



Mentor



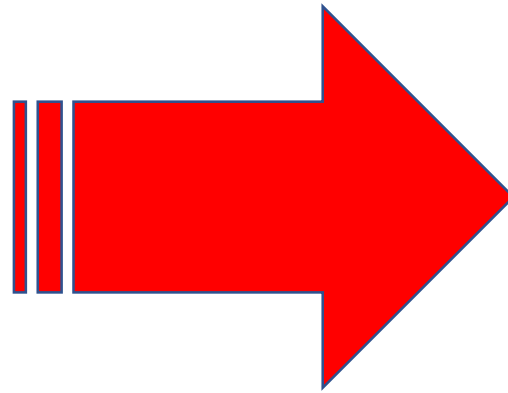
A J Sanchez

Mustafa KADIOGLU Women ECommerce

Natural Language Processing Capstone Project
Springboard Data Science Career Track April-2018 Cohort
github.com/mustafakadioglu



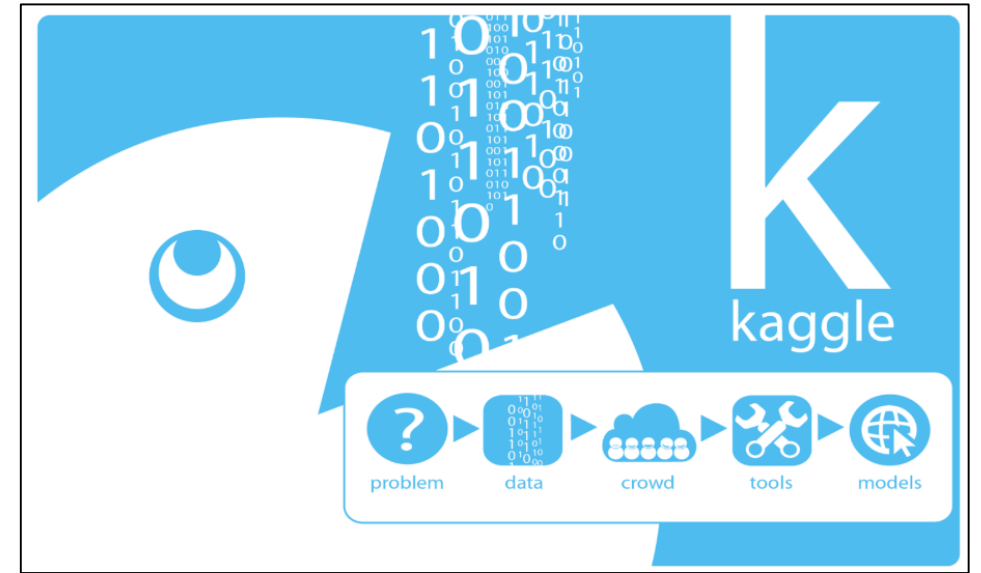
Problem Definition



The selling rates mostly depend on the reviews and ratings left by the customers which shows how they are satisfied with the product. That is the reason why it becomes crucial to predict whether customers will leave a good, neutral or bad rating based on their reviews.

Data Information

Women's Clothing E-Commerce dataset revolving around the reviews written by customers. Its nine supportive features offer a great environment to parse out the text through its multiple dimensions. Because this is real commercial data, it has been anonymized, and references to the company in the review text and body have been replaced with "retailer".



<https://www.kaggle.com/nicapotato/women-s-e-commerce-clothing-reviews/home>

23486 rows and 10 features

Data Information

	Unnamed: 0	Clothing ID	Age	Title	Review Text	Rating	Recommended IND	Positive Feedback Count	Division Name	Department Name	Class Name
0	0	767	33	NaN	Absolutely wonderful - silky and sexy and comf...	4	1	0	Initmates	Intimate	Intimates
1	1	1080	34	NaN	Love this dress! it's sooo pretty. i happene...	5	1	4	General	Dresses	Dresses
2	2	1077	60	Some major design flaws	I had such high hopes for this dress and reall...	3	0	0	General	Dresses	Dresses
3	3	1049	50	My favorite buy!	I love, love, love this jumpsuit. it's fun, fl...	5	1	0	General Petite	Bottoms	Pants
4	4	847	47	Flattering shirt	This shirt is very flattering to all due to th...	5	1	6	General	Tops	Blouses

Basic Information of the Data Set

```
1 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 23486 entries, 0 to 23485
Data columns (total 10 columns):
Clothing ID      23486 non-null int64
Age              23486 non-null int64
Title            19676 non-null object
Review Text      22641 non-null object
Rating           23486 non-null int64
Recommended IND  23486 non-null int64
Positive Feedback Count  23486 non-null int64
Division Name    23472 non-null object
Department Name  23472 non-null object
Class Name       23472 non-null object
dtypes: int64(5), object(5)
memory usage: 1.8+ MB
```

The data set has 10 columns and 23486 data points. There are some null values under Title, Review Text, Division Name, Department Name and Class Name columns.

