# DATA612 Project 2

Michael O'Donnell

June 14, 2019

# Overview:

In the following R code, two recommender systems are implemented on MovieLense data First, an Item-Based Collaborative Filtering recommender model Second, a User-Based Collaborative Filtering recommender model After both models are implemented, both models are evaluated This code was written with obvious help from "Building a Recommender System with R" chapters 3 and 4

## import libraries

```
library(recommenderlab)
```

```
## Loading required package: Matrix
```

```
## Loading required package: arules
```

```
##
## Attaching package: 'arules'
```

```
## The following objects are masked from 'package:base':
##
##     abbreviate, write
```

```
## Loading required package: proxy
```

```
##
## Attaching package: 'proxy'
```

```
## The following object is masked from 'package:Matrix':
##
##     as.matrix
```

```
## The following objects are masked from 'package:stats':
##
##     as.dist, dist
```

```
## The following object is masked from 'package:base':
##
##     as.matrix
```

```
## Loading required package: registry
```

```
library(ggplot2)
set.seed(1)
```

# import the MovieLense data

```
data(MovieLense)
MovieLense
```

```
## 943 x 1664 rating matrix of class 'realRatingMatrix' with 99392 ratings.
```

# View the size of the MovieLense data

```
object.size(MovieLense)
```

```
## 1409432 bytes
```

```
object.size(as(MovieLense, "matrix"))
```

```
## 12761360 bytes
```

# converting the matrix into vector to see values

```
vector_ratings <- as.vector(MovieLense@data)
unique(vector_ratings)
```
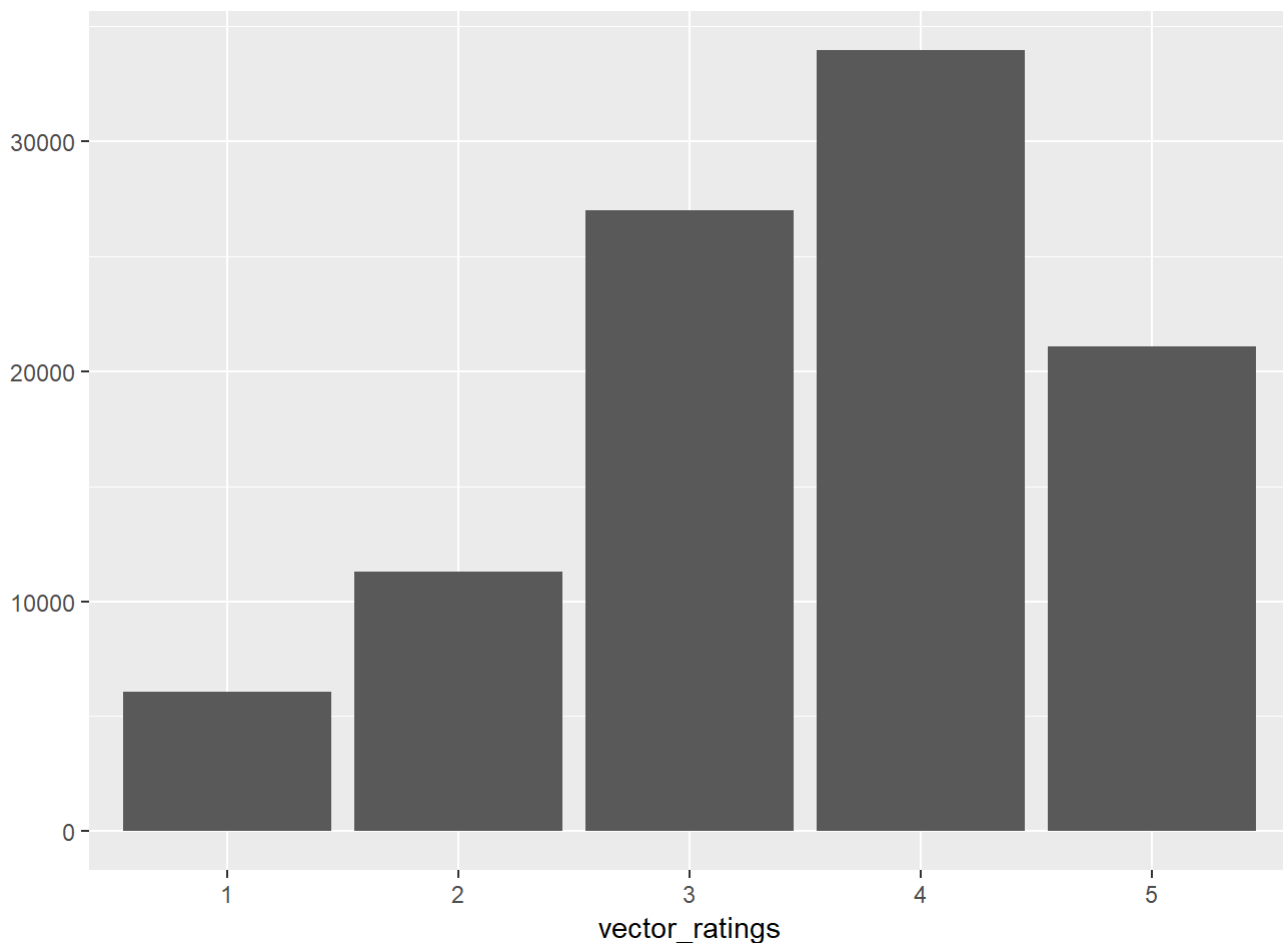
```
## [1] 5 4 0 3 1 2
```

```
table(vector_ratings)
```

```
## vector_ratings
##       0       1       2       3       4       5
## 1469760    6059   11307   27002   33947   21077
```

