**Inferential Report**
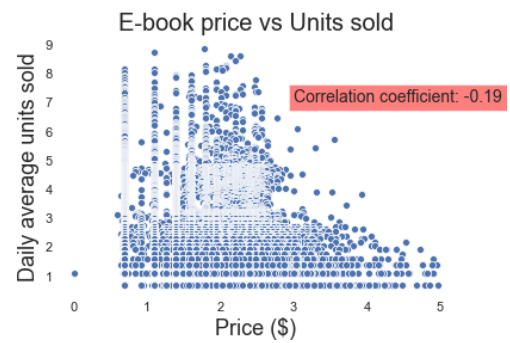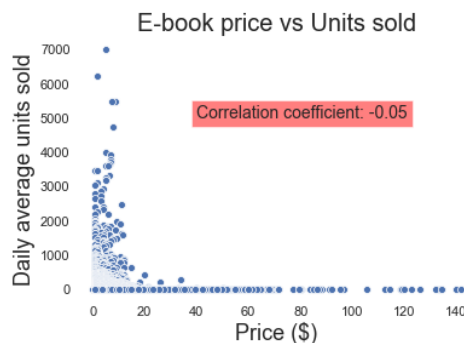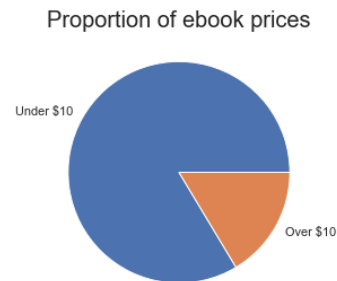
The initial findings from exploratory analysis shows 4 main features effecting e-book sales: price, ratings and reviews, publishers and genre. Most e-book titles average less than 100 units sold daily. However, the data is very skewed with about 100 titles that sell over 1000 copies daily. This is important thing to note before building a predictive model. Nevertheless, the findings are as followed:
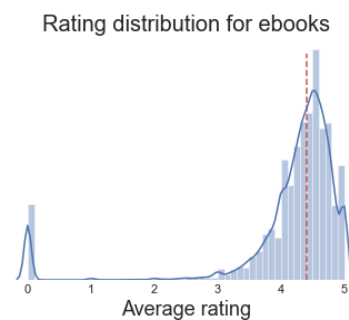
1. Price:
   a. sale_price has a large variance. Almost all books with price higher than $30 is non-fiction. One reason for this is non-fiction includes textbooks, academic books, collection books that are very expensive to produce with high markup.

   
   Proportion of ebook prices

   b. There is a very minor relationship between price and sales without transformation. However, with the right transformation (log), the relationship is better shown.

   

   c. Median price of 5.12 is very closed to average book price of 6.78. Over 75% of e-books has the price tag under 10 dollars.

2. Ratings and reviews
   a. Rating is scored from 1 to 5. If the score is missing, it is replaced by 0. The distribution of rating is skewed right with most of the e-books have 4+ star rating. A small set of e-books that does not have any rating and thus are rated at 0.
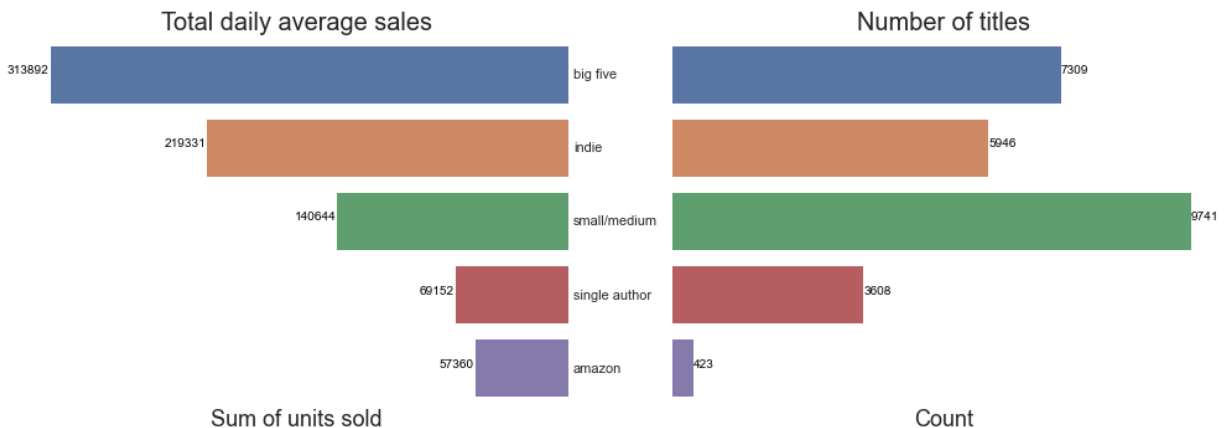
   