Feature Engineering

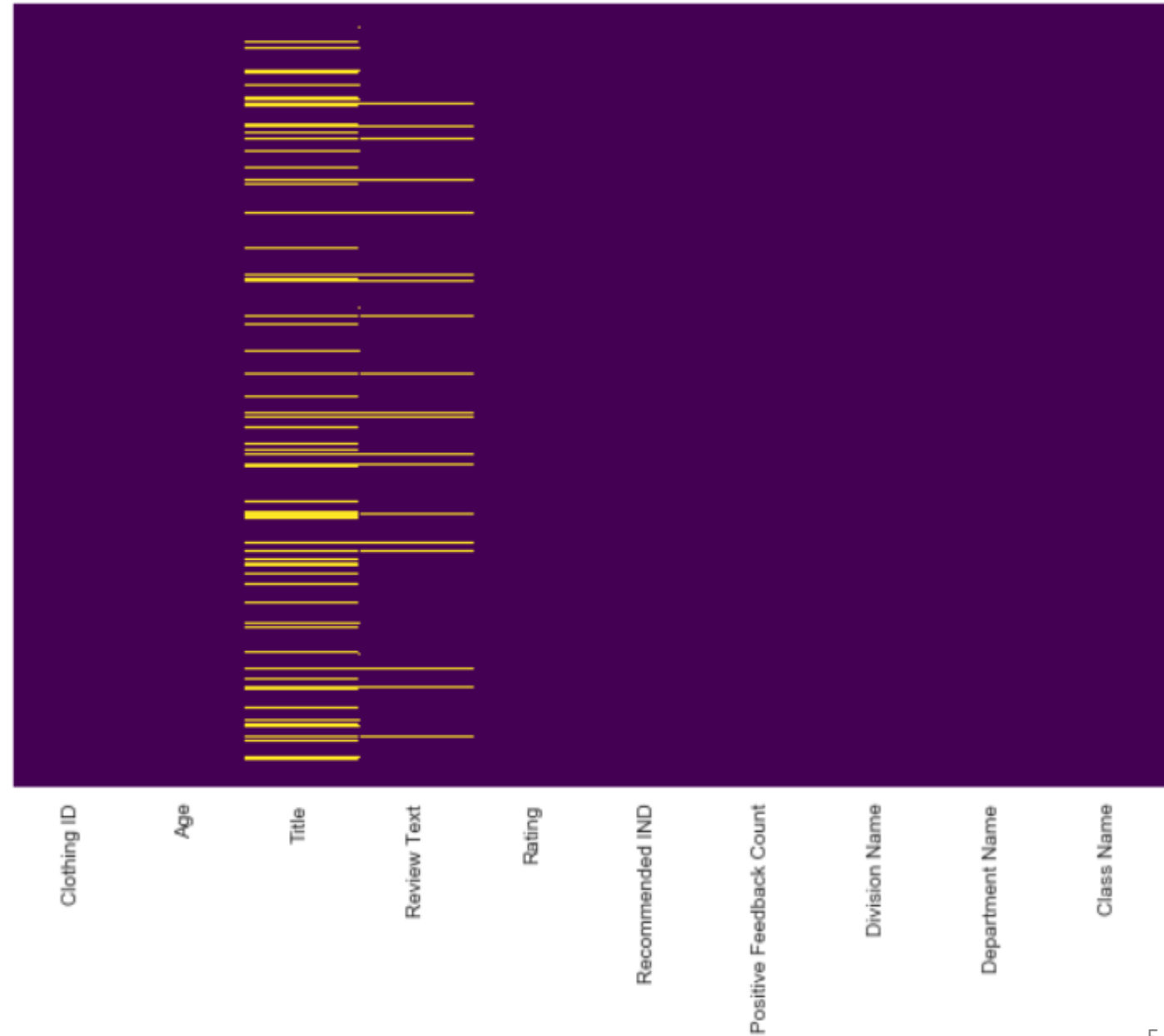
Data Set:
Features & Missing Values



Either Non-Null Values

Title

Review Text



Feature Engineering

Renaming Columns Based on _ and Lowercase Rules

```
1 def cleanup_column_names(df, rename_dict={}, do_inplace=True):
2     if not rename_dict:
3         return df.rename(columns={col: col.lower().replace(' ', '_')
4                                for col in df.columns.values.tolist()},
5                             inplace=do_inplace)
6     else:
7         return df.rename(columns=rename_dict, inplace=do_inplace)
```

```
1 cleanup_column_names(df)
```

Concatenating the Title and Review Text Columns (Based on Either Non-Null Values)

```
1 df2 = df[df.title.notnull() | df.review_text.notnull()]
2 df2.review_text.astype(str)
3 df2.title.astype(str)
4 df2['new_text'] = df2[['title', 'review_text']].apply(lambda x: ' '.join(str(y) for y in x if str(y) != 'nan'), axis=1)
5 df2.drop('title', axis = 1, inplace = True)
6 df2.head()
```


Feature Engineering

Text Cleaning

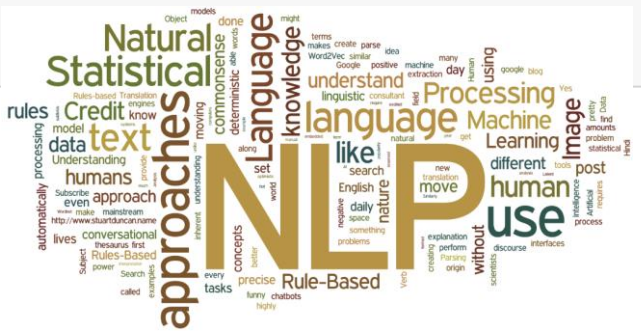
```
1 #nltk.download("wordnet", "C:\Users\Mike\nltk_data/")
2 df2['clean_text'] = df2['new_text'].map(lambda text: normalize_corpus(text))
```

After applying the function, we have a new column named as clean_text

The Length of the Clean_Text

```
1 df2['review_length'] = df2['clean_text'].map(len)
```

```
1 df2.head(1)
```



ng_id	age	review_text	rating	recommended_ind	positive_feedback_count	division_name	department_name	class_name	new_text	clean_text	review_length
767	33	Absolutely wonderful - silky and sexy and comf...	4	1	0	Initmates	Intimate	Intimates	Absolutely wonderful - silky and sexy and comf...	absolutely wonderful silky sexy comfortable	45

We have a new column named as review_length which shows the length of the reviews by counting the letters

Feature Engineering

Classifying the Ratings as Good, Neutral and Bad

```
1 bad_rat = len(df2[df2.rating < 3])
2 neut_rat = len(df2[df2.rating == 3])
3 good_rat = len(df2[df2.rating > 3])
4
5 print ('Bad ratings : {}'.format(bad_rat))
6 print ('Neutral ratings : {}'.format(neut_rat))
7 print ('Good ratings : {}'.format(good_rat))
```

```
Bad ratings : 2370
Neutral ratings : 2823
Good ratings : 17449
```

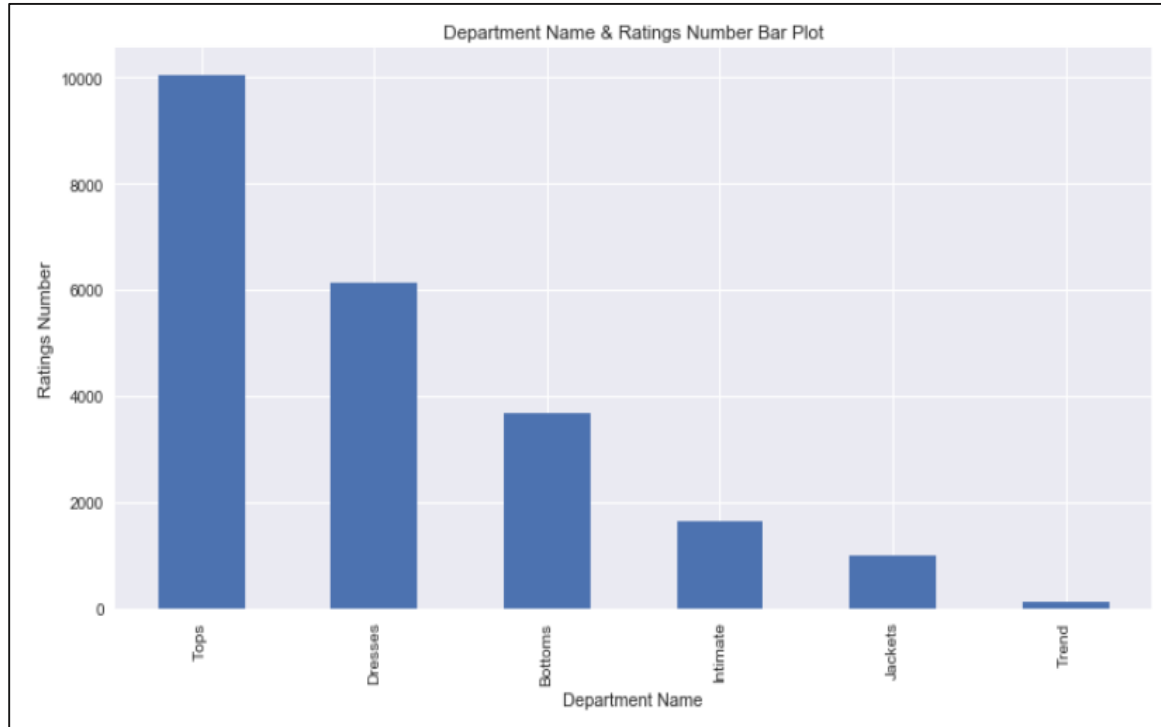
We classified the ratings as **Good** which is bigger than 3, **Neutral** which equals 3 and **Bad** which is less than 3.

