

Comparing approaches to analyze sentiment

Luka Žontar

Faculty of Computer and Information Science

University of Ljubljana

Ljubljana, Slovenia

Email: lz3057@student.uni-lj.si

Abstract—Sentiment analysis can help us automatically determine writer’s attitude towards the text on enormous datasets. In practice, it helps companies make important business decisions and evaluate public opinions on their products. This article should help readers understand different machine learning approaches to sentiment analysis. We also explain how textual data should be preprocessed generally and for different models.

This article focuses on deep learning machine learning models. At first, we take the Naive Bayesian model to compare it with more complex machine learning models. In particular, we try to evaluate a deep neural network with LSTM architecture and a BERT model that was fine-tuned on data that was used in training and validation of other models. As our dataset, we used Sentiment140, which contains pre-classified Twitter posts. In the end, we combine the aforementioned models to construct a majority voting ensemble model.

Index Terms—Sentiment analysis; classification; deep learning; neural networks.

I. INTRODUCTION

Sentiment analysis is becoming an important technique used in many companies worldwide. With it, we use natural language processing to systematically identify subjectiveness in texts. In other words, in sentiment analysis we try to determine, whether the writer is positively or negatively affected by the written content. We try to evaluate writer’s inclination towards the text using natural language processing and machine learning techniques to assign sentiment scores to phrases and sentences.

In practice, sentiment analysis helps us automatically process enormous amounts of textual data that customers write on social media and in survey responses. Customer feedback analysis can help us determine whether or not customers are happy with our product. This can improve our company’s strategic decisions, such as deciding in which product to invest or how to make customers happier.

Furthermore, it can help better understand public opinions about a particular subject. This can help politics predicting voting results and consequently form alliances between parties. Understanding people’s opinions can also make us better investors, since stock market stability depends on other investors.

which can help us predicting voting results or even stock market state.

In this article, we focus on comparing different sentiment analysis techniques. The goal of this article is to help readers understand different types of sentiment analysis approaches

and when they should be used. We explain different preprocessing techniques and compare a baseline Naive Bayesian model with more complex deep learning based models LSTM and BERT. In the end, we compare results with a majority voting ensemble model that is combined of all three aforementioned models.

For learning and validation we use the Sentiment140 dataset, which contains pre-classified Twitter posts. Tweets represent the general public and thus all sorts of opinions can be found in the dataset we used. Some opinions are stronger and in some it is difficult to determine the writer’s inclination towards the text topic. This makes the given dataset an appropriate training and validation set.

In next section, we make an overview of related work that was relevant for our research. After that, we explain the methodology of our work, that is, how data was preprocessed and how models were evaluated. In the forth and fifth section, we present results and discuss them. Lastly, we conclude this article by presenting the main findings.

II. RELATED WORK

In this section, we present related work that was relevant for our article. Sentiment analysis was proven to be useful in different areas. For example, in medicine, sentiment analysis is used to evaluate healthcare textual data. In medical sentiment, we use medical information to achieve the best result to increase healthcare quality [1]. As mentioned before, social media holds loads of textual data of different sorts. Muzafar et al. discussed how SARS-Cov-2 affected people on social media [2]. At first, the most commonly searched terms were Coronavirus related. And lots of negativity could be seen textual data from social media.

Feldman [3] discusses several approaches of sentiment analysis based on the basic unit that we will be classifying. We can classify whole documents, sentences or even particular aspects. Furthermore, we have to consider comparative statements such as less and more. As the most crucial field of sentiment analysis Feldman states sentiment lexicon acquisition, where words can be expanded in a set of words using corpus-based approaches or WordNet.

In past works, lots of researchers also worked on textual data preprocessing. Magliani and his colleagues [4] evaluated lots of preprocessing techniques and found that preprocessing vastly increases model’s accuracy. However, they found that basic preprocessing with stemming is the best way to improve

accuracy with preprocessing. They also found that processing negations, stopwords and emoticons also slightly increases accuracy.

