

MovieLens Report

Introduction

I set out to make a movie recommendation system based on the MovieLens dataset, which provides data on many users and how they rated many different movies. For my analysis I am using a smaller subset of the dataset as my computer does not have the memory or processing power to handle the larger dataset. You can find both datasets at

<http://files.grouplens.org/datasets/movielens/ml-latest-small.zip>

and

<http://files.grouplens.org/datasets/movielens/ml-latest.zip>

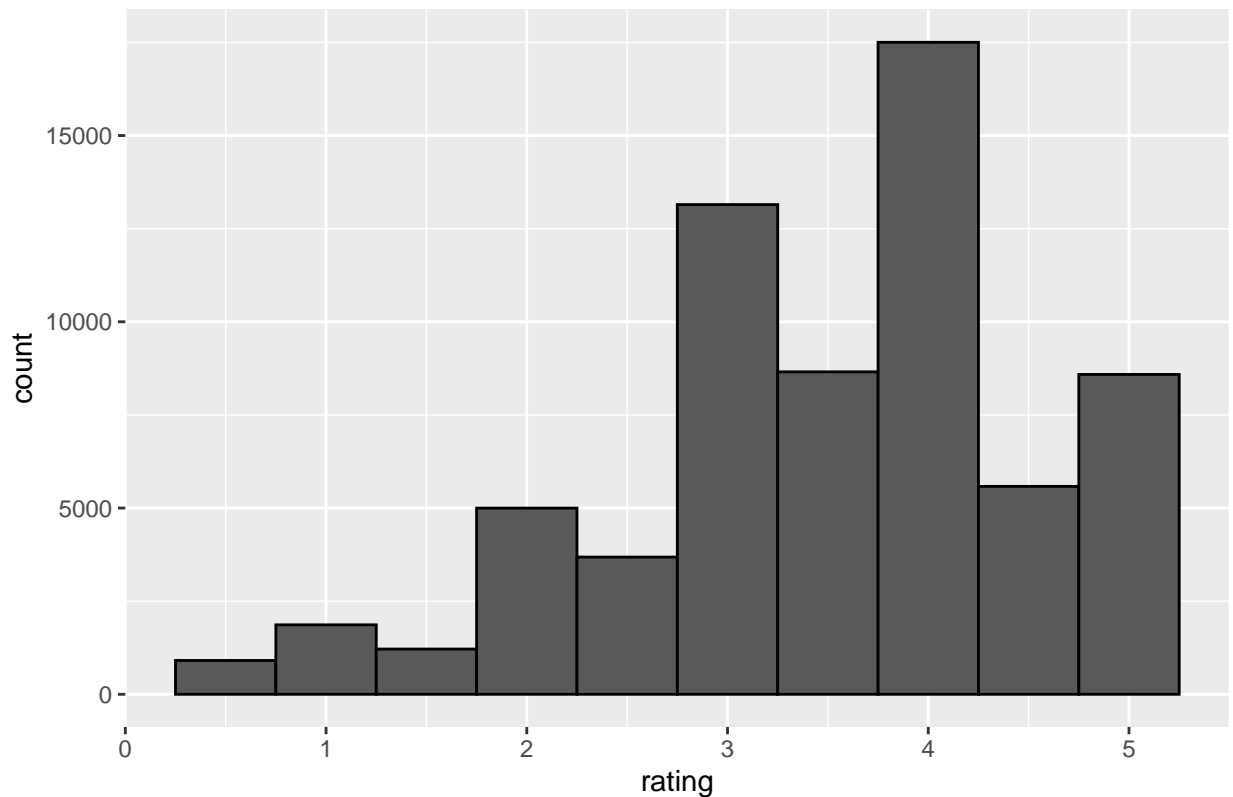
You may also recreate my results by running the 3 scripts download-data, train-validation-split, and model in that order. I have included these scripts with this report. If the code in these reports is not working, try running at least the first 2 scripts and running the report again. I have also uploaded the datasets with this report in case the website ever stops working.

The last thing I did before beginning work with the data set was to set aside 20% of the data set at random to use for validating my final model to prevent overfitting.

Modeling

The first thing I did split out another test set of 20% of the data to test various different models on. I then proceeded to graph a histogram of all the ratings in the training set to see the distribution of the data.