Applying the New Classification to the Ratings Column

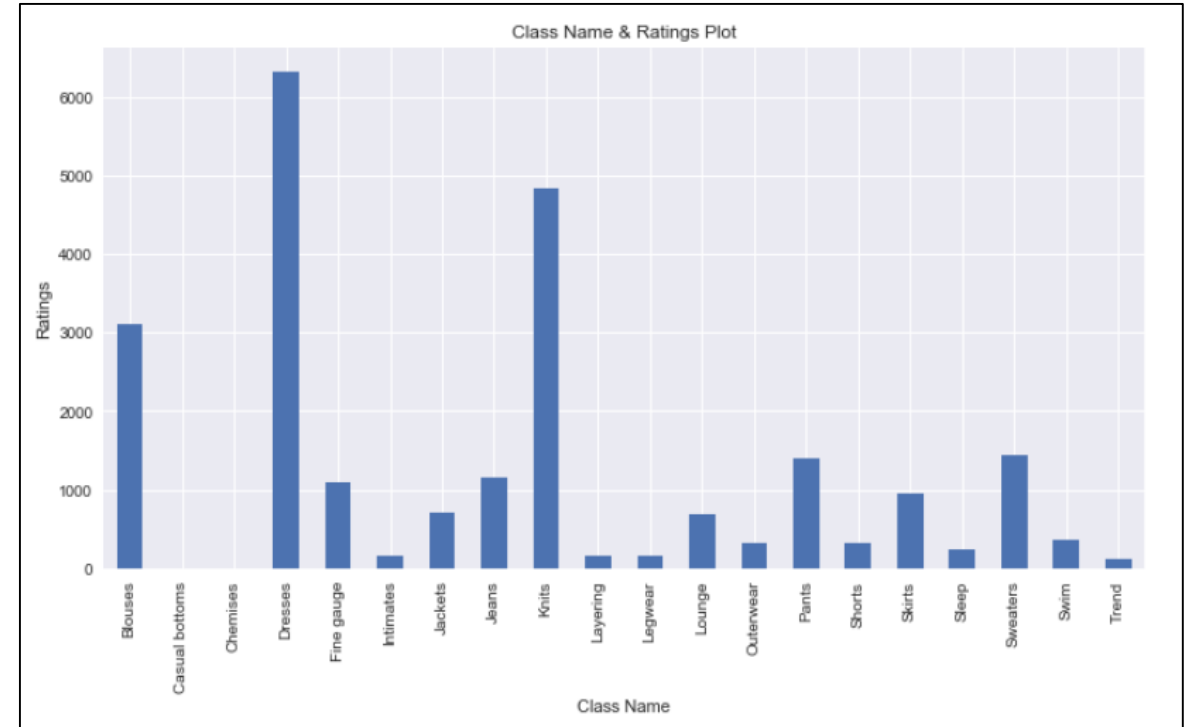
```
1 df2['rating_class'] = df2['rating'].apply(lambda x: 'bad' if x < 3 else ('good' if x > 3 else 'neutral'))
```

After applying the new classification we have a new column named as `rating_class` consists of three classes, '**Good**, **Neutral** and **Bad**'

Data Exploration Analysis

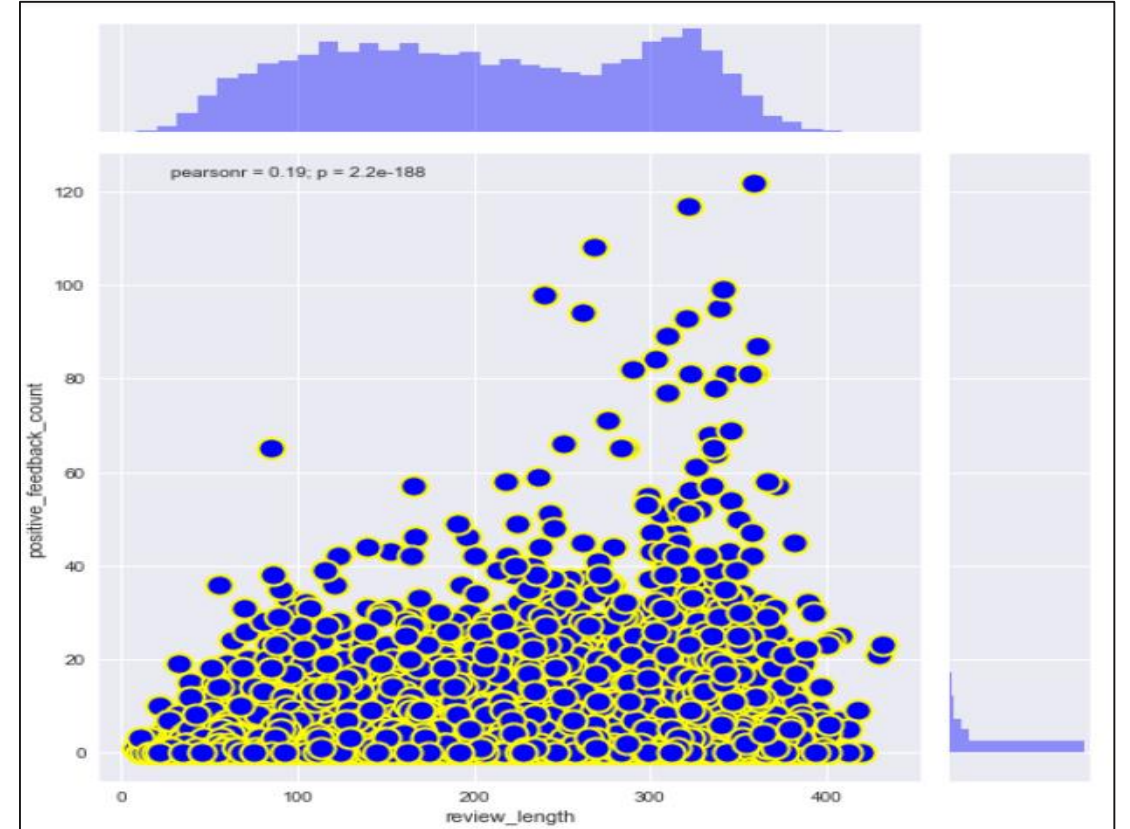
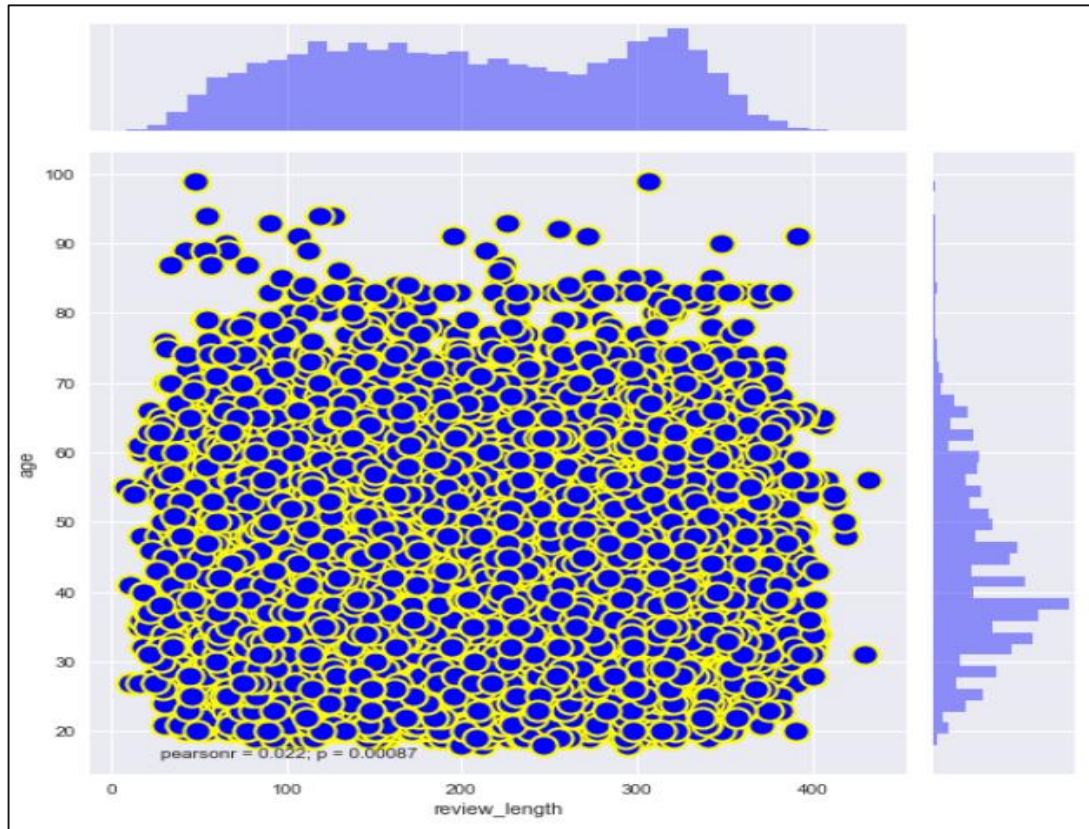


