

# Predicting Amazon Book Reviews Ratings Using Customers Review

Jin Park

# Objectives

- Data Cleaning / Exploratory Data Analysis.
- Predict Ratings using the sentiment analysis.



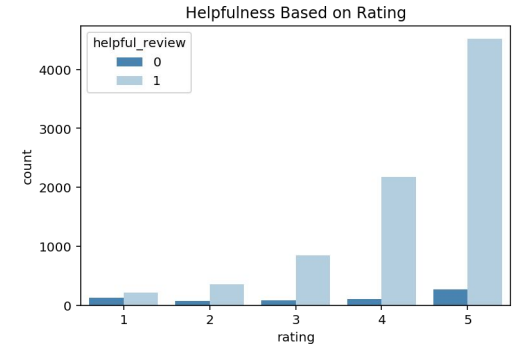
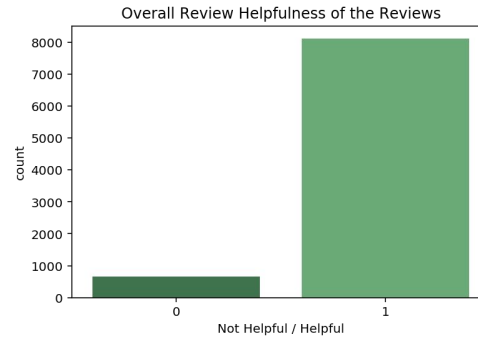
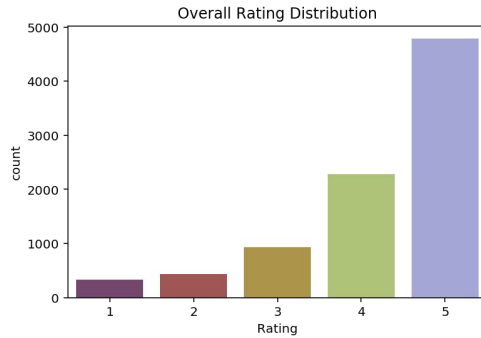
# Data

- The dataset for this project was collected from the University of California San Diego (UCSD) website, <http://jmcauley.ucsd.edu/data/amazon>.
- Dataset contains Amazon customer book reviews and metadata from May 1996 - July 2014.
- Web scraping Amazon.com to gather book categories by using Selenium.



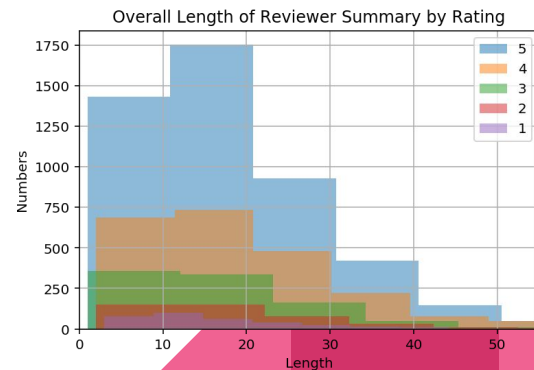
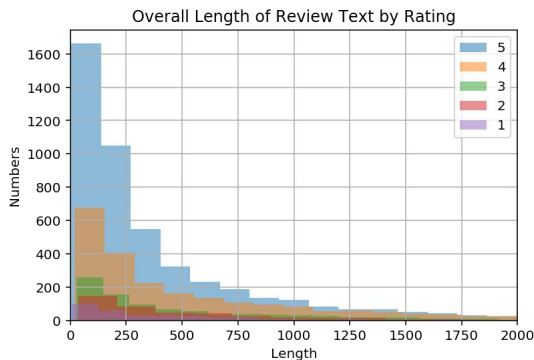
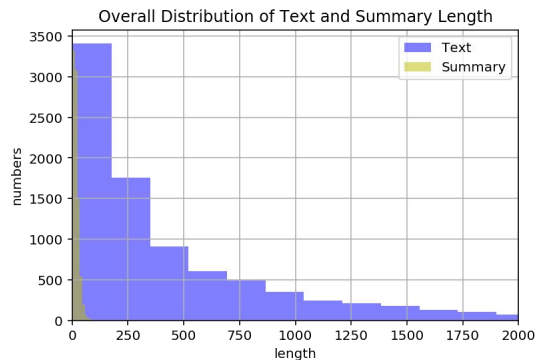
# Exploratory Data Analysis (EDA)

- Distributions of ratings, helpful counts, and helpful counts by ratings



# Reviews Length

- Length of review text and summary by ratings.