Histogram of Movie Ratings

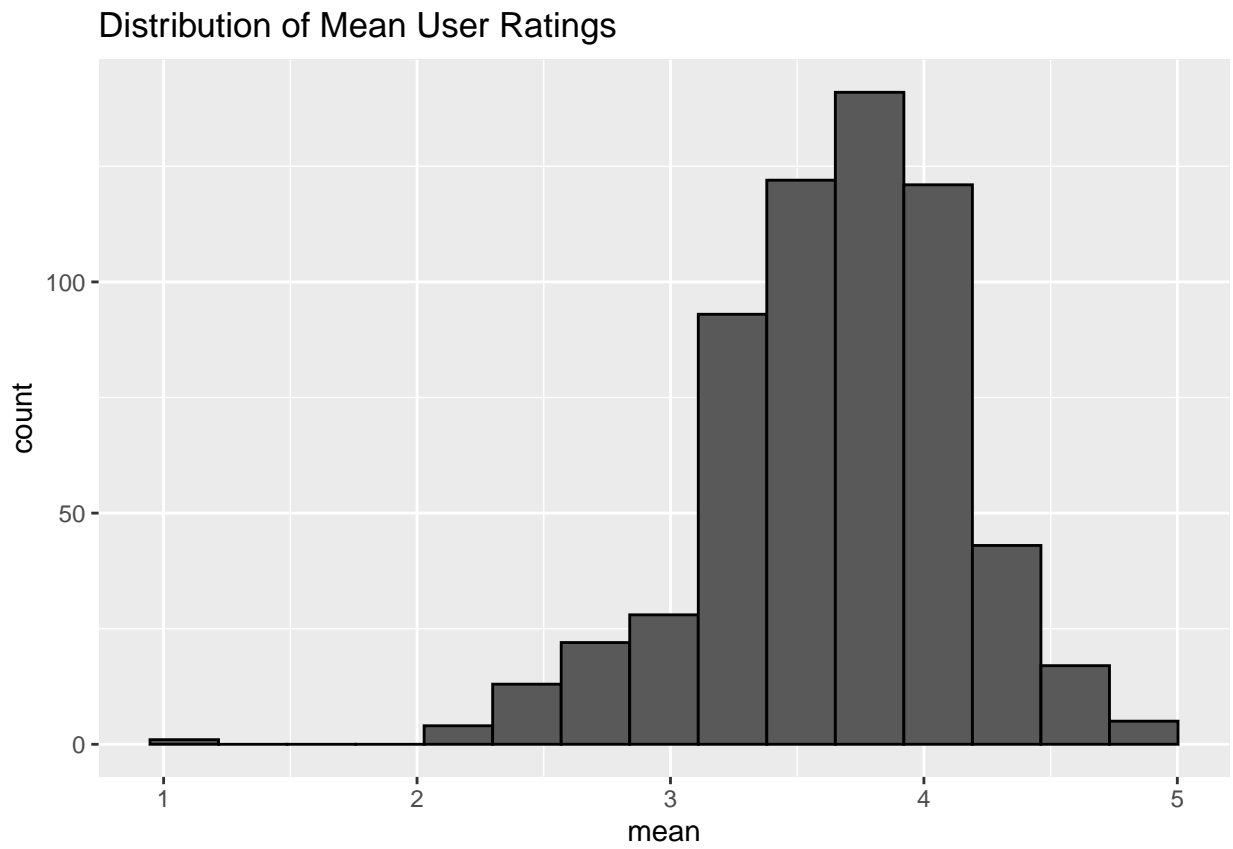


From here I proceeded to calculate the summary statistics of the ratings.

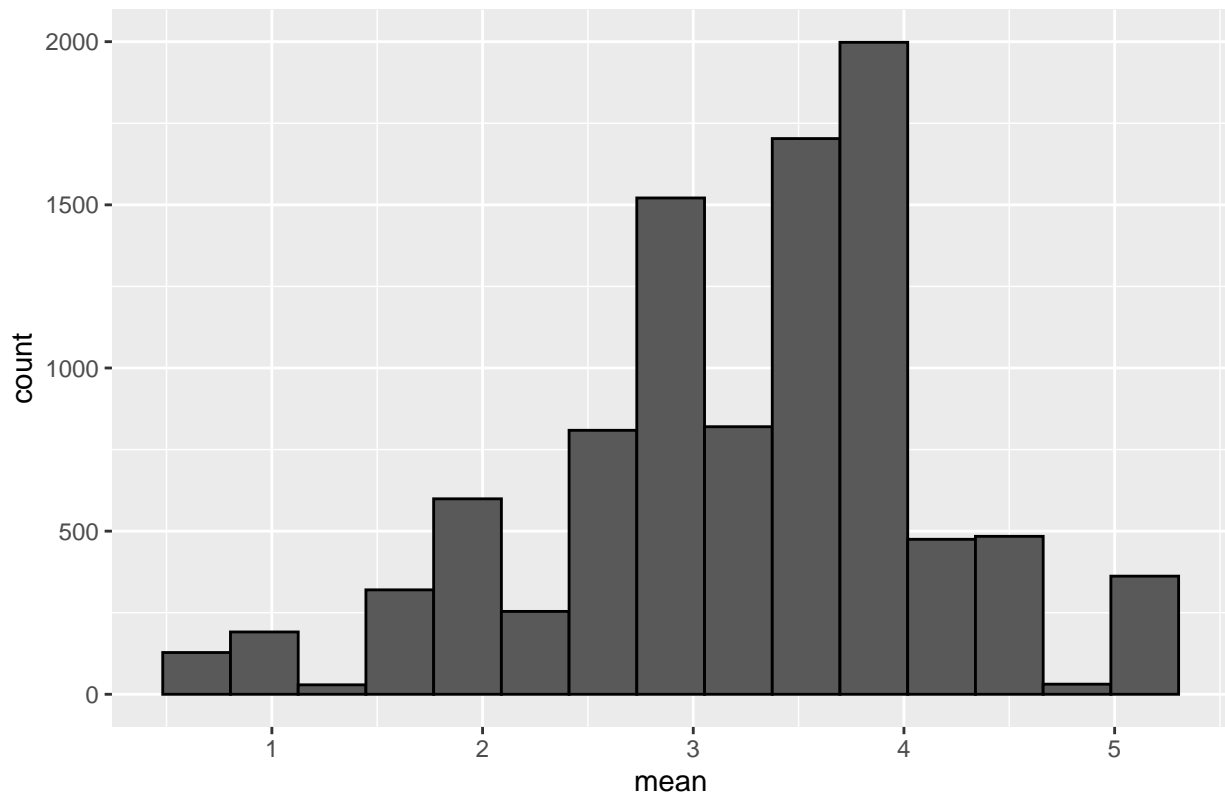
```
##      mean median      sd
## 1 3.494633    3.5 1.044586
```

We can see that the mean value of the ratings is 3.49. I used this as a default prediction to see how accurate a model is that simply guesses the mean value for any and all ratings. I evaluated the models based on root mean squared error (RMSE) which is a measurement of how far my predictions are from the actual ratings. Simply guessing the mean rating for all ratings yields a RMSE of 1.036, which is my benchmark to beat.