On Department basis, Tops and Dresses are sold mostly. Trend is the weakest sold as seen.



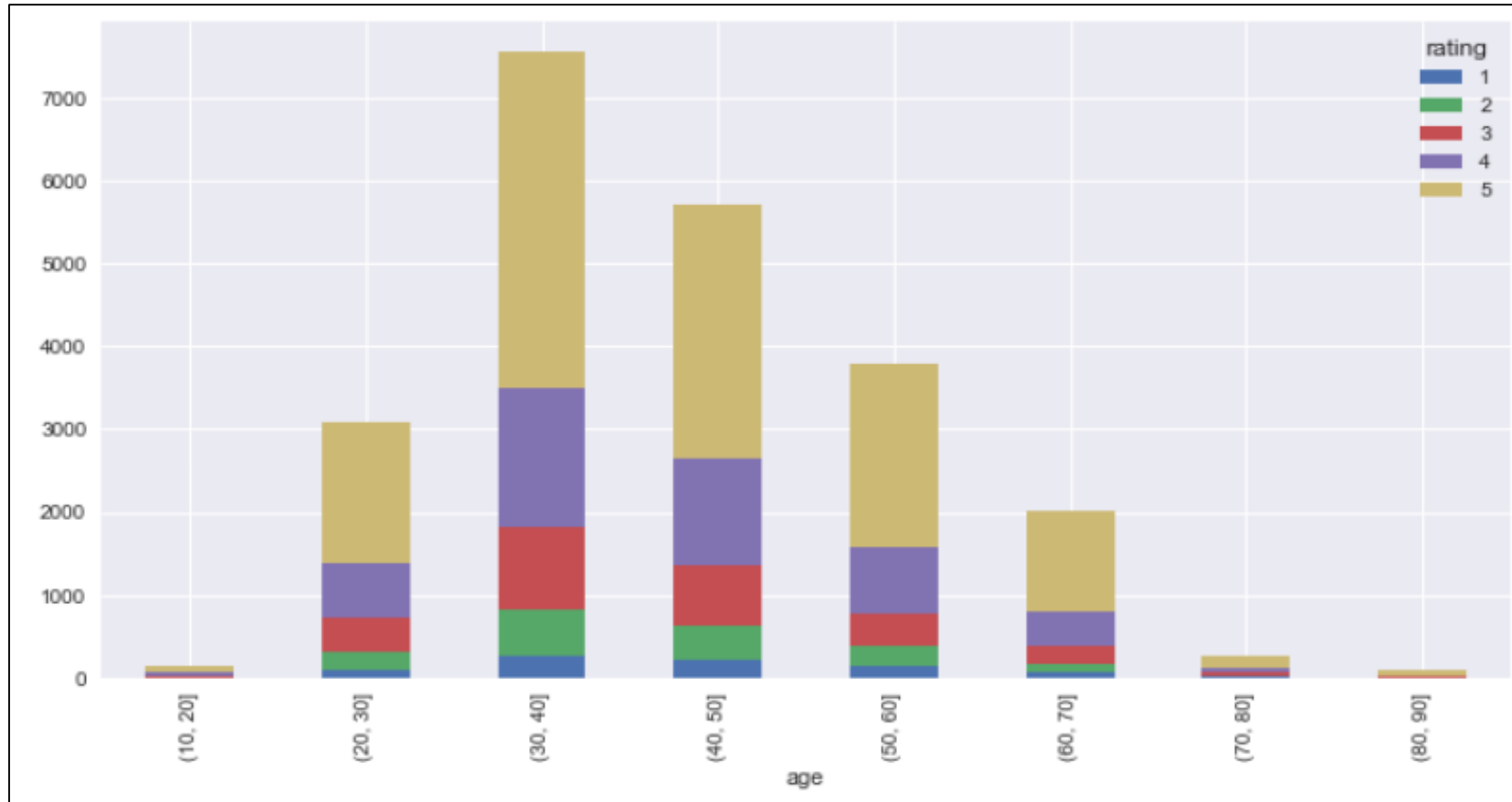
Most ratings were given to Dresses, Knits and Blouses.

Data Exploration Analysis



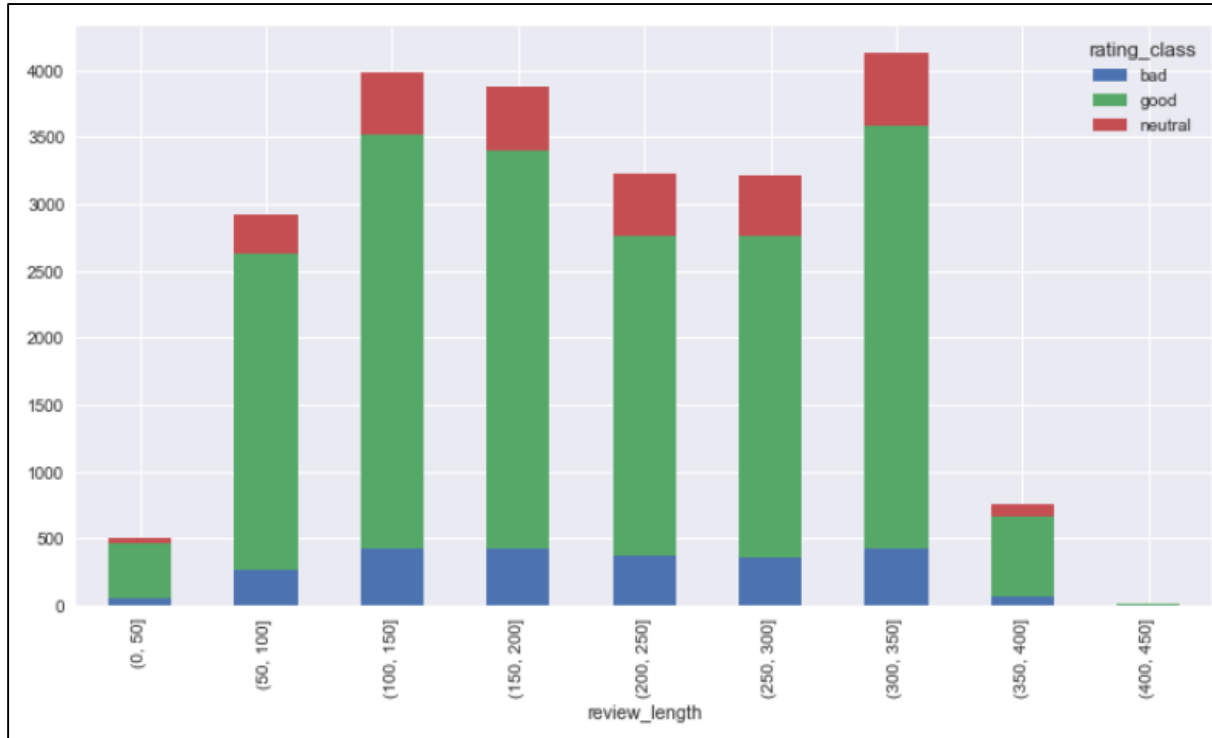
Mostly the customers between 20 and 65 ages left reviews and the review length is between 50 to 350 characters.

