

Sentiment Analysis of Yelp Reviews: A Comparison of Techniques and Models

Siqi Liu
University of Waterloo
sq2liu@uwaterloo.ca

Abstract—We use over 350,000 Yelp reviews on 5,000 restaurants to perform an ablation study on text preprocessing techniques. We also compare the effectiveness of several machine learning and deep learning models on predicting user sentiment (negative, neutral, or positive). For machine learning models, we find that using binary bag-of-word representation, adding bi-grams, imposing minimum frequency constraints and normalizing texts have positive effects on model performance. For deep learning models, we find that using pre-trained word embeddings and capping maximum length often boost model performance. Finally, using macro F1 score as our comparison metric, we find simpler models such as Logistic Regression and Support Vector Machine to be more effective at predicting sentiments than more complex models such as Gradient Boosting, LSTM and BERT.

Index Terms—Sentiment Analysis, Text Preprocessing

I. INTRODUCTION

A. Background

Sentiment analysis is the process of analyzing a piece of text and determining its author's attitude. Using modern natural language processing (NLP) techniques, one can predict the most likely emotional tone that the author carries based on what's presented in the text. Sentiment analysis is widely used by companies around the world to evaluate customer satisfaction, conduct market research and monitor brand reputation.

Yelp is a popular crowd-sourced review platform with millions of active users who rate and review hundreds of thousands of businesses across the globe. Since 2016, Yelp has been releasing and updating subsets of its enormous database to the general public.

B. Goal

While there is no one-fits-all approach to performing sentiment analysis, the goal of this project is to provide some guidance for practitioners to consider when developing models that fit their needs.

There are two main steps to building a model: data preparation and model selection. First, we perform an ablation study on several text preprocessing techniques (e.g., stop word removal, normalization) using a simple multinomial Naive Bayes model [15]. Afterwards, using the preprocessed data, we train several machine learning (e.g., Logistic Regression, Support Vector Machine) and deep learning (e.g., LSTM, BERT) models and compare their effectiveness at predicting sentiments.

II. DATA

A. Preparation

We download our data from the Yelp Open Data website (<https://www.yelp.com/dataset>). For the scope of this project, we only use the information on reviews, ratings and businesses, and forgo other information such as tips, pictures and user data that could also help predict user sentiment.

The raw dataset contains 6.7 million reviews and ratings on 192,609 businesses. Due to limited computational power, we narrow the scope of our study to businesses that are:

- In the restaurant/food servicing industry
- Located in the Greater Toronto Area
- With at least 10 reviews

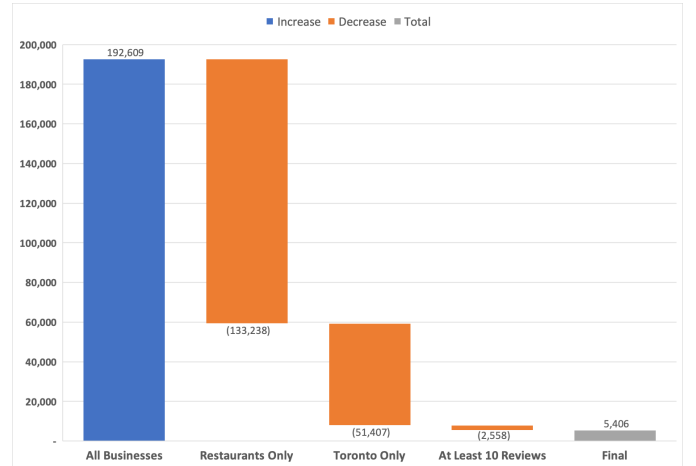


Fig. 1. Study Population Selection Waterfall

Fig. 1 shows the impacts of each of our narrowing criteria. As a result, our final study population consists of 5,406 businesses. We then select all 362,554 reviews and ratings on these businesses for our dataset.

There are two columns in our dataset - *review* contains the text reviews, and *rating* contains the star ratings that accompany the reviews. The ratings range from 1 to 5, with 1 being the lowest and 5 being the highest. This is a multi-class classification problem, and our goal is to build classification models that can accurately predict the ratings based on the reviews.

B. Distribution

Before jumping into model building, it is important to understand the distribution of our target variable (i.e., ratings), as having an imbalanced dataset could lead to sub-optimal model performance.