The next step in the modeling process was to examine different user and movie averages. It stands to reason that certain movies are better than other movies and should be receiving higher ratings, and also that certain users are more critical than others and will give lower ratings. We can see the distribution of both users and movies here.



Distribution of Mean Movie Ratings



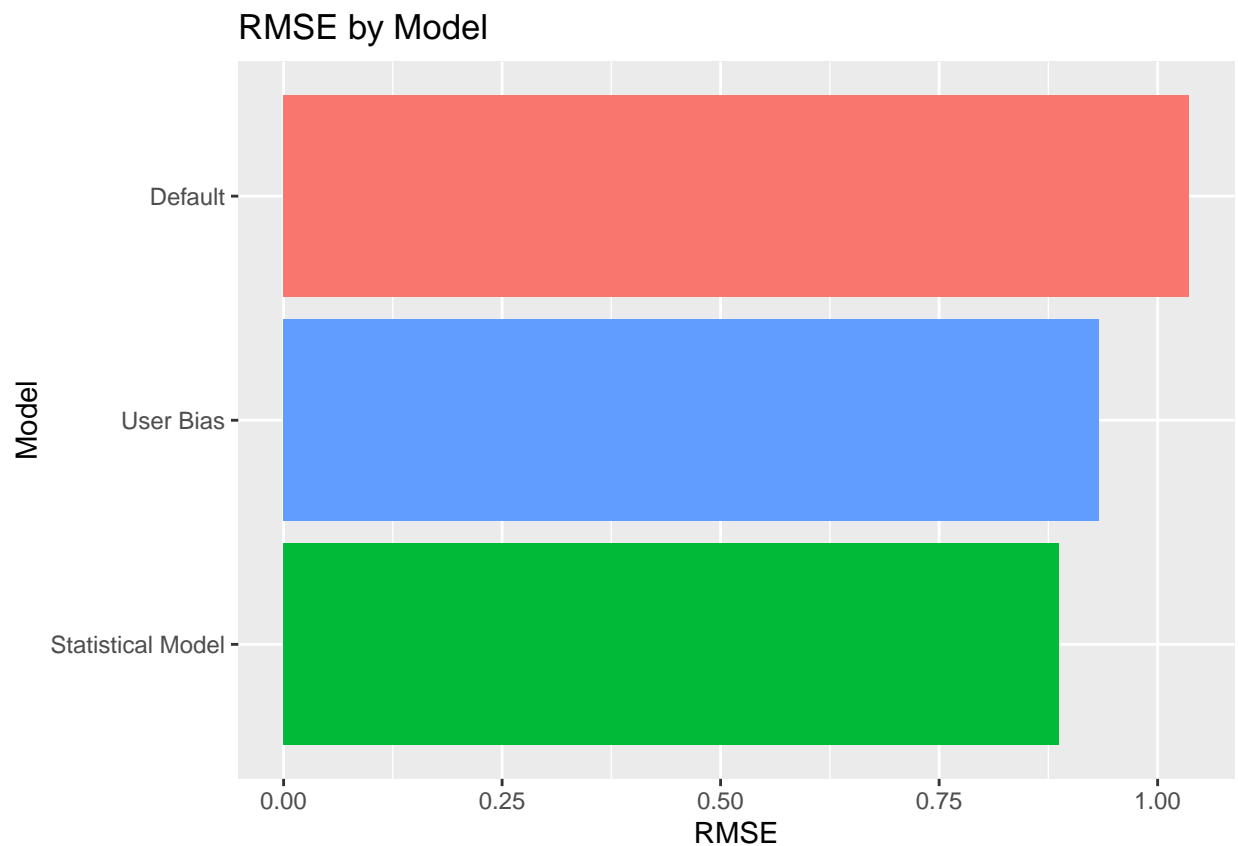
We see that as expected, there is variety in the mean rating across movies and users. My next step is to factor this variance into my model to hopefully improve my predictions. I chose to first include the difference in users, which I termed “user bias”. User bias was calculated by subtracting the total mean rating of 3.49 from all ratings and then taking the average rating for each user, which shows the users average difference from the mean rating.

```
## # A tibble: 610 x 2
##   userId user_bias
##   <int>    <dbl>
## 1     1      0.858
## 2     2      0.388
## 3     3     -0.668
## 4     4      0.179
## 5     5      0.398
## 6     6     -0.0405
## 7     7     -0.161
## 8     8     -0.110
## 9     9     -0.379
## 10    10     -0.245
## # ... with 600 more rows
```

If I add the user bias to the default prediction of 3.49, my predictions should more accurately reflect the true ratings. Incorporating user bias into the model decreases the RMSE to .93, which is a significant improvement. I now look to add the movie bias into the model, but the movie bias must account for the fact that I have already added the user bias. I will calculate the movie bias by subtracting both the overall mean rating and the user bias from each rating, and the mean of the resulting ratings will show whether certain movies are above or below average.

```
## # A tibble: 9,724 x 2
##   movieId movie_bias
##   <dbl>    <dbl>
## 1      1      0.368
## 2      2     -0.0235
## 3      3     -0.193
## 4      4     -1.20
## 5      5     -0.540
## 6      6      0.370
## 7      7     -0.316
## 8      8     -0.716
## 9      9     -0.348
## 10     10      0.0246
## # ... with 9,714 more rows
```

We can see that movie bias also plays a role as some movies are above average and some are below average. As such, a new model given by the overall mean + the user bias + the movie bias should improve our predictions. This newest model gives an RMSE of .89, another improvement.