# removing the null values and turning vector into factors

```
vector_ratings <- vector_ratings[vector_ratings != 0]
vector_ratings <- factor(vector_ratings)

qplot(vector_ratings)
```



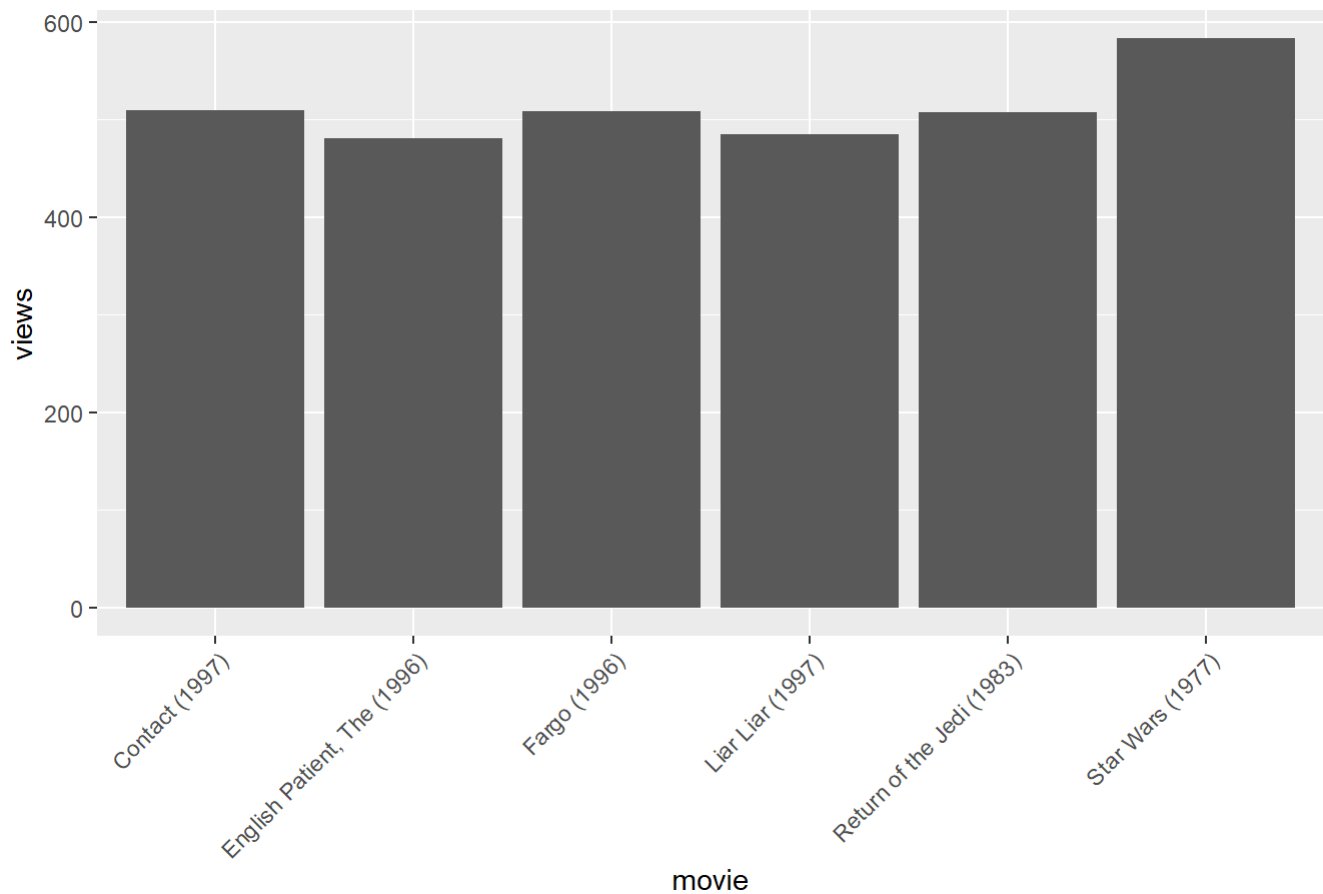## calculating and visualizing which movies have been viewed

```
views_per_movie <- colCounts(MovieLense)

table_views <- data.frame(
  movie = names(views_per_movie),