# Top Least Occuring Words

review_text	count	review_text	count	review_text	count
0	00	1	0	00 farina	1
1	mitrokhin2	1	1	peopleat fascinating	1
2	mitrokhins	1	2	peopleaside overused	1
3	mitt	1	3	peopleas story	1
4	mittenburg	1	4	peopleas actual	1
5	mittens	1	5	people34 mentality	1
6	mitts	1	6	people34 martin	1
7	mitzi	1	7	people34 dumb	1
8	mixedup	1	8	people34 could	1
9	mixeswhile	1	9	people zombies	1

review_text	count
0	00 farina flour
1	point view well
2	point view world
3	point view worthwhile
4	point view5 question
5	point viewand please
6	point viewcommentators political
7	point view various
8	point viewi would
9	point views concerning

review_summary	count
0	lara
1	nazism
2	nazi
3	nazarea
4	navarro
5	naughty
6	nativity
7	native
8	neal
9	nate

review_summary	count
0	10 best
1	pace star
2	paced adventure
3	paced brutal
4	paced historical
5	paced romance
6	paced suspenseful
7	paced thrill
8	paced unique
9	pacedkept attention

review_summary	count
0	10 best novels
1	pageturner sad sad
2	pages turning rapid
3	pages keep coming
4	page turner memoir
5	page bronies rule
6	page amazonspecifichtml codes
7	paen 80s nerd
8	padding short story
9	packed well ordered

