

UNIVERSITY OF CALIFORNIA

Los Angeles

Predicting New York Times Bestselling Fiction Books
Using Text & Image-Based Machine Learning

A thesis submitted in partial satisfaction
of the requirements for the degree
Master of Applied Statistics and Data Science

by

Anum Iqbal Damani

2025

PREVIEW

ABSTRACT OF THE THESIS

Predicting New York Times Bestselling Fiction Books Using Text & Image-Based Machine Learning

by

Anum Iqbal Damani

Master of Applied Statistics and Data Science

University of California, Los Angeles, 2025

Professor Yingnian Wu, Chair

This analysis focuses on predicting New York Times (NYT) bestselling fiction books, specifically on the Combined Print & E-Book Fiction list, using book descriptions and book covers. Approximately ten years of data from the NYT Books API and Open Library API is collected, and an exploratory data analysis is performed. All machine learning models are evaluated using performance metrics in the classification report, including accuracy, weighted average precision, and weighted average recall. In addition, the model performances are assessed using the precision-recall curve, ROC curve, and confusion matrix. Although the dataset is fairly balanced, NYT bestsellers are rare in reality, so the classification threshold is tuned for each model and precision is prioritized. Seven text-based models are implemented by utilizing cleaned book descriptions as input, which are transformed using either TF-IDF or BERT. Based on the overall performances, the text-based models with the best performances are Logistic Regression with BERT and XGBoost with TF-IDF, with strong accuracies at 92% and 91%, respectively. Next, five image-based models are implemented by extracting image features from book covers. The overall performances of the image-based models indicated that the XGBoost model with CLIP embeddings and the ResNet50 model are the strongest, with accuracies 87% and 83%, respectively. The best text-based and image-based models are then considered for multimodal modeling. Four multimodal models are implemented and compared. The Logistic Regression model with BERT & CLIP features was the top-performing multimodal model, achieving a strong accuracy of 97%.

The thesis of Anum Iqbal Damani is approved.

David Anthony Zes

Vivian Lew

Frederic R. Paik Schoenberg

Yingnian Wu, Committee Chair

University of California, Los Angeles

2025

PREVIEW

PREVIEW

To my supporters

Table of Contents

| | | |
|----------|---|-----------|
| 1 | Introduction | 1 |
| 2 | Data | 3 |
| 2.1 | Data Collection Methods for Tabular Dataset | 3 |
| 2.1.1 | Bestsellers | 3 |
| 2.1.2 | Non-Bestsellers | 4 |
| 2.1.3 | Combining the NYT Bestsellers & Non-Bestsellers | 5 |
| 2.2 | Data Collection Methods for Image Dataset | 6 |
| 3 | Text Preprocessing | 7 |
| 4 | Exploratory Data Analysis (EDA) | 8 |
| 4.1 | EDA for Tabular Dataset | 8 |
| 4.2 | EDA for Image Dataset | 11 |
| 5 | Machine Learning Methodology | 13 |
| 6 | Text-Based Models | 15 |
| 6.1 | Vectorization | 15 |
| 6.2 | Logistic Regression | 17 |
| 6.2.1 | Background | 17 |
| 6.2.2 | Model 1: Logistic Regression (TF-IDF) | 17 |
| 6.2.3 | Model 2: Logistic Regression (BERT) | 19 |
| 6.3 | Naive Bayes | 21 |
| 6.3.1 | Background | 21 |
| 6.3.2 | Model 3: Naive Bayes (TF-IDF) | 22 |
| 6.4 | Random Forest | 23 |
| 6.4.1 | Background | 23 |