Data Exploration Analysis



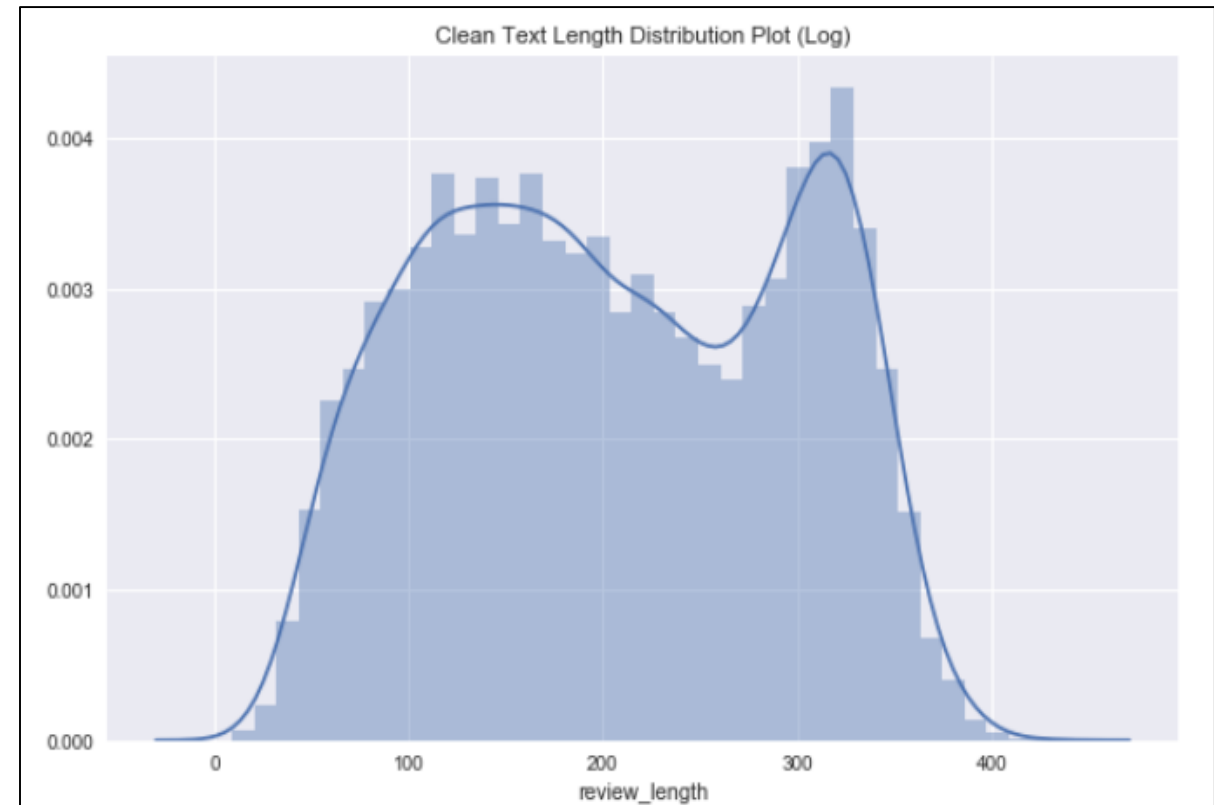
The most satisfied age group is between 30 and 50. Customers mostly left positive reviews.

Data Exploration Analysis



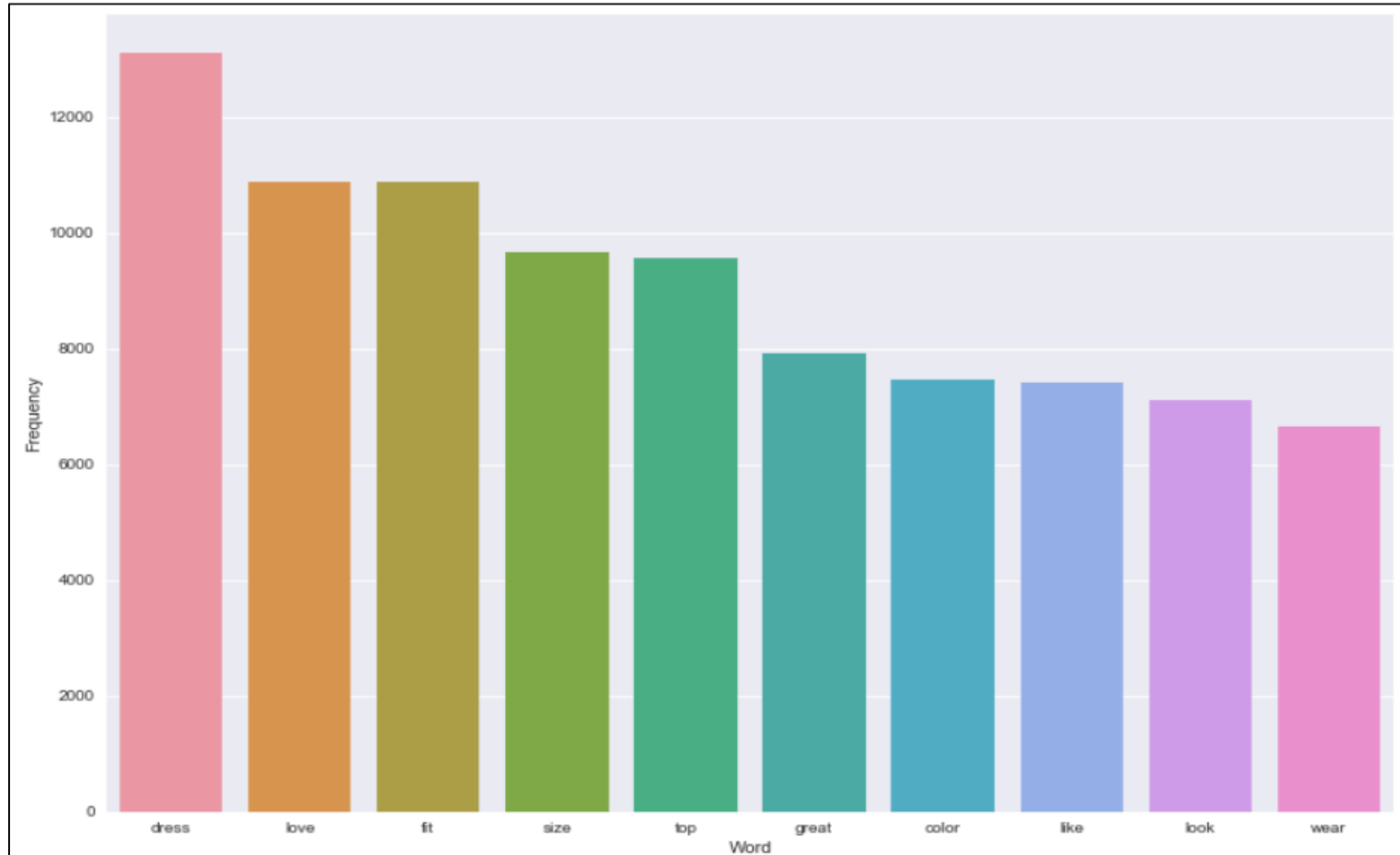
Class Based Review Length

Cleaned Text Length Distribution Plot(Log)