# Most Frequent Words in Reviews Word Clouds

### Overall Most Frequent Words in Review Text

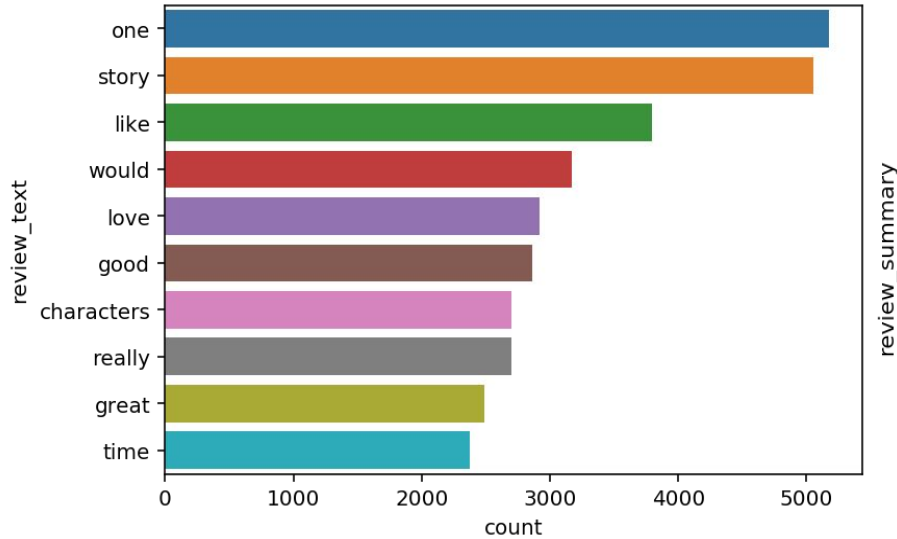


### Overall Most Frequent Words in Review Summary

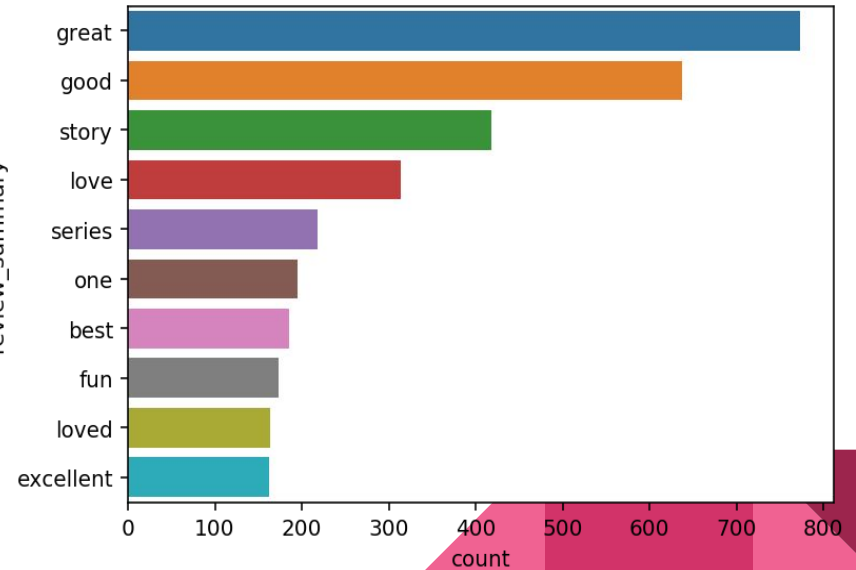


# Unigram

Top 10 Words in Review Text (ngram = 1)

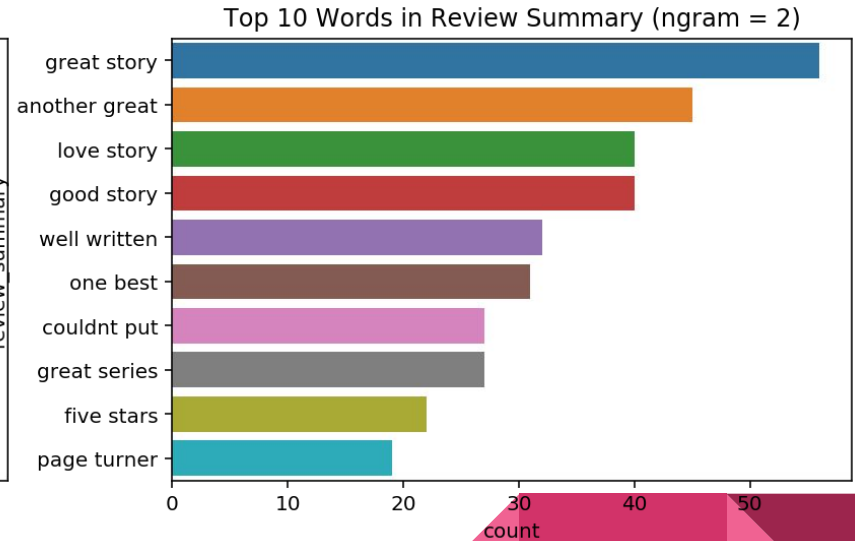
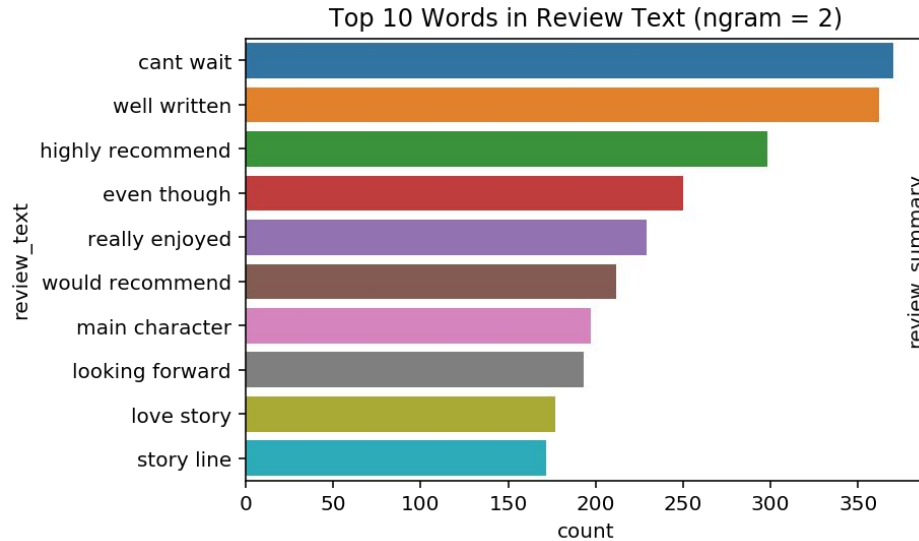


Top 10 Words in Review Summary (ngram = 1)

