

**Data Science Career Track
Capstone Project-2
Final Report**

**NATURAL LANGUAGE PROCESSING
Sentiment Analysis
Amazon Home and Kitchen Product Reviews**

**Gokmen Oran
June 2019**

CAPSTONE PROJECT 2- FINAL REPORT

1. INTRODUCTION	2
1.1. General	2
1.2. Problem Statement	2
1.3. Data Set Description	2
2. DATA WRANGLING	3
2.1. Inspecting the Dataset	3
2.2. Descriptive Statistics	4
2.3. Preprocessing the Text	5
3. EXPLORATORY DATA ANALYSIS	5
3.1. Target Variable : "rating_class" Feature	5
3.2. Other Features	7
3.2.1. "year" Feature	7
3.2.2. "customer" Feature	7
3.2.3. "product" Feature	8
3.2.4. "review_length" Feature	8
3.2.5. "Text Review" Feature	9
4. NATURAL LANGUAGE PROCESSING	11
4.1 Feature Engineering and Selection	11
4.2 Data Preprocessing	11
4.3 Selecting the Right Evaluation Metric	11
4.4 Modeling	11
4.4.1 Count Vectorizer	12
4.4.2 TF-IDF Vectorizer	13
4.4.3 Hash Vectorizer	14
4.4.4 Adding Most Common and Lest Common Words to Stopwords List (Count Vectorizer)	15
4.4.5 Synthetic Minority Oversampling Technique (SMOTE)	16
4.4.6 Applying PCA to Decrease the Linear Dimensionality + SMOTE	17
4.4.7 Truncated SVD + SMOTE	18
4.4.8 Applying Word2Vec and Simple Neural Network	19
5. CONCLUSION	20
5.1 Recommendations	20
5.2 Future Study	20

CAPSTONE PROJECT 2- FINAL REPORT

1. INTRODUCTION

1.1. General

Sentiment analysis ,which is a subtopics of Natural Language Processing (NLP), has been gradually becoming more and more popular. It is a contextual mining of text which identifies and extracts subjective information in source material and helping a business to understand the social sentiment of their brand, product or service while monitoring online conversations.

Sentiment Analysis has many applications ranging from ecommerce, marketing, to politics and any other research to tackle with text or unstructured text data. Companies, especially in e-commerce, also do sentiment analysis to collect and analyze customer feedback about their products. Besides that, potential customers prefer to review the opinions of existing customers before they purchase a product or use a service of a company. As seen here, there are two parts in e-commerce; one is the online retailer, which wants to maximize e-commerce sales or services, and the other is the consumers, who want to have the best product or service over alternatives.

1.2. Problem Statement

In this project, Amazon is our client. The company wants to develop a software tool that will identify the positive and negative words which customers use when they write reviews for the home and kitchen products as their purchase inclination. For that, they gave their 14 years home and kitchen products' reviews between 2000-2014 and asked us to develop a model which will identify positive and negative words used in the reviews as a component of customer's sentiment towards to the company's home and kitchen products.

According to the customer request, we will build a sentiment analysis model as part of natural language processing, based on their reviews on the home and kitchen product online purchases. Our dataset consists mainly of customers' reviews and ratings.

1.3. Data Set Description

Home and kitchen dataset revolving around the reviews written by customers. This is a real commercial data.

	reviewerID	asin	reviewerName	helpful	reviewText	overall	summary	unixReviewTime	reviewTime
0	A1115ST6F5CWYP	B00000JGRT	Amalfi Coast Girl	[29, 33]	I have had one of these for about 10 years. I...	4.0	good for a first ice cream machine	1148256000	05 22, 2006
1	A188JQXWF4EY1R	B00000JGRT	Ann B. Hibbard "anbee"	[4, 4]	We actually found this product on clearance sa...	4.0	Wonderful Product!	1282176000	08 19, 2010
2	AUAX1QWUCYKSX	B00000JGRT	Ashley S	[1, 1]	This product works great, if the unit kept in ...	5.0	Works as expected	1243555200	05 29, 2009
3	A2C271QUH9N1Z	B00000JGRT	audrey	[12, 13]	After trying other ice cream makers with mixed...	5.0	this will be one of your favorite small applia...	1043712000	01 28, 2003
4	A2PN65B6BSTIYZ	B00000JGRT	B. A. Chaney	[1, 1]	I bought this ice cream maker last summer and ...	5.0	You'll be addicted to homemade ice cream!	1214179200	06 23, 2008

CAPSTONE PROJECT 2- FINAL REPORT

Each row corresponds to a customer review, and includes the variables:

reviewerID : ID of the reviewer, e.g. A2SUAM1J3GNN3B - type: object

asin : ID of the product , e.g. 0000013714 – type: object

reviewerName : name of the reviewer – type: object

helpful : helpfulness of the review, e.g. 2/3 – type: object

reviewText : text of the review – type: object

overall : Rating (1,2,3,4,5)– type: float64

summary : summary of the review – type: object

unixReviewTime : time of the review (unix time) – type: int64

reviewTime : time of the review (raw) – type: object

The data was in Stanford Analysis Project webpage. The original data was in a JSON format there. In order to analyze the data, I should change the data format. For that, I import JSON and decode JSON file with using query in order to convert JSON file to csv file format.

Data Source:

http://seotest.ciberius.info/seo--snap.stanford.edu/data/amazon/productGraph/categoryFiles/reviews_Home_and_Kitchen_10.json.gz

