

Recommender Systems

Signal Data Science

In this assignment, we'll explore one way to make a [recommender system](#), something which predicts the rating a user would give to some item. Specifically, we'll be using [collaborative filtering](#) on the [MovieLens 1M Dataset](#), a set of one million different movie ratings. Collaborative filtering operates on the assumption that if one person A has the same opinion as another person B on item X , A is *also* more likely to have the same opinion as B on a *different* item Y than to have the opinion of a randomly chosen person on Y .

Collaborative filtering is a type of [unsupervised learning](#) and can serve as a *prelude* to dimensionality reduction (*e.g.*, with PCA or factor analysis) because filling in missing values is typically required for such methods. Specifically, we will be working with an [imputation](#)-based method of collaborative filtering, which infers *all* of the missing values from the given data.

In the following, write up your work in an R Markdown file with elaboration about *what* you're doing at each step and *why* you're doing it. Include interpretation of results as well whenever appropriate. Your goal should be to produce, at the end, an HTML (or PDF) file from the R Markdown writeup that gives a coherent and reasonably accessible description of the process you followed, the reasoning behind each step, and the results attained at the end.

Getting started

We'll first need to spend some time preparing the data before we can use any collaborative filtering methods.

- Download the [MovieLens 1M Dataset](#). Read the associated [README.txt](#), which describes the contents of the dataset.
- The first 5 lines of `ratings.dat` are:

```
1::1193::5::978300760
1::661::3::978302109
1::914::3::978301968
1::3408::4::978300275
1::2355::5::978824291
```

Use `read.csv()` with the appropriate options to load the file into R. The resulting data frame should have **1000209 rows** and **7 columns**.

- Restrict to the columns containing user IDs, movie IDs, and movie ratings. Name the columns appropriately.
- Compute the sets of `unique()` user IDs and movie IDs as well as the mean rating given. Compare the numbers of different user IDs and movie IDs with the *maximum* user ID and movie ID.
- Set the seed for consistency. Generate a training set using 80% of the data and a test set with the remaining 20%.

Because there are some movies which are rated by very few people and some people who rated very few movies, we have two corresponding problems: (1) there will be movies in the test set which were not rated by any people in the training set and (2) there will be people in the test set who do not show up in the training set. As such, we need to add to the training set a fake movie rated by every user and a fake user who rated every movie.

- Create two data frames corresponding to the above fake data using the previously calculated mean rating. (For the fake movie and user respectively, use a movie ID and user ID which are both 1 greater than their respective maximum values in the entire dataset.) When creating the fake user who has rated every movie, allow the movie IDs to range from 1 to the maximum movie ID in the dataset (which will include movie IDs not present in the dataset). The fake user should not have a rating for the fake movie.
- Perturb the ratings of the fake data slightly by adding a normally distributed noise term with mean 0 and standard deviation 0.01. Add your fake data to the training data frame, which should increase in size by **9994 rows**.

Next, we need to create a matrix containing rating data for (user, movie) pairs. We can store this as a *sparse* matrix, which is a special data structure designed for handling matrices where only a minority of the entries are filled in (because each user has only rated a small number of movies).

- Use `Incomplete()` to generate a sparse ratings matrix with one row per user ID and one column per movie ID. The resulting matrix should have **6041 rows** and **3953 columns**.

Alternating least squares imputation

We will proceed to use the method of alternating least squares (ALS) to impute the missing entries in the sparse ratings matrix.¹

Formally, we would like to fill in the missing entries of the matrix \mathbf{X} , and we do so by (1) expressing the *filled-in matrix* \mathbf{Z} as a function of two different matrices \mathbf{A} and \mathbf{B} and (2) minimizing a function of \mathbf{A} and \mathbf{B} . From the form of the minimization problem we can characterize the minimal solution for each of the two matrices as a function of the other one. As such, this suggests an iterative strategy where we use \mathbf{A} to estimate the optimal \mathbf{B} and vice versa, repeating until the values stabilize. We do so via L^2 regularized (“ridge”) regression, controlled by a regularization parameter λ . This is the method of alternating least squares.

The final filled-in ratings matrix can be expressed in terms of a *reduced-rank singular value decomposition*, where $\mathbf{Z} = \mathbf{U}\mathbf{D}\mathbf{V}^T$ and \mathbf{D} is a diagonal matrix containing the *singular values* of \mathbf{Z} . We say that the solution has reduced rank because as we increase λ , the *rank* of \mathbf{D} (i.e., the number of nonzero values on the diagonal of \mathbf{D}) decreases. In addition, we can think of the *rows* of \mathbf{U} and the *columns* of \mathbf{V}^T as being “factor scores” for users and movies. As such, we can (loosely) think of ALS as simultaneously performing imputation *and* dimensionality reduction on the ratings matrix.

First, we need to calculate what values of the regularization parameter λ we’ll search over.