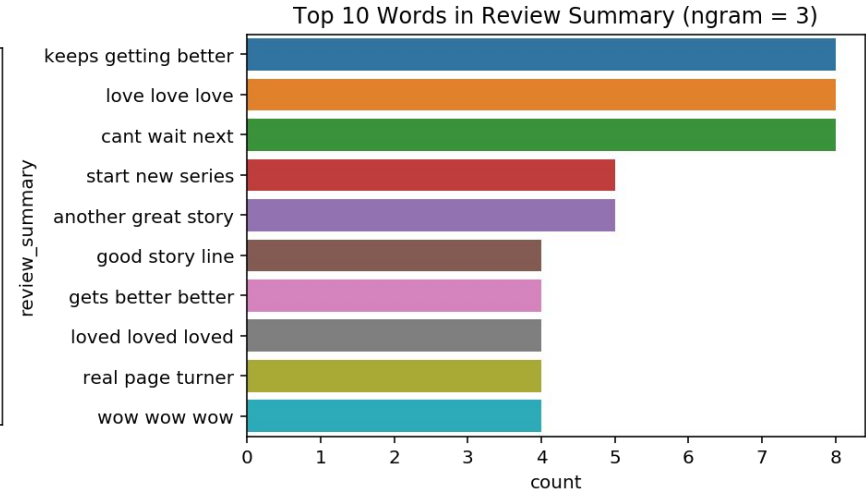
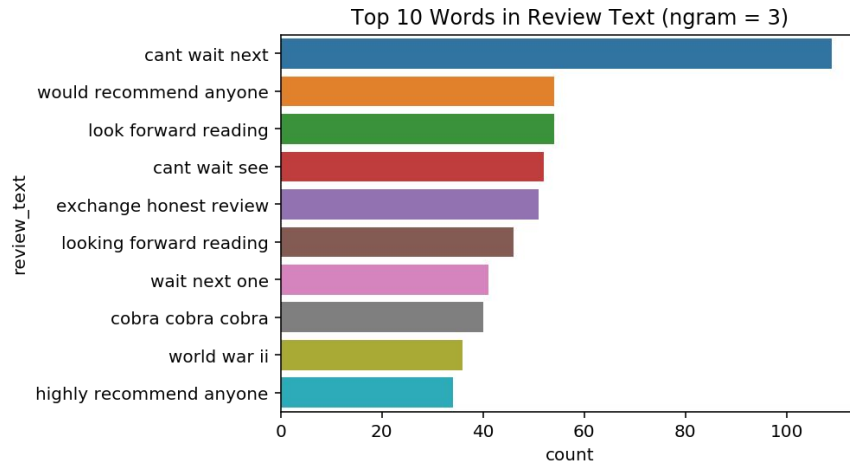




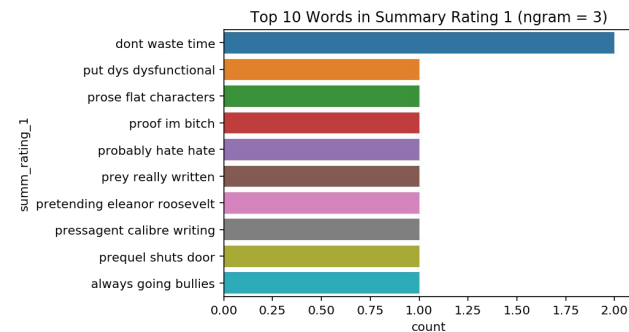
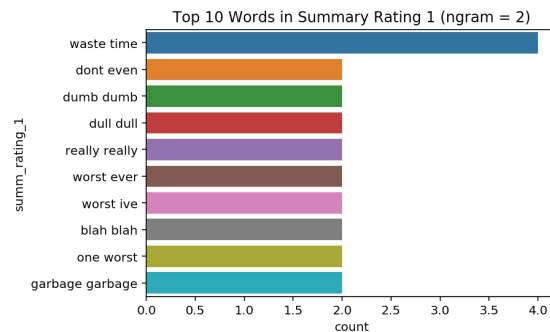
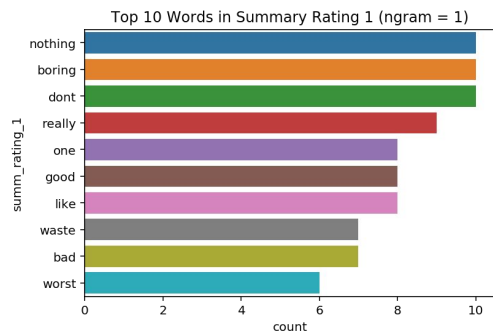
# Bigrams



