

# Project 4

*Max Wagner*

*July 7, 2016*

---

## Libraries

```
library(reshape2)
library(dplyr)
library(recommenderlab)
library(NMF)
library(ggplot2)
library(stringr)
library(knitr)
```

## Rec Sys

The MovieLens data was put into a smaller, more compact version that includes only user and movies with sufficient information by the University of Minnesota. One of the problems with the data set is that time is not taken into consideration when looking at user's recommendations. The set provides a timestamp for each recommendation, albeit in a strange format. The timestamp is in number of seconds from midnight on January 1, 1970. Using this information, it will be possible to weight ratings that are closer to present day higher than older ratings.

## Data

The first step is to load in the data from the small data set, and then separate the dates into a separate data frame, and convert the seconds to years. The year of the rating, and the year of the movie will be added to the ratings table. We can then preview the movies and ratings tables to make sure they appear normal.

```
#load in data from github
movies <- read.csv("https://raw.githubusercontent.com/maxwagner/643/master/ml-latest-small/movies.csv",
ratings <- read.csv("https://raw.githubusercontent.com/maxwagner/643/master/ml-latest-small/ratings.csv")

#get the rating dates
dates <- as.vector(ratings[,4])
ratings <- ratings[,-4]
ratings$r.year <- as.integer(dates/60/60/24/365 + 1970)

#get the movie dates
ratings <- inner_join(ratings, movies[,1:2],by = "movieId")
colnames(ratings)[5] <- "m.year"
ratings$m.year <- as.integer(str_sub(ratings$m.year, -5, -2))

kable(head(movies))
```

movieId	title	genres
1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy
2	Jumanji (1995)	Adventure Children Fantasy
3	Grumpier Old Men (1995)	Comedy Romance
4	Waiting to Exhale (1995)	Comedy Drama Romance
5	Father of the Bride Part II (1995)	Comedy
6	Heat (1995)	Action Crime Thriller

```
kable(head(ratings))
```

userId	movieId	rating	r.year	m.year
1	16	4.0	2008	1995
1	24	1.5	2008	1995
1	32	4.0	2008	1995
1	47	4.0	2008	1995
1	50	4.0	2008	1995
1	110	4.0	2008	1995

The above all seems reasonable for the information we have so far. One interesting approach to take when weighting based on time is to scale the rating itself, instead of scaling the recommender values. I'll take this approach for the purpose of this project. To decide how to weight the rating, we should include both the year of the rating, and the year of the movie. The most straightforward method is to scale based on where the most recent ratings and movies receive the least penalty to their score. For instance, a rating made in 2016 will be multiplied by 1, where a rating made in 1997 might be multiplied by 0.2. Scaling in this method keeps all values above 0, and below 5.

```
# first replace the few rows out the 100000 that had issues with the movie year
ratings$m.year[is.na(ratings$m.year)] <- mean(ratings$m.year, na.rm = TRUE)
```

```
# scale it
```

```
ratings$scale <- (ratings$r.year^25 / max(ratings$r.year)^25) * (ratings$m.year^25 / max(ratings$m.year)^25)
ratings$rating.scaled <- ratings$rating * ratings$scale
summary(ratings$scale)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.2248  0.5997  0.6630  0.6719  0.7606  1.0000
```

```
summary(ratings$rating)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.500   3.000   3.500   3.517   4.000   5.000
```

```
summary(ratings$rating.scaled)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.132   1.799   2.369   2.353   2.931   5.000
```

The summary of the scale show that on average a rating will be scaled by 0.67, while the lowest multiplier is 0.22, and the highest is 1.0. The overall ratings dropped dramatically for both mean and median, but the min and max remain within 0 and 5, which is desired.

Next, the ratings table is put into a real ratings matrix format using `recommenderlab`. In partial thanks to Sreejaya Nair and Suman K Polavarapu for the cleaner way to get to a ratings matrix with `acast`. I'll do this twice. Once for the original ratings, and once for the scaled ratings.

```
ratings_hori <- acast(ratings, userId ~ movieId, value.var="rating")
ratings_rm <- as(ratings_hori, "realRatingMatrix")
ratings_rm_r <- ratings_rm[, colCounts(ratings_rm) > 20]
```

## UBCF - Original Ratings

First, the original ratings will be placed into a UBCF recommender to get a baseline without any context weighting.

```
movies.recom <- Recommender(ratings_rm_r[1:400], method = "UBCF")
movies.predict <- predict(movies.recom, ratings_rm_r[c(500, 525, 600)],n=3)
movies.predict <- as.data.frame(as(movies.predict, "list"))
kable(movies.predict)
```

X500	X525	X600
318	296	541
356	318	1641
593	356	1213

In this case, for user 500 the top movies suggested are: Shawshank Redemption, The (1994), Wes Craven's New Nightmare (Nightmare on Elm Street Part 7: Freddy's Finale, A) (1994), and Silence of the Lambs, The (1991). We can see that each movie recommended is older in this selection.

## UBCF - Scaled Ratings

Now to try the same recommender with the weighted ratings.

```
ratings_hori <- acast(ratings, userId ~ movieId, value.var="rating.scaled")
ratings_rm <- as(ratings_hori, "realRatingMatrix")
ratings_rm_r <- ratings_rm[, colCounts(ratings_rm) > 20]
movies.recom <- Recommender(ratings_rm_r[1:400], method = "UBCF")
movies.predict <- predict(movies.recom, ratings_rm_r[c(500, 525, 600)],n=3)
movies.predict <- as.data.frame(as(movies.predict, "list"))
kable(movies.predict)
```

X500	X525	X600
318	296	50
48516	318	4886
7153	2028	4973

We can see that *Shawshank Redemption, The* (1994) is still the top ranked movie because of the original high rating it was given. But now *Departed, The* (2006) is 2nd, and *Lord of the Rings: The Return of the King, The* (2003) is now 3rd. Both of these are more recent movies, which makes sense given the weighting system.

---

The above is a basic version of what could be done. It is also possible to weight by something like genre. I had stated in a earlier project that each user tends to prefer certain genres, and some genres appear more than others. For instance, in this set the following frequencies can be found. It may be viable to scale based on the rarity of the genre as well as the year of the rating and movie.

```
movies.o <- movies[,3]
genres.split <- unlist(strsplit(movies.o, split = "|", fixed = TRUE))
genres.count <- data.frame(table(genres.split))
kable(head(genres.count[order(-genres.count$Freq),],10))
```

	genres.split	Freq
9	Drama	5220
6	Comedy	3515
18	Thriller	2187
16	Romance	1788
2	Action	1737
7	Crime	1440
3	Adventure	1164
12	Horror	1001
17	Sci-Fi	860
15	Mystery	675

The point of the above is that there are any number of scaling values, and scaling by just the year isn't a realistic method for getting the best results. Many different aspects should be considered, not just temporal ones.

## Citations

[https://github.com/srajeev1/MSDA-IS643/blob/master/projects/Project2/DATA643\\_Project\\_2.Rmd](https://github.com/srajeev1/MSDA-IS643/blob/master/projects/Project2/DATA643_Project_2.Rmd)