2. DATA WRANGLING

2.1. Inspecting the Dataset

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 25445 entries, 0 to 25444
Data columns (total 9 columns):
reviewerID      25445 non-null object
asin            25445 non-null object
reviewerName     25276 non-null object
helpful         25445 non-null object
reviewText      25445 non-null object
overall         25445 non-null float64
summary         25445 non-null object
unixReviewTime  25445 non-null int64
reviewTime      25445 non-null object
dtypes: float64(1), int64(1), object(7)
```

CAPSTONE PROJECT 2- FINAL REPORT

memory usage: 1.9+ MB

Amazon home and kitchen products data includes 25445 rows (observations) and 9 columns(feature variables) and its memory usage is 1.9+ MB. In the dataset, we have 7 object, 1 float64 and 1 int64 data types.

169 'reviewerName' information is missing in the dataset. Since customers don't give their identity, it may not be reliable to make an analysis on their reviews and ratings. I would prefer to drop the missing values from dataset since we have enough observations to conclude a prediction for sentiment analysis.

I concatenated 'reviewText' and 'summary' since both gave the approximately same type of information about product in text format, and later dropped both 'reviewText' and 'summary' columns.

'helpful' feature was dropped since I didn't need that column for our model.

I classified the 'overall' (ratings) as good (rating 3,4, and 5) and bad (rating 1 and 2) in order to make sentiment analysis. I created a new column named as 'rate_class' from 'overall' column and converted its' values as 'good' and 'bad'. Later, we dropped 'overall' column.

In the dataset, 'reviewerID' and 'reviewerName' were used both for identification of customers. I dropped one of them from the dataset. Preferably, I dropped 'reviewerName' since customer names were not standardized and there were lots of different style to represent them in it.

'unixReviewTime' was dropped since it has already been represented in 'reviewTime' feature in a more understandable format. Also, 'reviewTime' was converted to datetime data type and a new 'year' column was created to make analysis between other variables in the future work. After that, 'reviewTime' column was also dropped.

I renamed the columns in order to improve practicality/readability of coding:

reviewerID : "customer"

asin : "product"

reviewText: This will be concatenated with "summary" and renamed as "review_text"

overall: "rating_class"

reviewTime: "year"

2.2. Descriptive Statistics

In our dataset, we have 1276 reviews, which have bad ratings whereas 24000 reviews which have good ratings.

We have 1395 unique customers and 1171 products in this dataset. Each customer averagely gives 18 reviews for products and on the other hand, there is averagely 22 reviews for each product in the website.

CAPSTONE PROJECT 2- FINAL REPORT

2.3. Preprocessing the Text

Since, text is the most unstructured form of all the available data, various types of noise are present in it and the data is not readily analyzable without any pre-processing. The entire process of cleaning and standardization of text, making it noise-free and ready for analysis is known as text preprocessing. In this section, I apply the following text preprocessing respectively.

Removing HTML tags

We wrote a function to remove the HTML tags which typically does not add much value towards understanding and analyzing text.

Removing accented characters

We wrote a function to convert and standardize accented characters/letters into ASCII characters.

Expanding Contractions

We wrote a function to convert each contraction to its expanded, original form in order to help with text standardization.