  views = views_per_movie
)

table_views <- table_views[order(table_views$views,
                                 decreasing = TRUE), ]

ggplot(table_views[1:6, ], aes(x=movie, y=views)) +
  geom_bar(stat = "identity") +
  theme(axis.text.x = element_text(angle = 45, hjust = 1)) +
  ggtitle("Number of Views of Top 6 Movies")
```
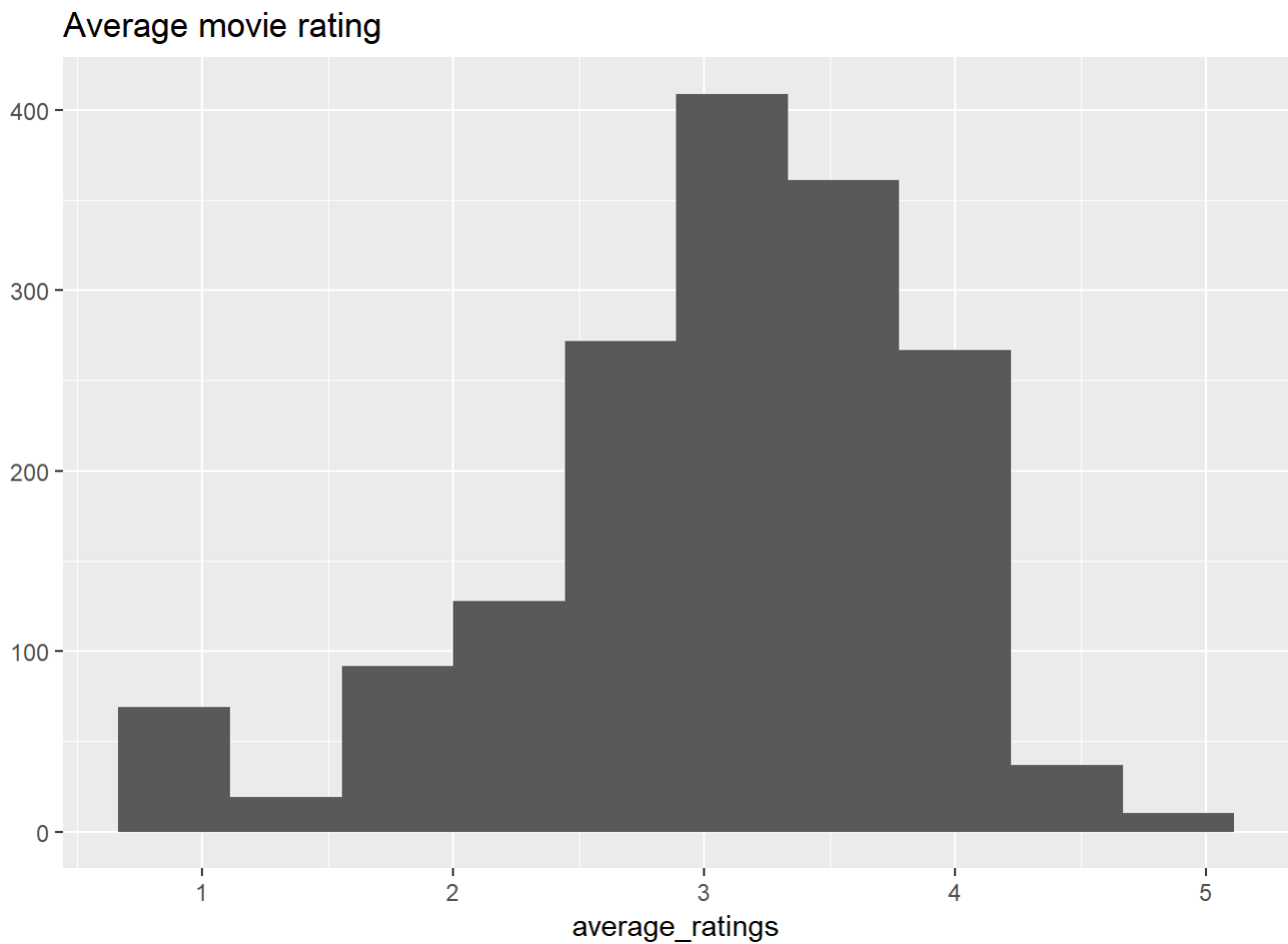
## Number of Views of Top 6 Movies



