



BARCELONA SCHOOL OF ECONOMICS

DATA SCIENCE METHODOLOGY PROGRAM

22DM014 - INTRO TO NATURAL LANGUAGE PROCESSING

Analysis of Amazon Books and Reviews

Authors:

CHEN, Joshua
EL DAOU, Mahmoud
GALLEGOS, Rafael

Professors:

MUELLER, Hannes
GARCIA, Luis

Contents

1	Introduction	2
2	Literature Reviews	2
3	Dataset	2
4	Exploratory Data Analysis	3
4.1	Descriptions	4
4.1.1	New York Times Bestsellers	4
4.2	Reviews	4
4.2.1	Sentiment Analysis	4
5	Predicting Review Score from Text	4
6	Review Description Alignment	4
7	Conclusions	4

List of Figures

4.1	Percent of Books by Average Score	3
4.2	Percentage of Books by Genre	3
4.3	Percentage of Books by Publisher	4
4.4	Percentage of Books by Publisher	4
4.5	Quantile Regression of Average Reviews per Book over Count of Reviews . .	5
4.6	Bag of Words - Book Descriptions - MinDF = 20, MaxDF = 0.01	5

List of Tables

1 Introduction

On July 16th, 1995, Amazon opened as an online book retailer, promoting itself as "the world's biggest bookstore". Since then, Amazon has clearly expanded into other ventures but still remains the largest seller of books in the world. Amazon generates \$28 billion worldwide every year with a selection of over 32 million books. With Kindle and Audible, Amazon also has majority control of the ebook and audiobook markets as well. While Amazon is no longer primarily known as predominantly a bookseller, they have succeeded in becoming the most dominant bookstore in the world.

In 2013, Amazon acquired Goodreads, a social cataloging platform for books. Goodreads now has more than 150 million users and has become a cornerstone (for better or worse) in the book community. However, the partnership between world's largest bookseller and the largest community of readers, gives us the opportunity to acquire and use lots of book data to identify book trends and reader preferences.

2 Literature Reviews

3 Dataset

In this project, we explore the Amazon Books Reviews dataset, which can be found on Kaggle at the following link: [Amazon Books Reviews Dataset](#)

Contained here are two large data sets contain information about Books and information about Reviews.

Feature	Description
Title	Book Title
Describe	Description of the book
authors	Name of book authors
image	URL for book cover
previewLink	Link to access this book on Google Books
publisher	Name of the publisher
publishedDate	The date of publish
infoLink	Link to get more information about the book on Google Books
categories	Genres of books
ratingsCount	Averaging rating for the book

(a) Book Dataset

Feature	Description
id	The Id of the Book
Title	Book Title
Price	The Price of the Book
User_id	Id of the user who rates the book
profileName	Name of user who rates the book
helpfulness	Helpfulness rating of the review, e.g., 2/3
score	Rating from 0 to 5 for the book
time	Time of giving the review
summary	The summary of a text review
text	The full text of a review

(b) Review Dataset

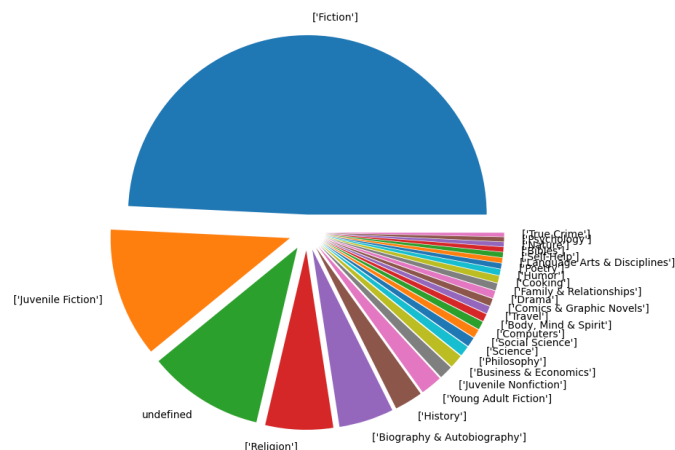
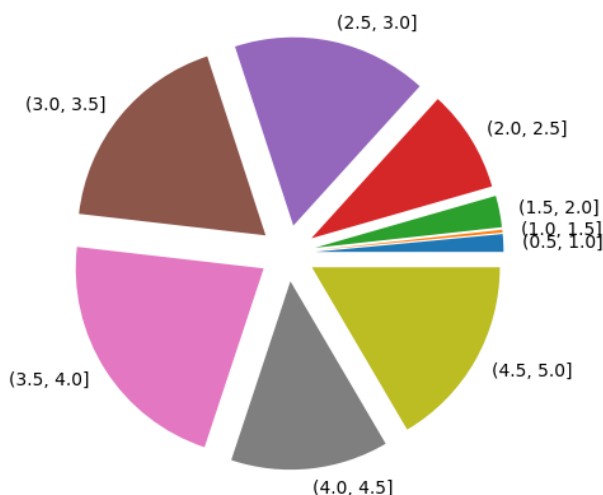


Figure 4.1: Percent of Books by Average Score Figure 4.2: Percentage of Books by Genre

There are 214,404 books in the book dataset and 3,000,000 reviews in the review dataset. Clearly, two very large datasets, we subsetting the datasets and reviews to only the most relevant for our purposes.