If we examine our model, we can see something slightly odd. Let us see the top 10 movies as predicted by our model.

```
## # A tibble: 10 x 2
##   title          mean
##   <chr>         <dbl>
## 1 'Salem's Lot (2004) 5
```

##	2	12 Angry Men (1997)	5
##	3	12 Chairs (1976)	5
##	4	20 Million Miles to Earth (1957)	5
##	5	3-Iron (Bin-jip) (2004)	5
##	6	42nd Street (1933)	5
##	7	61* (2001)	5
##	8	7 Faces of Dr. Lao (1964)	5
##	9	9/11 (2002)	5
##	10	A Detective Story (2003)	5

We get back some blockbusters we would expect like 12 Angry Men, alongside movies that we would not expect to be the top 10 movies such as 3-Iron (Bin-jip). This is because the model is not accounting for the number of ratings each movie has received, and cannot tell the difference between movies that have received 1000 ratings and movies that have received 1 rating. It becomes clear we need to add a regularization term that makes our model less confident in predictions with low sample sizes, causing these predictions to be closer to the mean value.

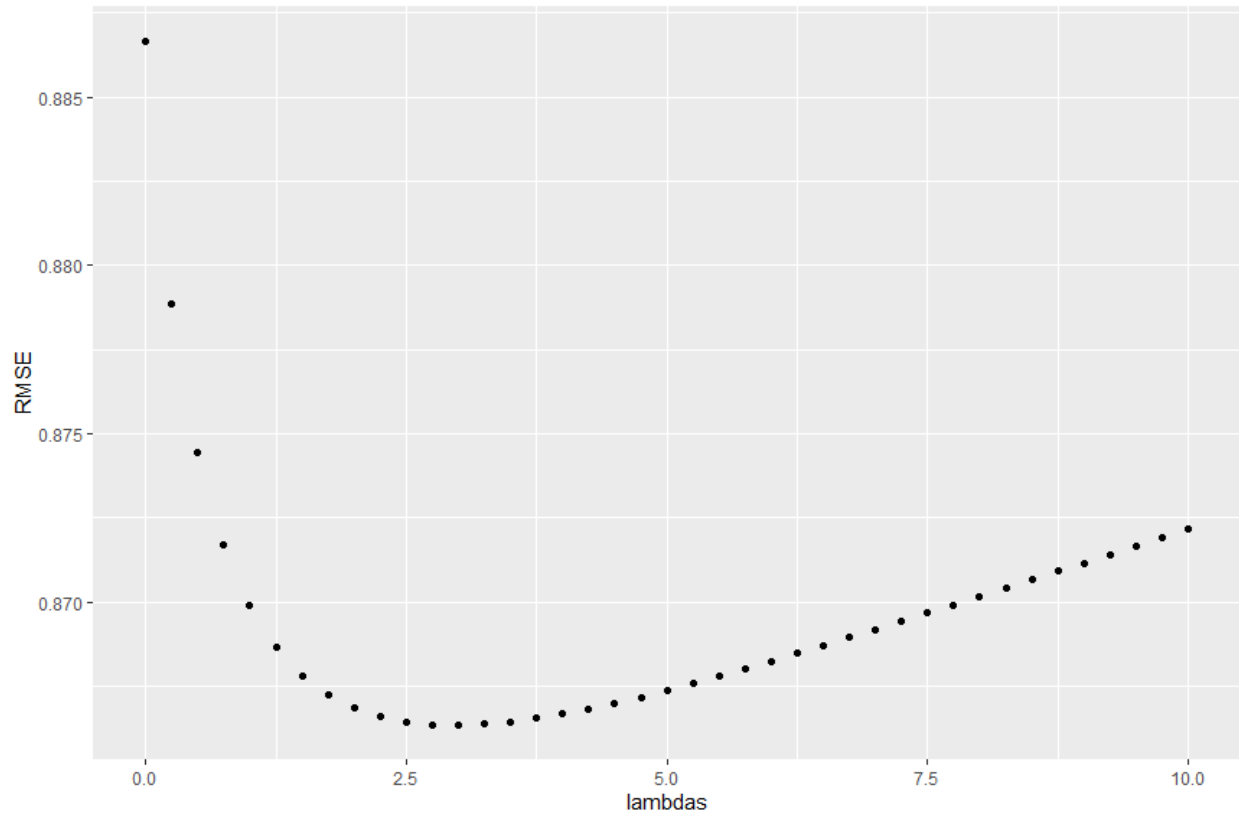
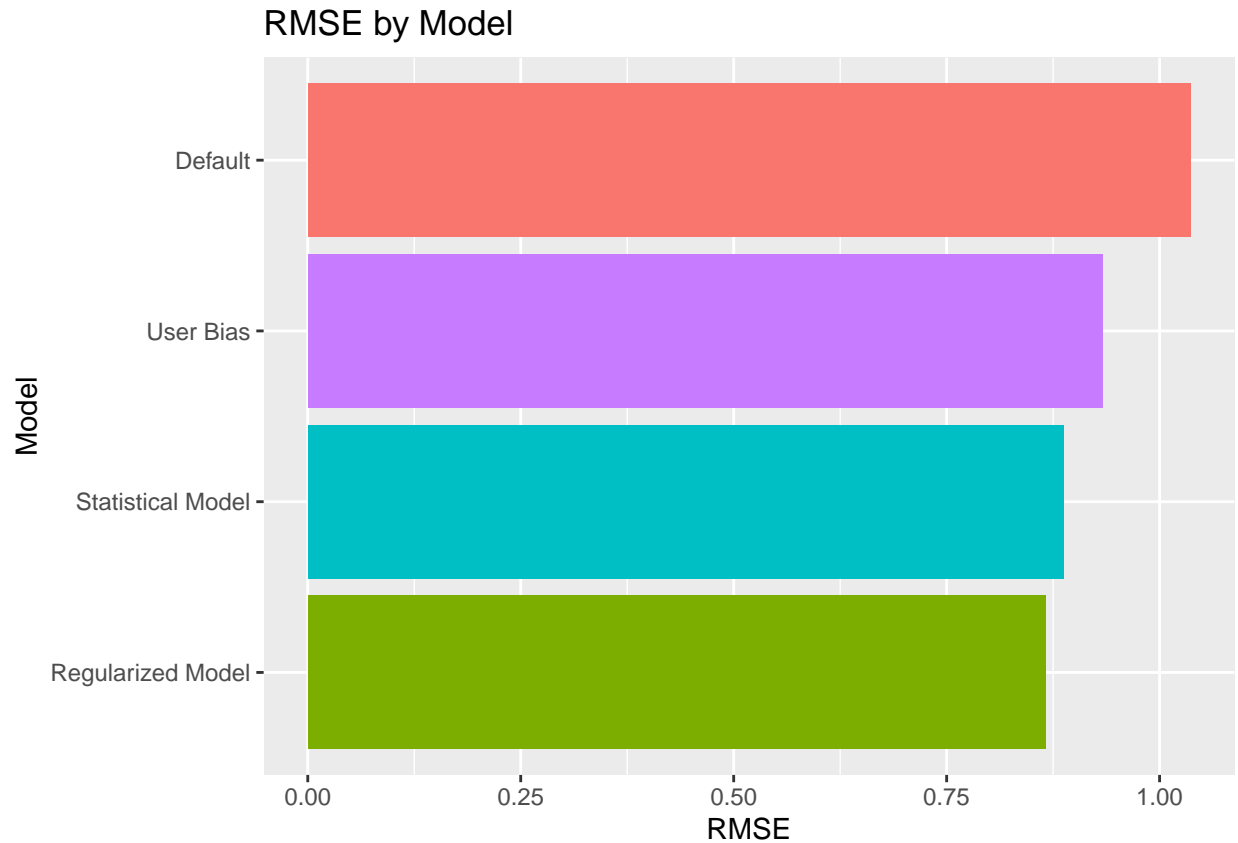


