

Business Problem

A recent increase in shipping prices is jeopardizing BigOnlineBookStore's trademark \$1 shipping costs. They are looking to increase sales to keep shipping at just \$1 as their analysis shows the shipping cost is a huge driver for their loyal customers and is responsible for them overtaking the market share of physical books sold online. They are looking to implement a recommendation engine to try and engage users to consider more books. They already have a repository of user information as they track ratings users give, along with a data on the books themselves. BigOnlineBookStore is also extremely concerned about operating cost, so they prefer a simple, fast algorithm with decent results as opposed to a super accurate one which would cost a lot in compute power.

Data

The data used is broken into two sections: book data which includes book_id, author information, series information, format, and others along with user data which includes books read, the rating given, text of review (if present), and similar data.

Data Cleaning & EDA

First, data was cleaned replacing NaN values with averages of the same books, plugging in the minimum year of multiple books when an original year was unavailable, and basic statistics were taken. It is of note that Most books do not have much for ratings:

and most users generally have a lower number of ratings:

It was decided that the best book_id (the one with the most ratings), would be used for each work (as a book can be published in different formats by different publishers). We took into account the data only for the best_book_id. This id was then joined with the consumer interactions table using an inner_join so each entry had a user id, book id, and a wealth of other data (number pages, number reviews, ratings, distribution of ratings by count, etc.).

When looking at this data, it became clear that the user prescribed rating was really correlated with nothing else; if anything, it mostly relied on the book's author's total average rating, but even that was weak:

Therefore, it was decided to use the Surprise package, which takes into account only user id, book id, and rating. The skinny dataset helps ensure the algorithms are simpler and hopefully will address both the needs of BigOnlineBookStore and keep their operating costs low.

Algorithms and ML

5 algorithms were initially tested: SVD, SlopeOne, Baseline AGS, Baseline SGD, CoClustering.

SVD and the 2 Baseline performed the best, so those 3 were further analyzed. To much surprise, the baseline SGD algorithm outperformed SVD at a fraction of the cost. See Metrics below:

Predictions

Analysis of prediction here.

Future Work

While user rating does not correlate much with book information, the book metrics do correlate with each other. Two things could happen in the future: creating a dummy rating for Surprise that encompasses more than just the user rating as well as taking the books themselves and clustering. The algorithms could work in tandem to find the best books and create a hierarchy (promoting the books most like that the user has already identified that they like).