

# BOOK RECOMMENDATION

BY FARHANA AFROZE

# What is recommender system?

- Recommender system is information filtering system that try to predict ratings or choices of a users and recommend product that are interesting to them.
- To build a recommender system there are two most popular approaches which are-
  1. Content based approach
  2. Collaborative approach

**Content Based Approach:** This approach requires lots of information rather than just knowing user interests.

**Collaborative Approach:** On the other hand, this approach doesn't need that much information. It needs user past experiences on a set of items so it can recommend product based on user historical preference.

# Book Recommendation using Recommender system

The goal of this project is-

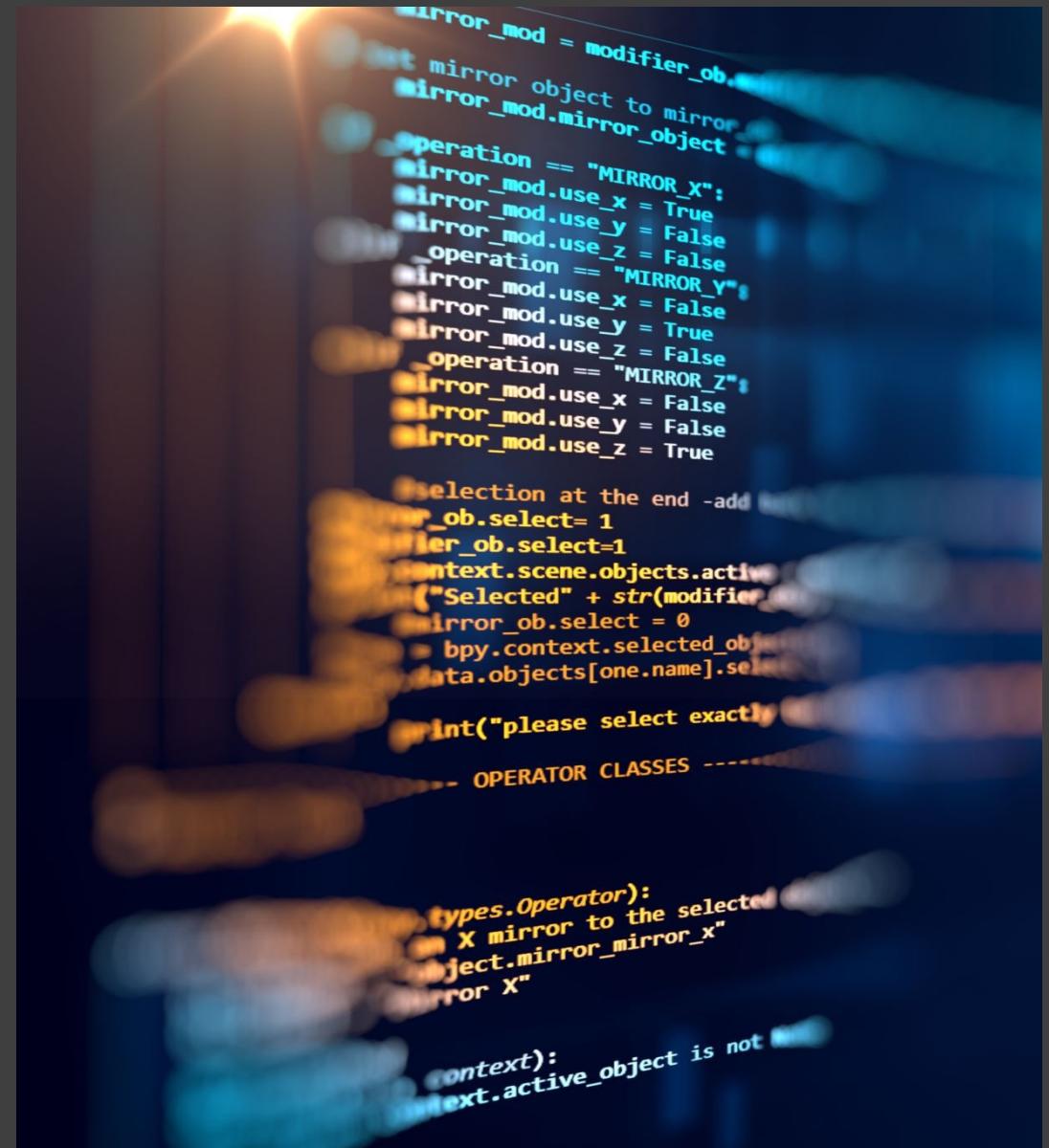
To find out similar users so we can recommend the book to new users based on the new user's similar interest of other users.

Finding out popular books.

Finding out new book that will interest the users based on their past buying of book.

# Project Implementation Tools

- For project implementation I used ‘R’ programming language.
- I used recommender lab package for implementation.
- Used various algorithms from recommender lab



A hand is holding a smartphone, which displays a Python script for Blender operators. The script defines several functions related to mirror modifiers:

```
mirror_mod = modifier_obj
# set mirror object to mirror
mirror_mod.mirror_object = mirror_object
if operation == "MIRROR_X":
    mirror_mod.use_x = True
    mirror_mod.use_y = False
    mirror_mod.use_z = False
elif operation == "MIRROR_Y":
    mirror_mod.use_x = False
    mirror_mod.use_y = True
    mirror_mod.use_z = False
elif operation == "MIRROR_Z":
    mirror_mod.use_x = False
    mirror_mod.use_y = False
    mirror_mod.use_z = True

selection at the end -add
mirror_mod.select=1
mirror_mod.select=1
context.scene.objects.active = eval("Selected" + str(modifier))
mirror_mod.select = 0
bpy.context.selected_objects = eval("Selected" + str(modifier))
data.objects[one.name].select = 1
print("please select exactly one object")

- OPERATOR CLASSES -
types.Operator:
    X mirror to the selected object.mirror_mirror_x"
    or X"
context):
    ext.active_object is not None
```

# Data set for project

- The dataset I chose for this project is from Kaggle. Here is the link-

<https://www.kaggle.com/arashnic/book-recommendation-dataset?select=Books.csv>

The data set has three csv files.

1. Books.csv,
2. Ratings.csv and
3. Users.csv

# Exploring the dataset

- In the Book.csv file it has 8 columns and 271,360 entries

Columns are- ISBN, Book.Title, Book.Author, Year.Of.Publication, Publisher, Image.URLS, Image.URL.M, Image.URL.L

- In the Rating.csv file it has 3 columns and 1,149,780 entries

Columns are- User.ID, ISBN, Book.Rating

- In the Users.csv file it has 3 columns and 278,858 entries

Columns are- User.ID, Location, Age

# Preprocessing the dataset

- At first, I would like to merge the book and rating dataset using the column ISBN
- Then I get the new dataset and again we would like to merge the new dataset with user dataset by using the column User.ID and the new dataset got 12 columns with 1,031,136 entries
- Next, I would like to see if there are any missing data in the dataset. There are like 277835 missing values. Then I omit missing values by deleting the rows. Now we got 73,301 entries.
- Next, I removed some columns that I don't need it for building recommendation system.