Data Exploration Analysis

Most Common Words

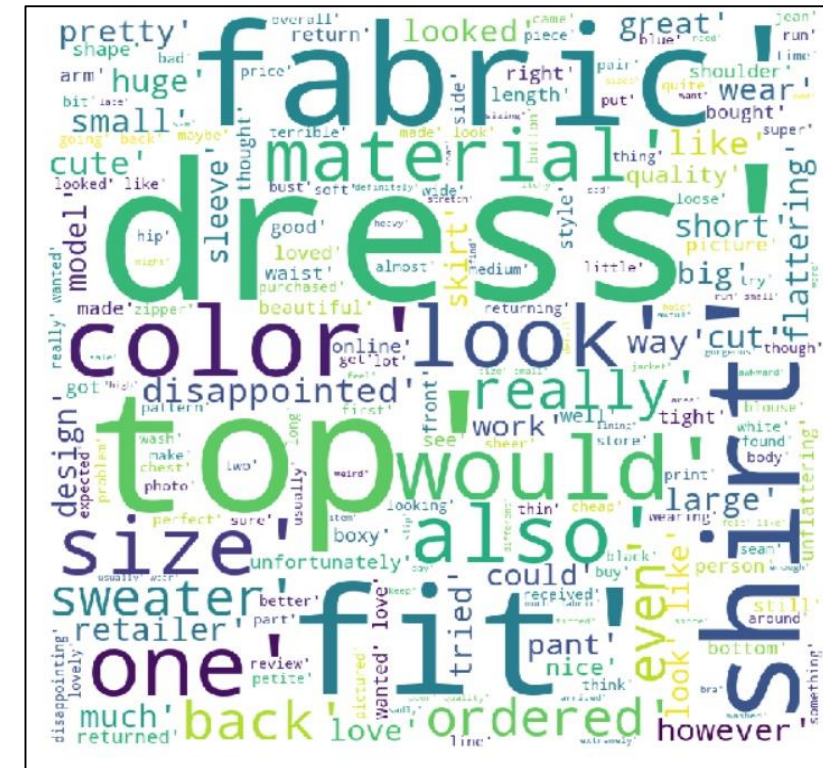


Word Clouds



Good

Neutral



Bad

Machine Learning

3 Rating Classes, Count-Vectorizing and the Algorithms

Proportion (Train/Test)	CLASS AMOUNT	BOW	Model	Accuracy Score			Classification Report				
							Category	precision	recall	f1-score	Support
0.75/0.25	3 (BAD, GOOD, NEUTRAL)	Count Vectorizer	Default Logistic Regression	Accuracy Score	Train	0.8247	Bad	1.00	0.26	0.41	1770
					Test	0.7555	Good	0.81	1.00	0.90	13059
			Random Forest Classifier N-Estimator = 200	Accuracy Score	Test	0.7622	Neutral	1.00	0.22	0.36	2142
							Bad	0.20	0.04	0.06	600
							Good	0.78	0.97	0.86	4376
			Linear SVM	Accuracy Score	Test	0.7613	Neutral	0.15	0.03	0.05	681
							Bad	0.20	0.02	0.04	600
							Good	0.78	0.98	0.87	4376
			Gradient Boosting	Accuracy Score	Test	0.7735	Neutral	0.16	0.02	0.03	681
							Bad	0.21	0.02	0.04	600
							Good	0.78	0.98	0.87	4376
			Xg Boosting	Accuracy Score	Test	0.7735	Neutral	0.18	0.02	0.04	681
							Bad	0.00	0.00	0.00	600
							Good	0.77	1.00	0.87	4376
			Naïve Bayes	Accuracy Score	Test	0.7723	Neutral	0.00	0.00	0.00	681
							Bad	0.00	0.00	0.00	600
							Good	0.77	1.00	0.87	4376
							Neutral	0.11	0.00	0.00	681

Machine Learning

3 Rating Classes, TF-IDF and the Same Algorithms

							Classification Report				
Proportion (Train/Test)	CLASS AMOUNT	BOW	Model	Accuracy Score			Category	precision	recall	f1-score	Support
0.75/0.25	3 (BAD, GOOD, NEUTRAL)	TF-IDF	Default Logistic Regression	Accuracy Score	Train	0.8247	Bad	1.00	0.26	0.41	1770
							Good	0.81	1.00	0.90	13059
							Neutral	1.00	0.22	0.36	2142
					Test	0.76	Bad	0.19	0.02	0.03	600
							Good	0.78	0.98	0.87	4376
							Neutral	0.14	0.01	0.68	681
			Random Forest Classifier N-Estimator = 200	Accuracy Score	Test	0.7618	Bad	0.20	0.02	0.04	600
							Good	0.78	0.98	0.87	4376
							Neutral	0.16	0.02	0.03	681
			Linear SVM	Accuracy Score	Test	0.7624	Bad	0.19	0.02	0.03	600
							Good	0.78	0.98	0.87	4376
							Neutral	0.17	0.02	0.04	681
			Gradient Boosting	Accuracy Score	Test	0.7735	Bad	0.00	0.00	0.00	600
							Good	0.77	1.00	0.87	4376
							Neutral	0.00	0.00	0.00	681
			Xg Boosting	Accuracy Score	Test	0.7735	Bad	0.00	0.00	0.00	600
							Good	0.77	1.00	0.87	4376
							Neutral	0.00	0.00	0.00	681
			Naïve Bayes	Accuracy Score	Test	0.7735	Bad	0.00	0.00	0.00	600
							Good	0.77	1.00	0.87	4376
							Neutral	0.00	0.00	0.00	681