Removing Special Characters

We used simple regular expressions (regexes) to remove special characters and symbols which are usually non-alphanumeric characters or even occasional numeric characters.

Lemmatization

We removed word affixes to get to the base form of a word, known as root word.

Removing stopwords

We wrote a function to remove stopwords, which have little or no significance in the text.

Building a Text Normalizer

Based on the functions which we have written above and with additional text correction techniques (such as lowercase the text, and remove the extra newlines, white spaces, apostrophes), we built a text normalizer in order to help us to preprocess the new_text document.

After applying text normalizer to 'the review_text' document, we applied tokenizer to create tokens for the clean text. As a result of that, we had 3070479 words in total.

Eventually, after completing all data wrangling and preprocessing phases, we save the dataframe to csv file as a 'Cleaned_Reviews_Home_and_Kitchen.csv'. After cleaning, we have 25276 observations.

A clean dataset will allow a model to learn meaningful features and not overfit on irrelevant noise. After following these steps and checking for additional errors, we can start using the clean, labelled data to train models in modeling section.

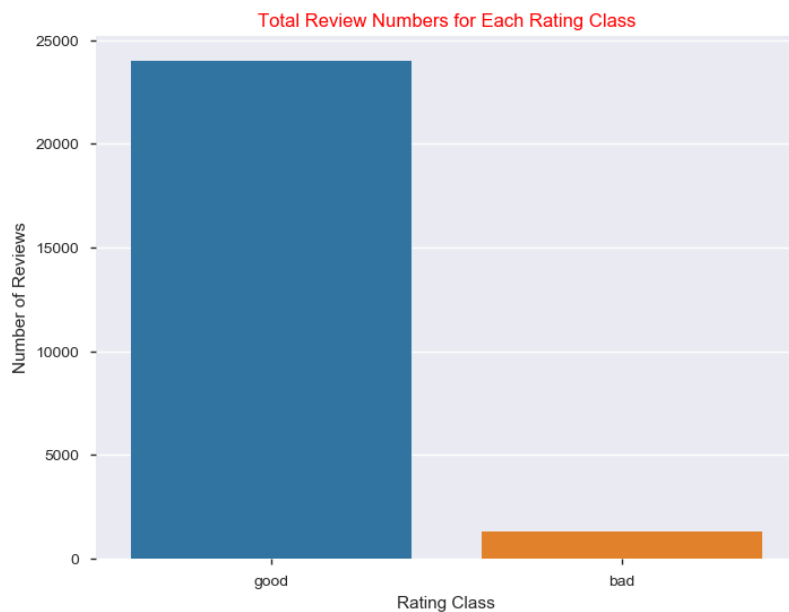
3. EXPLORATORY DATA ANALYSIS

3.1. Target Variable : "rating_class" Feature

CAPSTONE PROJECT 2- FINAL REPORT

Customers wrote reviews and gave ratings, which ranged between 1 to 5, for each home and kitchen product they bought in the Amazon online market between 2000 and 2014. In overall, customers were seemed to be averagely satisfied with the products they purchased.

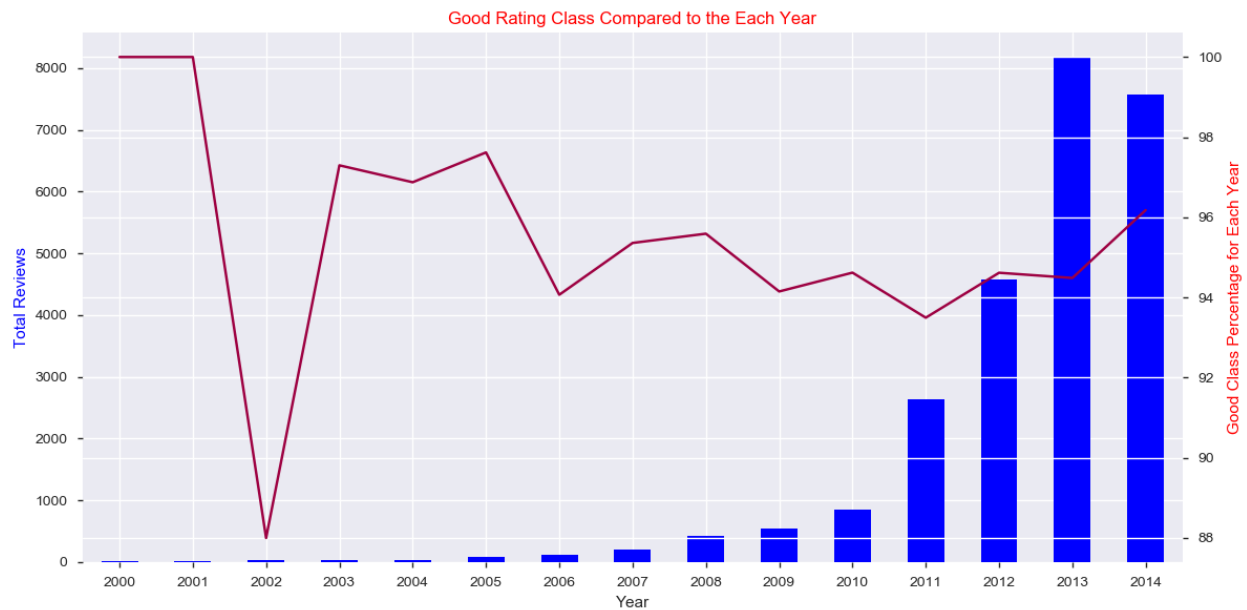
We diminished those 5 rating categories into two categories such as 'good' and 'bad' in order to develop a sentiment analysis model based on their reviews. According to those reviews, 95% of them (24000) are classified as good, whereas 5% of them (1276) are bad.



CAPSTONE PROJECT 2- FINAL REPORT

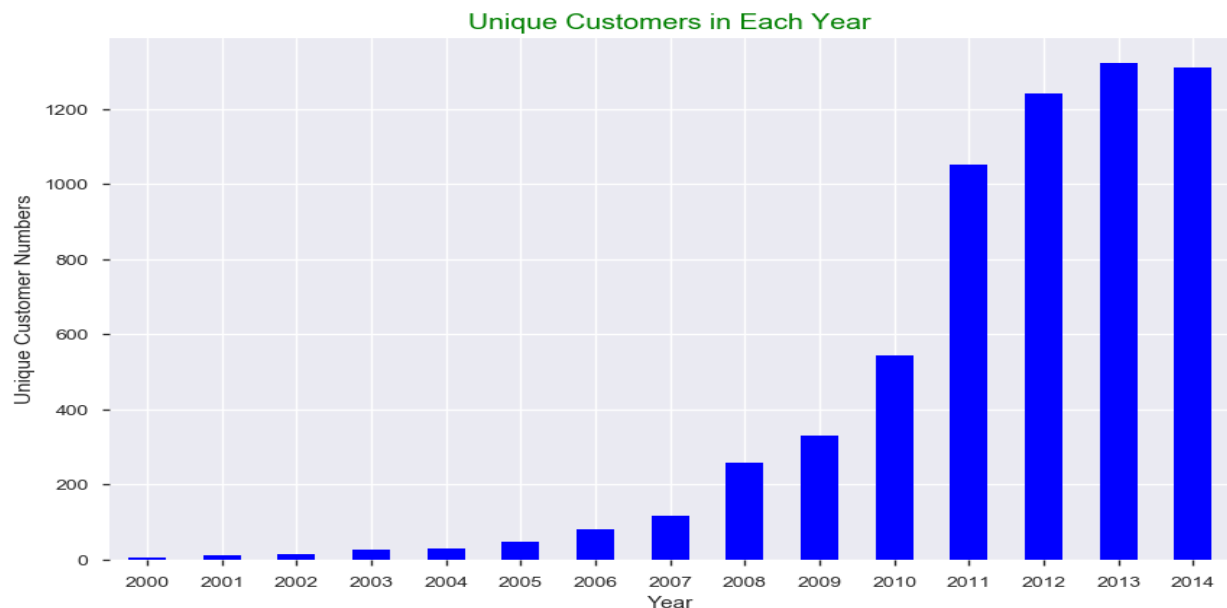
3.2. Other Features

3.2.1. "year" Feature



Except 2002, 'good ratings' percentage is progressing over 92%. 2002 has the lowest good ratings with 88% overall (There are only 25 reviews). 'good ratings' percentage is 100% in 2000 (10 reviews) and 2001 (16 reviews). As it might be seen in the graph, the overall good rating is progressing between 93% and 97% in home and kitchen products.

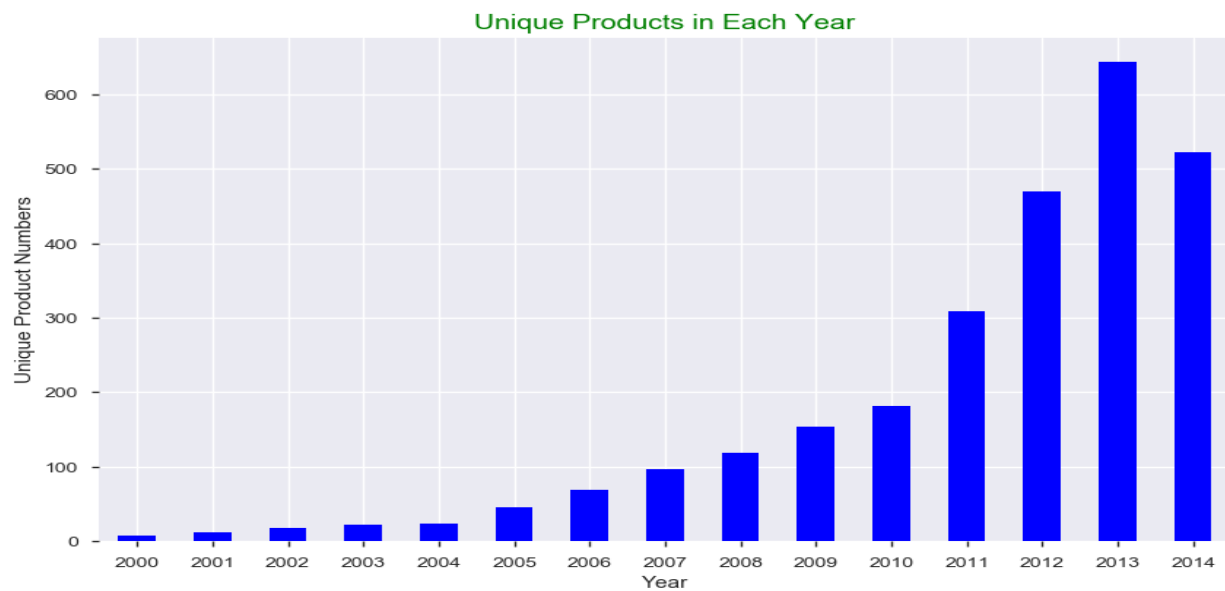
3.2.2. "customer" Feature



We have total 1395 unique customers who gave good reviews and 699 customers who gave bad reviews in the dataset. As it may be observed in the chart and table, the number of unique customers for each year has increased with the progress of the year.

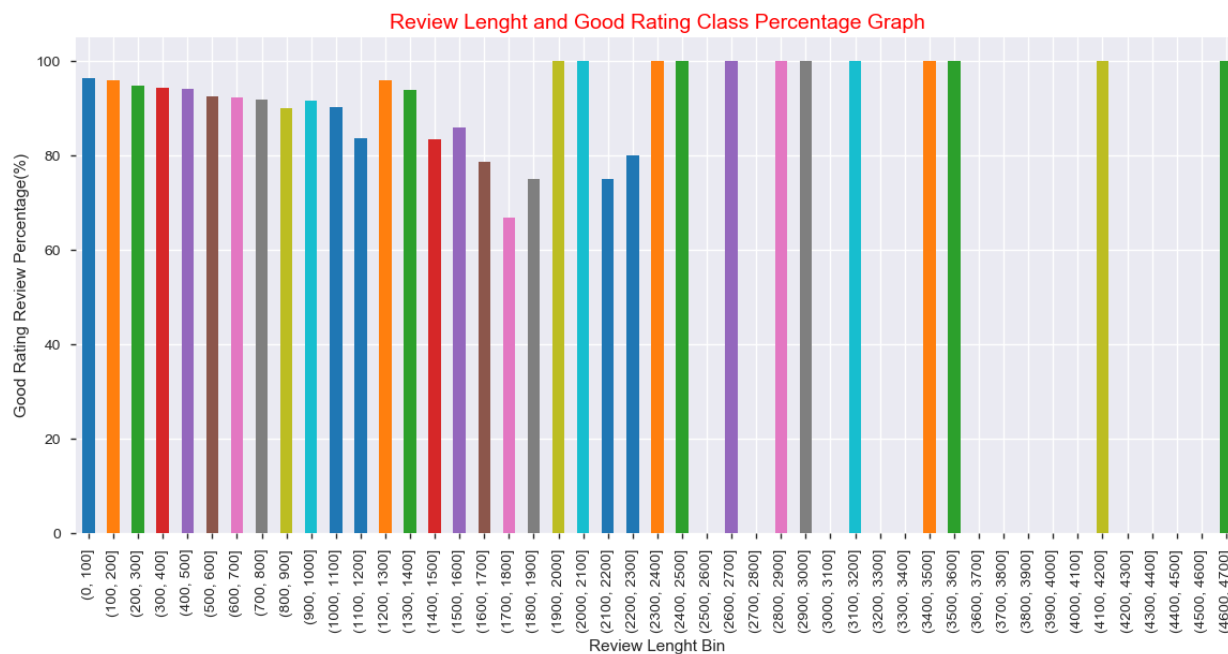
CAPSTONE PROJECT 2- FINAL REPORT

3.2.3. "product" Feature



We have total 1171 unique products in the dataset which belongs to year between 2000 and 2014. As it may be observed in the chart and table, the number of unique products for each year has increased generally with the progress of the year except 2014. There is a slight decrease in 2014 but we have only data until June in 2014.

3.2.4. "review_length" Feature



As it might be seen the graph, the highest percentage of good rating reviews lies between 0-1000 words with 96.2% whereas lowest percentage of good rating reviews

CAPSTONE PROJECT 2- FINAL REPORT

lies between 1700-1800 words with 66.6%. As the review length extends, the good rating tends to increase. Generally, the customers who have write longer reviews (more than 1900 words) tends to give good ratings.