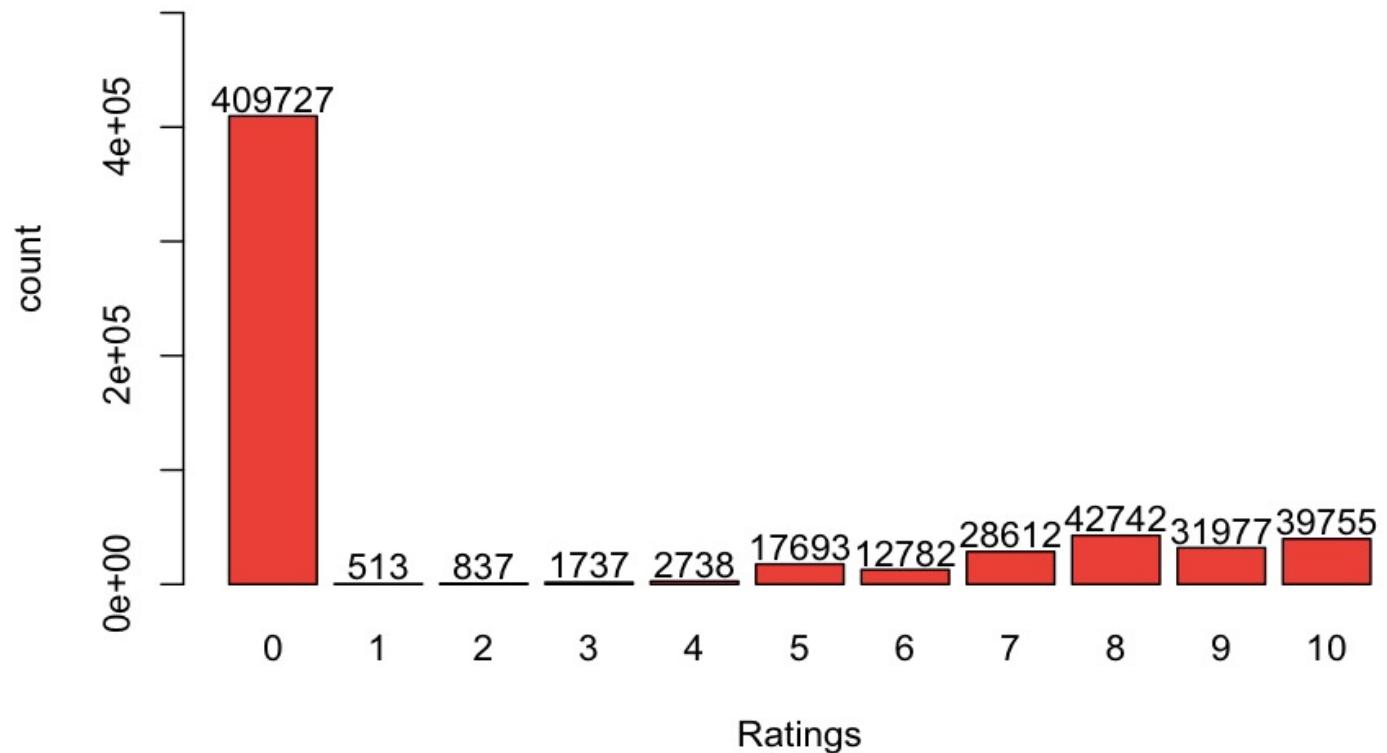
# Preprocessing the dataset

- Removed columns are-
- ISBN, Book.Author, Year.Of.Publication, Publisher, Location, Age, Image.URL.S, Image.URL.M, Image.URL.L
- Now, we have three columns which are User.ID, Book.Title and Book.Rating and has 753,301 entries.
- Next, I would remove duplicate rows where user rated same book twice. And there should be unique rows where user rated book once.
- After removing duplicate rows, we got 750095 entries in it.
- Then we would like to remove users who rated less than 25 books. After removing users we got 589,113 entries with 3 columns

# Visualization of the Dataset

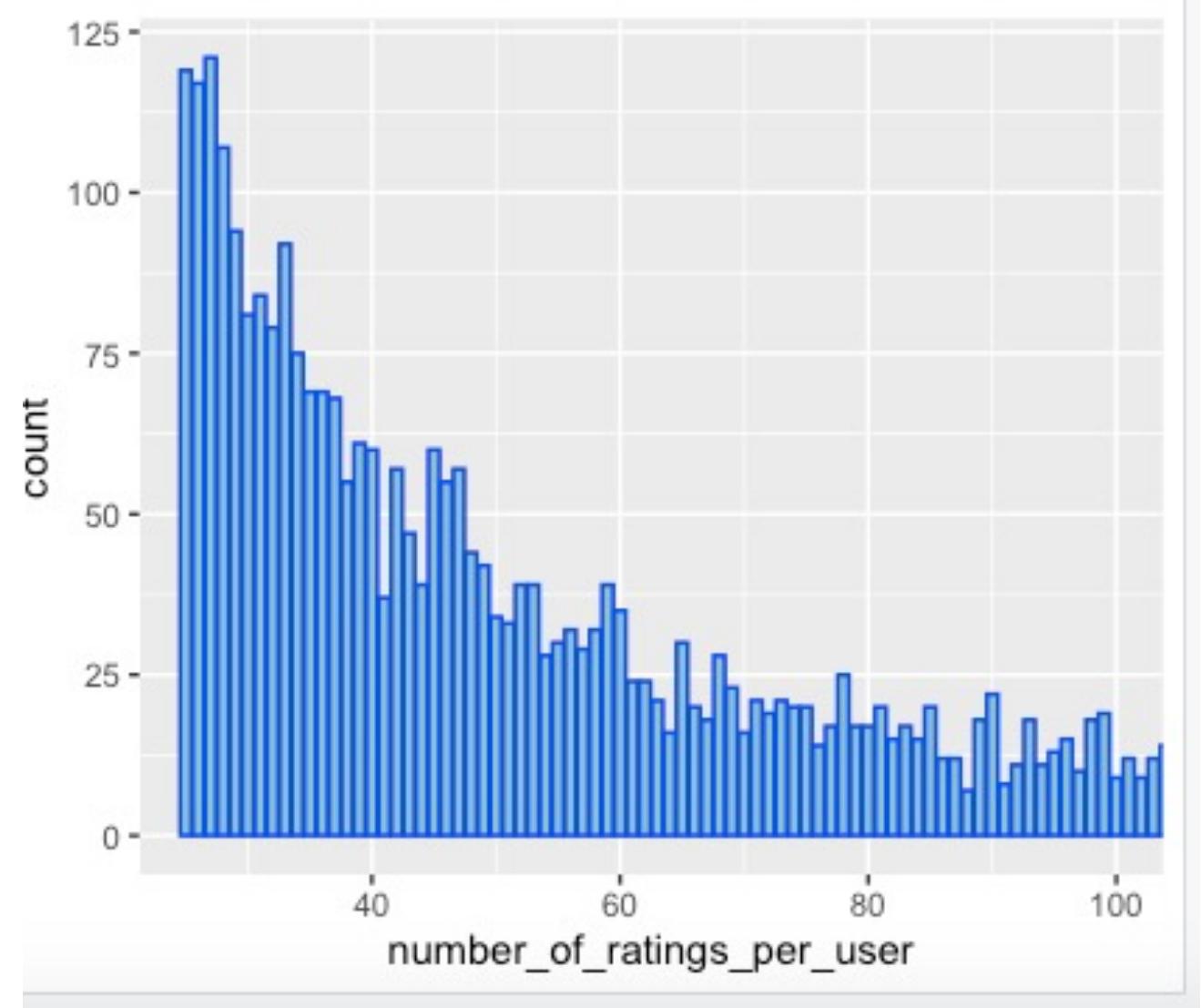
First, We would like to see ratings frequency of the dataset-