Latha [5] focused on how to process informal social media texts and emphasized the importance of removing URLs, special characters and repeated letters from texts. Furthermore, she proposed that question words should be removed from texts. We also noticed an interesting approach by Krouska et al. [6], who focused on the length of n-grams passed to classifiers. N-grams of length greater than 1 hold a lot more information about the text polarity than a single word, which is why they are used to get better accuracy on sentiment analysis classification.

To better understand sentiment analysis, we also need to get acquainted with the top performing machine learning models. Zainudin [7] split the sentiment analysis models into three major categories, lexicon based, machine learning based and hybrid approaches. Machine learning based approaches are further divided into supervised, semi-supervised and unsupervised approaches. In this article, we will be focusing on supervised machine learning approaches.

Starting with simpler methods, Dey [8] et al. used statistical models to capture elements of subjective style and the sentence polarity. They managed to get above 80% accuracy on a movie review dataset that is comprised of 4500 samples. In this study, Naive Bayesian classifier was compared with k-NN classifier. However, k-NN classifier failed to achieve comparable results on this dataset. Researchers also noticed high discrepancies amongst accuracies on different datasets.

In recent years, there have been several breakthroughs in the area of sentiment analysis. Recurrent neural network architectures, such as long short-term memory (LSTM), have proven to be very successful in learning to evaluate polarity of text. Besides the standard feedforward connections, LSTM adds feedback connections. Putra et al. [9] improved the basic LSTM architecture by adding exploiting the use of Word2Vec embeddings. They managed to achieve 85.96% accuracy. Word embeddings make the regular LSTM architecture more robust to overfitting. For Arabic texts, an additional improvement to LSTM architecture was proposed by Elfaik and Nfaoui [10], who used a deep bidirectional LSTM network to further increase the accuracy. While the complex model clearly overfitted on some datasets, it reached over 90% accuracy on others.

A novel approach was proposed by Devlin and colleagues [11], who used deep bidirectional transformers for language understanding. These transformers were trained on enormous datasets and finally fine-tuned on sentiment analysis tasks to produce state-of-the-art results.

III. METHODS

In this section, we explain different methods that were used in this comparative study. Firstly, we describe the data that we used to get our results. Next, we explain the preprocessing methodology and finally, we describe the machine learning models that were used in this research. We used four different

models: Naive Bayesian classifier, neural network with deep bidirectional LSTM architecture, BERT model fine-tuned with our data and a majority voting ensemble classifier.

A. Data

As already mentioned, we decided to use the Sentiment140 dataset that contains 1,600,000 tweets extracted using Twitter API. In Figure 1 we can see the most common words in the wordcloud representation. Due to time complexity of our research, we only use a randomly sampled subset of 40,000 tweets. While the dataset holds multiple features, we only use the tweet text and target variable from which we filter out the neutral tweets to make our problem a binary classification problem. Also, the positive tweet polarity, which is denoted with 4 is mapped to 1 due to convenience.

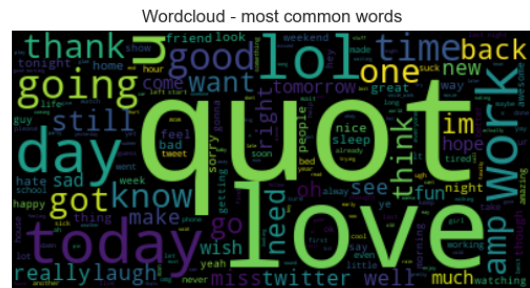


Fig. 1. **Wordcloud:** from figure we see the most common words in tweet. People often use quotes and words like day, today, love and word.

As can be seen from Figure 2, our classification problem is balanced. That is, the two classes are uniformly distributed. The next steps of preprocessing pipeline consists of the following:

- 1) Remove non-ASCII characters.
- 2) Remove links, hashtags, mentions and repetitive vowels (such as “hahaa”).
- 3) Replace laughing with a predefined word “laugh”.
- 4) Replace emojis with predefined combinations such as “smile”.
- 5) Switch to lower case.
- 6) Filter stopwords.
- 7) Tokenize words and turn them into appropriate form (depending on the algorithm).