3.2.5 “Text Review” Feature

Good Rating Words:

	words	Avg		words	Avg		words	Avg
1	oven	0.91509	18	end	0.87255	35	nicely	0.85841
2	light	0.91045	19	comfortable	0.87129	36	second	0.85833
3	ever	0.90291	20	happy	0.87075	37	set	0.85827
4	cut	0.90099	21	side	0.87059	38	need	0.85787
5	especially	0.89655	22	new	0.86957	39	every	0.85714
6	might	0.89381	23	highly	0.86932	40	definitely	0.85714
7	done	0.89216	24	sturdy	0.86722	41	something	0.85714
8	find	0.89011	25	cooking	0.86705	42	baking	0.85577
9	try	0.8882	26	know	0.86458	43	kitchen	0.85488
10	dish	0.88288	27	room	0.86441	44	way	0.85484
11	shape	0.87851	28	le	0.86184	45	glass	0.85417
12	high	0.87845	29	food	0.8617	46	gift	0.85402
13	recommended	0.87805	30	cup	0.8617	47	like	0.85342
14	getting	0.87681	31	home	0.86066	48	best	0.85326
15	old	0.87662	32	everything	0.86014	49	fun	0.85294
16	could	0.8764	33	needed	0.85981	50	store	0.85276
17	another	0.87562	34	safe	0.85976			

The most common 50 words, which belong to good rating class, are shown in the table above. Each of these words define which products what kind of good impression have on the customers.

Bad Rating Words:

	words	Avg		words	Avg		words	Avg
1	machine	0.752212	18	say	0.79835	35	made	0.80941
2	perfectly	0.754717	19	cleaning	0.8	36	received	0.80952
3	bottom	0.76506	20	right	0.80091	37	first	0.80969
4	fit	0.776224	21	actually	0.80272	38	may	0.81035
5	big	0.776786	22	using	0.80451	39	problem	0.81035
6	space	0.779874	23	issue	0.80451	40	back	0.81068
7	attractive	0.784314	24	piece	0.80473	41	water	0.81108
8	although	0.784314	25	inside	0.80556	42	take	0.81139

CAPSTONE PROJECT 2- FINAL REPORT

9	three	0.787037	26	item	0.80591	43	nice	0.81191
10	amount	0.790476	27	holder	0.80734	44	long	0.81197
11	though	0.790698	28	day	0.80745	45	grip	0.8125
12	counter	0.79085	29	going	0.80791	46	used	0.81426
13	simple	0.792453	30	floor	0.80833	47	small	0.81461
14	wash	0.792593	31	stick	0.8087	48	enough	0.81544
15	pot	0.793103	32	worth	0.80909	49	worked	0.81553
16	give	0.793548	33	press	0.80916	50	house	0.816
17	little	0.797849	34	pan	0.80928			

Same standards as above, the most common 50 words, which belong to bad rating class, are shown in this table. Likewise, in good ratings, each of these words define which products what kind of bad impression have on the customers.

Controversial Cases:

The controversial case such as "I was expecting better - negative meaning" or "it was better than my expectation - positive meaning " will be handled in the modelling section via using deep learning technique (Keras with Word2Vec).

CAPSTONE PROJECT 2- FINAL REPORT

4. NATURAL LANGUAGE PROCESSING

Due to computational considerations, I reduced the number of observations. I dropped the good_rating_class_reviews longer than 150 words, and I dropped all observations earlier than year 2010. After reducing, there are 8933 observations (7731 "good" reviews and 1202 "bad" reviews).

4.1 Feature Engineering and Selection

Machine Learning models take numerical values as input. The reviews are made of sentences, so in order to extract patterns from the data; we need to find a way to represent it in a way that machine learning algorithm can understand, i.e. as a list of numbers.

We implemented CounterVectorizer, TF-IDF, Hashing Vectorizer, Word2Vec, adding most common words into the stopwords list, SMOTE, PCA, and Truncated SVD techniques into classification models in the following sections as a part of feature engineering and selection.

4.2 Data Preprocessing

Separate response variable and features:

First, we separated features (clean text) and response variable(rating class) and as X and y respectively.

Splitting the dataset into the Training set and Test set:

