

DSCI 726: Project 4 – Predictive Analysis for Book Recommendations

Code

Becca Halford

2025-10-03

- Executive Summary
- Data Preparation
 - 1.1 Load & Clean
 - 1.2 Filter by Author Productivity (fs2|books removed)
- 2. Feature Engineering
 - 2.1 Transformations
 - 2.2 Bin author_num_books (quartiles)
 - 2.3 Scale numeric features (train params)
- 3. Baselines (for context)
- 4. Correlation & Collinearity Checks
- 5. Models
 - 5.1 Simple OLS (baseline comparator; not useful)
 - 5.2 Full OLS
 - 5.4.5-fold Cross-Validation (Full OLS)
- 6. Evaluation & Comparison
 - 6.1 Diagnostics (Full OLS)
- 7. Brief Findings & Limitations (in plain language)
- Appendix A – Supplemental EDA (kept minimal for rubric compliance)
 - A1. Post-Transformation Distributions
 - A2. Feature-Target Scatter with Linear Trend (sampled)
- 8. Save Key Outputs

Executive Summary

This project predicts the star rating of new books by authors who have published at least three titles on Amazon. I used a random 80/20 split by book title to create the training and test sets.

The modeling methods include a baseline model, a simple linear model, a full linear model, and a pruned decision tree. I evaluated model performance using RMSE, R², and the train-test performance gap. I also performed 5-fold cross-validation on the full linear model to assess its stability and generalizability.

1. Data Preparation

The column "save" was renamed to "savings" for clarity. The column "size" was removed due to redundancy. Duplicate rows were removed. Convert applicable columns to integer.

1.1 Load & Clean

```
kindle_data <- read_excel('cleaned_kindle_books.xlsx') >
  select(-size) |>
  mutate(across(c(pages, author_num_books, customer_reviews), as.integer)) |>
  distinct(title, .keep_all = TRUE) |>
  rename(savings = save)

cat("Dataset:", nrow(kindle_data), "books x", ncol(kindle_data), "variables\n")
```

Dataset: 44739 books x 11 variables

```
glimpse(kindle_data)
```

Rows: 44,739
Columns: 11
\$ title <chr> "Middle School Confidential: The Alexa Black Diaries ...
\$ author <chr> "#NororityProblem", "#NororityProblem", "#NororityPro...
\$ price <dbl> 4.99, 4.99, 4.99, 4.99, 2.99, 2.99, 15.49, 26.8...
\$ savings <dbl> 0.00, 0.00, 0.00, 0.00, 0.00, 0.00, 0.00, 0.00, 0.00, 0.00, 0.00
\$ pages <int> 117, 65, 72, 63, 115, 41, 172, 338, 336, 40...
\$ publisher <chr> "Self-Published", "Self-Published", "Self-Published", "Self-Publis...
\$ customer_reviews <int> 4, 3, 7, 2, 30, 6, 4, 1, 106, 219, 754, 359, 11, 7...
\$ stars <dbl> 3.0, 4.7, 3.0, 3.0, 2.8, 4.7, 5.0, 3.0, 4.6, 3.2, 4.3...
\$ author_num_books <int> 4, 4, 4, 4, 1, 1, 1, 2, 1, 3, 3, 3, 1, 1, 2, 1, 1...
\$ avg_stars <dbl> 3.0, 3.9, 3.9, 3.9, 2.8, 4.7, 5.0, 3.0, 3.4, 3.2, 4.4...
\$ discount_pct <dbl> 1.00, 1.00, 1.00, 1.00, 0.43, 1.00, 0.55, 0.67, -

1.2 Filter by author productivity (≤2 books removed)

Authors with 2 or less books were removed from the dataset to more accurately represent stars and avg_stars.

```
original_n <- nrow(kindle_data)
author_counts <- kindle_data |> count(author, name = "n_books")
few_book_authors <- author_counts |> filter(n_books <= 2) |> pull(author)

kindle_data <- kindle_data |> filter(!author %in% few_book_authors)

cat("Removed:", original_n - nrow(kindle_data), "rows (",
  round(100*(original_n - nrow(kindle_data))/original_n), "%\n", sep = "")
```

Removed: 2419 rows (54.1%)