**Users ratings of books**



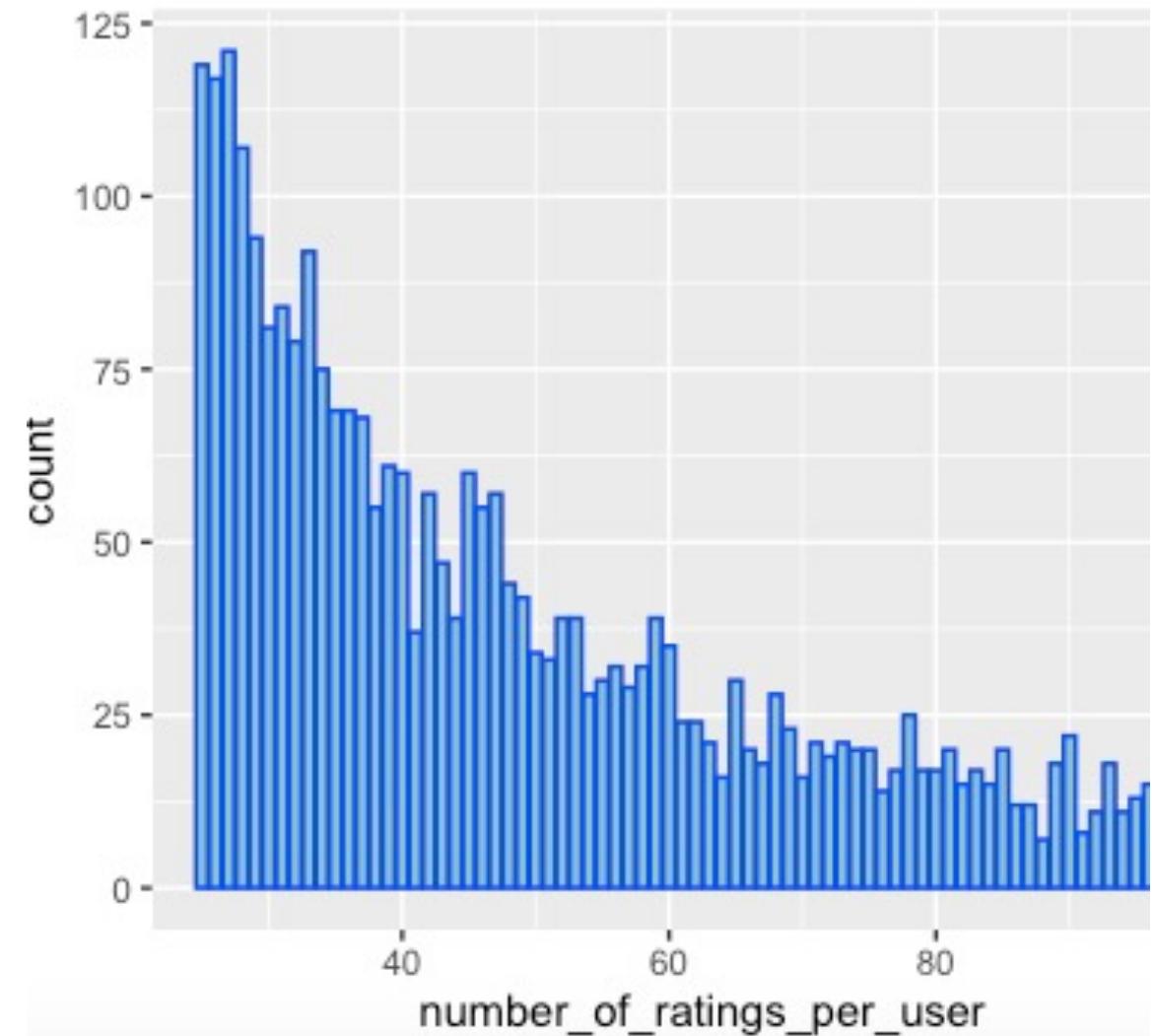
# Visualization of the Dataset

- Here we would like to see number of book ratings per user-



# Subset of the dataset

- In case we get memory issue while run that big dataset we would like to subset our dataset, so our recommender system doesn't halt while running.
- we chose 10% of the data using `sample_frac()` and then again we removed the users who rated less than 20 books. Now we have 36,803 entries
- Here is the number of ratings per user plot looks now-



# Recommender System Modeling

- For Recommender modeling we are going to use recommender lab package which will allow us building recommender model and evaluate recommendation system. For recommender lab to use we need data in matrix form. The data will convert into matrix using reaRatingMatrix.
- After converted into matrix we got 670 users with 26192 books with 36803 ratings.

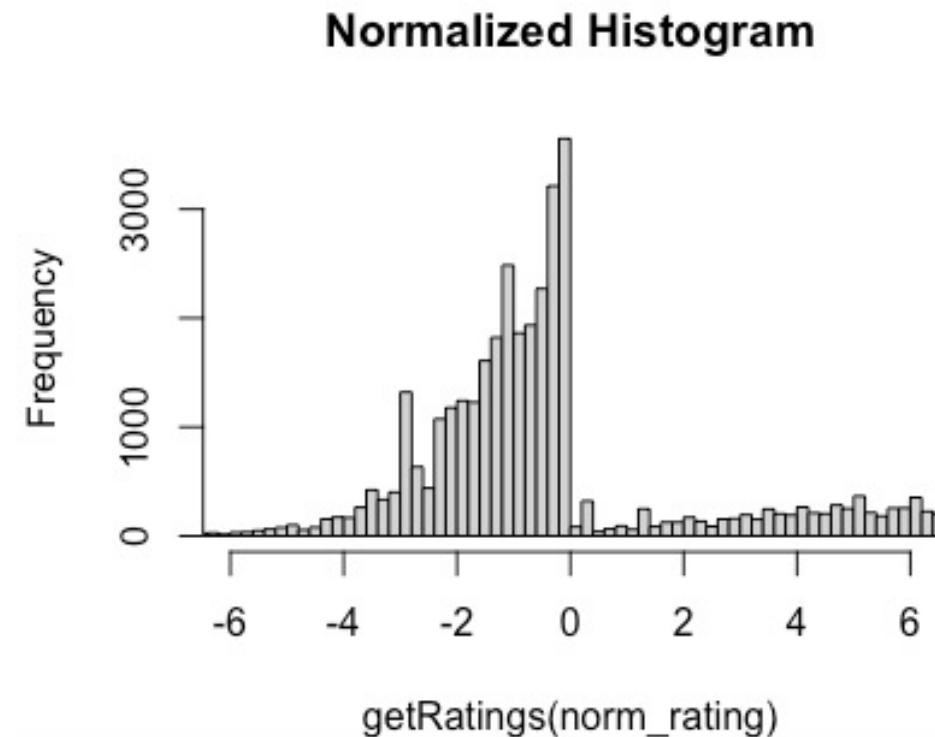
# Looking through matrix and visualization