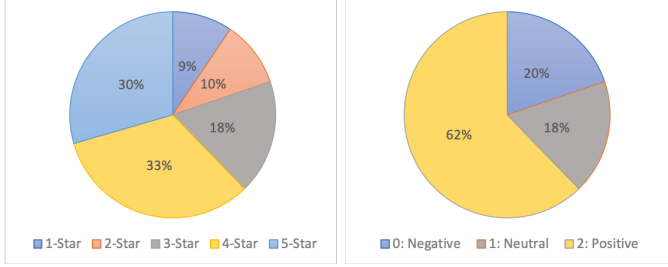


Fig. 2. Target Class Ratio - Before (left) and After (right) Grouping

Fig. 2 (left) shows the proportion of ratings in our dataset. It turns out that only 9% of our dataset has 1-star ratings, compared to 4-star ratings and 5-star ratings, which occupy a third of our dataset each. From an analytical perspective, businesses like to know more about the difference between a 1-star rating and a 5-star rating, rather than the difference between a 1-star rating and a 2-star rating. Therefore, we group ratings 1 and 2 as “negative” sentiments, group ratings 4 and 5 as “positive” sentiments, and consider rating 3 as “neutral”. Fig. 2 (right) shows the target class ratio after grouping. For modelling purposes, we use 0 for negative sentiments, 1 for neutral, and 2 for positive sentiments. Later on, we will examine the effect on model performance when the dataset is completely balanced.

C. Train Test Split

Since we will be comparing model performances, we need to establish a universal test set. Using *scikit-learn*, we split the dataset into 75% training and 25% testing, while preserving the target class ratio. In the end, we have 271,915 data points for training and 90,639 data points for testing.

III. EVALUATION METRIC

We will use macro F1 score defined in Sokolova, 2009 [1] as the evaluation metric for measuring and comparing model performances:

$$\text{Precision}_M = \frac{\sum_{i=1}^l \frac{tp_i}{tp_i + fp_i}}{l} \quad (1)$$

$$\text{Recall}_M = \frac{\sum_{i=1}^l \frac{tp_i}{tp_i + fn_i}}{l} \quad (2)$$

$$\text{F1}_M = \frac{2 \times \text{Precision}_M \times \text{Recall}_M}{\text{Precision}_M + \text{Recall}_M} \quad (3)$$

where l is the total number of classes and tp , fp , fn are the number of True Positives, False Positives and False Negatives respectively.

TABLE I
EXAMPLE CONFUSION MATRIX

		Predicted		
		Class 0	Class 1	Class 2
Actual	Class 0	1	0	1
	Class 1	0	1	1
	Class 2	1	0	3

Consider the example in Table I. In this example, macro precision, macro recall and macro F1 scores are computed as follows:

$$\text{Precision}_M = \frac{\frac{1}{1+1} + \frac{1}{1+0} + \frac{3}{3+2}}{3} = \frac{0.5 + 1.0 + 0.6}{3} = 0.7 \quad (4)$$

$$\text{Recall}_M = \frac{\frac{1}{1+1} + \frac{1}{1+1} + \frac{3}{3+1}}{3} = \frac{0.5 + 0.5 + 0.75}{3} = 0.5833 \quad (5)$$

$$\text{F1}_M = \frac{2 \times 0.7 \times 0.5833}{0.7 + 0.5833} = 0.63 \quad (6)$$

Note that this macro F1 score is slightly different from the macro F1 score implemented in *scikit-learn*, which is a simple average of within-class F1 scores.

IV. ABLATION STUDY RESULTS

A. Training Set Size

We want to examine the marginal improvements in the model performance of having a larger training set. Many studies (e.g., Renault, 2019 [2]) have shown the importance of having a large dataset on sentiment analysis tasks. Starting with a small training set size of 1,000, we train eight multinomial Naive Bayes models using default parameters and record the macro F1 scores.