Then, we divided the dataset into the training and test sets. We have used 25% of dataset as testing set and 75% of the dataset as the training set. We also set the random state of the split to ensure consistent results.

4.3 Selecting the Right Evaluation Metric

Since we have a data imbalance in our case, the evaluation of the classifier performance must be carried out using adequate metrics in order to consider the class distribution and to pay more attention to the minority class. Based on this thought we used f1 score which is harmonic average of precision and recall as my evaluation metric.

Understanding the types of errors our model makes are important. A good way to visualize that information is using a Confusion Matrix, which compares the predictions our model makes with the true label. With that in mind, we used confusion matrix besides our evaluation metric (f1 score).

4.4 Modeling

This is a supervised binary classification problem. We are trying to predict the sentiment based on the reviews left by customers who bought home and kitchen products in Amazon e-commerce online platform. We used Python's Scikit Learn libraries to solve the problem. In this context, we implemented Logistic Regression, Random Forest, Naive Bayes, XGBOOST, CatBoost algorithms and Simple Neural Network as well. Since the ratings of the reviews were not distributed normally, we decided to decrease rating

CAPSTONE PROJECT 2- FINAL REPORT

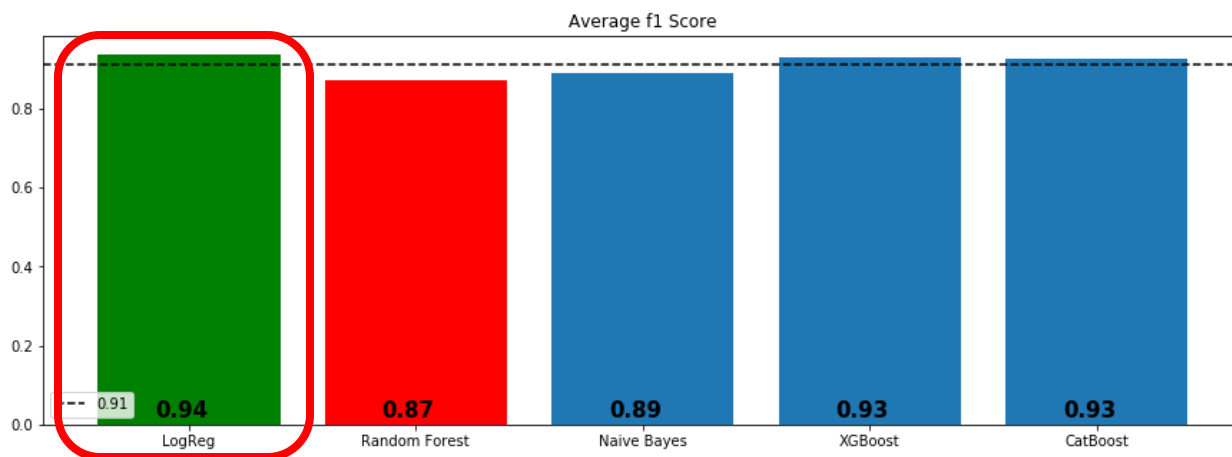
classes from 5 to 2 by merging Rating 1-2 as 'Bad' and Rating 3-4-5 as 'Good' . For feature selection, we applied threshold for word occurrence with using min_df/max_df, PCA and Singular Value Decomposition. For feature engineering, we applied CountVectorizer, TF-IDF, Hashing Vectorizer and Word2Vec to the text data in order to turn a collection of text documents into numerical feature vectors.

Before start modeling, we looked the dummy classifier f1 score. This classifier is useful as a simple baseline to compare with other real classifiers. It is equal to 0.77. It means that randomly selection will yield this score.

4.4.1 Count Vectorizer

			precision	recall	f1-score	support
model	accuracy	class				
LogReg	0.936437	bad	0.745161	0.785714	0.764901	294.0
		good	0.967256	0.959278	0.963251	1940.0
		average	0.938028	0.936437	0.937147	2234.0
Random Forest	0.900627	bad	1.000000	0.244898	0.393443	294.0
		good	0.897317	1.000000	0.945880	1940.0
		average	0.910831	0.900627	0.873178	2234.0
Naive Bayes	0.886750	bad	0.559420	0.656463	0.604069	294.0
		good	0.946533	0.921649	0.933925	1940.0
		average	0.895588	0.886750	0.890515	2234.0
XGBoost	0.935542	bad	0.921348	0.557823	0.694915	294.0
		good	0.936770	0.992784	0.963964	1940.0
		average	0.934741	0.935542	0.928556	2234.0
CatBoost	0.931513	bad	0.840580	0.591837	0.694611	294.0
		good	0.940799	0.982990	0.961432	1940.0
		average	0.927610	0.931513	0.926317	2234.0

We call vectorization the general process of turning a collection of text documents into numerical feature vectors. This specific strategy (tokenization, counting and normalization) is called the Bag of Words or “Bag of n-grams” representation. Documents are described by word occurrences while completely ignoring the relative position information of the words in the document. "CountVectorizer" implements both tokenization and occurrence counting in a single class. **Logistic Regression is the winner with 0.937147.**

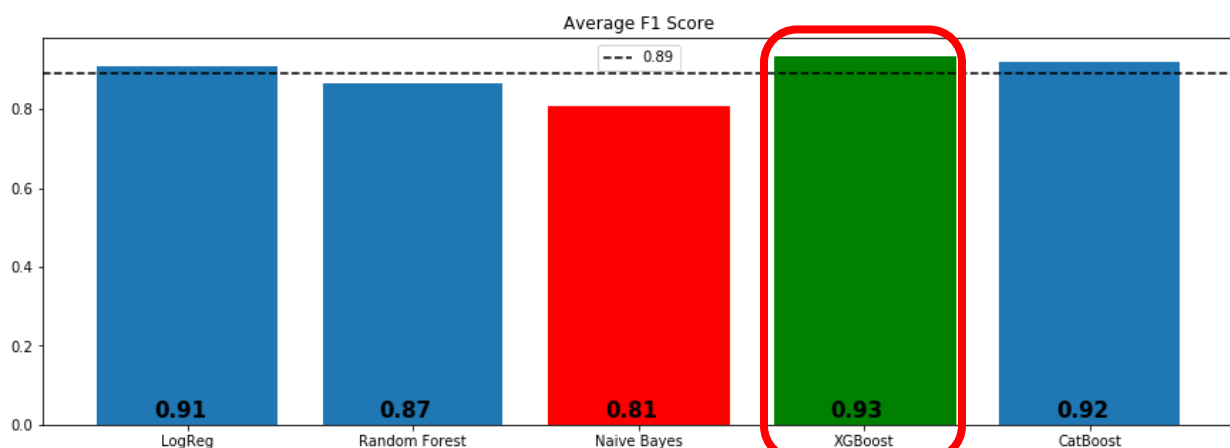


CAPSTONE PROJECT 2- FINAL REPORT

4.4.2 TF-IDF Vectorizer