## visualizing the average movie score

```
average_ratings <- colMeans(MovieLense)

qplot(average_ratings) +
  stat_bin(bins = 10) +
  ggtitle("Average movie rating")
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```
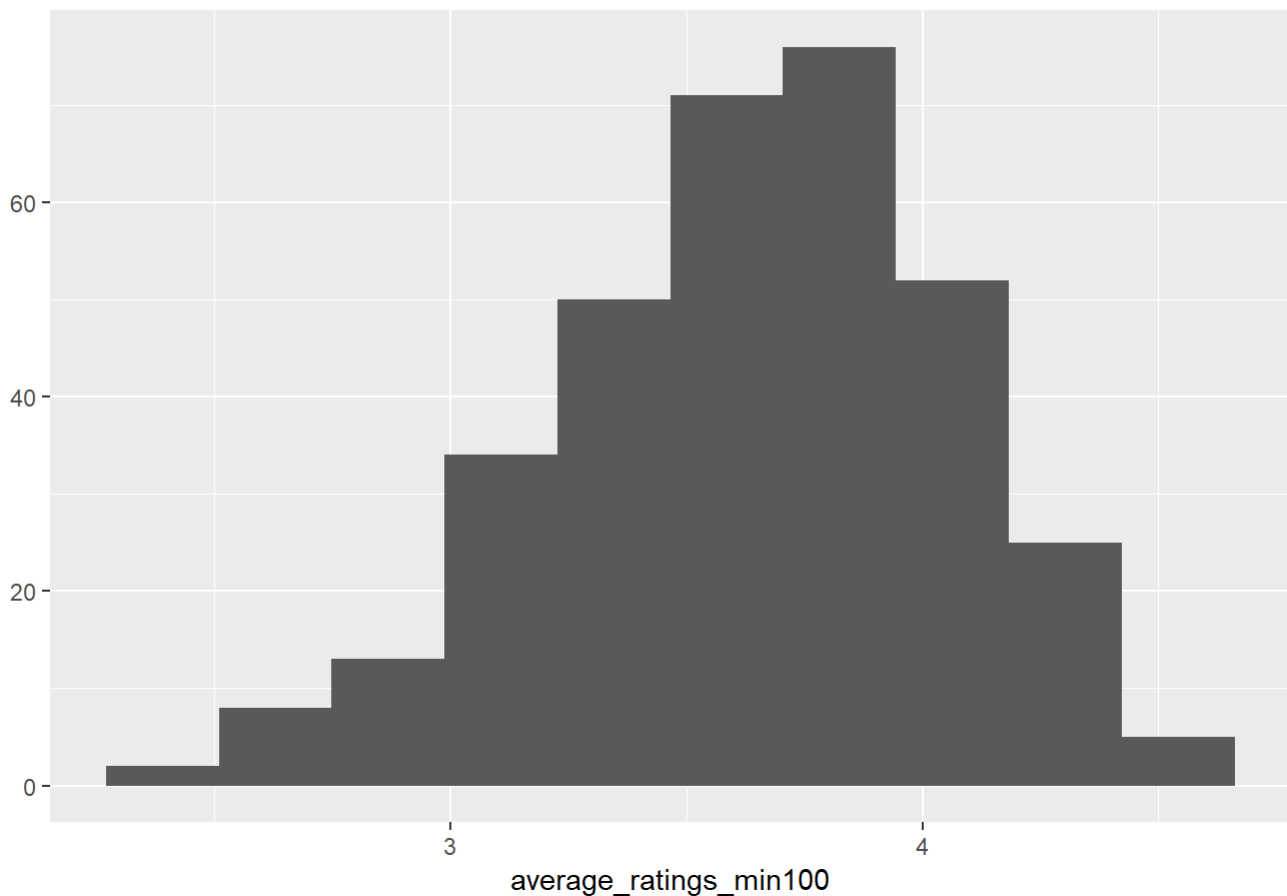
Average movie rating

## view the average ratings of only movies with 100 views minimum