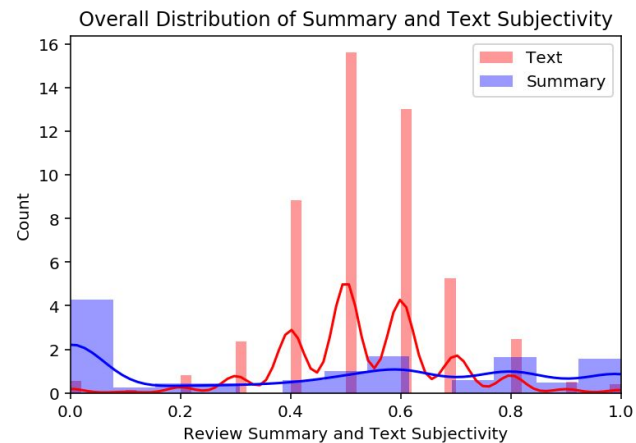
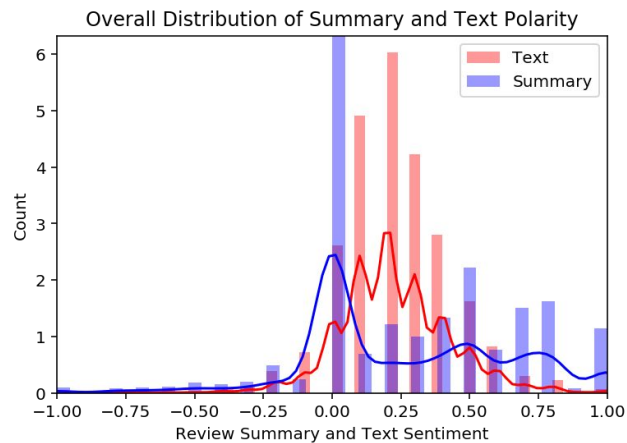
# Trigrams



