

# Data 612 - Project 1

Anna Moy & Natalie Kalukeerthie

2025-06-08

## Overview

**Briefly describe the recommender system that you're going to build out from a business perspective, e.g. "This system recommends data science books to readers."**

For this project, we are building a recommender system for Yelp-style restaurant reviews. Users rate restaurants from 1 to 5, but not every user rates every restaurant (in this project our item will be 'restaurants'). Our goal is to predict missing ratings, using two simple models:

- A **Raw Average Model** is using the overall average rating across all users (Yelp users) and items (Restaurants)
- A **Baseline Model** that adjusts predictions based on both user and restaurant (item) biases. It uses the raw average model and then adjusts the users (Yelp users) and item (Restaurants) biases.

Our core idea is that evaluating how well each model performs (using RMSE), we aim to learn how bias-aware predictions improve recommendation accuracy. We simulate reviews from 6 Yelp users (Alice, Ben, Cindy, David, Ella, Frank) for 5 restaurants (Pasta Place, Sushi Spot, Burger Barn, Curry Corner, Taco Town). Keep in mind that some users have not rated every restaurant.

```
# Load library
library(reshape2)
library(knitr)
library(dplyr)
```

```
##
## Attaching package: 'dplyr'

## The following objects are masked from 'package:stats':
##
##   filter, lag

## The following objects are masked from 'package:base':
##
##   intersect, setdiff, setequal, union
```

```
# Create Dataset with 6 users and 5 restaurants
ratings_df <- data.frame(
  user = c(
    rep("Alice", 5),
```

```

    rep("Ben", 5),
    rep("Cindy", 5),
    rep("David", 5),
    rep("Ella", 5),
    rep("Frank", 5)
  ),
  restaurant = rep(c("Pasta_Place", "Sushi_Spot", "Burger_Barn", "Curry_Corner", "Taco_Town"), 6),
  rating = c(
    5, NA, 4, NA, 4,      # Alice
    4, 3, 5, 3, 4,      # Ben
    4, 2, NA, NA, 3,     # Cindy
    2, 2, 3, 1, 2,      # David
    4, NA, 5, 4, 5,      # Ella
    4, 2, 5, 4, 4        # Frank
  )
)

```

We pivot the data into a matrix format (User-Item Matrix), which helps visualize missing ratings.

```

#Transform data frame into a user-item matrix

#pivot the long format into a matrix for user, item and ratings
matrix <- dcast(ratings_df, user ~ restaurant, value.var = "rating")

# row labels will be the based on user
rownames(matrix) <- matrix$user

# remove the first column since it is repeating duplicate value
matrix <- matrix[, -1]

#print the matrix
print(matrix)

```

```

##      Burger_Barn Curry_Corner Pasta_Place Sushi_Spot Taco_Town
## Alice           4           NA           5           NA           4
## Ben             5            3           4            3           4
## Cindy          NA           NA           4            2           3
## David           3            1           2            2           2
## Ella           5            4           4            NA           5
## Frank          5            4           4            2           4

```

Next, we will split the data into a test (20% of data) and training (80% of data) dataset

```

set.seed(42)
# Take the data and split into training 80% and testing 20% dataframe
train_indices <- sample(1:nrow(ratings_df), size = 0.8 * nrow(ratings_df))

train_df <- ratings_df[train_indices, ]
test_df <- ratings_df[-train_indices, ]

```

Our first model will be the raw average predictor, this model takes the average of all the known ratings in the training dataset and ignores any missing ratings. The RMSE calculates the difference between the Ratings and the raw average for the training and test data.

**RSME (Root Mean Squared Error)** - measures the average difference between the predicted values and actual values

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2}$$

Where:

- $\hat{y}_i$  = predicted value for the  $i$ -th observation
- $y_i$  = actual value for the  $i$ -th observation
- $n$  = total number of observations

```
# Find the Raw average and calculate RMSE

# Raw Average on Training set and ignore NA
global_mean <- mean(train_df$rating, na.rm = TRUE)

# Add the raw avg onto the training and test dataset
train_df$raw_avg <- global_mean
test_df$raw_avg <- global_mean

# RMSE function
rmse <- function(actual, predicted) {
  sqrt(mean((actual - predicted)^2))
}

train_fil <- train_df[!is.na(train_df$rating), ]
test_fil <- test_df[!is.na(test_df$rating), ]

# Calculate RSME which takes the difference of ratings - avg mean for training and test set
rmse_train_raw <- rmse(train_fil$rating, train_fil$raw_avg)
rmse_test_raw <- rmse(test_fil$rating, test_fil$raw_avg)

rmse_train_raw

## [1] 1.133719

rmse_test_raw

## [1] 1.182748
```

Our second model 2 is the **Baseline Model** assumes that each user and item (restaurant) has a tendency to rate higher or lower than average. This model adjusts predictions using:

- **User bias:** How much a user rates higher/lower than average.
- **Item bias:** How popular (or unpopular) a restaurant is compared to others.