We also tried to improve evaluation metric scores with stemming words, fixing misspells, handling negations and removing retweets. However, these preprocessing techniques made the models overfit faster and effectively reduced evaluation metric scores. While, Naive Bayesian classifier was not really effected by those, LSTM and BERT performed worse. We assume that performance drop occurred due to

different pretraining in BERT and failing to find appropriate embeddings in LSTM with word embeddings. Rows with null values were dropped and finally, we used 80% of data for training and the rest for testing our models.

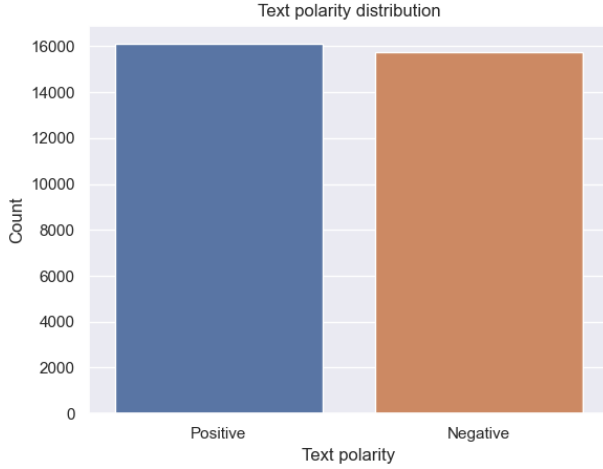


Fig. 2. **Text polarity distribution:** from figure we denote that text polarity is almost uniformly distributed. Both classes are representable.

In Figure I, we can see the distribution of tweet length after preprocessing. This later helps us define parameters and model complexities in neural networks.

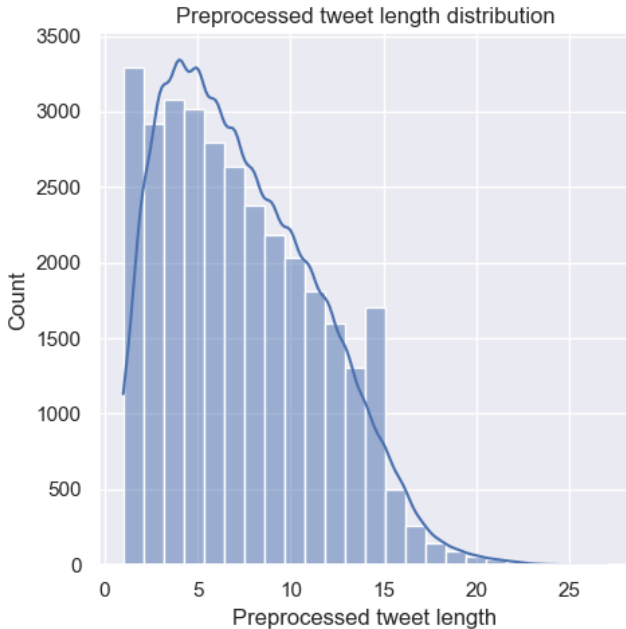


Fig. 3. **Preprocessed tweet length distribution:** as can be seen from the figure, the distribution is right-skewed. This means that more tweets are short. The distribution does not have a long tail, which is why, we do not have to handle outliers in this situation.

B. Naive Bayes

Naive Bayesian classifier is one of the fastest and easiest methods that can also perform very well. The basic idea of this classifier is to find the probabilities of classes based on the input text. Textual data is passed to the model with word count vectors. With bigger datasets this model overfits, however, it tends to be used as the baseline model.

C. LSTM

As already mentioned, LSTM has been proven to be very succesful in sentiment analysis. We used the aforementioned bidirectional LSTM architecture with GloVe word embeddings instead of matrix of tokenized word encodings. We used GloVe with 200 dimensional representations. Next, we describe the architecture of the LSTM model in more detail:

- 1) Model accepts a matrix of tokenized tweets padded to the pre-defined maximum length. This was set to 60 due to tweet length distribution.
- 2) To decrease the chance of overfitting, we transform this matrix to a matrix of word embeddings and only use the top 4000 most common words.
- 3) Learning rate is set to $5e-5$ and we train the model for maximum 100 epochs. Early stopping is implemented if the model starts to overfit.
- 4) The model consists of two bidirectional LSTM layers with dropout (once again, to prevent overfitting).
- 5) In the end results are densed using ReLU activation function and on the last layer a sigmoid function.
- 6) We use binary crossentropy loss function and Adam optimizer.