# Lowest Rating Words

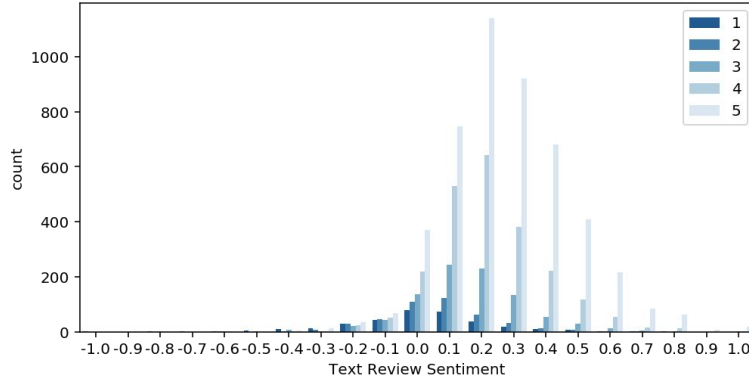


# Sentiment Analysis

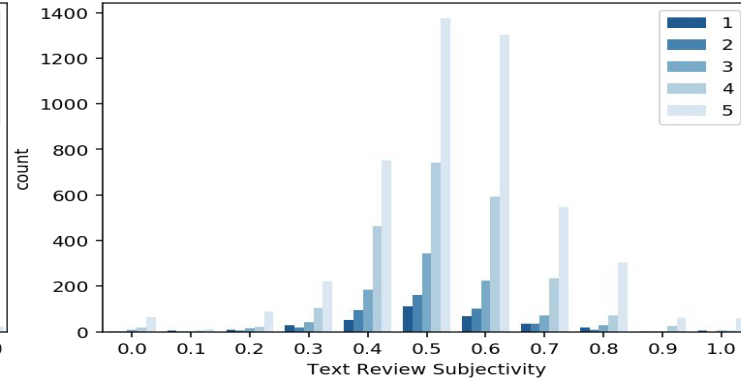


# Sentiment Analysis by Ratings

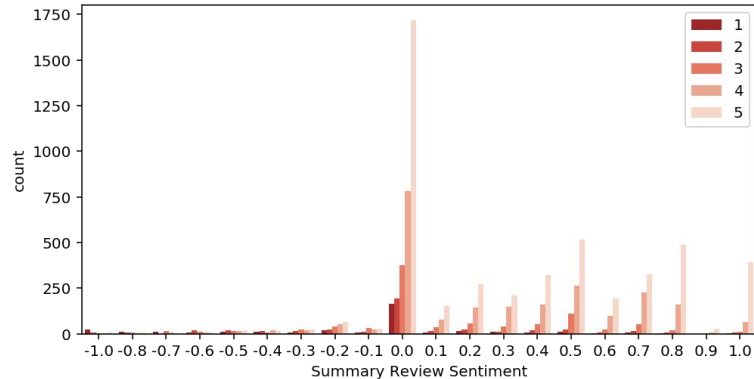
Text Review Sentiment by Ratings



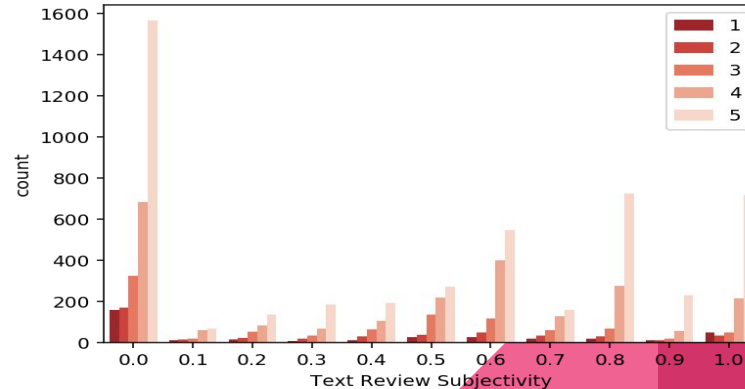
Text Review Subjectivity by Ratings



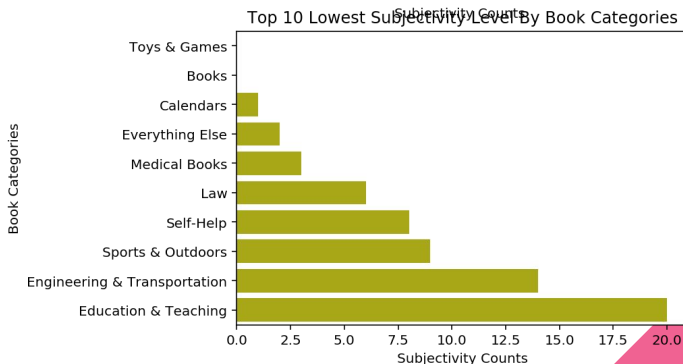
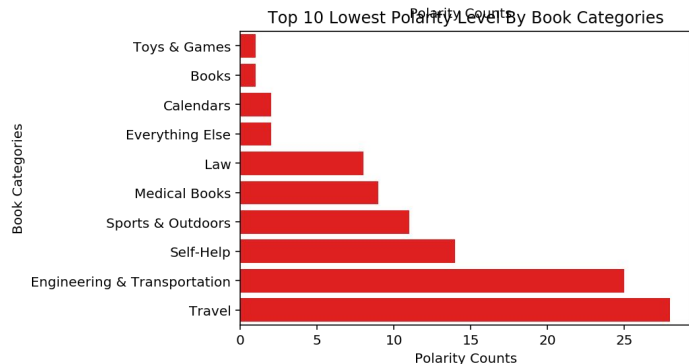
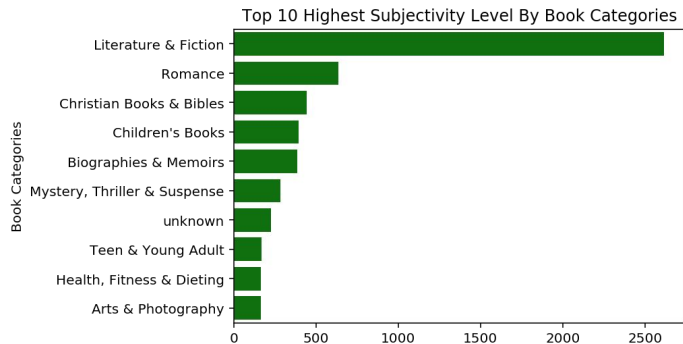
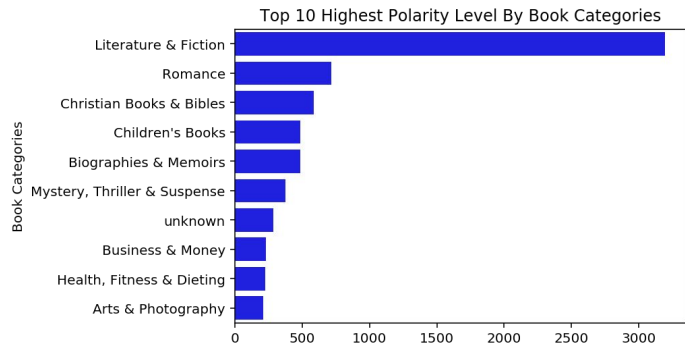
Summary Review Sentiment by Rating