- First, we would like to see ratings of the user given.

essing 26151 columns in show(); maybe adjust

# Looking through matrix and visualization

- As we saw from previous page that matrix is very very sparse. We saw lots of users didn't rate books. And this is the reason matrix is very sparse.
- Next, we would like to normalize the matrix so if some of users gives very high ratings and some of user gives low ratings it won't effect. Let's see normalized rating graph-



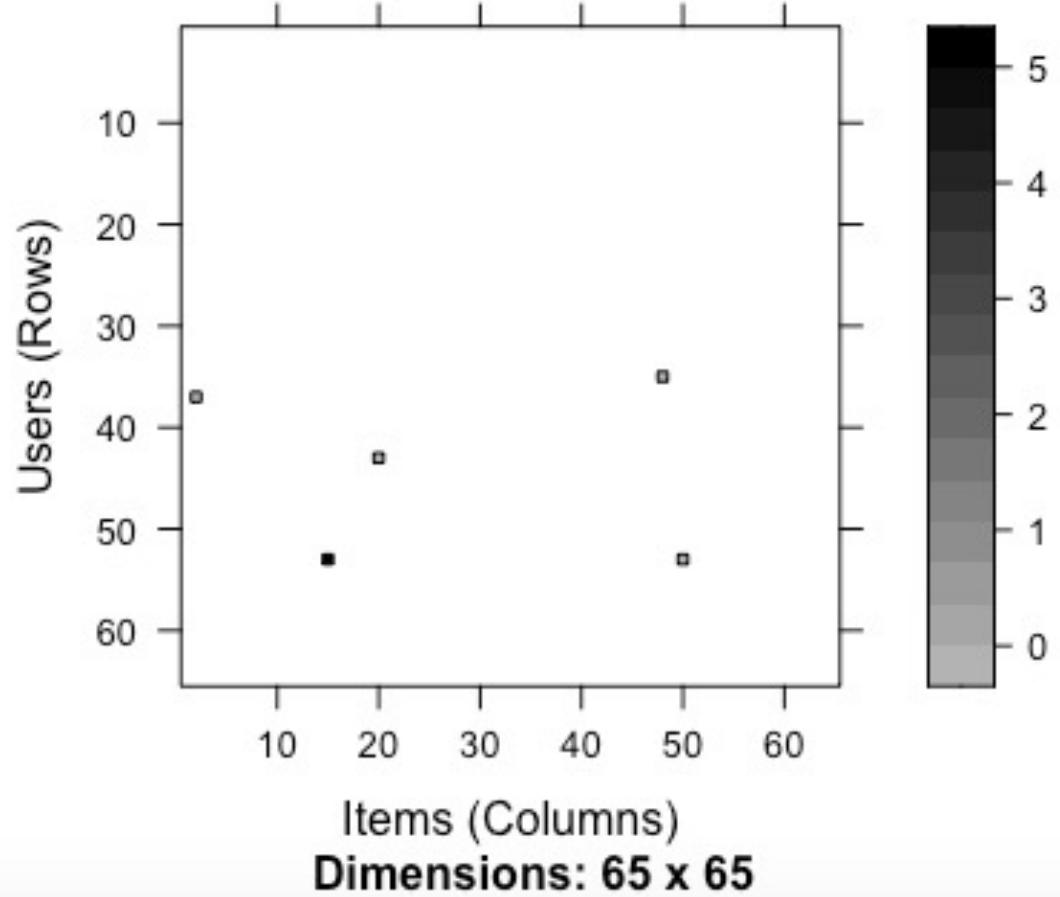
# Looking through matrix and visualization

- Let's see now how many ratings each rating has-

# Looking through matrix and visualization

---

- Here we would like to visualize overall ratings between users and books using 65 users and 65 books. And we would like to observe this through heatmap.



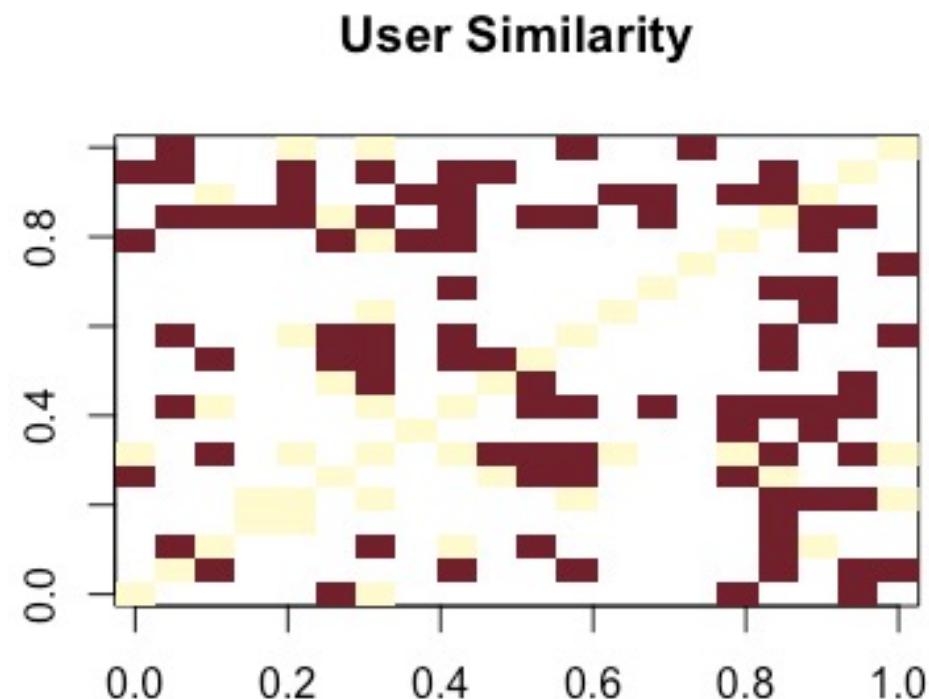


# Data preparation for Recommender system

- As our data matrix is big we would like to subset the data matrix.
- We selected data where user rated books more than 7 times plus book has been rated more than 5 times.
- We got 670 users with 448 movies with 3797 ratings. And again we selected another dataset from this matrix. Where we selected users who rated more than 7 times. And the final data has  $162 \times 448$  rating matrix of class 'realRatingMatrix' with 2026 ratings.