| | | |
|----------|--|-----------|
| 6.4.2 | Model 4: Random Forest (TF-IDF) | 23 |
| 6.4.3 | Model 5: Random Forest (BERT) | 25 |
| 6.5 | XGBoost | 27 |
| 6.5.1 | Background | 27 |
| 6.5.2 | Model 6: XGBoost (TF-IDF) | 28 |
| 6.5.3 | Model 7: XGBoost (BERT) | 29 |
| 6.5.4 | Comparison of Performances for Text-Based Models | 32 |
| 7 | Image-Based Models | 33 |
| 7.1 | EfficientNetB0 | 33 |
| 7.1.1 | Background | 33 |
| 7.1.2 | Model 1: EfficientNetB0 | 33 |
| 7.2 | MobileNet | 35 |
| 7.2.1 | Background | 35 |
| 7.2.2 | Model 2: MobileNetV3Small | 36 |
| 7.3 | ResNet50 | 38 |
| 7.3.1 | Background | 38 |
| 7.3.2 | Model 3: ResNet50 | 38 |
| 7.4 | VGG16 | 40 |
| 7.4.1 | Background | 40 |
| 7.4.2 | Model 4: VGG16 | 41 |
| 7.5 | XGBoost with CLIP Embeddings | 43 |
| 7.5.1 | Background | 43 |
| 7.5.2 | Model 5: XGBoost with CLIP Embeddings | 43 |
| 7.6 | Comparison of Performances for Image-Based Models | 46 |
| 8 | Multimodal Models | 47 |
| 8.1 | Background | 47 |
| 8.2 | Model 1: XGBoost with TF-IDF & CLIP Features | 47 |
| 8.3 | Model 2: XGBoost with TF-IDF & ResNet50 Features | 49 |
| 8.4 | Model 3: Logistic Regression with BERT & CLIP Features | 51 |
| 8.5 | Model 4: Logistic Regression with BERT & ResNet50 Features | 52 |
| 8.6 | Comparison of Performances for Multimodal Models | 55 |

| | |
|----------------------------------|-----------|
| 9 Conclusion | 57 |
| 10 Bibliography | 59 |

PREVIEW

List of Figures

| | | |
|------|--|----|
| 4.1 | Word Clouds of Important Words in Book Titles | 9 |
| 4.2 | Word Clouds of Important Words in Book Descriptions | 9 |
| 4.3 | Distribution of Sentiment Polarity for Bestsellers and Non-Bestsellers | 10 |
| 4.4 | Box Plot of Compound Sentiment | 10 |
| 4.5 | Examples of Book Covers in Dataset | 11 |
| 4.6 | Color Histograms for Bestsellers vs. Non-Bestsellers | 12 |
| 6.1 | Logistic Regression (TF-IDF) Precision-Recall Plots | 18 |
| 6.2 | Logistic Regression (TF-IDF) Performance Plots | 19 |
| 6.3 | Logistic Regression (BERT) Precision-Recall Plots | 20 |
| 6.4 | Logistic Regression (BERT) Performance Plots | 21 |
| 6.5 | Naive Bayes (TF-IDF) Precision-Recall Plots | 22 |
| 6.6 | Naive Bayes (TF-IDF) Performance Plots | 23 |
| 6.7 | Random Forest (TF-IDF) Precision-Recall Plots | 24 |
| 6.8 | Random Forest (TF-IDF) Performance Plots | 25 |
| 6.9 | Random Forest (BERT) Precision-Recall Plots | 26 |
| 6.10 | Random Forest (BERT) Performance Plots | 27 |
| 6.11 | XGBoost (TF-IDF) Precision-Recall Plots | 28 |
| 6.12 | XGBoost (TF-IDF) Performance Plots | 29 |
| 6.13 | XGBoost (BERT) Precision-Recall Plots | 30 |
| 6.14 | XGBoost (BERT) Performance Plots | 31 |
| 7.1 | EfficientNetB0 Model Architecture | 33 |
| 7.2 | EfficientNetB0 Precision-Recall Plots | 34 |
| 7.3 | EfficientNetB0 Performance Plots | 35 |
| 7.4 | MobileNet Convolutional Block | 36 |
| 7.5 | MobileNetV3Small Precision-Recall Plots | 37 |