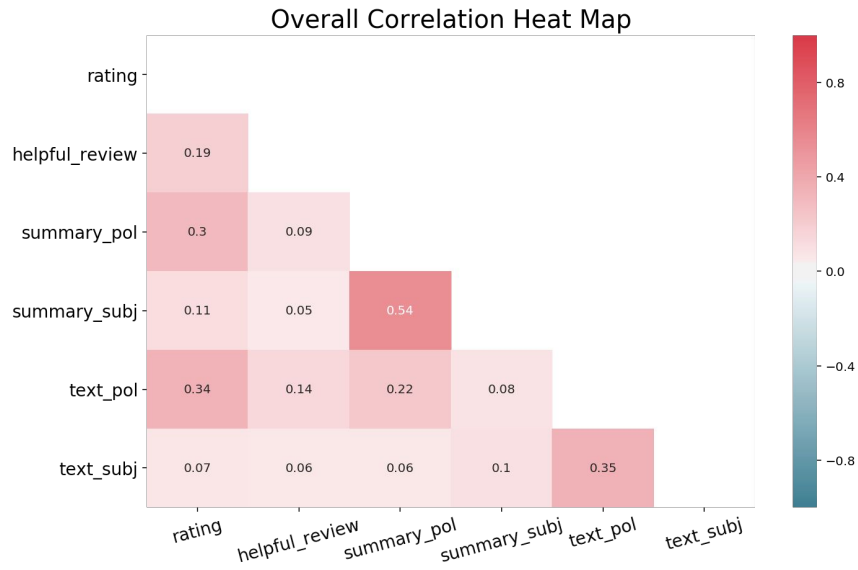
Text Review Subjectivity by Ratings



# Polarity and Subjectivity Level by Book Categories

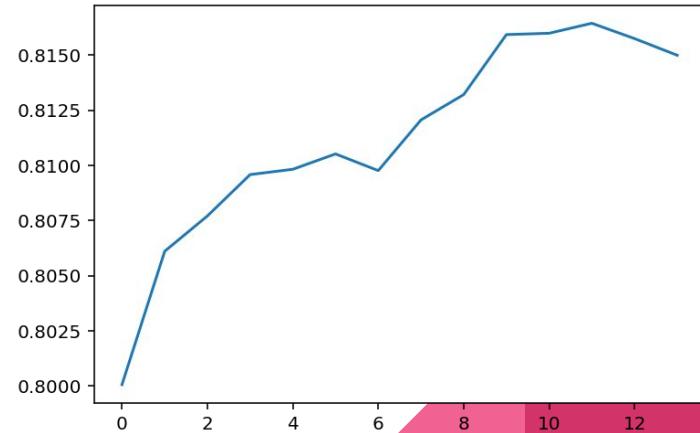
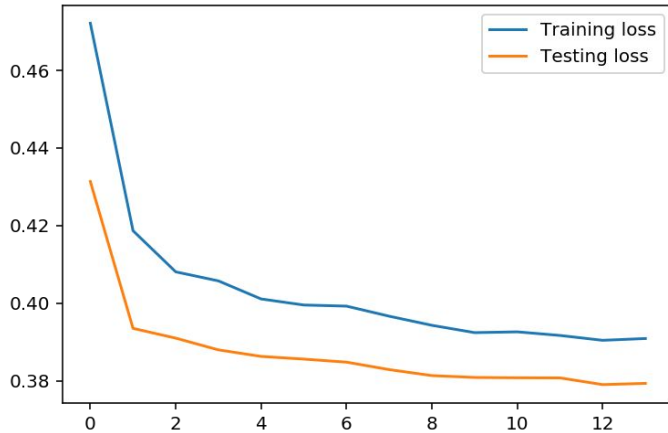


# Correlation Heatmap



# Modeling

- Baseline: 55%
- RandomForest Classifier, Logistic Regression, Neural Network using Keras.
- Neural Network: 82% Accuracy on test set.





# Conclusion

- The product ratings can be predicted just by using customers sentiment analysis.
- Text data itself could not outperform the neural network that uses sentiment analysis as features.



# Any Questions?

This is the end of my presentation and thank you for listening.