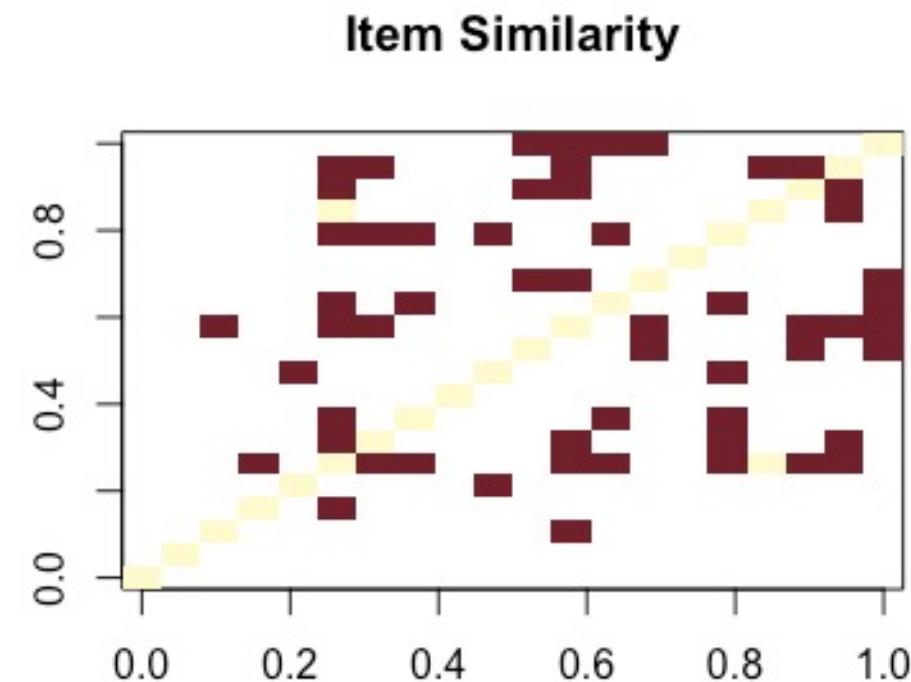
# Similarity Matrix

- Here we would like to see similarity matrix between users and books
- Similarity between users calculated using cosine distance.



# Similarity Matrix

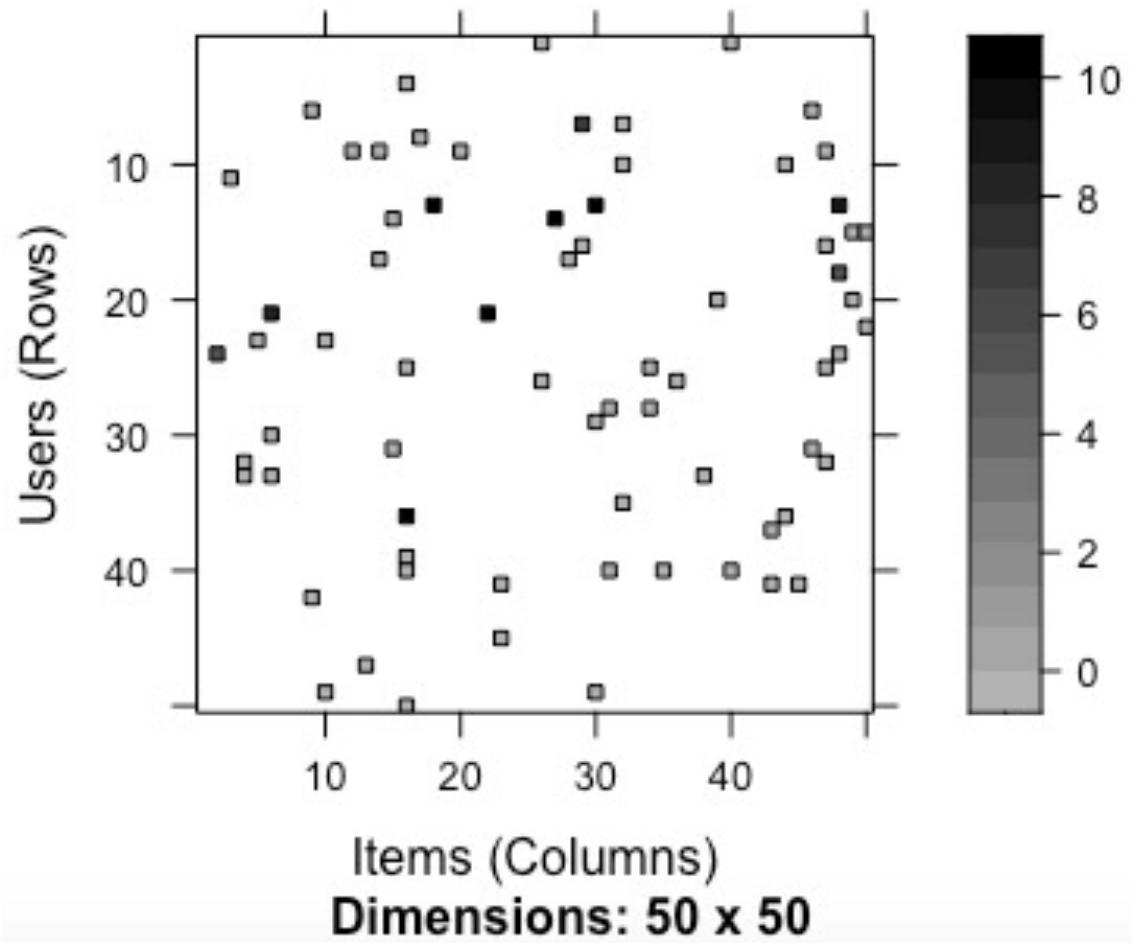
- Similarity between items is calculated using cosine distance.



# Heatmap to see overall ratings

- Here we would like to see overall ratings between users and items through heatmap and it's based on the final dataset where we had 2026 ratings on 162 users and 448 books. We would like to see heatmap on first 50 users and first 50 books