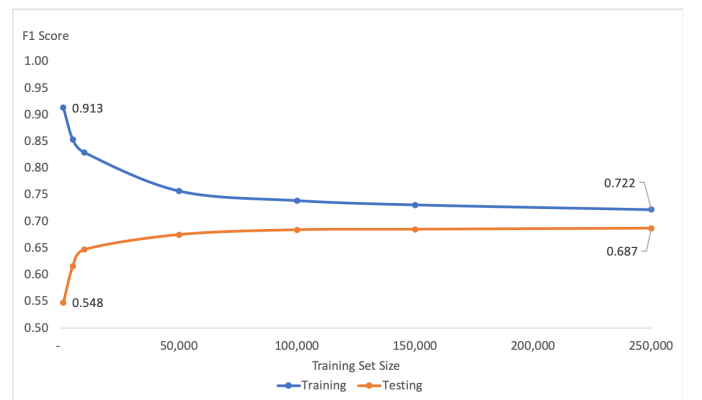


Fig. 3. F1 Score by Training Set Size

Fig. 3 shows the model performance on the training set and the test set. Fig. 4 shows the normalized confusion matrices on then test set when training set sizes are 1,000 (left) and 250,000 (right). As we can see, although the marginal benefit

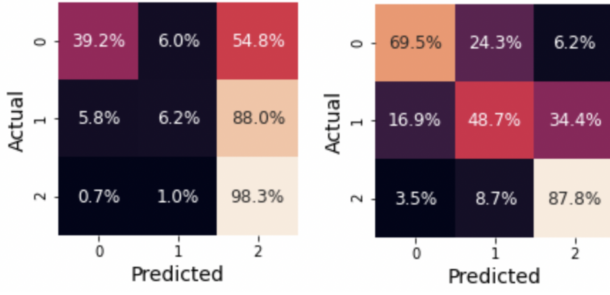


Fig. 4. Size = 1,000 (left) vs. Size = 25,000 (right)

of having a large dataset diminishes over time, it prevents overfitting effectively and improves overall model performance significantly.

B. Balanced Target Class Ratio

Currently, our target class ratio is at 20/20/60. While this is not considered severely imbalanced, we want to examine if having a completely balanced dataset could improve model performance.

Since our smallest class has a training set size of 49,000, we set up our balanced training set by down-sampling all three classes to 45,000 data points each. For our benchmark case, we down-sample the entire training set to $45,000 * 3 = 135,000$ data points, preserving the 20/20/60 imbalanced target class ratio.

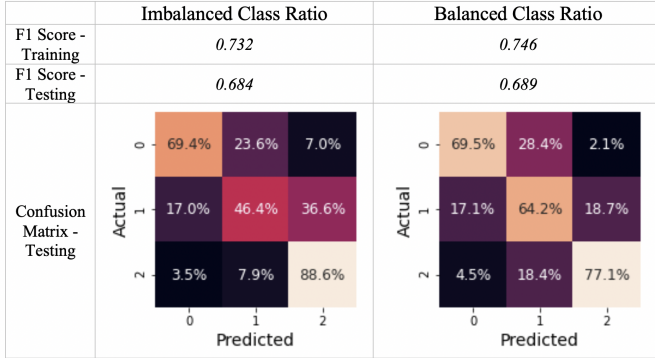


Fig. 5. Imbalanced (left) vs. Balanced (right)

Fig. 5 shows the results when models are trained with imbalanced (left) and balanced (right) training sets. As we can see, having a balanced training set not only improves the model performance, it also prevents the model from always predicting the majority class (e.g., actual = 1, predicted = 2 goes down from 36.6% to 18.7%). Although the True Positive rate for class 2 decreases from 88.6% to 77.1%, this is a reasonable trade-off for lower False Negative rates in the other two classes.

Moreover, further experiments show that the model trained on the down-sampled balanced dataset outperforms the model trained on the entire training set without any down-sampling, despite the difference in sample size. Therefore, we update our

training set for all future ablation studies and model builds to be the down-sampled balanced training set. We do not, however, modify the test set, since the class imbalance is an actual representation of the distribution in our target population.

C. Word Representation

Now that we have our final dataset, we can start building our text preprocessing pipeline. First, for each text corpus, we need to extract features for our models to use. Two of the most common representation methods are bag-of-words and term frequency-inverse document frequency (TF-IDF).

Under the bag-of-words representation, each text corpus is converted into a vector of word counts, so the feature space becomes all unique words among all training texts. One variation of this model is to use binary indicators rather than word counts. The motivation behind this variation is that sometimes, word occurrences are more important than word frequency.

The TF-IDF representation also converts each text corpus into a vector. However, instead of word counts, it is

$$\text{tf-idf}(t, d, D) = \text{tf}(t, d) \times \text{idf}(t, D) \quad (7)$$

where $\text{tf}(t, d)$ is the term frequency of the term t in document d , and $\text{idf}(t, D)$ is the inverse (log) document frequency of t in the entire training set of documents D . $\text{tf-idf}(t, d, D)$ is high if the term has high term frequency and low document frequency. Essentially, this representation emphasizes terms that are rare and discounts terms that are less valuable because they appear in a lot of documents.