```
average_ratings_min100 <- average_ratings[views_per_movie >= 100]

qplot(average_ratings_min100) +
  stat_bin(bins = 10) +
  ggtitle("Average movie rating (minimum 100)")
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

## Average movie rating (minimum 100)



## selecting only data with enough ratings and power users

```
# greater than 100 views
# only accounting for users that have rated at least 50 movies
ratings_movies <- MovieLense[rowCounts(MovieLense) > 50,
                             colCounts(MovieLense) > 100]
ratings_movies
```

```
## 560 x 332 rating matrix of class 'realRatingMatrix' with 55298 ratings.
```

```
#average ratings per user
avg_ratings_user <- rowMeans(ratings_movies)
```

## normalize the user ratings to zero

```
ratings_movies_normalize <- normalize(ratings_movies)
```

## splitting the data into training and testing sets

```
which_train <- sample(x = c(TRUE, FALSE), size = nrow(ratings_movies),
                       replace = TRUE, prob = c(0.8, 0.2))

train <- ratings_movies[which_train, ]
test <- ratings_movies[!which_train, ]
```

## use k-fold to split the users into 5 groups

```
which_set <- sample(x = 1:5, size=nrow(ratings_movies),
                     replace = TRUE)
for(i in 1:5) {
  which_train <- which_set == i
  train <- ratings_movies[which_train, ]
  test <- ratings_movies[!which_train, ]
}
```