We predict the rating as is: Predicted Rating = Global Mean + User Bias + Item Bias

$$\hat{r}_{ui} = \mu + b_u + b_i$$

Where:

- $\hat{r}_{ui}$ : predicted rating for user  $u$  and item  $i$
- $\mu$ : global average rating
- $b_u$ : user bias
- $b_i$ : item bias

This is considered a collaborative filtering approach, where bias correction leads to more personalized predictions.

To prevent overfitting, we applied regularization to user and item bias estimates. This penalizes users or restaurants with very few ratings, reducing the risk of exaggerated bias values. For example, if a user has rated only one item extremely high, regularization ensures that this does not overly skew their predicted ratings across the system.

Regularization is a technique that is used to prevent overfitting and to improve the model generalization.

```
#Finds out the bias values for all columns and rows in grid

# Set regularization strength
lambda <- 10

# Compute regularized user bias
user_bias_df <- train_df[!is.na(train_df$rating), ] %>%
  group_by(user) %>%
  summarise(
    num_ratings = n(),
    sum_diff = sum(rating - global_mean),
    user_bias = sum_diff / (num_ratings + lambda)
  )

# Compute regularized restaurant bias
restaurant_bias_df <- train_df[!is.na(train_df$rating), ] %>%
  group_by(restaurant) %>%
  summarise(
    num_ratings = n(),
    sum_diff = sum(rating - global_mean),
    restaurant_bias = sum_diff / (num_ratings + lambda)
  )

# Find the mean for all bia users in training data
#user_avg <- aggregate(rating ~ user, data = train_df, mean)
# Take the user bias values and the difference from the raw mean
#user_avg$user_bias <- user_avg$rating - global_mean

#user_avg

# Find the mean for all items in the training data
#restaurant_avg <- aggregate(rating ~ restaurant, data = train_df, mean)

# Take the item values and the difference frm the raw mean
#restaurant_avg$restaurant_bias <- restaurant_avg$rating - global_mean

#restaurant_avg
```

```
#Baseline Predictor
```

```
# Merge regularized user and restaurant bias into training set
```

```
train_df <- merge(train_df, user_bias_df[, c("user", "user_bias")], by = "user", all.x = TRUE)
```

```
train_df <- merge(train_df, restaurant_bias_df[, c("restaurant", "restaurant_bias")], by = "restaurant", all.x = TRUE)
```

```
# Predict using regularized baseline model
```

```
train_df$pred_baseline <- global_mean + train_df$user_bias + train_df$restaurant_bias
```

```
# Repeat for test set
```

```
test_df <- merge(test_df, user_bias_df[, c("user", "user_bias")], by = "user", all.x = TRUE)
```

```
test_df <- merge(test_df, restaurant_bias_df[, c("restaurant", "restaurant_bias")], by = "restaurant", all.x = TRUE)
```

```
test_df$pred_baseline <- global_mean + test_df$user_bias + test_df$restaurant_bias
```

```
# Merge user bias into the training data set by user
```

```
#train_df <- merge(train_df, user_avg[, c("user", "user_bias")], by = "user", all.x = TRUE)
```

```
#Merge item into the training data set by item
```

```
#train_df <- merge(train_df, restaurant_avg[, c("restaurant", "restaurant_bias")], by = "restaurant", all.x = TRUE)
```

```
#train_df
```

```
#Baseline Predictor for Training data avg + user bias + item bias
```

```
#train_df$pred_baseline <- global_mean + train_df$user_bias + train_df$restaurant_bias
```

```
# Repeat same steps for testing dataset
```

```
#test_df <- merge(test_df, user_avg[, c("user", "user_bias")], by = "user", all.x = TRUE)
```

```
#test_df <- merge(test_df, restaurant_avg[, c("restaurant", "restaurant_bias")], by = "restaurant", all.x = TRUE)
```

```
#test_df$pred_baseline <- global_mean + test_df$user_bias + test_df$restaurant_bias
```

```
# Baseline Predictor RMSE
```

```
#remove the NA in the row
```

```
train_filtered <- train_df[!is.na(train_df$rating), ]
```

```
test_filtered <- test_df[!is.na(test_df$rating), ]
```

```
# Using the RMSE function taking the difference between rating and baseline predictor for training and testing
```

```
rmse_train_base <- rmse(train_filtered$rating, train_filtered$pred_baseline)
```

```
rmse_test_base <- rmse(test_filtered$rating, test_filtered$pred_baseline)
```

```
rmse_train_base
```

```
## [1] 0.8429098
```

```
rmse_test_base
```

```
## [1] 1.020685
```

```

# Prepare the data
model_names <- c("Raw Average Model", "Baseline Predictor Model")
rmse_train <- c(round(rmse_train_raw, 3), round(rmse_train_base, 3))
rmse_test <- c(round(rmse_test_raw, 3), round(rmse_test_base, 3))
rmse_test_change <- c(NA, round((1 - (rmse_test_base / rmse_test_raw)) * 100, 3))

results_df <- data.frame(
  Model = model_names,
  RMSE_Train = rmse_train,
  RMSE_Test = rmse_test,
  RMSE_Test_Change_Percent = rmse_test_change
)

# Print the table
kable(results_df, caption = "RMSE Comparison of Models")