| | | |
|------|---|----|
| 7.6 | MobileNetV3Small Performance Plots | 38 |
| 7.7 | ResNet50 Precision-Recall Plots | 39 |
| 7.8 | ResNet50 Performance Plots | 40 |
| 7.9 | VGGNet Model Architecture | 41 |
| 7.10 | VGG16 Precision-Recall Plots | 42 |
| 7.11 | VGG16 Performance Plots | 43 |
| 7.12 | XGBoost (CLIP) Precision-Recall Plots | 44 |
| 7.13 | XGBoost (CLIP) Performance Plots | 45 |
| 8.1 | XGBoost with TF-IDF & CLIP Precision-Recall Plots | 48 |
| 8.2 | XGBoost with TF-IDF & CLIP Performance Plots | 49 |
| 8.3 | XGBoost with TF-IDF & ResNet50 Precision-Recall Plots | 49 |
| 8.4 | XGBoost with TF-IDF & ResNet50 Performance Plots | 50 |
| 8.5 | Logistic Regression with BERT & CLIP Precision-Recall Plots | 51 |
| 8.6 | Logistic Regression with BERT & CLIP Performance Plots | 52 |
| 8.7 | Logistic Regression with BERT & ResNet50 Precision-Recall Plots | 53 |
| 8.8 | Logistic Regression with BERT & ResNet50 Performance Plots | 54 |

List of Tables

| | | |
|-----|---|----|
| 2.1 | Variable Descriptions for Tabular Dataset | 6 |
| 6.1 | Performance Report of Logistic Regression (TF-IDF) Model | 18 |
| 6.2 | Performance Report of Logistic Regression (BERT) Model | 20 |
| 6.3 | Performance Report of Naive Bayes (TF-IDF) Model | 22 |
| 6.4 | Performance Report of Random Forest (TF-IDF) Model | 24 |
| 6.5 | Performance Report of Random Forest (BERT) Model | 26 |
| 6.6 | Performance Report of XGBoost (TF-IDF) Model | 28 |
| 6.7 | Performance Report of XGBoost (BERT) Model | 30 |
| 6.8 | Comparison of Text-Based Model Performances | 32 |
| 7.1 | Performance Report of EfficientNetB0 Model | 35 |
| 7.2 | Performance Report of MobileNetV3Small Model | 37 |
| 7.3 | Performance Report of ResNet50 Model | 40 |
| 7.4 | Performance Report of VGG16 Model | 42 |
| 7.5 | Performance Report of XGBoost Model with CLIP Embeddings | 44 |
| 7.6 | Comparison of Image-Based Model Performances | 46 |
| 8.1 | Performance Report of XGBoost Model with TF-IDF & CLIP Features | 48 |
| 8.2 | Performance Report of XGBoost Model with TF-IDF & ResNet50 Features | 50 |
| 8.3 | Performance Report of Logistic Regression Model with BERT & CLIP Features | 51 |
| 8.4 | Performance Report of Logistic Regression Model with BERT & ResNet50 Features | 53 |
| 8.5 | Comparison of Multimodal Model Performances | 55 |

Chapter 1

Introduction

For decades, The New York Times (NYT) bestseller lists have served as indicators for book popularity, reflecting and influencing readers' preferences. After browsing through the shelves of bookstores or libraries, one will quickly notice that several books receive praise by NYT. Although NYT bestselling books are prominent throughout bookstores and libraries, the books that receive this special recognition are actually quite rare. Every year, around three million books are published, but the number of books that are featured on the NYT bestselling lists is under five hundred books [23].

My motivation for this analysis is to gain a better understanding of the books that receive this special recognition. In order to do this, it is crucial to investigate the relationship between book features and NYT bestseller status. The descriptions and covers of books can be influential in a reader's decision to select a book to read. This analysis aims to address the following research questions. To what extent can text-based features, specifically extracted from text descriptions, predict NYT bestsellers? Additionally, how accurately can image-based features, specifically extracted from book covers, predict NYT bestsellers? Finally, how effective are multimodal features, which are combined textual and visual features, in predicting NYT bestsellers?

While the NYT has many bestseller lists, this analysis focuses solely on the Combined Print & E-Book Fiction list, which considers print and e-book sales of fiction books. My focus on the Combined Print & E-Book Fiction list stems from the widespread popularity of both print and e-book formats among readers. Moreover, the Combined Print & E-Book Fiction list provides a holistic perspective on fiction book success. For the purpose of this analysis, I have defined the term "bestseller" to refer to books that have been featured on the NYT Combined Print & E-Book Fiction list between January 4, 2015 to January 19, 2025. Furthermore, the

books that have not been featured on the NYT Combined Print & E-Book Fiction list during this date range are labeled as “non-bestseller.” A book that is labeled as “non-bestseller” in the dataset could have received recognition on a different NYT bestselling list, such as “Hardcover Fiction.”