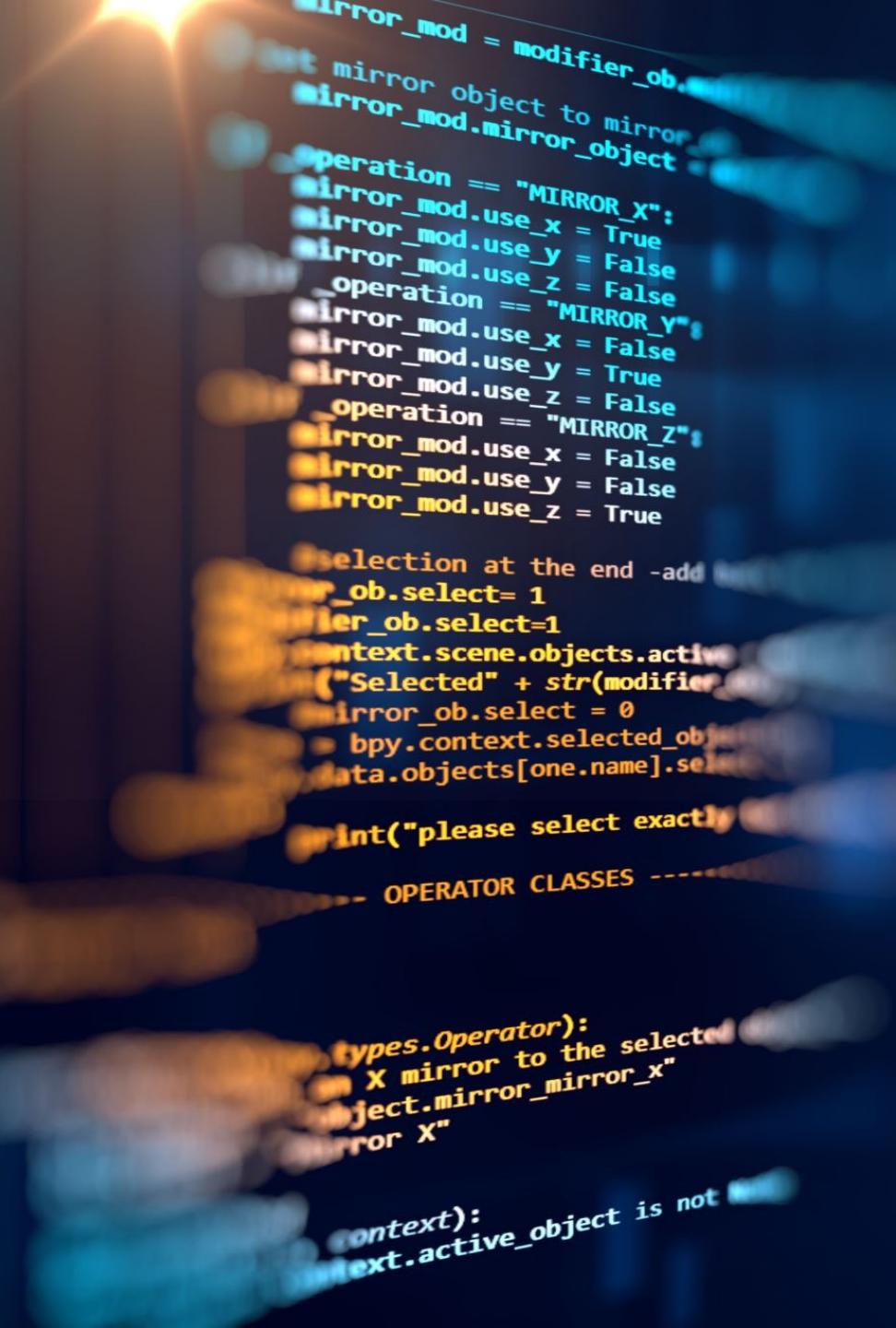
**First 50 users and books**



# Looking through some algorithm in recommender lab

- In the recommender lab there are various algorithm and supported algorithms are:

```
> names(model)
[1] "HYBRID_realRatingMatrix"      "ALS_realRatingMatrix"
[3] "ALS_implicit_realRatingMatrix" "IBCF_realRatingMatrix"
[5] "LIBMF_realRatingMatrix"       "POPULAR_realRatingMatrix"
[7] "RANDOM_realRatingMatrix"     "RERECOMMEND_realRatingMatrix"
[9] "SVD_realRatingMatrix"         "SVDF_realRatingMatrix"
[11] "UBCF_realRatingMatrix"
> 14|
```



# Algorithms of recommender lab

---

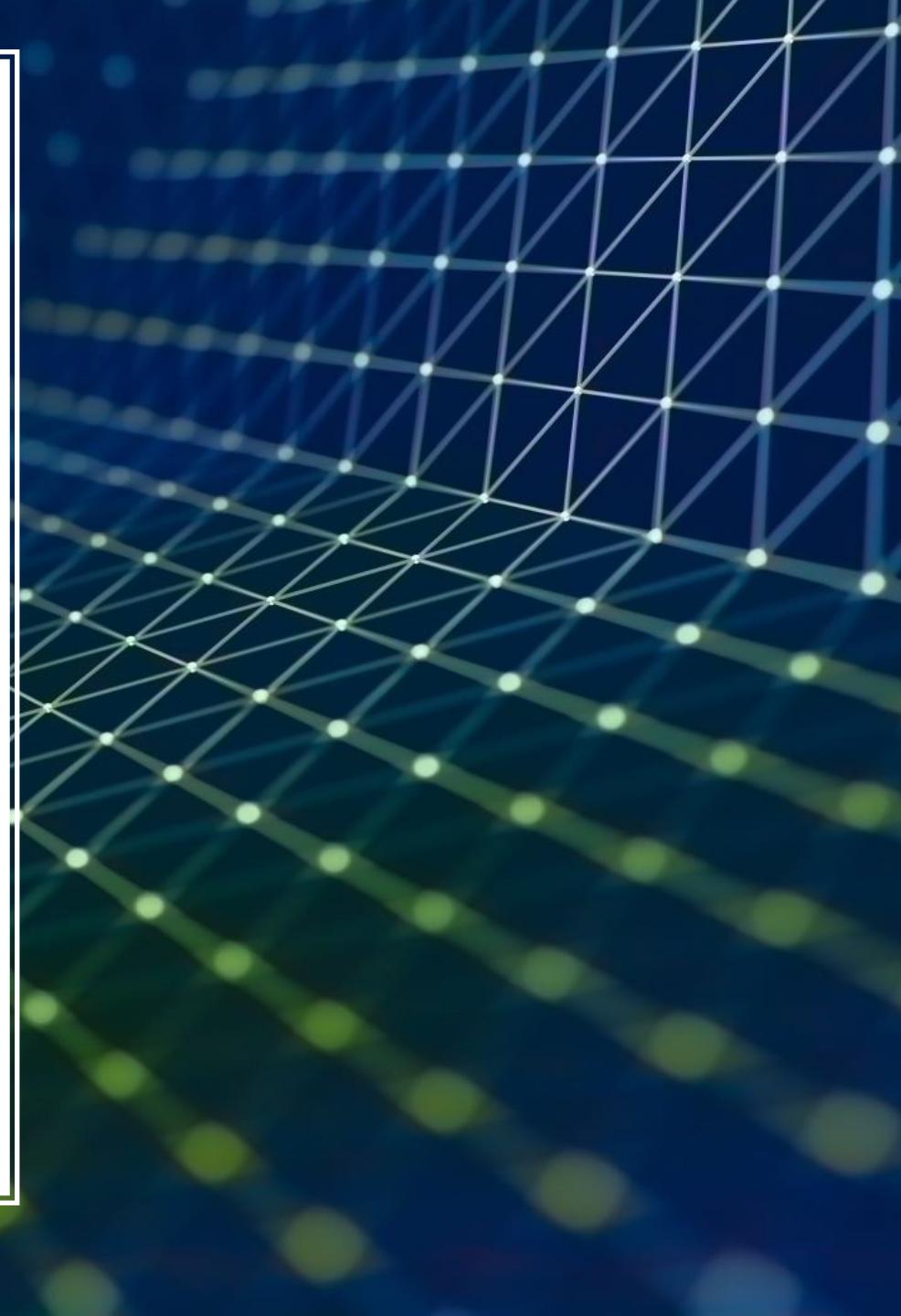
From those algorithm I implemented Three of those-

- UBCF algorithm
- Popular algorithm
- IBCF algorithm

# Algorithm description

Let's see the description of those three algorithms-

- Popular Algorithm: This algorithm is based on the item popularity. You choose top items based on the most rated items
- UBCF algorithm: UBCF algorithm recommends items by finding similar users to the user who recommends the movie. It used k-nearest neighbors to measure the distance between each pair of users.
- IBCF algorithm: In IBCF we try to find out item similar items to the user. We select the items based on the past buying history of users. IBCF uses cosine based similarity, jaccard distance and other types of distance measurement.



# Implementation of Algorithms

# Popular Algorithm



Here we implementing algorithm on the first 100 users.