Figure 1: Optimization of the regularization parameter

Here we see that adding a regularization term λ has a positive effect on the RMSE up until approximately $\lambda = 3$, at which point we have reached the optimal amount of regularization and any increase also increases the RMSE of our model. So for our model, we will add the regularization term of $\lambda = 3$ which decreases our RMSE down to .87, another improvement.



Machine Learning with Matrix Factorization

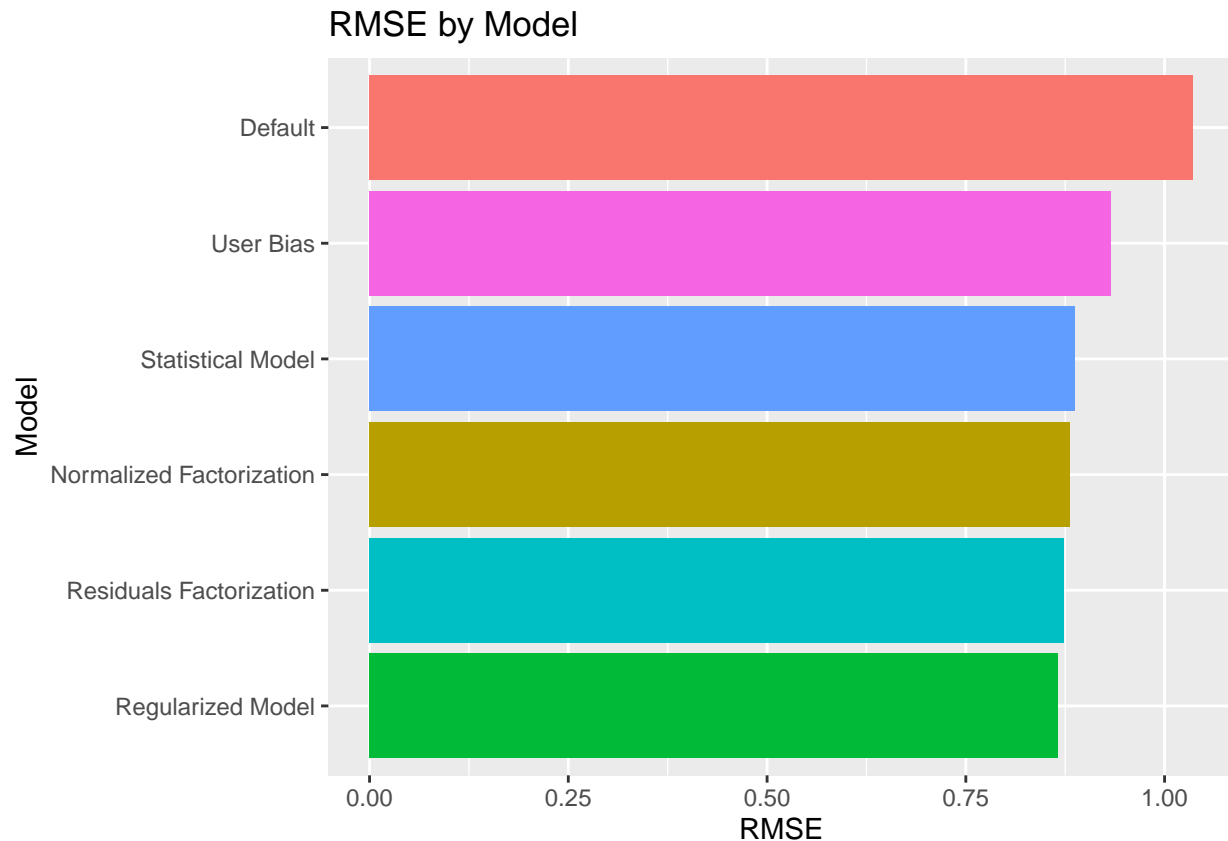
At this point I felt the model was not going to improve any further without using machine learning, so I attempted to improve the model using recommenderlab. At first I attempted to create a machine learning algorithm using Matrix Factorization. I create a matrix with the users representing the rows and the movies representing the columns. Said matrix had many missing values as not all users had rated every movie. Here is a small snippet of said matrix.

```
##      1  2  3  4  5  6  7  8  9 10
## [1,]  4 NA  4 NA NA  4 NA NA NA NA
## [2,] NA NA NA NA NA NA NA NA NA NA
## [3,] NA NA NA NA NA NA NA NA NA NA
## [4,] NA NA NA NA NA NA NA NA NA NA
## [5,]  4 NA NA NA NA NA NA NA NA NA
## [6,] NA NA  5  3 NA NA  4  3 NA  3
## [7,] NA NA NA NA NA NA NA NA NA NA
## [8,] NA NA NA NA NA NA NA NA NA  2
## [9,] NA NA NA NA NA NA NA NA NA NA
## [10,] NA NA NA NA NA NA NA NA NA NA
```

