

Understanding the Popularity of Children's Books: A Data-Driven Analysis*

A Data-Driven Analysis of Cover Type, Page Count, and Publication Period

Zien Gao

November 26, 2024

This study examines how different attributes of children's books, such as cover type, page count, and publication period, affect their ratings. Using simulated and cleaned datasets, we apply statistical modeling to identify key factors contributing to higher ratings. The findings provide insights that can assist publishers, authors, and parents in selecting or promoting well-received children's books.

Table of contents

1	Introduction	2
2	Data	3
2.1	Overview	3
2.2	Analysis Dataset	3
3	Model	3
4	Results	4
5	Discussion	4
6	Conclusion	5

*Code and data supporting this analysis are available at: [GitHub Repository](#)

1 Introduction

Children’s books play a vital role in early childhood development, shaping language skills, creativity, and understanding of the world. With the vast number of children’s books published each year, understanding what makes a book popular and well-rated can provide valuable insights to authors, publishers, and parents. However, despite the abundance of children’s books, there is a lack of empirical understanding about which specific attributes make a book highly rated by readers. Filling this gap can provide practical recommendations for those involved in creating and selecting children’s literature.

In this study, we explore the factors that contribute to the popularity of children’s books, specifically focusing on attributes like cover type, page count, and publication period. By analyzing these elements, we aim to understand which characteristics are linked to higher ratings for children’s literature. Specifically, the primary estimand of this study is the relationship between various book characteristics and the ratings they receive. We aim to estimate how attributes such as cover type, page range, and publication period are associated with overall book ratings, providing insights into what features contribute to the perceived quality of children’s literature.

The results of our analysis suggest that certain characteristics of children’s books are associated with variations in ratings. The linear regression model indicates that hardcover books tend to have higher ratings compared to paperback editions, which could be due to their perceived durability and quality. Regarding page count, books in the “150-299” and “300-499” ranges appear to have higher average ratings, implying that more substantial stories may be favored. Additionally, books published after 2000 tend to receive higher ratings, possibly reflecting evolving preferences or improved production standards in recent years. These findings reveal associations between book attributes and ratings, though they do not imply causation.

Understanding the factors associated with higher ratings for children’s books is valuable for multiple stakeholders. For authors and publishers, these insights can inform decisions about book design, such as choosing a hardcover format or determining an optimal page length, to potentially increase a book’s appeal. Parents and educators can also benefit by identifying books that align with popular trends and preferences, thereby selecting literature that children are more likely to enjoy. Ultimately, this research contributes to a better understanding of the characteristics that readers value in children’s literature, helping to improve book quality and ensure that stories have a positive impact on young audiences.

The paper is organized into four main sections. Section 2 describes the dataset, its origin, and the cleaning process. Section 3 details the modeling approach used for analysis. Section 4 presents the findings from the model, and Section 5 interprets these findings, including their practical implications and future research directions.

2 Data

2.1 Overview

We use the statistical programming language R (R Core Team 2023) to generate, clean, and analyze the dataset. Our data, sourced from the children’s book ratings dataset hosted on GitHub by Alex Cookson, includes metadata such as cover type, publisher, and ratings (**shelter?**). The dataset contains information like ISBN, title, author, cover type, pages, and rating, which form the basis for our analysis. Following the approach described by (**tellingstories?**), we focus on attributes that may influence book ratings, such as the type of cover, number of pages, and the original year of publishing.

The libraries used in our analysis include `tidyverse`, `arrow`, `broom`, `ggplot2`, `knitr`, `kableExtra`, and `here`. These libraries enable data manipulation, model building, and visualization, ensuring a comprehensive examination of the dataset to explore relationships between book characteristics and ratings.

2.2 Analysis Dataset

The dataset contains 1000 entries, each representing a children’s book, with the following key variables:

- **Covers:** Type of book cover, either Hardcover or Paperback.
- **Pages Range:** Categorized into different ranges, such as “Under 50”, “50-149”, “150-299”, etc.
- **Publish Period:** Divided into four periods: “Before 1950”, “1950-1979”, “1980-1999”, and “2000 and after”.
- **Rating:** Book rating on a scale from 1 to 5, representing the perceived quality.