Machine Learning

Expanded Stop Words, 3 Rating Classes, Count Vectorizer and the Same Algorithms

Proportion (Train/Test)	CLASS AMOUNT	BOW	Model	Accuracy Score			Classification Report				
							Category	precision	recall	f1-score	Support
0.75/0.25	3 (BAD, GOOD, NEUTRAL)	Count Vectorizer	Default Logistic Regression	Accuracy Score	Train	0.8429	Bad	0.63	0.81	0.71	1770
							Good	0.97	0.86	0.92	13059
							Neutral	0.51	0.74	0.61	2142
					Test	0.7689	Bad	0.44	0.57	0.50	600
							Good	0.95	0.84	0.89	4376
							Neutral	0.32	0.46	0.38	681
			Random Forest Classifier N-Estimator = 200	Accuracy Score	Test	0.7898	Bad	0.49	0.36	0.42	600
							Good	0.84	0.96	0.89	4376
							Neutral	0.31	0.10	0.15	681
			Linear SVM	Accuracy Score	Test	0.8041	Bad	0.51	0.45	0.48	600
							Good	0.88	0.94	0.91	4376
							Neutral	0.36	0.24	0.28	681
			Gradient Boosting	Accuracy Score	Test	0.7977	Bad	0.62	0.22	0.33	600
							Good	0.82	0.99	0.89	4376
							Neutral	0.43	0.09	0.15	681
			Xg Boosting	Accuracy Score	Test	0.7931	Bad	0.62	0.18	0.28	600
							Good	0.81	0.99	0.89	4376
							Neutral	0.42	0.06	0.10	681
			Naïve Bayes	Accuracy Score	Test	0.8037	Bad	0.55	0.52	0.53	600
							Good	0.92	0.90	0.91	4376
							Neutral	0.36	0.44	0.40	681

Machine Learning

Expanded Stop Words, 3 Rating Classes, Count Vectorizer and the Same Algorithms and SMOTE

Proportion (Train/Test)	CLASS AMOUNT	BOW	Model	Accuracy Score		Classification Report				
						Category	precision	recall	f1-score	Support
0.75/0.25	3 (BAD, GOOD, NEUTRAL)	(SMOTE) Count Vectorizer	Default Logistic Regression	Train	0.8533	Bad	0.86	0.87	0.87	13059
						Good	0.91	0.85	0.88	13059
				Test	0.735	Neutral	0.79	0.84	0.82	2142
						Bad	0.42	0.53	0.47	600
						Good	0.94	0.81	0.87	4376
						Neutral	0.27	0.45	0.34	681
			Random Forest Classifier N-Estimator = 200	Test	0.7983	Bad	0.51	0.31	0.39	600
						Good	0.84	0.96	0.89	4376
						Neutral	0.32	0.12	0.17	681
			Linear SVM	Test	0.7254	Bad	0.38	0.48	0.42	600
						Good	0.92	0.81	0.86	4376
						Neutral	0.26	0.39	0.31	681
			Gradient Boosting	Test	0.7542	Bad	0.51	0.34	0.41	600
						Good	0.87	0.88	0.87	4376
						Neutral	0.26	0.31	0.28	681
			Xg Boosting	Test	0.7429	Bad	0.49	0.32	0.38	600
						Good	0.87	0.87	0.87	4376
						Neutral	0.24	0.32	0.28	681
			K Neighbors Classifier	Test	0.4581	Bad	0.20	0.40	0.26	600
						Good	0.89	0.46	0.61	4376
						Neutral	0.15	0.47	0.23	681
			Naïve Bayes	Test	0.7546	Bad	0.46	0.55	0.50	600
						Good	0.95	0.82	0.88	4376
						Neutral	0.31	0.53	0.39	681

Machine Learning

Expanded Stop Words, 3 Rating Classes, Count Vectorizer, Logistic Regression & Random Forest, SMOTE and Linear Dimensionality Reduction (PCA and Truncated SVD)

								Classification Report				
Proportion (Train/Test)	CLASS AMOUNT	BOW		Model	Accuracy Score			Category	precision	recall	f1-score	Support
0.75/0.25	3 (BAD, GOOD, NEUTRAL)	(SMOTE) Count Vectorizer	PCA Linear Dimentionality Reduction	Default Logistic Regression	Accuracy Score	Test	0.482	Bad	0.13	0.18	0.15	600
								Good	0.76	0.57	0.65	4376
								Neutral	0.09	0.19	0.12	681
			Truncated SVD Linear Dimentionality Reduction	Random Forest Classifier N-Estimator = 200	Accuracy Score	Test	0.6955	Bad	0.08	0.07	0.08	600
								Good	0.77	0.89	0.82	4376
								Neutral	0.00	0.00	0.00	681
				Default Logistic Regression	Accuracy Score	Test	0.7931	Bad	0.12	0.23	0.16	600
								Good	0.81	0.61	0.70	4376
								Neutral	0.16	0.28	0.21	681
				Random Forest Classifier N-Estimator = 200	Accuracy Score	Test	0.7275	Bad	0.06	0.04	0.05	600
								Good	0.77	0.94	0.84	4376
								Neutral	0.00	0.00	0.00	681

Machine Learning

Expanded Stop Words, 2 Rating Classes, Count Vectorizer and All Algorithms

Proportion (Train/Test)	CLASS AMOUNT	BOW	Model	Accuracy Score			Classification Report				
							Category	precision	recall	f1-score	Support
0.75/0.25	2 (BAD, NOT BAD)	Count Vectorizer	Default Logistic Regression	Accuracy Score	Train	0.91	Bad	0.56	0.97	0.71	1178
							Not Bad	1.00	0.91	0.95	13059
					Test	0.85	Bad	0.40	0.68	0.50	592
							Not Bad	0.96	0.88	0.92	5062
			Random Forest Classifier N-Estimator = 200	Accuracy Score	Test	0.89	Bad	0.49	0.23	0.31	592
							Not Bad	0.92	0.97	0.94	5065
			Linear SVM	Accuracy Score	Test	0.8985	Bad	0.52	0.46	0.49	592
							Not Bad	0.94	0.95	0.94	5065
			Gradient Boosting	Accuracy Score	Test	0.9043	Bad	0.69	0.15	0.25	592
							Not Bad	0.91	0.99	0.95	5065
			Xg Boosting	Accuracy Score	Test	0.9022	Bad	0.71	0.11	0.19	592
							Not Bad	0.91	0.99	0.95	5065
			Naïve Bayes	Accuracy Score	Test	0.8914	Bad	0.49	0.68	0.57	592
							Not Bad	0.96	0.92	0.94	5065

Machine Learning

Expanded Stop Words, 2 Rating Classes, Count Vectorizer and All Algorithms and SMOTE

							Classification Report				
Proportion (Train/Test)	CLASS AMOUNT	BOW	Model	Accuracy Score			Category	precision	recall	f1-score	Support
0.75/0.25	2 (BAD, NOT BAD)	(SMOTE) Count Vectorizer	Default Logistic Regression	Accuracy Score	Train	0.93	Bad	0.91	0.97	0.94	15193
							Not Bad	0.97	0.90	0.94	15193
				Accuracy Score	Test	0.84	Bad	0.37	0.67	0.48	592
							Not Bad	0.96	0.87	0.91	5065
			Random Forest Classifier N-Estimator = 200	Accuracy Score	Test	0.89	Bad	0.47	0.23	0.30	592
							Not Bad	0.91	0.97	0.94	5065
			Linear SVM	Accuracy Score	Test	0.8439	Bad	0.36	0.62	0.46	592
							Not Bad	0.95	0.87	0.91	5065
			Gradient Boosting	Accuracy Score	Test	0.8955	Bad	0.50	0.32	0.39	592
							Not Bad	0.92	0.96	0.94	5065
			Xg Boosting	Accuracy Score	Test	0.8964	Bad	0.51	0.34	0.40	592
							Not Bad	0.93	0.96	0.94	5065
			Naïve Bayes	Accuracy Score	Test	0.8449	Bad	0.39	0.81	0.52	592
							Not Bad	0.97	0.85	0.91	5065

Machine Learning

Expanded Stop Words, 2 Rating Classes, Count Vectorizer and All Algorithms and N_Grams(1,2)