I proceeded to normalize the ratings in the matrix by subtracting the mean of each column from said column, and the mean of each row from said row. I then fit a recommender using the recommenderlab package on the normalized data, and predicted the ratings of the normalized test set. The recommender had a RMSE

of only .12 on the training set but of .88 on the test set, indicating that the model was seriously overfitting the training data.

I then tried a slightly different approach. Instead of giving the normalized ratings to the recommender, I gave the recommender the variance that was not explained by the statistical model I had created earlier. This was done by subtracting the model prediction from each actual rating, and whatever resulted was data not accounted for by the model. Using this data instead of the normalized ratings, we had a training RMSE of .11 and test RMSE of .87, still massively overfitting. At this point I decided to switch models to a model that could add a regularization term to reduce overfitting.

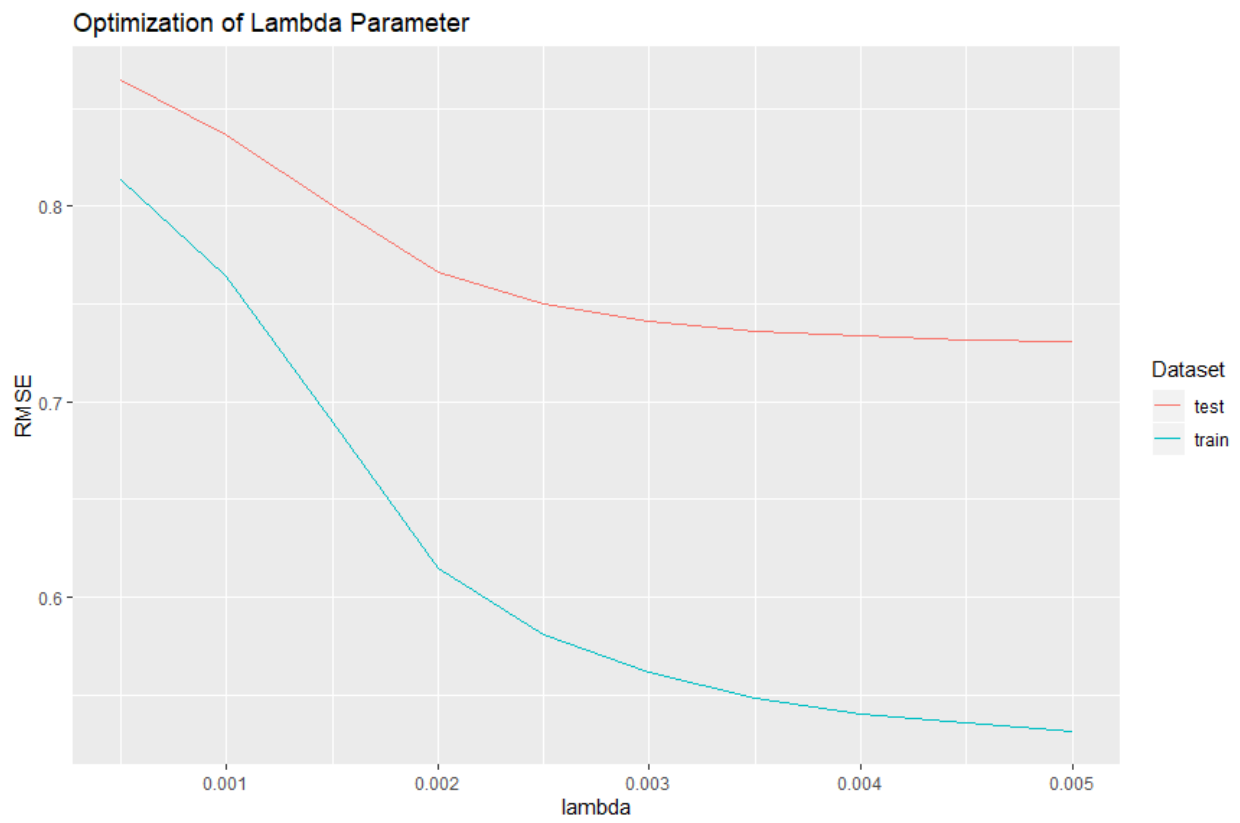
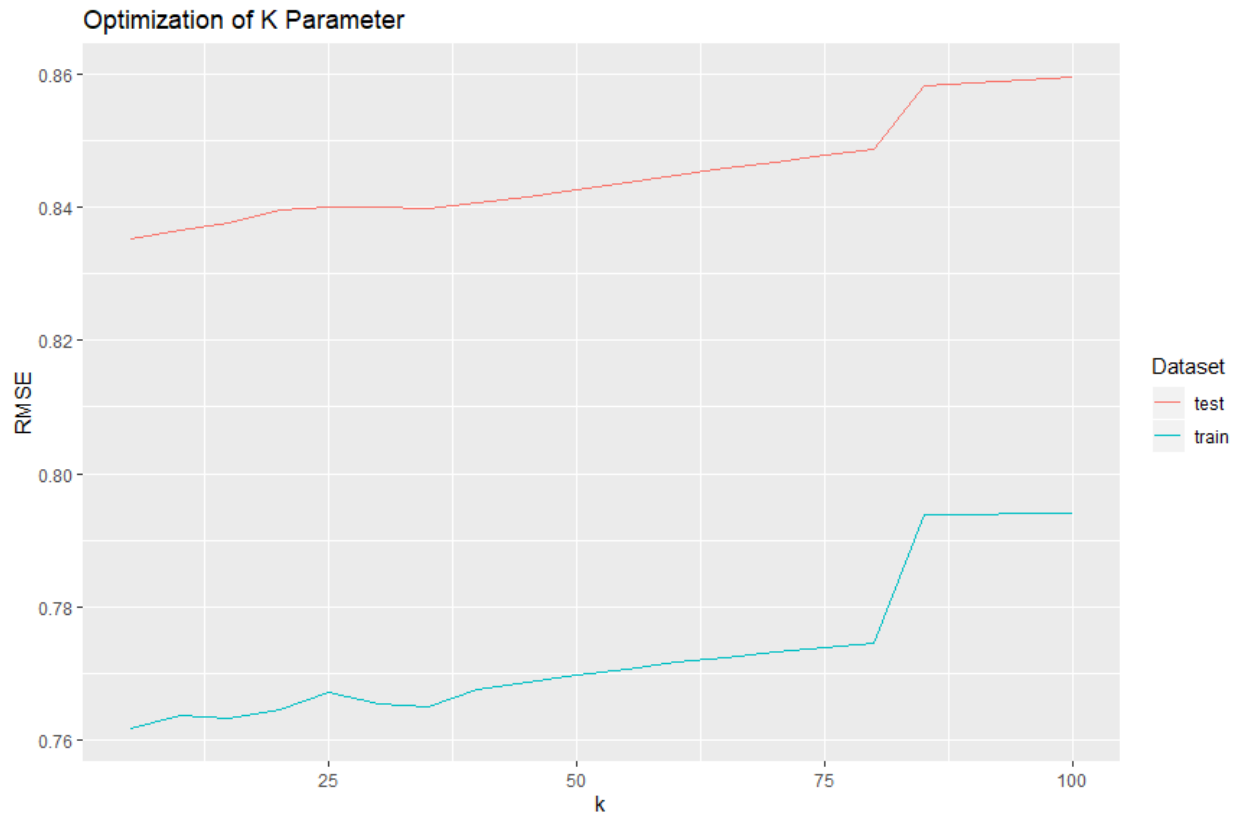