This analysis revolves around a binary classification problem with fairly balanced classes. In this dataset, there are 4,666 observations in total, where 2,216 observations are bestsellers and 2,450 observations are non-best sellers. To account for the rarity of NYT bestsellers in real life, my approach involved tuning the classification threshold for all models to above 0.5, if possible. Through this approach, I prioritized precision, while aiming to maintain a reasonable recall around 0.8. Although it is not ideal to miss bestsellers, it can be argued that the cost of false positives is greater. It can be risky to falsely predict a book as a bestseller due to the amount of resources that are allocated in the book publishing process. Through the implementation of text-based, image-based, and multimodal models in this manner, I hope that my exploration of this topic provides valuable and applicable insights about fiction book success across print and e-book formats.

Chapter 2

Data

2.1 Data Collection Methods for Tabular Dataset

2.1.1 Bestsellers

The New York Times bestselling fiction book dataset was collected using the New York Times Books API [22] [16]. This process involved obtaining a NYT API key from the NYT Developer portal and generating an API request for the specified bestseller list of interest for this analysis: “Combined Print & E-Book Fiction.” Next, the bestseller lists were downloaded using the desired date range. The date range selected for this data is January 4, 2015 to January 19, 2025. Specifically, the variables obtained are: title, author, list publication date, book image URL, and book description.

The NYT Books API also provides other information about bestsellers, such as rank, number of weeks on the list, ISBN numbers, and publisher name; however, these variables were excluded from the analysis. Note that the NYT Books API includes *list_publication_date*, which is the publication date of the list in which the book was featured. In order to obtain the book publication year, it was necessary to fetch this information from Open Library API [11]. Open Library API allows users to connect to a large database containing information pertaining to more than 39 million works [12]. A function was written to obtain the book publication date from Open Library API for every row in the dataset, specifically using the title and author information. The variable name for the resulting column for the book publication date is called *year*. Finally, the *list_publication_date* variable was dropped from the dataset.

In terms of data cleaning, the duplicate observations were dropped from the dataset. There were 76 observations with missing book descriptions, so these observations were dropped. There are 352 observations with missing values in the *year* column, but these observations were

kept in the dataset. A new column for the class label 1 was added to this dataset. The resulting NYT fiction bestseller dataset includes 2,216 observations. To summarize, the columns of the bestseller dataset are:

- ***title***: The title of the book
- ***author***: The author(s) of the book
- ***year***: The publication year of the book
- ***book_image***: The URL for the book cover image
- ***description***: The description of the book
- ***label***: The label column contains the value 1

2.1.2 Non-Bestsellers

Recall that the NYT bestseller dataset spans ten years of data, from 2015 to 2025. In the NYT bestseller dataset, there are several books published prior to 2015 that were featured on the NYT Combined Print & E-Book Fiction bestseller list in 2015. Specifically, *The Alchemist* by Paulo Coelho was originally published in 1988 but it was featured on the NYT bestselling fiction list in 2015. Another book featured on the NYT bestselling fiction list in 2015 is *Sharp Objects* by Gillian Flynn, which was originally published in 2006. Lastly, a book that was published in 2011 and featured on the NYT bestselling fiction list in 2015 is *The Martian* by Andy Weir. It is important to note that the film adaptation of *The Martian* was also released in 2015. The book *The Martian* by Andy Weir is an example of a “backlist revival,” a term that is prevalent in the book publishing industry which emphasizes surges in popularity, or prolonged popularity, of books years after their publication date. It is possible for the popularity of a book to be driven by an upcoming film adaptation.

The non-best-selling fiction book dataset was collected using Open Library API, specifically the Subjects API. The data spans 2010 to 2024, and this date range was carefully selected. To mitigate the issue of data sparsity in Open Library API, the non-bestseller dataset includes more than ten years of data. Given the backlist revivals present in the bestseller dataset, and to allow for a fairer comparison, the lower bound for the date range of the non-bestseller dataset is chosen to be 2010. The upper bound for the date range for the non-bestseller list is 2024, rather