- Use `biScale()` to scale both the columns and the rows of the sparse ratings matrix with `maxit=5` and `trace=TRUE`. You can ignore the resulting warnings (increasing the number of maximum iterations doesn’t improve the outcome, which you can verify for yourself).

`lambda0()` will calculate the lowest value of the regularization parameter which gives a zero matrix for \mathbf{D} , i.e., drives all rating estimates to zero.

- Use `lambda0()` on the scaled matrix and store the returned value as `lam0`.
- Create a vector of λ values to test by (1) generating a vector of 20 *decreasing* and uniformly spaced numbers from `log(lam0)` to 1 and then (2) calculating e^x with each of the previously generated values as x . You should obtain a vector where entries 1 and 5 are respectively 103.21 and 38.89.
- Initialize a data frame `results` with three columns: `lambda`, `rank`, and `rmse`, where the `lambda` column is equal to the previously generated sequence of values of λ to test. Initialize a list `fits` as well to store the results of alternating least squares for each value of λ .

¹Hastie *et al.* (2014), [Matrix Completion and Low-Rank SVD via Fast Alternating Least Squares](#).

- Write a RMSE function to calculate the [root-mean-square error](#) between two vectors.

Now, we're ready to try using ALS for varying values of λ . In order to reduce computation time and find a low-dimensionality solution, we will constrain the rank of \mathbf{D} to a maximum of 30.

- Iterate through the calculated values of λ . For each one, do the following:
 - Use `softImpute()` with the current value of λ to calculate a singular value decomposition of the scaled sparse ratings matrix. Set `rank.max=30` to restrict solutions to a maximum rank of 30 and `maxit=1000` to control the number of iterations allowed. For all but the first call of `softImpute()`, pass into the `warm.start` parameter the *previous* result of calling `softImpute()` to reduce the required computation time via a “warm start”. Read the documentation for details on what these parameters mean.
 - Calculate the *rank* of the solution by (1) rounding the values of the diagonal matrix of the resulting decomposition, stored in `$d`, to 4 decimal places and (2) calculating the number of nonzero entries.
 - Use `impute()` to fill in the entries of the test set with the calculated matrix decomposition. (Pass in to `impute()` the calculated matrix decomposition as well as the user and movie ID columns in the test set.) Calculate the corresponding RMSE between the predicted ratings and the actual ratings.
 - Store the calculated decomposition in the previously initialized list `fits` as well as the calculated rank and RMSE in the corresponding row of the `results` data frame. Print out the results of the current iteration as well.

You should find that the minimum RMSE is attained at approximately $\lambda \approx 20$ with an RMSE of approximately 0.858.

- Store the best-performing matrix decomposition into a variable called `best_svd`.

Analyzing the results

Now that we have a way to fill in missing entries, we can do some further analysis of the MovieLens dataset.

- As with the ratings dataset, load the movies dataset (in `movies.dat`) and name the columns appropriately.
- How many different genres are listed in the dataset? (You may find [strsplit\(\)](#) helpful.) There is a single genre which is obviously the result

of a data entry error. Add an appropriately named column for all of the *other* genres to serve as an *indicator variable* for whether each movie belongs to a particular genre. Fill in the entries of those columns accordingly.

- Restrict to movies which were listed at least once in the ratings dataset.

Examine the dimensions of the calculated matrix V in `best_svd`. The i th row corresponds to the movie with ID i and the j th column represents the “scores” for the j th “movies factor” (loosely speaking). We’re interested in analyzing these “factors”. To that end:

- Subset `best_svd$v` with the movie IDs in the movie dataset which remain after removing rows corresponding to movies not present in the ratings dataset. (Pay attention to the data type of the movie ID column, which is loaded in as a *factor*.) After doing so, add the factor columns to the data frame created from the movies dataset.²

Next, we’ll illustrate one possible path of analysis by looking at the “Drama” genre.

- Examine the correlation between the indicator variable for movies tagged as dramas and the factor columns. Using `glm()`, run an unregularized logistic regression of the indicator variables against the factors.
- Use `CVbinary()` (from the [DAAG](#) package) on the resulting model to generate *cross-validated probability predictions* for the whole dataset (stored in `CVbinary(fit)$cvhat`). Plot the associated ROC curve and calculate the AUC.

We now have a *probability* for each movie corresponding to how likely it is to be a drama or not given the information stored in the factor variables.

- Create a new data frame with (1) movie titles, (2) the indicator variable for dramas, and (3) the predicted probability for each movie. Order the rows from largest to smallest probability. Which movies are the most likely to be dramas and which movies are the most unlikely to be dramas? How well does this correspond with the actual genre labeling in the dataset?
- Repeat the above analysis for 3 other genres of your choice.

²Something like `movies = cbind(movies, best_svd$v[as.numeric(as.character(movies$mid))])`. (Be sure to understand what this code does!)