## establishing the Item Based Collaborative Filtering recommender model

```
model <- Recommender(data = train, method = "IBCF",
                     parameter = list(k=30))
model
```

```
## Recommender of type 'IBCF' for 'realRatingMatrix'
## learned using 111 users.
```

## apply model onto the test set (IBCF model)

```
# number of items to recommend
n_recommend <- 5

predicted <- predict(object = model, newdata = test, n = n_recommend)
predicted
```

```
## Recommendations as 'topNList' with n = 5 for 449 users.
```

## see the list of recommended movies for the first test user (IBCF model)

```
test_user_one <- predicted@items[[1]]
test_movies_one <- predicted@itemLabels[test_user_one]
test_movies_one
```

```
## [1] "Craft, The (1996)"
## [2] "Little Women (1994)"
## [3] "E.T. the Extra-Terrestrial (1982)"
## [4] "G.I. Jane (1997)"
## [5] "Star Trek: The Motion Picture (1979)"
```

## now, recommend movies for each user in the test set (IBCF model)

```
recommender_matrix <- sapply(predicted@items, function(x){
  colnames(ratings_movies)[x]
})

recommender_matrix[, 2:4]
```

```
##      2
## [1,] "Babe (1995)"
## [2,] "Natural Born Killers (1994)"
## [3,] "Mystery Science Theater 3000: The Movie (1996)"
## [4,] "Frighteners, The (1996)"
## [5,] "Event Horizon (1997)"
##      3                             5
## [1,] "Craft, The (1996)"           "Ace Ventura: Pet Detective (1994)"
## [2,] "Aladdin (1992)"              "Cape Fear (1991)"
## [3,] "While You Were Sleeping (1995)" "Dumbo (1941)"
## [4,] "Emma (1996)"                 "Little Women (1994)"
## [5,] "Michael (1996)"              "Peacemaker, The (1997)"
```

## Now, to view the most frequently recommended movies (IBCF model)

```
items <- factor(table(recommender_matrix))
items <- sort(items, decreasing = TRUE)
top_items <- data.frame(names(items), items)
head(top_items)
```

```
##                                                names.items.
## Spawn (1997)                                    Spawn (1997)
## Ace Ventura: Pet Detective (1994)   Ace Ventura: Pet Detective (1994)
## Natural Born Killers (1994)            Natural Born Killers (1994)
## Craft, The (1996)                               Craft, The (1996)
## Outbreak (1995)                                  Outbreak (1995)
## Ghost and the Darkness, The (1996) Ghost and the Darkness, The (1996)
##                                    items
## Spawn (1997)                          45
## Ace Ventura: Pet Detective (1994)     32
## Natural Born Killers (1994)           31
## Craft, The (1996)                     28
## Outbreak (1995)                       23
## Ghost and the Darkness, The (1996)    22
```

# We've implemented a IBCF model

# Now, we will implement a User Based Collaborative Filtering model

(on the same data)

```
model <- Recommender(data = train, method = "UBCF")
model
```

```
## Recommender of type 'UBCF' for 'realRatingMatrix'
## learned using 111 users.
```

# To view some more details of this model

```
names(getModel(model))
```

```
## [1] "description" "data"        "method"      "nn"          "sample"
## [6] "normalize"   "verbose"
```

# apply model onto the test set (UBCF model)

```
# number of items to recommend
n_recommend <- 5


predicted <- predict(object = model, newdata = test, n = n_recommend)
predicted
```

```
## Recommendations as 'topNList' with n = 5 for 449 users.
```

## see the list of recommended movies for the first test user (UBCF model)

```
test_user_one <- predicted@items[[1]]
test_movies_one <- predicted@itemLabels[test_user_one]
test_movies_one
```

```
## [1] "Leaving Las Vegas (1995)" "Schindler's List (1993)"
## [3] "L.A. Confidential (1997)" "Trainspotting (1996)"
## [5] "Titanic (1997)"
```