   Rating distribution for ebooks

b. Rating does not seem to correlate with sales. The correlation is only at 0.01. This is very understandable as most books have high rating and thus decrease its significance.



c. In contrast, number of reviews effect sales greatly with the correlation coefficient at 0.33. However, would more sales generate the number of reviews or the opposite? While it makes sense that a popular book has more people writing reviews about it, there are many ways that author can do to encourage more people writing the review about the book too. Thus, number of reviews is a legit factor in our analysis.

3. Publishers
    a. There are over 300 publishers in the dataset with smaller ones grouped together as indie, medium, single authors... The big five publishers are: Penguin Random House, HarperCollins, Macmillan, Hachette Book Group, and Simon & Schuster.
    b. With a more resources at hand, one would expect e-book by big 5 do better than the average e-book and it is true:
        i. big five has less titles count but a much better sale figure
        ii. big five has a higher median starting price
    c. However, trying to get published by the big 5 is impractical because they do not accept author manuscript without literary agency.

d. With the same log transform as before, the distribution of e-book price is better shown across different publishers. Big five has the highest median price with balanced variation. All other publishers have lower median price and a long tail toward the top (very high prices)



Book price by publishers in log scale

4. Genre
   a. non-fiction and genre fiction dominate the market in both number of titles and sales. Only these two have multiple titles that sales on average 1000+ copy daily.
   b. Trying to be popular in niche genre such as children (Dr.Seuss) and comics (Stan Lee) is much harder compared to other genres.



Daily average unit sold for each genre



Number of titles for each genre

## Total number of review and sales

Daily average units sold vs Number of reviews
Correlation: 0.33

## Daily average unit sold for each genre

Genre vs Daily average units sold
(nonfiction, genre fiction, children, fiction, comics, foreign language)

## Rating and sales

Daily average units sold vs Average rating
Correlation: 0.01

## Sales distribution for ebooks (< 1000)

Daily average units sold (< 1000)

## Proportion of ebook prices

Under $10
Over $10

## Rating distribution for ebooks

Average rating

### Total daily average sales

| | Sum of units sold |
|---|---|
| big five | 313892 |
| indie | 219331 |
| small/medium | 140644 |
| single author | 69152 |
| amazon | 57300 |

### Number of titles

| | Count |
|---|---|
| big five | 7509 |
| indie | 5946 |
| small/medium | 8741 |
| single author | 3608 |
| amazon | 423 |

## Price vs Units sold

Price ($) vs Daily average units sold
Correlation coefficient: -0.05

## Book price by publishers in log scale

log(price + 1)
(big five, small/medium, amazon, indie, single author)

## Number of titles for each genre

| Genre | Number of titles |
|---|---|
| nonfiction | 14161 |
| genre fiction | 8903 |
| children | 2541 |
| fiction | 733 |
| comics | 568 |
| foreign language | 121 |

## Random forest top 10 feature importance

- total_reviews
- genre_genre fiction
- sale_price
- average_rating
- genre_fiction
- genre_nonfiction
- publisher_type_small/medium
- sold_by_Random House
- sold_by_HarperCollins
- sold_by_Penguin Group (USA) LLC

0.0   0.2   0.4