First we dropped all the books without any ratings or descriptions as these are irrelevant for our purpose. This reduced the amount of books to 45,127. We then dropped any books with less than 10 reviews, considering them insignificant enough. Furthermore, using the *Langid* package in Python, we removed any books that were identified to be any language other than English. This leaves us with 6,399 books.

We only kept the reviews that were for one of our 6,399 books and used the same method to remove non-English reviews. Because Goodreads allows other users to review other people's reviews, we used this to further subset only the most relevant reviews. Any reviews with less than 10 responses were therefore removed. This brought our review count to 94,573.

4 Exploratory Data Analysis

First, to better understand the dataset we performed some exploratory data analysis. By calculating the average review scores for each book, we find that the largest representation of reviews are within the 3.0 - 4.0 range, with a significant majority of the books receiving a rating above 3.0. The median rating was a 3.6 and the mean was also around 3.6. While this in balance may be reflective of our selection bias of only the most reviewed books (i.e. the most popular books), it is typical for people to tend towards higher reviews as people selectively consume things they tend to enjoy rather than sampling randomly (ex. MovieLens dataset).

By genre we see that the majority of the books are fiction. However, this is only broken down into Juvenile Fiction, Young Adult Fiction, and Fiction whereas non-fiction books are broken down into many topics/subtopics. Genres with less than 15 books were grouped

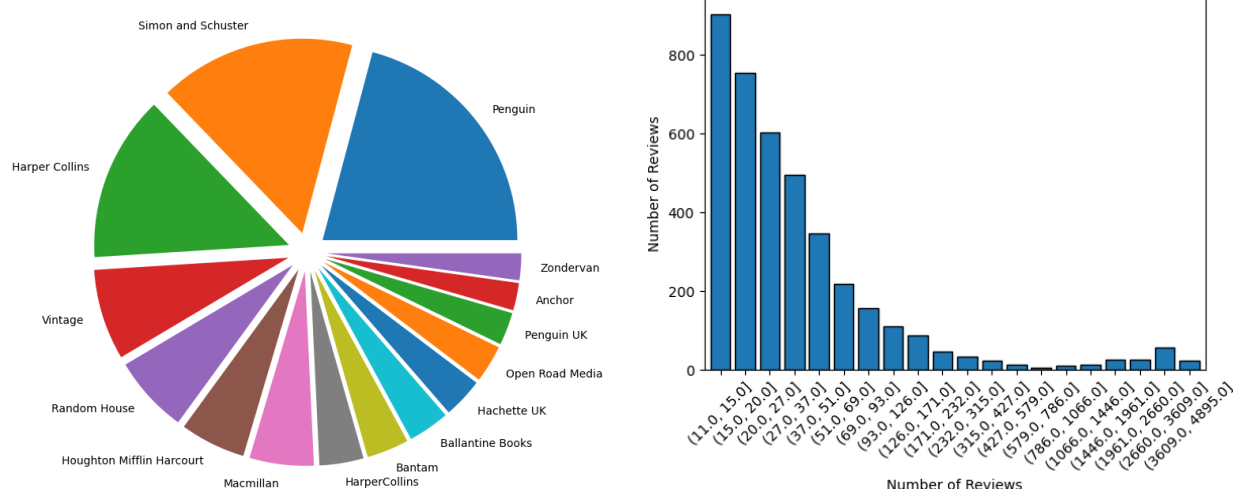


Figure 4.3: Percentage of Books by Publisher Figure 4.4: Percentage of Books by Publisher

into "Other" as there were a number of hyper-specific or uninformative genres. Particularly amusing ones include "Fairy", "Diners", "Bachelors", and "Human-computer Interaction".

We also have the breakdown by publishers. While this was not utilized in this paper, an interesting extension of this project would be to use the description data to see if there are any significant differences among how different publishers market their books.

Finally, we have the counts of books by the number of reviews they received (binned by a log 10 scale). Here we see that a significant majority of the books have received less than 100 reviews. However, there is a significant tail with the maximum review count reaching 4,895.

4.1 Descriptions

4.1.1 New York Times Bestsellers

4.2 Reviews

In the reviews, there are two important text features - the review summary and the full review itself.