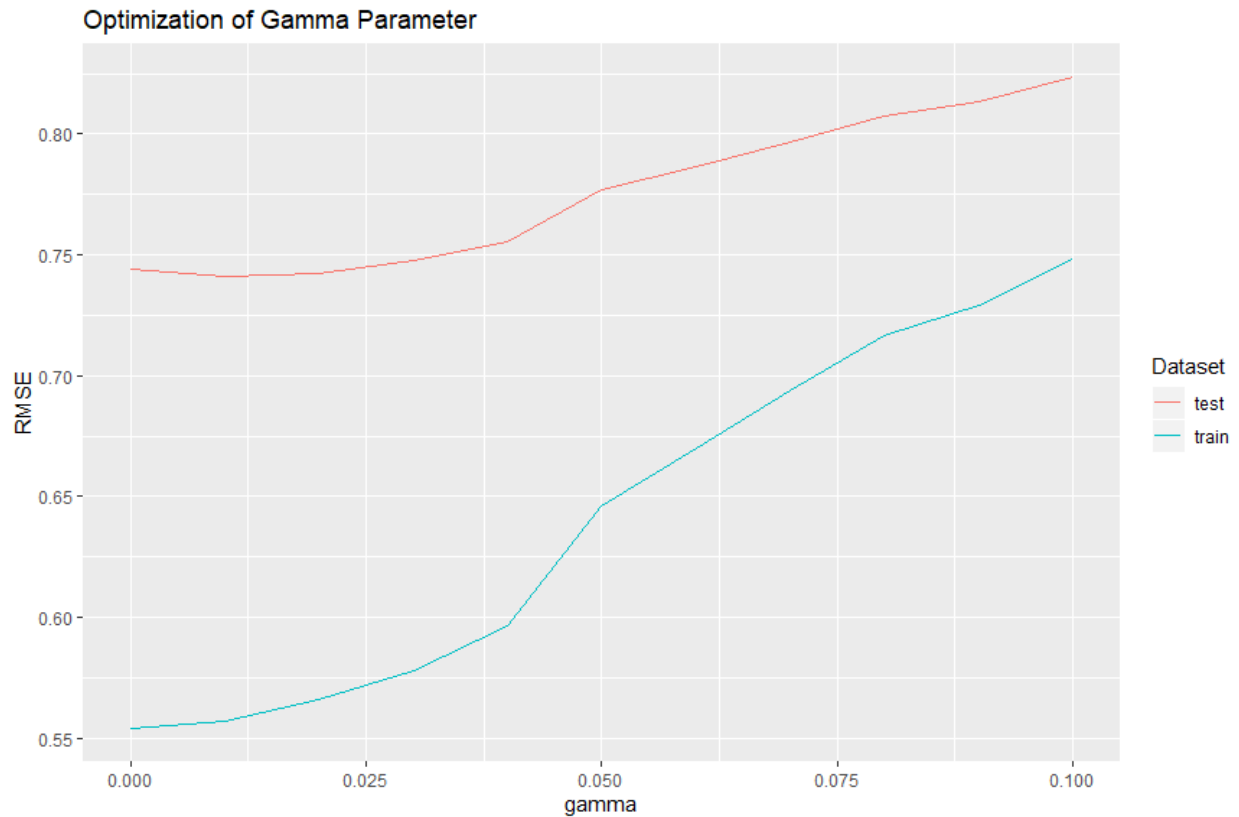
Machine Learning with Funk Singular Value Decomposition