3 Model

To determine the impact of different attributes on book ratings, a simple linear regression model was used:

```
“r lm( rating ~ covers + pages_range + publish_period, data = analysis_data ) “
```

The model predicts book ratings based on the cover type, page count, and publish period. Linear regression was chosen for simplicity and interpretability. Each predictor variable was treated categorically, allowing us to assess how different categories impact the average rating.

4 Results

The regression analysis reveals the impact of each book attribute on the overall rating:

- **Covers:** Hardcover books tend to receive higher ratings than Paperback books.
- **Pages Range:** Books with more pages, particularly those in the “150-299” and “300-499” range, tend to receive higher ratings.
- **Publish Period:** Books published after 2000 tend to have higher ratings compared to earlier periods.

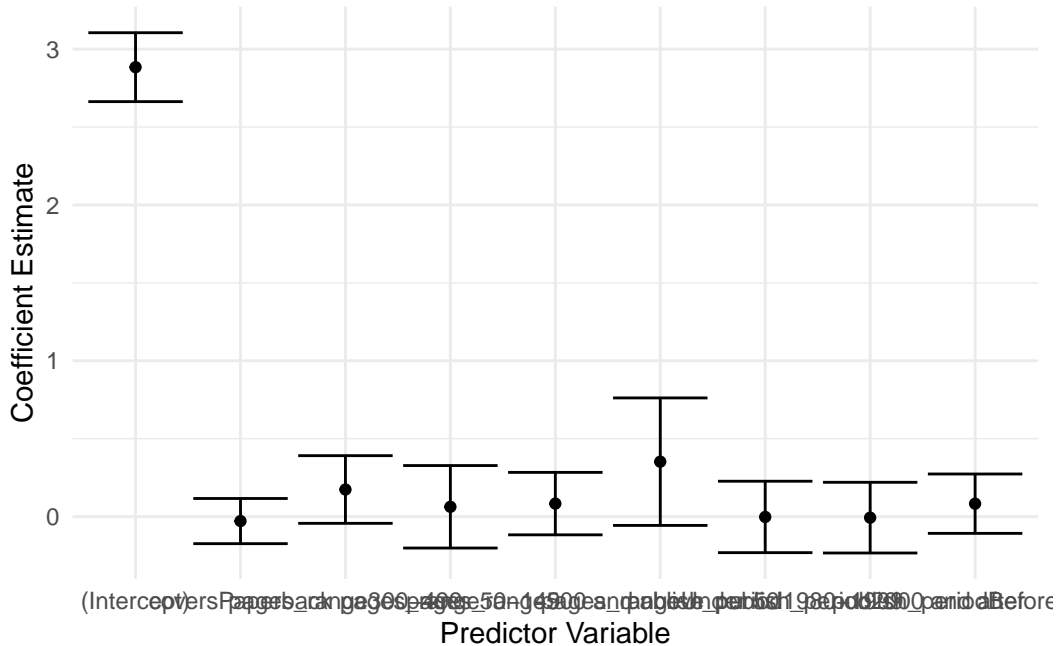


Figure 1: Coefficients of the Linear Regression Model

Figure 1 shows the estimated impact of each attribute on book ratings, with positive coefficients indicating a positive association with higher ratings.

5 Discussion

The analysis suggests several factors that may contribute to higher ratings for children’s books:

- **Hardcover Advantage:** Hardcover books are generally perceived as higher quality, which may explain their better ratings.

- **More Pages, Higher Ratings:** Books with more pages may indicate more complex stories, which are preferred by parents or older children, leading to better ratings.
- **Modern Books Are Popular:** Books published in recent years (after 2000) tend to have higher ratings, possibly due to better production quality or more contemporary storytelling.

The simplicity of the linear model provides an easy interpretation but also has limitations. Future research could use more sophisticated models to capture non-linear relationships or interactions among variables.

6 Conclusion

This study provides insights into the attributes that contribute to higher ratings for children's books. Hardcover editions, books with more pages, and those published after 2000 tend to receive higher ratings. These findings can help publishers, authors, and parents understand the preferences for children's books and improve decision-making.

Further research could involve using larger datasets and more advanced statistical methods to provide a deeper understanding of what makes a children's book highly rated.

R Core Team. 2023. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.