Despite all the precautions we did against overfitting, the model performance still becomes too dependent on the training set after first few epochs as we can see in Figures 4, 5.

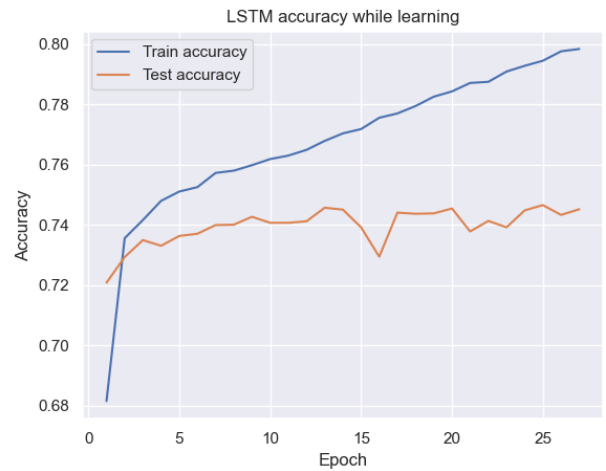


Fig. 4. **LSTM accuracy:** in first few epochs accuracy is rising quickly. However, it becomes more or less constant in median after 5 or 10 epochs, while varying up and down.

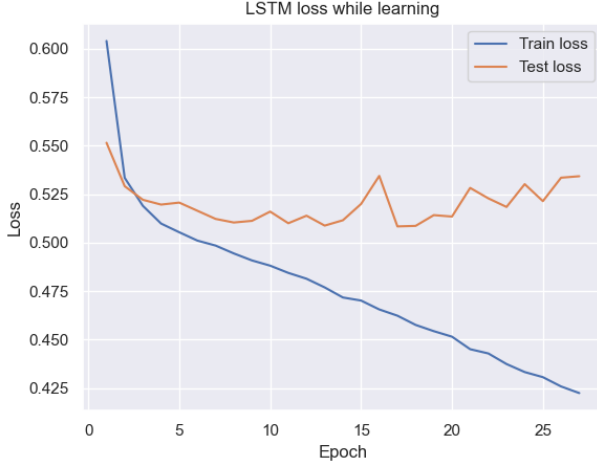


Fig. 5. **LSTM loss**: in first few epochs loss is dropping quickly. However, it becomes more or less constant in median after 5 or 10 epochs, while varying up and down.

D. BERT

E. Majority voting ensemble classifier

Lastly, we explain the ensemble classifier that exploits the aforementioned three models and predicts the class with majority voting. With such ensemble method, we should get a more robust model. Evaluation metric scores could also be increased with such methods.

IV. RESULTS

Finally, we present results using the methodology that we described above. To better understand, how well our models generalize, we used a 10-fold cross validation in all the models. Using local resources with 6GB GPU, we managed to train 33 different models.

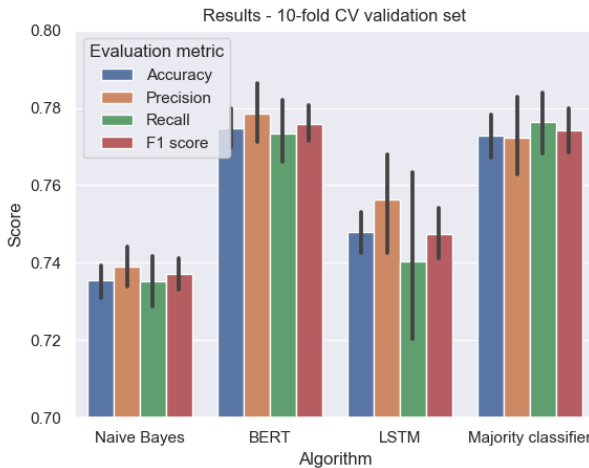


Fig. 6.