In order to help our model focus more on meaningful words, we can use a TF-IDF score (Term Frequency, Inverse Document Frequency) on top of our Bag of Words model. TF-IDF weighs words by how rare they are in our dataset, discounting words that are too frequent and just add to the noise. TF-IDF works by penalizing these common words by assigning them lower weights while giving importance to words which appear in a subset of a particular document. **XGboosting is the winner with 0.934794 score.**

			precision	recall	f1-score	support
model	accuracy	class				
LogReg	0.904208	bad	0.595694	0.846939	0.699438	294.0
		good	0.975220	0.912887	0.943024	1940.0
		average	0.925274	0.904208	0.910968	2234.0
Random Forest	0.897493	bad	1.000000	0.221088	0.362117	294.0
		good	0.894421	1.000000	0.944269	1940.0
		average	0.908316	0.897493	0.867656	2234.0
Naive Bayes	0.868397	bad	0.000000	0.000000	0.000000	294.0
		good	0.868397	1.000000	0.929564	1940.0
		average	0.754114	0.868397	0.807231	2234.0
XGBoost	0.941361	bad	0.965714	0.574830	0.720682	294.0
		good	0.939291	0.996907	0.967242	1940.0
		average	0.942768	0.941361	0.934794	2234.0
CatBoost	0.926141	bad	0.824121	0.557823	0.665314	294.0
		good	0.936118	0.981959	0.958491	1940.0
		average	0.921379	0.926141	0.919908	2234.0

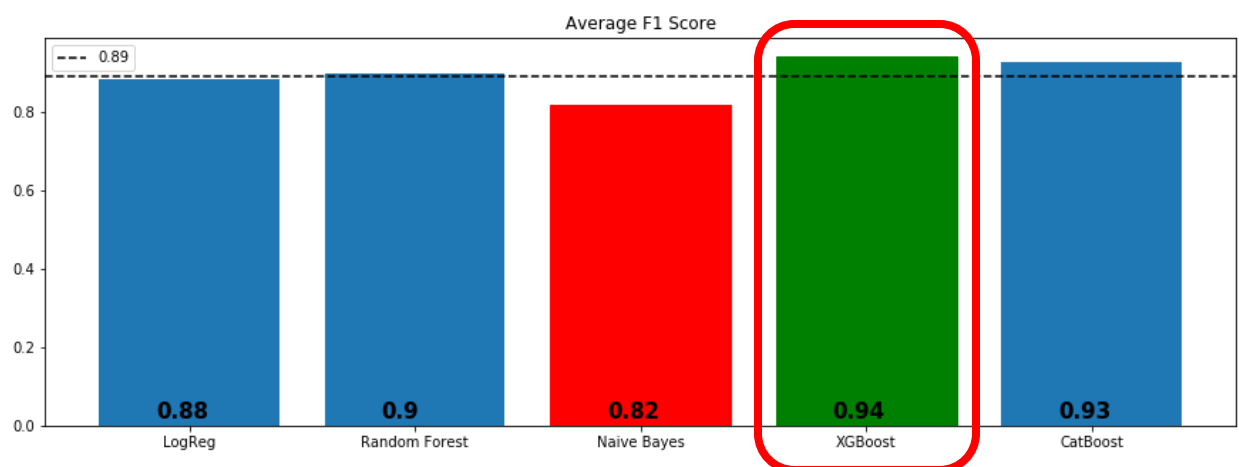


CAPSTONE PROJECT 2- FINAL REPORT

4.4.3 Hash Vectorizer

Hash Vectorizer is designed to be as memory efficient as possible. Instead of storing the tokens as strings, the vectorizer applies the hashing trick to encode them as numerical indexes. The downside of this method is that once vectorized, the features' names can no longer be retrieved. **XGboost is the winner with 0.941092 score.**

			precision	recall	f1-score	support
model	accuracy	class				
LogReg	0.870636	bad	0.505133	0.836735	0.629962	294.0
		good	0.972524	0.875773	0.921616	1940.0
		average	0.911015	0.870636	0.883234	2234.0
Random Forest	0.914951	bad	1.000000	0.353741	0.522613	294.0
		good	0.910798	1.000000	0.953317	1940.0
		average	0.922537	0.914951	0.896635	2234.0
Naive Bayes	0.871979	bad	1.000000	0.027211	0.052980	294.0
		good	0.871518	1.000000	0.931349	1940.0
		average	0.888427	0.871979	0.815753	2234.0
XGBoost	0.946285	bad	0.962766	0.615646	0.751037	294.0
		good	0.944770	0.996392	0.969895	1940.0
		average	0.947139	0.946285	0.941092	2234.0
CatBoost	0.931513	bad	0.847291	0.585034	0.692153	294.0
		good	0.939931	0.984021	0.961471	1940.0
		average	0.927739	0.931513	0.926028	2234.0

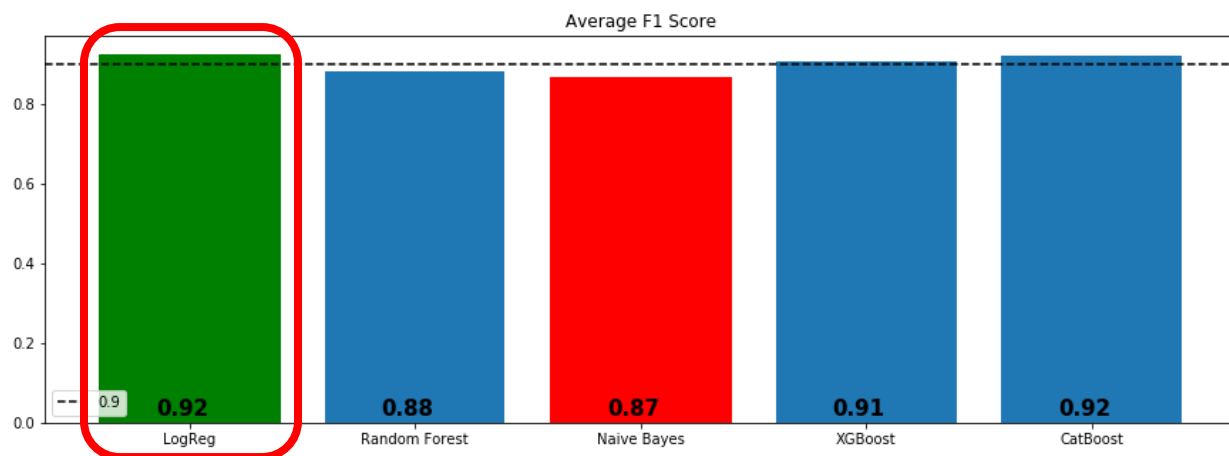


CAPSTONE PROJECT 2- FINAL REPORT

4.4.4 Adding Most Common and Lest Common Words to Stopwords List (Count Vectorizer)

Since there were not too many distinguisher words in different classes, the most and least common 70 words added to the stopwords list and models were applied in order to see any changes in evaluation metrics. **Logistic Regression is the winner with 0.924259 score.** Adding most and least common words to the stopword list didn't have impact on models' performance.

			precision	recall	f1-score	support
model	accuracy	class				
LogReg	0.923456	bad	0.699029	0.734694	0.716418	294.0
		good	0.959481	0.952062	0.955757	1940.0
		average	0.925204	0.923456	0.924259	2234.0
Random Forest	0.905998	bad	0.988372	0.289116	0.447368	294.0
		good	0.902700	0.999485	0.948630	1940.0
		average	0.913975	0.905998	0.882663	2234.0
Naive Bayes	0.854073	bad	0.463964	0.700680	0.558266	294.0
		good	0.950838	0.877320	0.912601	1940.0
		average	0.886764	0.854073	0.865969	2234.0
XGBoost	0.921218	bad	0.953846	0.421769	0.584906	294.0
		good	0.919202	0.996907	0.956479	1940.0
		average	0.923761	0.921218	0.907579	2234.0
CatBoost	0.928380	bad	0.860215	0.544218	0.666667	294.0
		good	0.934570	0.986598	0.959880	1940.0
		average	0.924785	0.928380	0.921292	2234.0

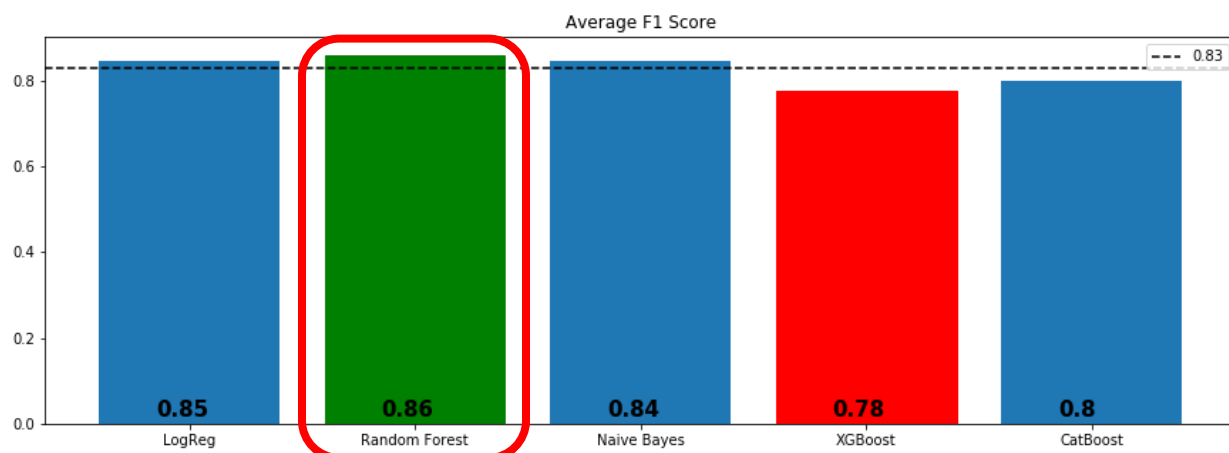


CAPSTONE PROJECT 2- FINAL REPORT

4.4.5 Synthetic Minority Oversampling Technique (SMOTE)