Then we are predicting top 5 books on the 15 users. Here is the screen shot of first 4 users and their top 5 books-

'7072'

"Little Altars Everywhere: A Novel"

"Stupid White Men ...and Other Sorry Excuses for the State of the N

"Interview with the Vampire"

"A Time to Kill"

"The Lovely Bones: A Novel"

'8667'

"Little Altars Everywhere: A Novel"

"Stupid White Men ...and Other Sorry Excuses for the State of the N

"Interview with the Vampire"

"A Time to Kill"

"The Lovely Bones: A Novel"

'4299'

"Little Altars Everywhere: A Novel"

"Stupid White Men ...and Other Sorry Excuses for the State of the N

"Interview with the Vampire"

"The Lovely Bones: A Novel"

"The Catcher in the Rye"

'5233'

"Little Altars Everywhere: A Novel"

"Stupid White Men ...and Other Sorry Excuses for the State of the N

"Interview with the Vampire"

"A Time to Kill"

"The Lovely Bones: A Novel"

# UBCF Algorithm

- So first we will implement the algorithm and then predict the algorithm on each user.
- We selected 5 most similar user for prediction
- So, what it will do is- UBCF will identify books in terms of having same interest by the same 5 users. And then it will predict books to a new user.

# UBCF Algorithm outcome

---

- So, for the first user we found 40 ratings of 40 books-
- For the second user we found 0 books-
- For third user we found 34 books with 34 ratings
- Snapshot of third user books and ratings-(first couple ratings and books)

\$`143175`	
	A Bend in the Road 9.600000
	Bleachers 5.8333333
	Circle of Friends 9.600000
	Confessions of an Ugly Stepsister : A Novel 0.8333333
	Daddy 6.8000000
	Dawn (Cutler) 6.8000000
E Is for Evidence: A Kinsey Millhone Mystery (Kinsey Millhone Mysteries (Paperback))	0.8333333
	Heart of the Sea (Irish Trilogy) 0.8333333
	Ladder of Years 0.8333333
	Lucy Sullivan Is Getting Married 6.8000000
	Midnight Bayou

# IBCF Algorithm

- Before implementing IBCF algorithm we splits our data matrix where 75% for train data and 25% for test data.
- In train data we got 112 users with 448 books with 1424 ratings
- In test data we got 50 users with 448 books with 602 ratings
- Next, we implemented IBCF algorithm. We have the parameter is  $k = 30$  which is k-most similar items.
- Then we predicted using test data and we have parameters  $n = 5$  which is number of books recommend to each user.
- Next page, Let's see the first 5 users with their 5 books recommended.

# IBCF Algorithm

```
$`7346`  
[1] "\\" Lamb to the Slaughter and Other Stories (Penguin 60s S.)\""  
[2] "\\"A\\\" is for Alibi : A Kinsey Millhone Mystery (A Kinsey Millhone Mystery)\""  
[3] "\\"It's Never Too Late to Have a Happy Childhood\\\"": Inspirations for Adult Children\""  
[4] "\\"Surely You're Joking, Mr. Feynman!\\\"": Adventures of a Curious Character\""  
[5] "\\"The Man with the Twisted Lip (Penguin 60s S.)\""  
  
$`12538`  
character(0)  
  
$`13552`  
[1] "101 Dinosaur Jokes"  
[2] "12 Steps to Living Without Fear"  
[3] "18Mm Blues"  
[4] "A 5th Portion of Chicken Soup for the Soul : 101 Stories to Open the Heart and Rekindle the Spirit"  
[5] "A Confederacy of Dunces (Evergreen Book)"  
  
$`16634`  
[1] "102 Haunted House Jokes"  
[2] "A Cast of Killers"  
[3] "A Deadly Little Christmas"  
[4] "A Deepness in the Sky : A Novel (Zones of Thought)"  
[5] "2001 Lippincott's Nursing Drug Guide (Book with Mini CD-ROM for Windows & Macintosh)"  
  
$`31315`  
[1] " Twenty Thousand Leagues Under the Sea (Illustrated Classics)"  
[2] "101 Hopelessly Hilarious Jokes"  
[3] "40 Tons Of Trouble (Women Who Dare) (Harlequin Super Romance, No 726)"  
[4] "A 4th Course of Chicken Soup for the Soul: 101 More Stories to Open the Heart and Rekindle the Spirit"  
[5] "A Breach of Promise (William Monk Novels (Paperback))"
```

# Comparing different Recommender Algorithms

# Comparing Different algorithms

- Here we would like see performance of IBCF, Popular and random methods and would like to compare them so we know which methods is predicting best.
- We like to evaluate these using recommender lab. We would like to use k-fold cross validation for that. The parameters for evaluation scheme are  $k = 3$  which means we chose 3-fold cross validation, then  $given = 4$  which is at least 4 books user rated and the  $goodRating = 4$  which is we consider ratings good at least 4.
- Next, we will evaluate this model with all 3 algorithms with parameters  $n = c(1,3,5,10,15,20,25,30)$  which is the top recommendation list

# Algorithm prediction time

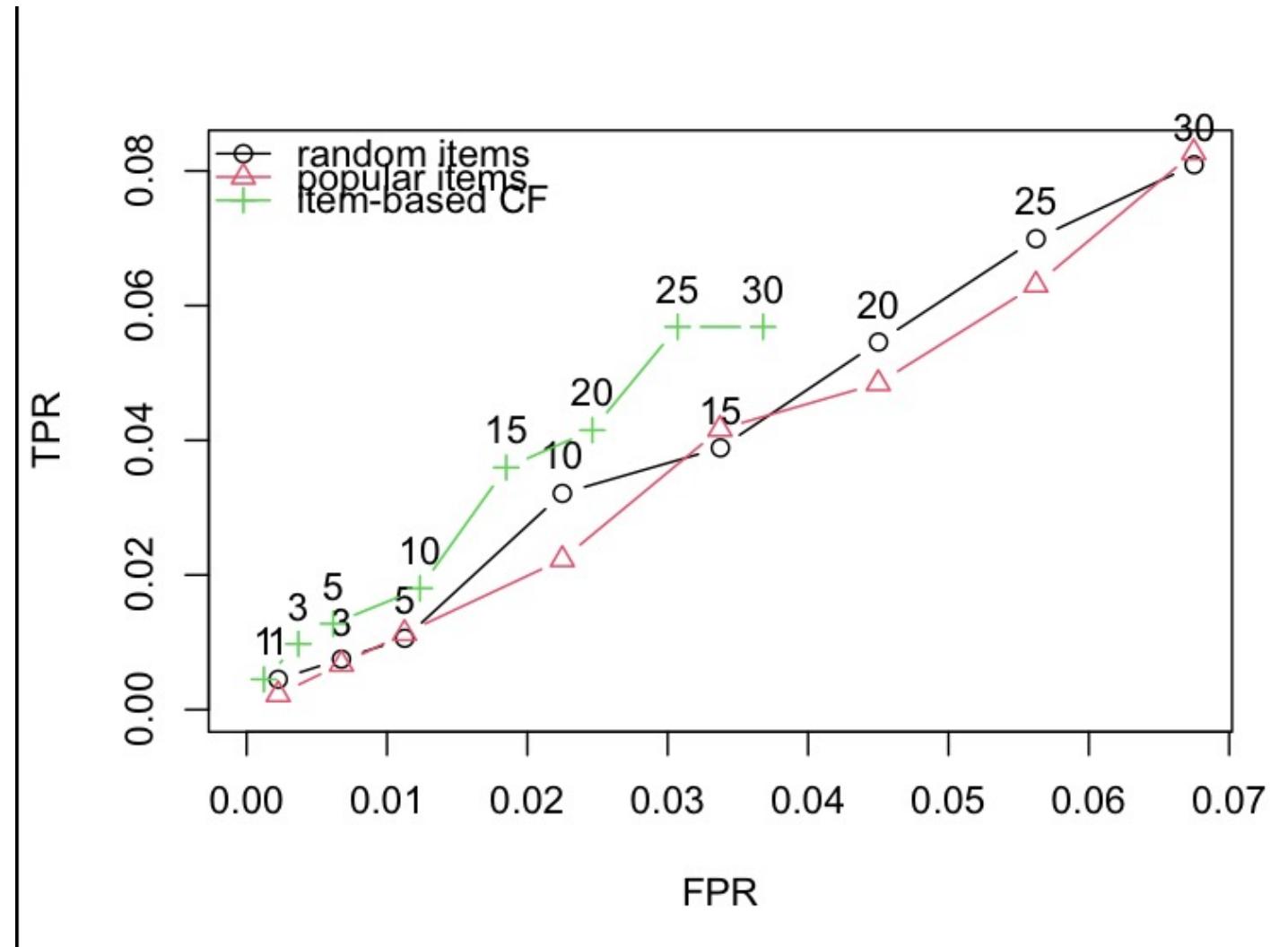
- Let's see all 3 methods prediction or modeling time-

```
RANDOM run fold/sample [model time/prediction time]
  1 [0.017sec/0.097sec]
  2 [0.005sec/0.066sec]
  3 [0.002sec/0.072sec]
POPULAR run fold/sample [model time/prediction time]
  1 [0.016sec/0.271sec]
  2 [0.01sec/0.257sec]
  3 [0.011sec/0.259sec]
IBCF run fold/sample [model time/prediction time]
  1 [0.19sec/0.049sec]
  2 [0.179sec/0.037sec]
  3 [1.272sec/0.036sec]
```