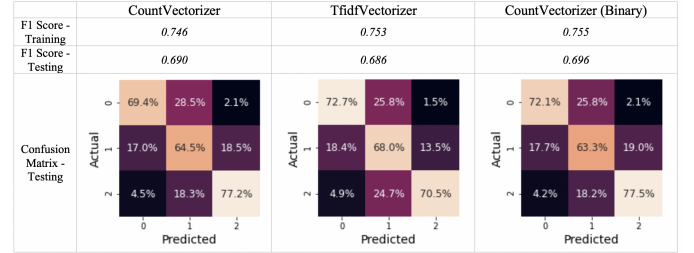


Fig. 6. CountVectorizer (left) vs. TfidfVectorizer (middle) vs. Binary CountVectorizer (right)

Fig. 6 shows the results when models are trained using each of the three representation methods through *CountVectorizer* and *TfidfVectorizer* in *scikit-learn*. Even though the *TfidfVectorizer* model has the lowest F1 score, it has higher True Positive rates for both class 0 and class 1 compared to both *CountVectorizer* models. Since the binary variation of *CountVectorizer* achieves the highest F1 score, we will use it for our future experiments.

D. Number of N-Grams

Let us now investigate the benefit of incorporating n-gram into our model. N-gram is simply a contiguous sequence of

n words. For example, the text corpus “the food is not good” has:

- Five 1-gram/unigram words: “the”, “food”, “is”, “not”, “good”
- Four 2-gram/bigram phrases: “the food”, “food is”, “is not”, “not good”

Adding bigram, trigram or even four-gram phrases into our feature space allows us to capture modified verbs and nouns, thus improve model performance (Wang, 2012 [3]).

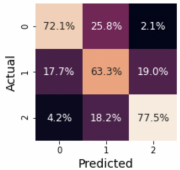
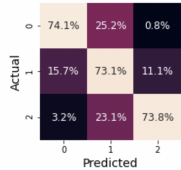
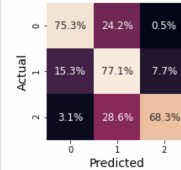
	CountVectorizer (Binary)	CountVectorizer (Binary) - Unigram + Bigram	CountVectorizer (Binary) - Unigram + Bigram + Trigram
F1 Score - Training	0.755	0.909	0.980
F1 Score - Testing	0.696	0.729	0.719
Confusion Matrix - Testing			

Fig. 7. Unigram (left) vs. Unigram + Bigram (middle) vs. Unigram + Bigram + Trigram (right)

For our experiment, we compare three n-gram models - unigram, bigram and trigram. Fig. 7 shows the results of our experiment. As we include more n-grams, the model tends to overfit since the vocabulary size increases. While the True Positive rates for class 0 and 1 both increase, it decreases for class 2. Consider one particular text corpus from the test set:

Food wasn't great. My gnocchi was incredibly greasy. Would not recommend anyone coming here.

This review has a ground-truth label of 0. While it is correctly predicted by our bigram model, our unigram model thinks that it belongs to class 2.

TABLE II
CONDITIONAL LOG-LIKELIHOODS

	Unigram			Unigram + Bigram		
	Class 0	Class 1	Class 2	Class 0	Class 1	Class 2
food	-4.878	-5.023	-5.019	-6.127	-6.184	-6.133
greasy	-8.485	-8.471	-8.792	-9.246	-9.244	-9.616
great	-6.531	-5.934	-5.279	-7.404	-6.880	-6.351
incredibly	-8.987	-9.464	-8.933	-9.722	-10.191	-9.714
incredibly greasy				-13.790	-14.882	-16.118
not	-4.665	-4.862	-5.296	-6.038	-6.139	-6.473
not recommend				-9.269	-10.747	-11.761
recommend	-7.536	-7.561	-6.719	-8.231	-8.314	-7.516
would	-5.748	-5.711	-5.921	-6.794	-6.734	-6.919
would not				-8.666	-9.404	10.417
Predicted Prob	0.343	0.129	0.528	0.999	0.001	0.000

Table II shows the conditional log-likelihoods ($\log \Pr(x|\text{class})$) of some of the features predicted by our unigram and bigram models. We can see that:

- For any given word (e.g., *food*, *great*), the conditional log-likelihood across all three classes is higher (less negative) in the unigram model than in the bigram model. This is because the vocabulary size is larger in the bigram model, resulting in lower relative-frequency for each word

- The “trend” for conditional log-likelihoods among the three classes is the same in the unigram model and the bigram model. For example, under both models, the word *incredibly* has higher log-likelihoods for classes 0 and 2 than class 1. This can be interpreted as: if someone expresses strong feelings in his/her review, he/she is more likely to be either negative or positive than neutral
- Under the bigram model, phrases like *incredibly greasy*, *not recommend*, *would not* are captured. All three phrases assign a much higher log-likelihood to class 0 than to class 2. Even though their subset words (e.g., *recommend*, *would*) favour class 2 over class 0, the bigram phrases’ high log-likelihoods more than offset these favours and tilt the prediction in the bigram model to a 99% probability toward class 0

This particular example illustrates the importance of including contiguous sequences of words in our model. However, trigram model suffers severely from overfitting, as the vocabulary size grows too large.

E. Remove Stopwords

Stopwords are commonly used words, such as “the”, “and”, “to”. They appear in almost every text corpus, so naturally, we would like to remove them from our models.

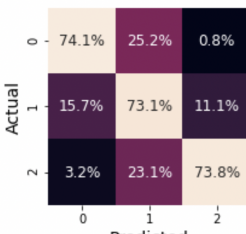
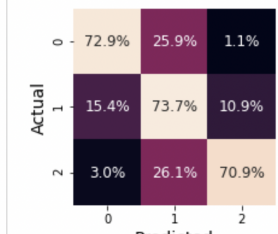
	CountVectorizer (Binary) - Unigram + Bigram	CountVectorizer (Binary) - Unigram + Bigram - Stopwords Removed
F1 Score - Training	0.909	0.958
F1 Score - Testing	0.729	0.711
Confusion Matrix - Testing		

Fig. 8. Stopwords Not Removed (left) vs. Removed (right)

TABLE III
TOP-10 MOST FREQUENT WORDS

	Not Removed	Removed
1	the	food
2	and	good
3	to	place
4	it	service
5	of	like
6	was	would
7	for	one
8	is	back
9	but	really
10	in	great

NLTK, a popular natural language toolkit, provides a list of 179 most common stopwords. Fig. 8 shows the results before

(left) and after (right) we remove these stopwords from our models. Table III shows the resulting top-10 most frequent words. Interestingly, the model performance decreases after we remove these stopwords. This corroborates with the findings in Saif, 2014 [4], which shows that using a pre-compiled list of stop words negatively impacts the performance of sentiment classification.

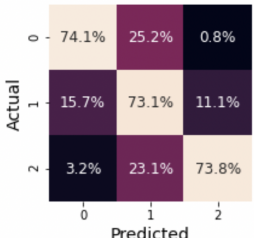
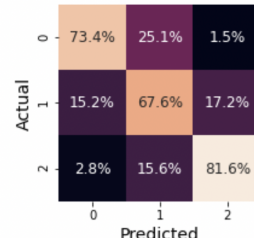
	CountVectorizer (Binary) - Unigram + Bigram	CountVectorizer (Binary) - Unigram + Bigram - Min DF = 6
F1 Score - Training	0.909	0.818
F1 Score - Testing	0.729	0.730
Confusion Matrix - Testing		

Fig. 9. Without (left) vs. With (right) Minimum Frequency Requirement

Instead, Saif, 2014 [4] suggests to eliminate infrequent terms. As such, we set up another experiment, in which we eliminate all words and phrases that appear in at most 5 reviews. Fig. 9 shows the result before (left) and after (right) we impose the minimum frequency requirement on our bigram model. Although the True Positive rates for class 0 and 1 decrease, there is a significant improvement on overfitting and the overall F1 score. More importantly, the feature space shrinks by over 85% from 2,130,891 to 264,977, resulting in much faster training and testing time.

F. Normalization

Normalization is a powerful and effective way to scale down feature space by reducing words to their roots. We will experiment with two popular normalization techniques - stemming and lemmatization.