1.3 Train/Test Split

I split the data using a simple random sample by book title after completing the filtering process. The final split allocated 80% of the data to the training set (16,438 rows) and 20% to the test set (4,109 rows).

```
set.seed(123)
test_titles <- sample(kindle_data$title, size = floor(0.2 * nrow(kindle_data)))
train_set <- kindle_data |> filter(!title %in% test_titles)
test_set <- kindle_data |> filter(title %in% test_titles)

cat("Train/Test sizes:", nrow(train_set), "/", nrow(test_set), "\n")
```

Train/Test sizes: 16438 / 4109

2. Feature Engineering

A log transformation was performed due to highly skewed variables.

2.1 Transformations

```
cols_to_log <- c("price", "customer_reviews", "pages", "savings")

train_set <- train_set |
  mutate(across(all_of(cols_to_log), ~log(.x + 1), .names = "log_(.col)"))
test_set <- test_set |
  mutate(across(all_of(cols_to_log), ~log(.x + 1), .names = "log_(.col)"))

cat("Log-transformed:", paste0(cols_to_log, collapse = ", "), "\n")
```

Log-transformed: price, customer_reviews, pages, savings

Log transformation was unsuccessful on 'author_num_books' so binning was applied.

2.2 Bin author_num_books (quartiles)

```
quartiles <- quantile(train_set$author_num_books, probs = c(.0, .25, .5, .75, 1), na.rm = TRUE)

train_set <- train_set |
  mutate(author_books_bin = cut(author_num_books, breaks = quartiles,
    labels = c("Few_Q1", "Some_Q2", "Many_Q3", "Prolific_Q4"),
    include.lowest = TRUE))

test_set <- test_set |
  mutate(author_books_bin = cut(author_num_books, breaks = quartiles,
    labels = c("Few_Q1", "Some_Q2", "Many_Q3", "Prolific_Q4"),
    include.lowest = TRUE))

table(train_set$author_books_bin)
```

Few_Q1 Some_Q2 Many_Q3 Prolific_Q4
5051 3250 4097 4030

2.3 Scale numeric features (train params)

```
scale_cols <- c(paste0("log_", cols_to_log), "discount_pct", "avg_stars")

scaling_params <- train_set |
  summarise(across(all_of(scale_cols), list(mean = -mean(.x, na.rm = TRUE),
    sd = -sd(.x, na.rm = TRUE))))
```

```
for (col in scale_cols) {
  m <- scaling_params[[paste0(col, "_mean")]]
  s <- scaling_params[[paste0(col, "_sd")]]
  train_set_scaled[[col]] <- (train_set_scaled[[col]] - m) / s
  test_set_scaled[[col]] <- (test_set_scaled[[col]] - m) / s
}
```

```
train_set_scaled <- train_set_scaled |> mutate(author_books_bin = as.factor(author_books_bin))
test_set_scaled <- test_set_scaled |> mutate(author_books_bin = factor(author_books_bin,
  levels = levels(train_set_scaled$author_books_bin)))
```

One-hot for binned variable

dummies <- dummyVars(~author_books_bin, data = train_set_scaled)

train_cat <- as.data.frame(predict(dummies, newdata = train_set_scaled))

test_cat <- as.data.frame(predict(dummies, newdata = test_set_scaled))

numeric_features <- c("log_price", "log_customer_reviews", "log_pages", "log_savings", "avg_stars", "discount_pct", "stars")

train_model <- cbind(train_set_scaled |> select(all_of(numeric_features)), train_cat)

test_model <- cbind(test_set_scaled |> select(all_of(numeric_features)), test_cat)

```
cat("Model matrices:", dim(train_model)[1], "*", dim(train_model)[2], "(train)\n",
  dim(test_model)[1], "*", dim(test_model)[2], "(test)\n")
```

Model matrices: 16438 x 11 (train);
4109 x 11 (test)

Best by Test RMSE: Full OLS with RMSE = 0.4553 ; R² = 0.4563

Baseline RMSE (Author-only): 0.3571

Baseline RMSE (Null Mean): 0.4797

Baseline RMSE (Author-only): 0.3571

Baseline RMSE (Null Mean): 0.4797

<