## now, recommend movies for each user in the test set (UBCF model)

```
recommender_matrix <- sapply(predicted@items, function(x){
  colnames(ratings_movies)[x]
})

recommender_matrix[, 2:4]
```

```
##       2
## [1,] "Princess Bride, The (1987)"
## [2,] "Raiders of the Lost Ark (1981)"
## [3,] "One Flew Over the Cuckoo's Nest (1975)"
## [4,] "Fugitive, The (1993)"
## [5,] "Terminator, The (1984)"
##       3                              5
## [1,] "Star Wars (1977)"             "Godfather, The (1972)"
## [2,] "Princess Bride, The (1987)"   "Leaving Las Vegas (1995)"
## [3,] "Fargo (1996)"                 "Boot, Das (1981)"
## [4,] "Empire Strikes Back, The (1980)" "Schindler's List (1993)"
## [5,] "Raiders of the Lost Ark (1981)"  "Nikita (La Femme Nikita) (1990)"
```

## View the most frequently recommended movies (UBCF model)

```
items <- factor(table(recommender_matrix))
items <- sort(items, decreasing = TRUE)
top_items <- data.frame(names(items), items)
head(top_items)
```

```
##                                                         names.items.
## One Flew Over the Cuckoo's Nest (1975) One Flew Over the Cuckoo's Nest (1975)
## Usual Suspects, The (1995)                        Usual Suspects, The (1995)
## Princess Bride, The (1987)                        Princess Bride, The (1987)
## Schindler's List (1993)                              Schindler's List (1993)
## Godfather, The (1972)                                  Godfather, The (1972)
## Silence of the Lambs, The (1991)          Silence of the Lambs, The (1991)
##                                          items
## One Flew Over the Cuckoo's Nest (1975)    123
## Usual Suspects, The (1995)                119
## Princess Bride, The (1987)                101
## Schindler's List (1993)                    98
## Godfather, The (1972)                       97
## Silence of the Lambs, The (1991)           97
```

Since we have now implemented both IBCF and UBCF systems, let's evaluate the models!

# First, evaluating the IBCF model

```
folds <- 4
items_keep <- 15
rating_threshold <- 3

eval_sets <- evaluationScheme(data = ratings_movies, method =
                          "cross-validation", k = folds,
                      given = items_keep, goodRating = rating_threshold)

eval_model <- "IBCF"
parameters <- NULL

eval_recommender <- Recommender(data = getData(eval_sets, "train"),
                          method = eval_model, parameter = parameters)


n_recommend <- 5

eval_predicted <- predict(object = eval_recommender, newdata =
                      getData(eval_sets, "known"), n=n_recommend,
                  type = "ratings")

eval_accuracy <- calcPredictionAccuracy(x = eval_predicted,
                              data = getData(eval_sets,
                                          "unknown"),
                          byUser = FALSE)
eval_accuracy
```

```
##     RMSE      MSE      MAE
## 1.434249 2.057070 1.092191
```

# Now, evaluating the UBCF model

```
folds <- 4
items_keep <- 15
rating_threshold <- 3

eval_sets <- evaluationScheme(data = ratings_movies, method =
                              "cross-validation", k = folds,
                      given = items_keep, goodRating = rating_threshold)

eval_model <- "UBCF"
parameters <- NULL

eval_recommender <- Recommender(data = getData(eval_sets, "train"),
                          method = eval_model, parameter = parameters)

n_recommend <- 5

eval_predicted <- predict(object = eval_recommender, newdata =
                      getData(eval_sets, "known"), n=n_recommend,
                  type = "ratings")

eval_accuracy <- calcPredictionAccuracy(x = eval_predicted,
                                  data = getData(eval_sets,
                                                "unknown"),
                              byUser = FALSE)
eval_accuracy
```

```
##      RMSE       MSE       MAE
## 0.9785090 0.9574799 0.7716670
```

# Analysis

The User-Based Collaborative Filtering recommender system outperformed the IBCF