TABLE IV
EXAMPLE TEXT NORMALIZATION

	Original	Stemmed	Lemmatized
1	what a trouble	what a troubl	what a trouble
2	this is very troubling	thi is veri troubl	this be very troubling
3	i am troubled	i am troubl	i be trouble

Table IV shows three sentences after passing through *PorterStemmer* and *WordNetLemmatizer* in *NLTK*. With stemming, “trouble”, “troubling” and “troubled” all get reduced to “troubl”, which is not a word. This is common in stemming. On the other hand, with lemmatization, the results are still words, and “troubling”, which is an adjective, remains to be “troubling” instead of “trouble”. This is because, with lemmatization, we can also do part-of-speech tagging, which

determines the identification of a word (e.g., noun, verb, adjective) depending on the context. However, this system is not perfect, as “troubled” is also an adjective in this context but still gets reduced to “trouble”. Overall, lemmatization with part-of-speech tagging attempts to minimize the loss in the original meaning, at the expense of lower feature space reduction than stemming.

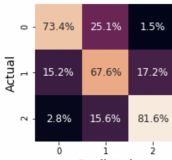
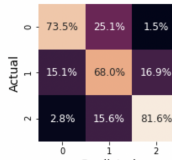
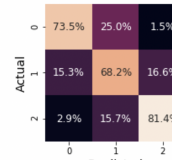
	CountVectorizer (Binary) - Unigram + Bigram - Min DF = 6	CountVectorizer (Binary) - Unigram + Bigram - Min DF = 6 - Stemmed	CountVectorizer (Binary) - Unigram + Bigram - Min DF = 6 - Lemmatized
F1 Score - Training	0.818	0.813	0.813
F1 Score - Testing	0.730	0.731	0.731
Confusion Matrix - Testing			

Fig. 10. Original (left) vs. Stemmed (middle) vs. Lemmatized (right)

Fig. 10 shows the results after we perform text normalization. As we can see, there is very little improvement in model performance. However, considering the reduction in overfitting and feature space, we will continue to lemmatize our texts for our model builds.

Finally, our text preprocessing pipeline is as follows:

- Use the binary variation of *CountVectorizer*
- Unigram + bigram
- Minimum document frequency of 6
- Lemmatized, with part-of-speech tagging

When trained on a multinomial Naive Bayes model, we obtain a macro F1 score of 0.731 on the test set. In the next section, we will experiment with other types of models to see if we can outperform this score.

V. MODEL RESULTS

We consider three popular machine learning models for text classification - Logistic Regression [9], Support Vector Machine (SVM) [10] and Gradient Boosting (XGBoost) [13]. We also consider two deep learning architectures that are widely used for NLP tasks - Long Short Term Memory (LSTM) [11] and Bidirectional Encoder Representations (BERT) [12].

TABLE V
MODEL RESULTS

	Macro F1 Score (Test)	Training Time
Multinomial Naive Bayes	0.731	341 ms
Logistic Regression	0.757	2 min 5 sec
Support Vector Machine	0.757	11 sec
Gradient Boosting (XGBoost)	0.750	1 hr 25 min
LSTM	0.701	18 min 20 sec
BERT	0.723	3 hr 15 min

Table V shows the model performances, as well as the training times for all of our trained models. All models are

trained on Google Colab (some using GPU/TPU backends). Note that training time does not account for the time for hyperparameter tuning, which could take hours or even days depending on the number of hyperparameters. As it turns out, more complex models (Gradient Boosting, LSTM, BERT) do not necessarily improve model performance. This corroborates with the findings in many other papers related to sentiment analysis in social media (e.g., Salinca, 2015 [5] and Xu, 2014 [6]).

A. Logistic Regression, SVM

Logistic Regression models offer great interpretability and transparency, something that other machine learning and deep learning models often lack. Coefficients in a logic regression model are log odds. Since our problem is a multi-class one, each feature has three coefficients (one for each class).

TABLE VI
MOST NEGATIVE (LEFT), NEUTRAL (MIDDLE) AND POSITIVE (RIGHT)
WORDS/PHRASES

	Negative	Neutral	Positive
1	horrible	3 star	delicious
2	terrible	3 5	excellent
3	bland	three star	amazing
4	disappointing	ok	awesome
5	2 star	a ok	perfect
6	disappointed	average	fantastic
7	not worth	decent	amaze
8	rude	not bad	best
9	poor	okay	4 star
10	awful	however	incredible

Table VI shows the most negative, neutral and positive words and phrases based on our trained logistic regression model. For example, the most negatively word is “horrible”, with fitted coefficients of (0.593, -0.11, -0.483) for classes 0, 1 and 2 respectively.