Proportion (Train/Test)	CLASS AMOUNT	BOW	Model	Accuracy Score			Classification Report				
							Category	precision	recall	f1-score	Support
0.75/0.25	2 (BAD, NOT BAD)	Ngrams (1,2) Count Vectorizer	Default Logistic Regression	Accuracy Score	Train	0.9321	Bad	0.61	0.98	0.75	1178
							Not Bad	1.00	0.93	0.96	15193
					Test	0.867	Bad	0.41	0.65	0.50	592
							Not Bad	0.96	0.89	0.92	5065
			Random Forest Classifier N-Estimator = 200	Accuracy Score	Test	0.8972	Bad	0.52	0.26	0.34	592
							Not Bad	0.92	0.97	0.94	5065
			Linear SVM	Accuracy Score	Test	0.8948	Bad	0.50	0.47	0.48	592
							Not Bad	0.94	0.94	0.94	5065
			Gradient Boosting	Accuracy Score	Test	0.9041	Bad	0.69	0.15	0.25	592
							Not Bad	0.91	0.99	0.95	5065
			Xg Boosting	Accuracy Score	Test	0.9027	Bad	0.74	0.11	0.19	592
							Not Bad	0.91	1.00	0.95	5065
			Naïve Bayes	Accuracy Score	Test	0.8895	Bad	0.48	0.70	0.57	592
							Not Bad	0.96	0.91	0.94	5065

Machine Learning

Expanded Stop Words, 2 Rating Classes, Count Vectorizer and All Algorithms, N_Grams(1,2) and SMOTE

						Classification Report						
Proportion (Train/Test)	CLASS AMOUNT	BOW	Model	Accuracy Score		Category	precision	recall	f1-score	Support		
0.75/0.25	2 (BAD, NOT BAD)	(SMOTE)	Default Logistic Regression		Train	0.9539	Bad	0.93	0.98	0.96	15193	
							Not Bad	0.98	0.92	0.95	15193	
				Accuracy Score	Test	0.8589	Bad	0.39	0.64	0.49	592	
							Not Bad	0.95	0.89	0.92	5065	
		Ngrams (1,2)	Random Forest Classifier N-Estimator = 200	Accuracy Score	Test	0.8904	Bad	0.45	0.22	0.29	592	
							Not Bad	0.91	0.97	0.94	5065	
				Linear SVM	Accuracy Score	Test	0.8523	Bad	0.37	0.58	0.45	592
								Not Bad	0.95	0.88	0.91	5065
		Count Vectorizer	Gradient Boosting	Accuracy Score	Test	0.8953	Bad	0.50	0.33	0.39	592	
							Not Bad	0.92	0.96	0.94	5065	
				Xg Boosting	Accuracy Score	Test	0.8951	Bad	0.50	0.33	0.40	592
								Not Bad	0.92	0.96	0.94	5065
			Naïve Bayes	Accuracy Score	Test	0.8497	Bad	0.40	0.82	0.53	592	
							Not Bad	0.98	0.85	0.91	5065	

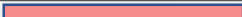
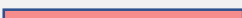
Machine Learning

Expanded Stop Words, 2 Rating Classes, Count Vectorizer, All Algorithms and N_Grams(1,3)

							Classification Report				
Proportion (Train/Test)	CLASS AMOUNT	BOW	Model	Accuracy Score			Category	precision	recall	f1-score	Support
0.75/0.25	2 (BAD, NOT BAD)	Ngrams (1,3) Count Vectorizer	Default Logistic Regression	Accuracy Score	Train	0.9013	Bad	0.52	0.96	0.67	1178
					Not Bad	0.99	0.89	0.94	15193		
					Test	0.8564	Bad	0.40	0.72	0.51	592
					Not Bad	0.96	0.87	0.92	5065		
			Random Forest Classifier N-Estimator = 200	Accuracy Score	Test	0.8893	Bad	0.44	0.23	0.30	592
							Not Bad	0.91	0.97	0.94	5065
			Linear SVM	Accuracy Score	Test	0.904	Bad	0.55	0.44	0.49	592
							Not Bad	0.94	0.96	0.95	5065
			Gradient Boosting	Accuracy Score	Test	0.9043	Bad	0.69	0.16	0.26	592
							Not Bad	0.91	0.99	0.95	5065
			Xg Boosting	Accuracy Score	Test	0.9025	Bad	0.71	0.12	0.20	592
							Not Bad	0.91	0.99	0.95	5065
			Naïve Bayes	Accuracy Score	Test	0.8845	Bad	0.46	0.68	0.55	592
							Not Bad	0.96	0.91	0.93	5065

Machine Learning

Expanded Stop Words, 2 Rating Classes, Count Vectorizer, All Algorithms, N_Grams(1,3) and SMOTE

							Classification Report				
Proportion (Train/Test)	CLASS AMOUNT	BOW	Model	Accuracy Score			Category	precision	recall	f1-score	Support
0.75/0.25	2 (BAD, NOT BAD)	(SMOTE) Ngrams (1,3) Count Vectorizer	Default Logistic Regression		Train	0.9266	Bad	0.90	0.96	0.93	15193
							Not Bad	0.96	0.89	0.92	15193
				Accuracy Score	Test	0.8465	Bad	0.37	0.69	0.49	592
							Not Bad	0.96	0.86	0.91	5065
			Random Forest Classifier N-Estimator = 200	Accuracy Score	Test	0.8919	Bad	0.47	0.24	0.31	592
							Not Bad	0.92	0.97	0.94	5065
			Linear SVM	Accuracy Score	Test	0.8424	Bad	0.37	0.69	0.48	592
							Not Bad	0.96	0.86	0.91	5065
			Gradient Boosting	Accuracy Score	Test	0.898	Bad	0.52	0.36	0.43	592
							Not Bad	0.93	0.96	0.94	5065
			Xg Boosting	Accuracy Score	Test	0.899	Bad	0.52	0.39	0.44	592
							Not Bad	0.93	0.96	0.94	5065
			Naïve Bayes	Accuracy Score	Test	0.8387	Bad	0.38	0.83	0.52	592
							Not Bad	0.98	0.84	0.90	5065

Machine Learning

Expanded Stop Words, 2 Rating Classes, Count Vectorizer, Genetic Algorithms

							Classification Report				
Proportion (Train/Test)	CLASS AMOUNT	BOW	Model	Accuracy Score			Category	precision	recall	f1-score	Support
0.75/0.25	2 (BAD, NOT BAD)	Ngrams (1,2) Count Vectorizer	CatBoosting	Accuracy Score	Test	0.9052	Bad	0.59	0.30	0.40	592
							Not Bad	0.92	0.98	0.95	5065
			TPOT (Linear SVC with CV)	Accuracy Score	Test	0.9116	Bad	0.69	0.28	0.40	592
							Not Bad	0.92	0.91	0.89	5065

Conclusions

In this study, we tried to predict the ratings based on the reviews left by the female E-customers. We applied Count Vectorizing, TF-IDF, Classification Algorithms, Synthetic Minority Oversampling Technique (SMOTE) and Linear Dimensionality Reduction as well as Genetic Algorithms.

Since in the study the most important thing to predict the classes correctly, precision scores were more important than the recalls. Therefore we took in to account accuracy and precision scores. Thus the boosting algorithms were the winners almost each combinations.

SMOTE and Linear Dimensionality Reduction techniques decreased the accuracy, precision and recall scores drastically.

Recommendations to the Client

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