4.2.1 Sentiment Analysis

5 Predicting Review Score from Text

6 Review Description Alignment

7 Conclusions

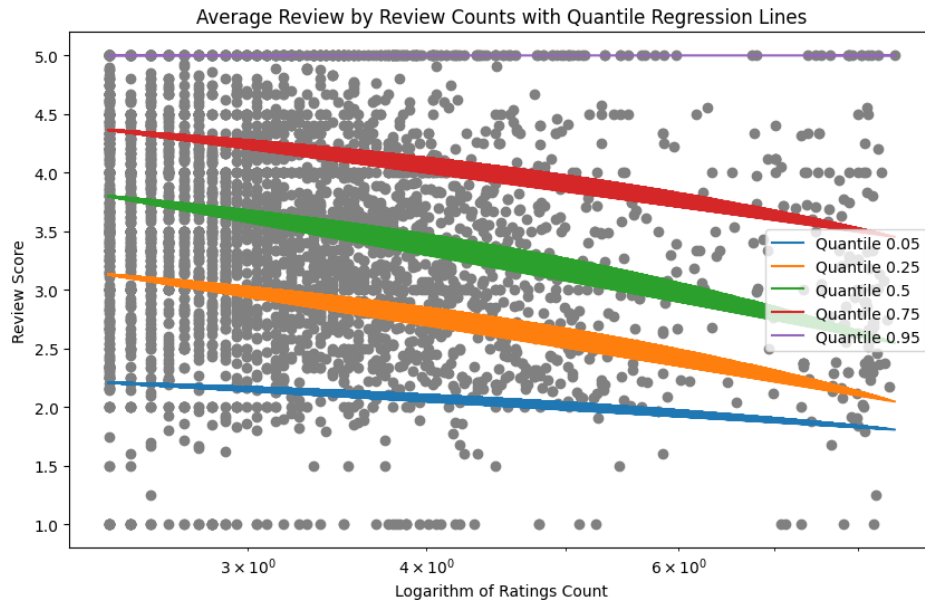


Figure 4.5: Quantile Regression of Average Reviews per Book over Count of Reviews

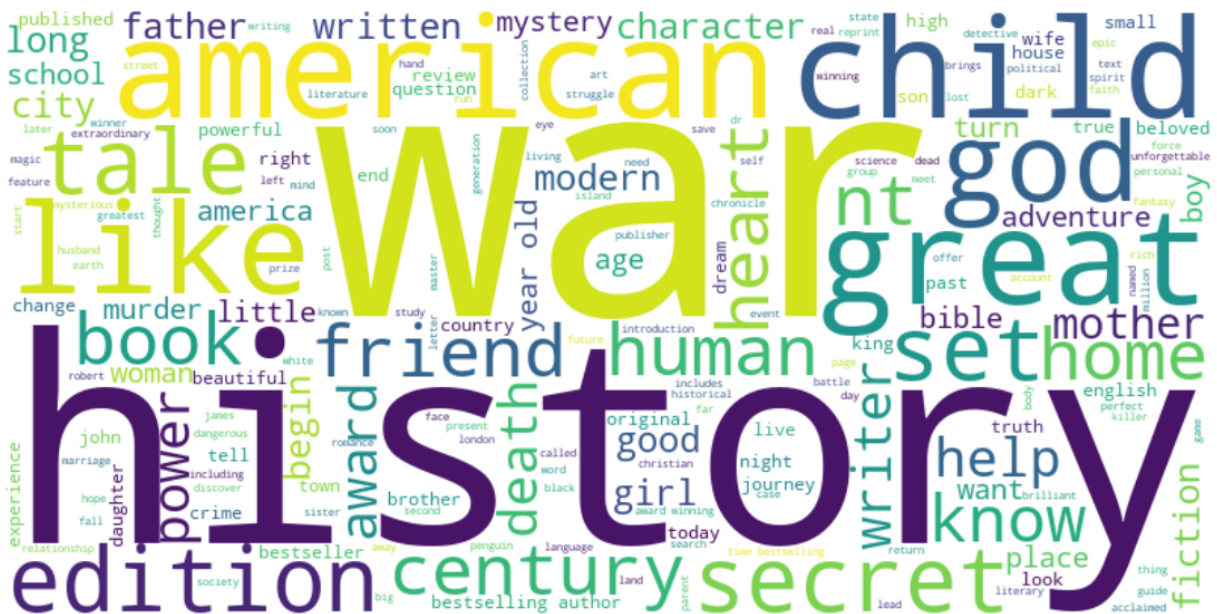


Figure 4.6: Bag of Words - Book Descriptions - MinDF = 20, MaxDF = 0.01