TABLE VII
MOST (LEFT) VS. LEAST (RIGHT) DISCRIMINATIVE WORDS/PHRASES

	Most	Least
1	3 star	even pull
2	delicious	our condo
3	3 5	decor that
4	three star	9 friend
5	terrible	that fine
6	horrible	crispy perfection
7	not worth	anyone because
8	bland	canned corn
9	excellent	mouth open
10	disappointing	rice either

Table VII shows the most and least discriminative words. This is defined as the standard deviation among the fitted coefficients. For example, the most discriminative phrase is “3 star”, with fitted coefficients of (-0.525, 0.868, -0.343), unsurprisingly high for class 1 and low for the other two classes. On the other hand, all of the least discriminating words/phrases have coefficients that are essentially 0, meaning

that they don’t contribute much to our model prediction and could be eliminated from our feature space in model builds.

For each list, we can get almost identical results from our trained SVM model. Similar results are seen in Yu, 2017 [7], which also uses SVM.

B. LSTM

Deep learning models, particularly Recurrent Neural Networks (RNN), have been proven successful in many NLP related tasks due to its nature of allowing information to persist, much like how humans think. LSTM is a special kind of RNN that solves the “vanishing gradient” problem when the dependency is long (Olah, 2015 [8]).

After experimenting with many different techniques, we find that:

- Text normalization (e.g., lemmatization with part-of-speech tagging) does not improve model performance
- Using pre-trained word embeddings such as GloVe (<https://nlp.stanford.edu/projects/glove>) helps with minimizing variance during the training phase
- If pre-trained word embeddings are used, it is better to not cap the vocabulary size. Otherwise, capping the vocabulary size to say 5,000 words yields better results
- Even though 20% of our reviews have more than 200 tokens, while only 5% of our reviews have more than 300 tokens, capping maximum length to 200 reduces overfitting and yields better overall performance
- Use dropouts in between stacked LSTM layers to reduce overfitting

TABLE VIII
LSTM ARCHITECTURE SUMMARY

Layer	Output Shape	Param #
Embedding	(None, 200, 200)	16,381,400
LSTM	(None, 200, 64)	67,840
Dropout	(None, 200, 64)	0
LSTM	(None, 64)	33,024
Dense	(None, 64)	4,160
Dense	(None, 3)	195
Total params: 16,486,619		
Trainable params: 105,219		
Non-trainable params: 16,381,400		

Table VIII shows the architecture of our final LSTM model. Since we are using the word embedding from GloVe pre-trained on Wikipedia, only a fraction of the parameters in our model are trainable, greatly reduces the computation load.

Fig. 11 shows the accuracy and loss per epoch during the training phase, where we use 20% of our training set as validation set for early-stopping. As we can see, the model learns very quickly and most of the biases disappear before epoch 5.

Fig. 12 shows the confusion matrix for our trained LSTM model on the test set. Most of the errors come from class 2, in which the model incorrectly assigns 34% of the reviews to class 1. This is because positive ratings tend to have shorter reviews than negative and neutral ratings – in our

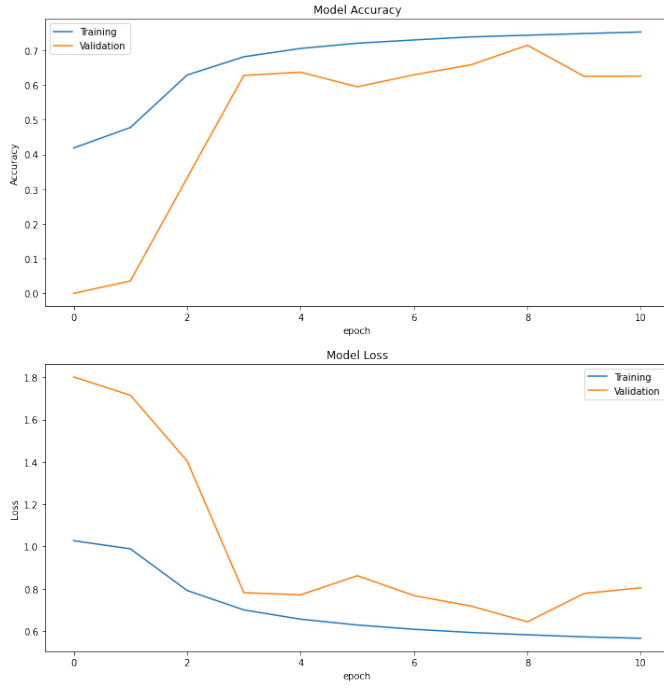


Fig. 11. LSTM Training History

Actual	0	79.8%	19.3%	0.9%
	1	18.7%	72.3%	9.0%
	2	2.8%	34.0%	63.2%
		Predicted		
		0	1	2

Fig. 12. LSTM Confusion Matrix

dataset, only 15% of positive ratings exceed the 200 tokens maximum length cap, while 23% of negative and neutral ratings exceed this cap. Our experiment shows that, if we increase the maximum length cap to 300, errors become much more symmetrical. However, the model overfits and the overall performance deteriorates.