```

Table 1: RMSE Comparison of Models

Model	RMSE_Train	RMSE_Test	RMSE_Test_Change_Percent
Raw Average Model	1.134	1.183	NA
Baseline Predictor Model	0.843	1.021	13.702

```

### RAW AVERAGE MODEL PREDICTIONS
# Fill missing values with the global average
ratings_raw <- ratings_df %>%
  mutate(pred_raw = ifelse(is.na(rating), global_mean, rating))

# Reshape to matrix
ratings_matrix_raw <- dcast(ratings_raw, user ~ restaurant, value.var = "pred_raw")

### BASELINE MODEL PREDICTIONS
# Compute biases
user_avg <- ratings_df %>%
  group_by(user) %>%
  summarize(u_avg = mean(rating, na.rm = TRUE))

item_avg <- ratings_df %>%
  group_by(restaurant) %>%
  summarize(i_avg = mean(rating, na.rm = TRUE))

# Merge with original data
baseline_df <- ratings_df %>%
  left_join(user_avg, by = "user") %>%
  left_join(item_avg, by = "restaurant") %>%
  mutate(pred_baseline = ifelse(
    is.na(rating),
    global_mean + (u_avg - global_mean) + (i_avg - global_mean),
    rating
  ))

# Reshape to matrix

```

```
ratings_matrix_base <- dcast(baseline_df, user ~ restaurant, value.var = "pred_baseline")

#Print both Predictions
print("Raw Average Model Prediction:")
```

```
## [1] "Raw Average Model Prediction:"
```

```
print(ratings_matrix_raw)
```

```
##      user Burger_Barn Curry_Corner Pasta_Place Sushi_Spot Taco_Town
## 1 Alice      4.000000      3.368421          5  3.368421          4
## 2 Ben        5.000000      3.000000          4  3.000000          4
## 3 Cindy      3.368421      3.368421          4  2.000000          3
## 4 David      3.000000      1.000000          2  2.000000          2
## 5 Ella       5.000000      4.000000          4  3.368421          5
## 6 Frank      5.000000      4.000000          4  2.000000          4
```

```
print("Baseline Model Prediction:")
```

```
## [1] "Baseline Model Prediction:"
```

```
print(ratings_matrix_base)
```

```
##      user Burger_Barn Curry_Corner Pasta_Place Sushi_Spot Taco_Town
## 1 Alice      4.000000      3.964912          5  3.214912          4
## 2 Ben        5.000000      3.000000          4  3.000000          4
## 3 Cindy      4.031579      2.631579          4  2.000000          3
## 4 David      3.000000      1.000000          2  2.000000          2
## 5 Ella       5.000000      4.000000          4  3.381579          5
## 6 Frank      5.000000      4.000000          4  2.000000          4
```

The Baseline Predictor Model improves test RMSE by about 13.7% compared to the Raw Average Model, indicating better generalization on unseen data. We incorporated regularization into our model to reduce over-fitting of the data and this process actually decreased our baseline predictor model's test and training RMSE since we shrunk extreme bias values toward the global average.

RMSE is used to measure how close our predictions are compared to the actual ratings. The lower RMSE of the baseline predictor confirms it outperforms the raw average model, indicating meaningful improvement. Sushi Spot and Curry Corner have negative item biases, suggesting they are consistently rated below average by users. Burger Barn and Taco Town on the other hand received higher average ratings from users which indicates customer satisfaction.

The users Alice and Ella will provide higher ratings compared to others on average which are positive user biases. The negative user biases would be David and Cindy. The prediction for the baseline model takes the bias into consideration compared to the raw average predictions which is the same for all of them.

Potential errors is not conduction regularization which would cause overfitting on the data. Another issue would be the handling of NAs in the data set. If we replace the NAs with zero it could consider it as a datapoint and it would be best to leave it as NA. But if we have too many NAs in the dataset it could potentially be an issue.

One of the issues we may encounter is cold start which is not having enough ratings from users or items. We can try to assign zeros biasto new users and items and fallback on the global average.