Since we have already noted the imbalance in the values within the target variable, let us implement the SMOTE method in the dealing with this skewed value in order to see whether we may improve our accuracy score. After SMOTE mechanism to improve target class imbalance and identifying best hyper-parameters, f1 score of our model did not show an improvement and decreased to 0.86. We should keep in mind that oversampling will generate artificial observations which may be tricky to evaluate the accuracy of the model. **Random Forest is the winner with 0.858826.**

			precision	recall	f1-score	support	
	model	accuracy	class				
	LogReg	0.834825	bad	0.411765	0.595238	0.486787	294.0
			good	0.934218	0.871134	0.901574	1940.0
			average	0.865462	0.834825	0.846987	2234.0
	Random Forest	0.873321	bad	0.531792	0.312925	0.394004	294.0
			good	0.901989	0.958247	0.929268	1940.0
			average	0.853270	0.873321	0.858826	2234.0
	Naive Bayes	0.829902	bad	0.405702	0.629252	0.493333	294.0
			good	0.938695	0.860309	0.897795	1940.0
			average	0.868552	0.829902	0.844566	2234.0
	XGBoost	0.740824	bad	0.285068	0.642857	0.394984	294.0
			good	0.933164	0.755670	0.835090	1940.0
			average	0.847873	0.740824	0.777171	2234.0
	CatBoost	0.771710	bad	0.310526	0.602041	0.409722	294.0
			good	0.929688	0.797423	0.858491	1940.0
			average	0.848204	0.771710	0.799432	2234.0

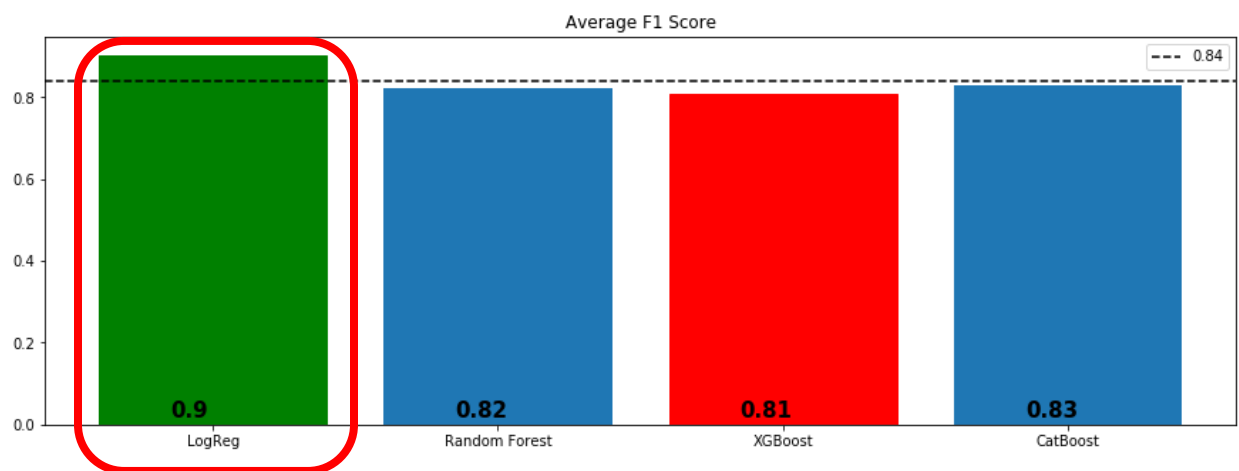


CAPSTONE PROJECT 2- FINAL REPORT

4.4.6 Applying PCA to Decrease the Linear Dimensionality + SMOTE

We have used Principal Component Analysis (PCA) which is a dimension-reduction tool and reduces a large set of variables to a small set that still contains most of the information in the dataset. **Logistic Regression is the winner with 0.903021.**

			precision	recall	f1-score	support
model	accuracy	class				
LogReg	0.900179	bad	0.605341	0.693878	0.646593	294.0
		good	0.952557	0.931443	0.941882	1940.0
		average	0.906862	0.900179	0.903021	2234.0
Random Forest	0.870636	bad	0.592593	0.054422	0.099688	294.0
		good	0.874037	0.994330	0.930311	1940.0
		average	0.836998	0.870636	0.820999	2234.0
XGBoost	0.863921	bad	0.187500	0.010204	0.019355	294.0
		good	0.868801	0.993299	0.926888	1940.0
		average	0.779140	0.863921	0.807454	2234.0
CatBoost	0.864369	bad	0.447059	0.129252	0.200528	294.0
		good	0.880875	0.975773	0.925899	1940.0
		average	0.823784	0.864369	0.830438	2234.0



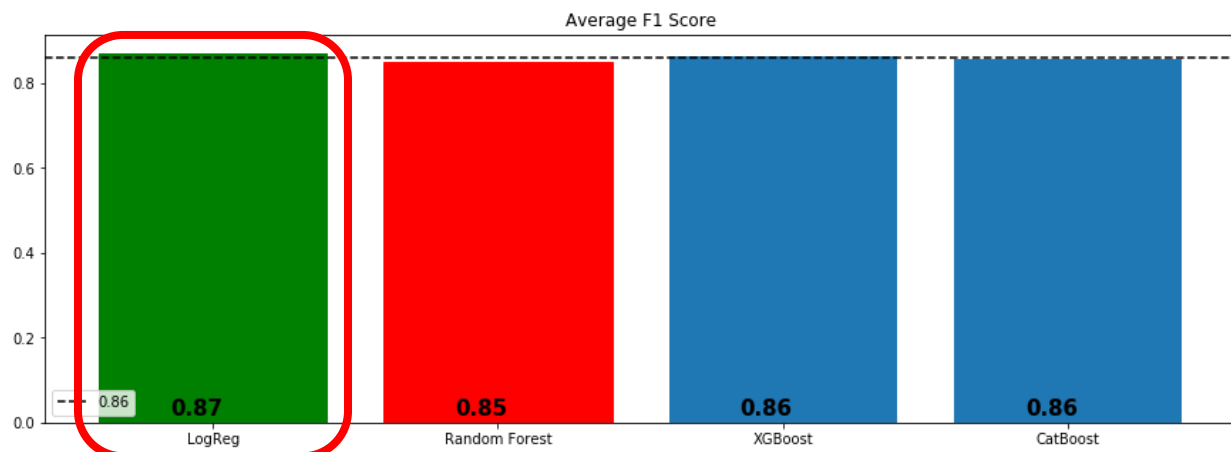
CAPSTONE PROJECT 2- FINAL REPORT

4.4.7 Truncated SVD + SMOTE

Since the dimension is reduced and some information is lost, the f1 score for truncated SVD models are worst of all.

Logistic Regression is the winner with 0.870563 score.

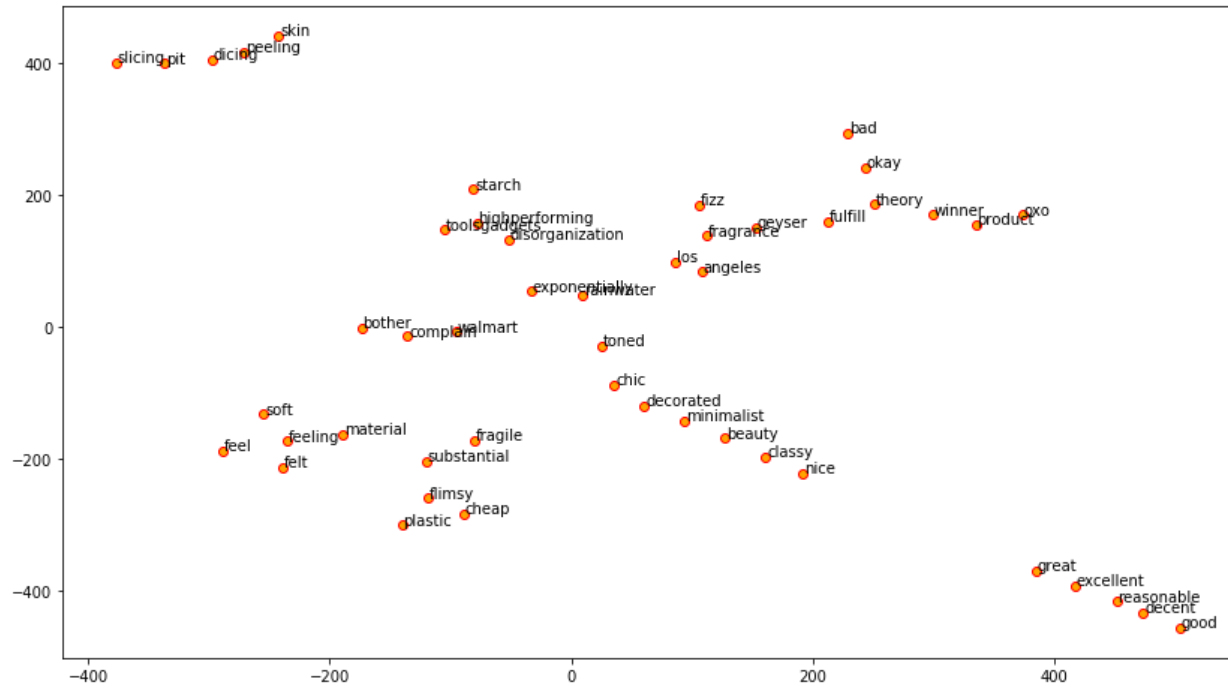
model	accuracy	class	precision	recall	f1-score	support
LogReg	0.860340	bad	0.478774	0.690476	0.565460	294.0
		good	0.949724	0.886082	0.916800	1940.0
		average	0.887746	0.860340	0.870563	2234.0
Random Forest	0.835273	bad	0.420259	0.663265	0.514512	294.0
		good	0.944068	0.861340	0.900809	1940.0
		average	0.875133	0.835273	0.849971	2234.0
XGBoost	0.851835	bad	0.458797	0.700680	0.554509	294.0
		good	0.950700	0.874742	0.911141	1940.0
		average	0.885965	0.851835	0.864207	2234.0
CatBoost	0.844673	bad	0.441501	0.680272	0.535475	294.0
		good	0.947221	0.869588	0.906745	1940.0
		average	0.880667	0.844673	0.857885	2234.0



CAPSTONE PROJECT 2- FINAL REPORT

4.4.8 Applying Word2Vec and Simple Neural Network

We created word vectors using Word2Vec and the model has 20149 unique words where each word has a vector length of 100. Then we used these dense vectors - word embeddings - in a simple neural network to predict. In training and validation accuracy graph, the model starts to overfit after 3th epoch. The accuracy for this simple neural network is 0.9161.