VI. CONCLUSION AND FUTURE WORKS

In this paper, we perform an ablation study on text preprocessing techniques and compare the effectiveness of several machine learning and deep learning models on predicting user sentiments.

For machine learning models, we find that using binary bag-of-word representation, adding bi-grams, imposing minimum frequency constraints and normalizing texts have positive effects on model performance. For deep learning models, we find that using pre-trained word embeddings and capping maximum length often boost model performance.

Overall, we find simpler models such as Logistic Regression and SVM to be more effective at predicting sentiments than more complex models, in terms of both model performance and training time. More importantly, simpler models offer great interpretability, while more complex models require other techniques (e.g., LIME analysis [14]) to understand its inner workings.

However, since deep learning models require much heavier computation than machine learning models, one possible improvement to our analysis is to run our models on better hardware setup, so we can have more relaxed constraints on maximum length and possibly use larger dataset.

Lastly, we re-emphasize that there is no one-fits-all solution to building sentiment analysis models, and the insights from this paper should only be considered as guidance. Deep learning models, albeit their lack of interpretability, tend to outperform machine learning models for more complex NLP tasks such as speech recognition, translation and automatic summarization.

REFERENCES

- [1] Sokolova, Marina, and Guy Lapalme. "A systematic analysis of performance measures for classification tasks." *Information processing management* 45.4 (2009): 427-437.
- [2] Renault, T. Sentiment analysis and machine learning in finance: a comparison of methods and models on one million messages. *Digit Finance* (2019).
- [3] Wang, Sida, and Christopher D. Manning. "Baselines and bigrams: Simple, good sentiment and topic classification." *Proceedings of the 50th annual meeting of the association for computational linguistics: Short papers-volume 2*. Association for Computational Linguistics, 2012.
- [4] Saif, Hassan, et al. "On stopwords, filtering and data sparsity for sentiment analysis of twitter." (2014): 810-817.
- [5] Salinca, Andreea. "Business reviews classification using sentiment analysis." *2015 17th International Symposium on Symbolic and Numeric Algorithms for Scientific Computing (SYNASC)*. IEEE, 2015.
- [6] Xu, Yun et al. "Sentiment Analysis of Yelp's Ratings Based on Text." (2014).
- [7] Yu, Boya, et al. "Identifying Restaurant Features via Sentiment Analysis on Yelp Reviews." *arXiv preprint arXiv:1709.08698* (2017).
- [8] Olah, Christopher. "Understanding LSTM Networks." *Understanding LSTM Networks*, colah.github.io/posts/2015-08-Understanding-LSTMs/.
- [9] Hosmer Jr, David W., Stanley Lemeshow, and Rodney X. Sturdivant. *Applied logistic regression*. Vol. 398. John Wiley Sons, 2013.
- [10] Suykens, Johan AK, and Joos Vandewalle. "Least squares support vector machine classifiers." *Neural processing letters* 9.3 (1999): 293-300.
- [11] Hochreiter, Sepp, and Jürgen Schmidhuber. "Long short-term memory." *Neural computation* 9.8 (1997): 1735-1780.
- [12] Devlin, Jacob, et al. "Bert: Pre-training of deep bidirectional transformers for language understanding." *arXiv preprint arXiv:1810.04805* (2018).
- [13] Chen, Tianqi, and Carlos Guestrin. "Xgboost: A scalable tree boosting system." *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*. 2016.
- [14] Ribeiro, Marco Tulio, Sameer Singh, and Carlos Guestrin. "Why should i trust you?" Explaining the predictions of any classifier." *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*. 2016.
- [15] Da Silva, Nadia FF, Eduardo R. Hruschka, and Estevam R. Hruschka Jr. "Tweet sentiment analysis with classifier ensembles." *Decision Support Systems* 66 (2014): 170-179.