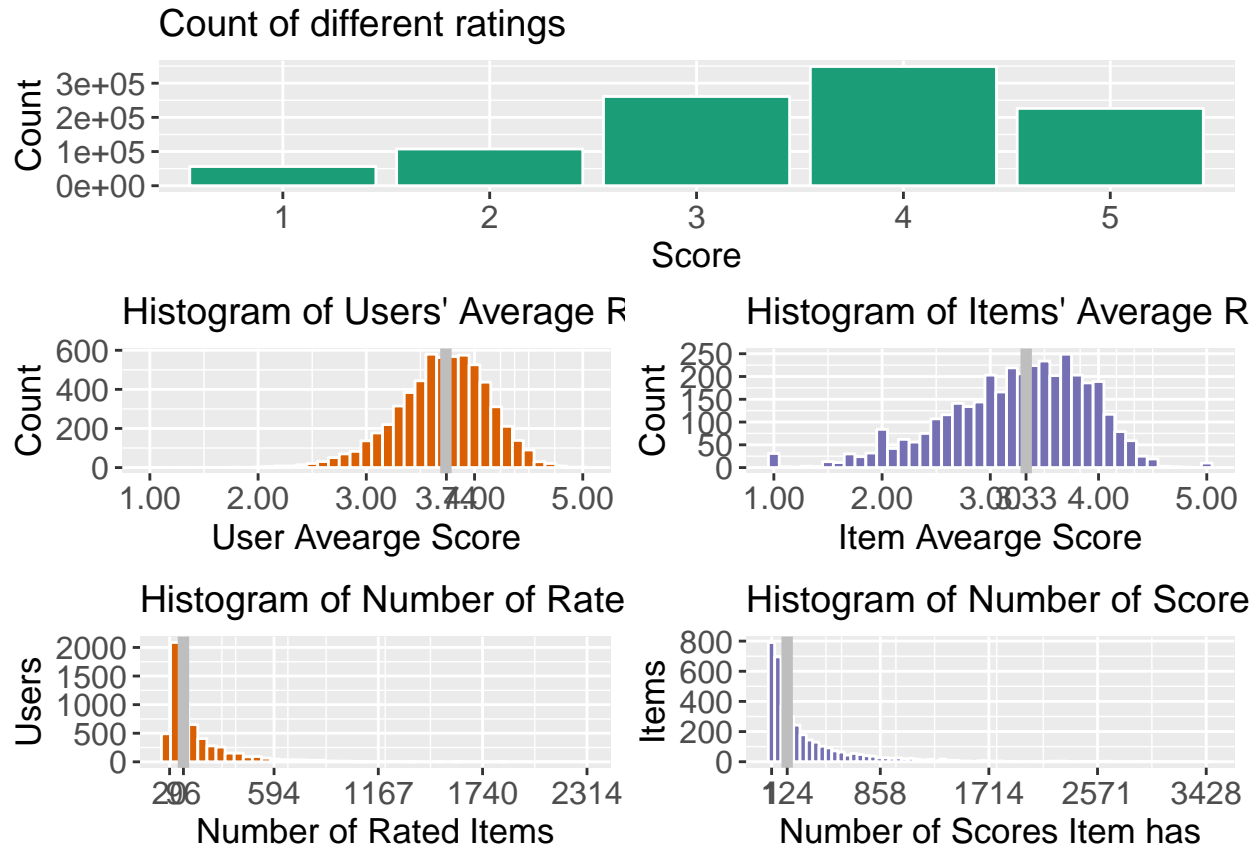
user_id	item_id	rating
Length:99392	Length:99392	Min. :1.00
Class :character	Class :character	1st Qu.:3.00
Mode :character	Mode :character	Median :4.00
NA	NA	Mean :3.53
NA	NA	3rd Qu.:4.00
NA	NA	Max. :5.00

Var1	Freq
1	6059
2	11307
3	27002
4	33947
5	21077



### 3.2 Data visualization

```
visualize_ratings(ratings_table = ratings, color_palett = "Dark2")
```



## 4 Model Evaluation

```
model <- svd_build(mtx)
```

```
model_tunes <- svd_tune(model, r = 2:50)
```

```
## Model is evaluated for r = 2
## Model is evaluated for r = 3
## Model is evaluated for r = 4
## Model is evaluated for r = 5
## Model is evaluated for r = 6
## Model is evaluated for r = 7
## Model is evaluated for r = 8
## Model is evaluated for r = 9
## Model is evaluated for r = 10
## Model is evaluated for r = 11
## Model is evaluated for r = 12
```

```
## Model is evaluated for r = 13
## Model is evaluated for r = 14
## Model is evaluated for r = 15
## Model is evaluated for r = 16
## Model is evaluated for r = 17
## Model is evaluated for r = 18
## Model is evaluated for r = 19
## Model is evaluated for r = 20
## Model is evaluated for r = 21
## Model is evaluated for r = 22
## Model is evaluated for r = 23
## Model is evaluated for r = 24
## Model is evaluated for r = 25
## Model is evaluated for r = 26
## Model is evaluated for r = 27
## Model is evaluated for r = 28
## Model is evaluated for r = 29
## Model is evaluated for r = 30
## Model is evaluated for r = 31
## Model is evaluated for r = 32
## Model is evaluated for r = 33
## Model is evaluated for r = 34
## Model is evaluated for r = 35
## Model is evaluated for r = 36
## Model is evaluated for r = 37
## Model is evaluated for r = 38
## Model is evaluated for r = 39
## Model is evaluated for r = 40
## Model is evaluated for r = 41
## Model is evaluated for r = 42
## Model is evaluated for r = 43
## Model is evaluated for r = 44
## Model is evaluated for r = 45
## Model is evaluated for r = 46
## Model is evaluated for r = 47
## Model is evaluated for r = 48
## Model is evaluated for r = 49
## Model is evaluated for r = 50
```

```
model_tunes$train_vs_valid
```



Train and Validation RMSE for SVD Approximation