```
{  
  'feel': ['feeling', 'felt', 'soft', 'substantial', 'material'],  
  'good': ['decent', 'great', 'nice', 'excellent', 'reasonable'],  
  'product': ['winner', 'toolsgadgets', 'disorganization', 'oxo', 'highperforming'],  
  'cheap': ['flimsy', 'fragile', 'plastic', 'walmart', 'exponentially'],  
  'beauty': ['chic', 'toned', 'minimalist', 'decorated', 'classy'],  
  'bad': ['complain', 'okay', 'theory', 'fulfill', 'bother'],  
  'skin': ['peeling', 'starch', 'slicing', 'dicing', 'pit'],  
  'fragrance': ['fizz', 'los', 'angeles', 'geyser', 'rainwater']  
}
```

Visualization of the words of interest and their similar words using their embedding vectors after reducing their dimensions to a 2-D space with t-SNE is presented above. Similar words based on gensim's model can be viewed as well.

CAPSTONE PROJECT 2- FINAL REPORT

5. CONCLUSION

In this project, we tried to predict the rating scores based on the reviews left by the customers. Before going through the model result, it is explicitly shown that data set preparation and feature engineering are as much as important as the model creation. We applied;

- Count Vector, TF-IDF, Hashing Vector, Word2Vec
- Classification Models and Simple Neural Network
- Adding most and least common words to CountVect,
- SMOTE,
- PCA + SMOTE
- Truncated SVD + SMOTE

We can go with XGboosting with Hash Vectorizing (f1 score is 0.941092) or Logistic Regression with Count Vectorizing (f1 score is 0.936437) in deploying section. Adding most and least common words to the stopword list didn't have impact on models' performance. Resampling technique and linear dimensionality reduction did not make a positive improvement of the model accuracy but decreased the model performance.

5.1 Recommendations

We recommend the client to use the model as it is and give us some time to develop neural network algorithms to get better scores. By the way, prediction time and the size of the data set are also important and need to be considered. Especially neural network algorithms may not be outperforming with the low size data sets.

Provide a system that user can write their reviews/feedbacks without mistakes such as not accepting the review if it does not satisfy the word/grammar correctness.

Encourage users to reflect their experience with products via giving incentives.

Use the reviews as a feedback mechanism, take actions towards them, and let the customers know about the conclusion.

5.2 Future Study

The future studies on the similar data may focus on;

- Using different methods in order to minimize the effect of the matching words
- Implementation of deep learning with different neural network types and different layer combinations.
- Using different AutoML tools.
- Implementation of Dask library for parallel processing to decrease run time.