> |

# Graph of the evaluation model

- Here we would like to see the TPR and FPR graph of three algorithms.



# Graph Interpretation

- So, from last page graph what's the TPR and FPR?
- TPR(True Positive Rate): These are recommended books that have been purchased.
- FPR(False Positive Rate): These are recommended books that haven't been purchased
- From the previous graph we see IBCF has FPR is less than other methods. It's showing that this algorithm finds out more recommending items that have been purchased and also it's not misinterpreting like other two methods where other two methods have more recommended items that haven't been purchased. But also other two recommendation system is doing better in TPR. Here IBCF is doing best as it is area under the curve an has less FPR.

# Confusion matrix result

Next page, let's see confusion matrix result for each methods



# IBCF confusion matrix result

---

TPI	FPI	FNI	TNI	NI	precision	recall	TPR	FPR	n
0.01851851	0.57407411	1.8518521	441.55561	4441	0.03125001	0.01351351	0.01351351	0.00130011	11
0.05555561	1.72222221	1.8148151	440.40741	4441	0.03125001	0.02927931	0.02927931	0.00389991	31
0.07407411	2.88888891	1.7962961	439.24071	4441	0.02500001	0.03828831	0.03828831	0.00654191	51
0.11111111	5.81481481	1.7592591	436.31481	4441	0.01875001	0.05405411	0.05405411	0.01316761	101
0.12962961	8.75925931	1.7407411	433.37041	4441	0.01458331	0.06756761	0.06756761	0.01983551	151
0.12962961	11.72222221	1.7407411	430.40741	4441	0.01093751	0.06756761	0.06756761	0.02654531	201
0.14814811	14.66666671	1.7222221	427.46301	4441	0.01000001	0.07657661	0.07657661	0.03321321	251
0.14814811	17.62962961	1.7222221	424.50001	4441	0.00833331	0.07657661	0.07657661	0.03992301	301

# POPULAR Confusion Matrix result

TPI	FPI	FNI	TNI	NI	precision	recall	TPR	FPR	nl
0.01851851	0.98148151	1.8518521	441.14811	4441	0.01851851	0.00675681	0.00675681	0.00221971	11
0.03703701	2.96296301	1.8333331	439.16671	4441	0.01234571	0.02027031	0.02027031	0.00670151	31
0.05555561	4.94444441	1.8148151	437.18521	4441	0.01111111	0.02702701	0.02702701	0.01118301	51
0.09259261	9.90740741	1.7777781	432.22221	4441	0.00925931	0.04144141	0.04144141	0.02240791	101
0.11111111	14.88888891	1.7592591	427.24071	4441	0.00740741	0.06846851	0.06846851	0.03367521	151
0.11111111	19.88888891	1.7592591	422.24071	4441	0.00555561	0.06846851	0.06846851	0.04498431	201
0.14814811	24.85185191	1.7222221	417.27781	4441	0.00592591	0.08873871	0.08873871	0.05620941	251
0.16666671	29.83333331	1.7037041	412.29631	4441	0.00555561	0.11576581	0.11576581	0.06747671	301

# RANDOM Confusion matrix result

---

	TP	FP	FN	TN	N	precision	recall	TPR	FPR	n
-----: -----: -----: -----: ---: -----: -----: -----: -----: -----: --:	0.0185185   0.9814815   1.851852   441.1481   444   0.0185185   0.0135135   0.0135135   0.0022199   1									
0.0370370   2.9629630   1.833333   439.1667   444   0.0123457   0.0225225   0.0225225   0.0067016   3										
0.0370370   4.9629630   1.833333   437.1667   444   0.0074074   0.0225225   0.0225225   0.0112252   5										
0.0555556   9.9444444   1.814815   432.1852   444   0.0055556   0.0315315   0.0315315   0.0224923   10										
0.1111111   14.8888889   1.759259   427.2407   444   0.0074074   0.0518018   0.0518018   0.0336751   15										
0.1481481   19.8518519   1.722222   422.2778   444   0.0074074   0.0878378   0.0878378   0.0449004   20										
0.1666667   24.8333333   1.703704   417.2963   444   0.0066667   0.0968468   0.0968468   0.0561675   25										
0.1851852   29.8148148   1.685185   412.3148   444   0.0061728   0.1103604   0.1103604   0.0674347   30										

> |

## Evaluating Metrics parameter-

- False Negative- These are not recommended books that have been purchased
- True Negative- These are not recommended books that haven't been purchased.
- Recall – This is the measure of how many true positives get predicted out of all the positives in the dataset.
- Precision – This is the measure of correctness of positive prediction

# Comparing the models- IBCF

IBCF metrics result we got-

Precision = 0.0083 for n=30 top items

Recall = 0.076 for n = 30 top items

TPR = 0.076 for n = 30 top items

FPR = 0.0399 for n = 30 top items

# Comparing the models- Popular

Popular method metrics result we got-

Precision = 0.0055 for n = 30 top items

Recall = 0.115 for n = 30 items

TPR = 0.115 for n = 30 items

FPR = 0.064 for n = 30 items

## Comparing the models- Random

Precision- 0.006 for n = 30 top books

Recall- 0.11 for n = 30 top books

TPR- 0.011 for n = 30 top books

FPR- 0.067 for n = 30 top book

# Conclusion

- Based on those previous models I would choose IBCF model as this model precision is higher and FPR is lower than other two models. So, I would like to have more correct positive prediction than false positive prediction. And also I want FPR is lower because I want the product that has been recommended wants them purchased.