This approach is similar to matrix factorization in that it requires a matrix, but instead of finding its principal components, it will break down the matrix into two vectors whose crossproduct returns as close to the original matrix as possible. Running singular value decomposition on the normalized ratings gave a training RMSE of .77 and a test RMSE of .87. There was much less overfitting occurring in this model, but unfortunately it also didn't outperform our simple statistical model.

I then ran singular value decomposition on the variance not explained by the statistical model. the training RMSE was .76 and the test RMSE was .84, an improvement and our best RMSE to date. I decided to select this model to move forward with and optimize the parameters for. There are several parameters to optimize for this algorithm: K, the minimum #of ratings a movie needs to have to be considered; gamma, the regularization parameter; and lambda, the learning rate.

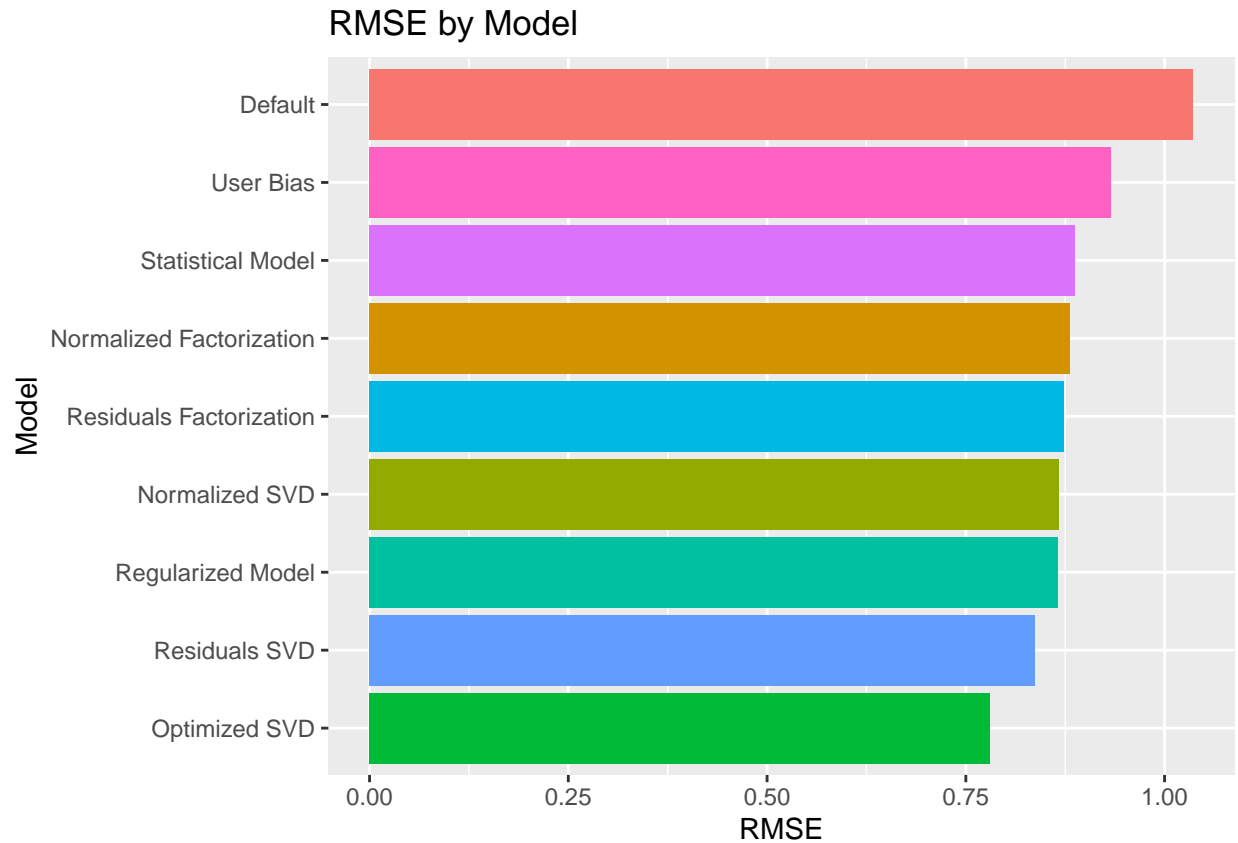
Here is the performance based on each optimization I ran.





We can see that low values of K are optimal, so I chose to use the default value of 10. For λ , I chose a value of .003 which is where the curve for the test set begins to flatten indicating no further improvement. However, this has introduced some overfit as there is significant distance between the training and test sets. To adjust for this, I chose a slightly high value of gamma of .05 which is not optimal on the test set, but partially closes the gap between training and test sets.

Here we have achieved an RMSE of .78, and the last step is to check this against the validation set to see if the model is overfit or if it will hold up on data it has never seen.



Validation Data Set

The validation dataset has not been used at all in the analysis. I transformed the validation set into variance not predicted by the statistical model and ran the funkSVD model on the result. The end result was an RMSE of .806, only slightly worse than the performance on the test set and still better than our previous best model. This indicates that the model is not badly overfit and likely has predictive power in the future. Here is a summary of all the models and their performance