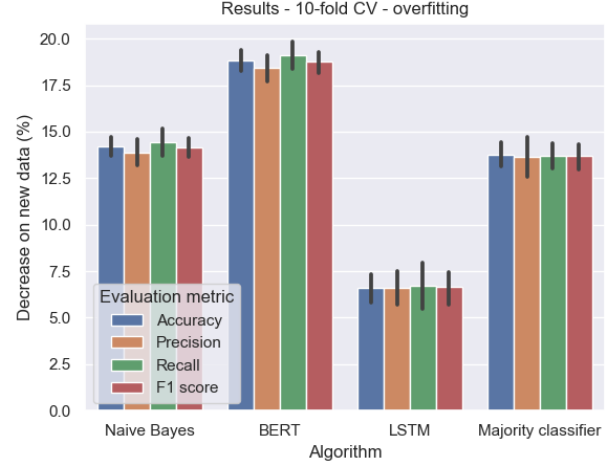


Fig. 7.

V. DISCUSSION

VI. CONCLUSION

REFERENCES

- [1] K. Denecke and Y. Deng, "Sentiment analysis in medical settings: New opportunities and challenges," *Artificial Intelligence in Medicine*, vol. 64, no. 1, pp. 17–27, 2015. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0933365715000299>
- [2] M. Bhat, M. Qadiri, N.-u.-A. Beg, M. Kundroo, A. N. Ahanger, and B. Agarwal, "Sentiment analysis of social media response on the covid19 outbreak," *Brain, Behavior, and Immunity*, vol. 87, 05 2020.
- [3] R. Feldman, "Techniques and applications for sentiment analysis," *Commun. ACM*, vol. 56, p. 82–89, 04 2013.
- [4] F. Magliani, T. Fontanini, P. Fornacciari, S. Manicardi, and E. Iotti, "A comparison between preprocessing techniques for sentiment analysis in twitter," 12 2016.
- [5] I. Latha, D. Varma, and D. Govardhan, "Preprocessing the informal text for efficient sentiment analysis," *International Journal of Emerging Trends & Technology in Computer Science*, vol. 1, pp. 58–61, 07 2012.
- [6] A. Krouska, C. Troussas, and M. Virvou, "The effect of preprocessing techniques on twitter sentiment analysis," 07 2016, pp. 1–5.
- [7] S. Zainudin, "Sentiment analysis techniques in recent works," 07 2015.
- [8] L. Dey, S. Chakraborty, A. Biswas, B. Bose, and S. Tiwari, "Sentiment analysis of review datasets using naïve bayes' and k-nn classifier," *International Journal of Information Engineering and Electronic Business*, vol. 8, no. 4, p. 54–62, Jul 2016. [Online]. Available: <http://dx.doi.org/10.5815/ijieeb.2016.04.07>
- [9] P. Muhammad, R. Kusumaningrum, and A. Wibowo, "Sentiment analysis using word2vec and long short-term memory (lstm) for indonesian hotel reviews," *Procedia Computer Science*, vol. 179, pp. 728–735, 01 2021.
- [10] H. Elfaik and E. H. Nfaoui, "Deep bidirectional lstm network learning-based sentiment analysis for arabic text," *Journal of Intelligent Systems*, vol. 30, no. 1, pp. 395–412, 2021. [Online]. Available: <https://doi.org/10.1515/jisys-2020-0021>
- [11] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," 2019.

Algorithm	Accuracy		Precision		Recall		F1 score	
	Train	Test	Train	Test	Train	Test	Train	Test
Naive Bayes	0.855	0.733	0.855	0.739	0.857	0.725	0.856	0.732
LSTM	0.810	0.745	0.790	0.730	0.848	0.781	0.818	0.755
BERT	0.951	0.775	0.953	0.783	0.951	0.763	0.952	0.773
Majority classifier	0.896	0.768	0.890	0.765	0.905	0.776	0.897	0.771

TABLE I