

Twitter Sentiment Analysis

Tijana Ilievska
Politecnico Di Torino
s301381
s301381@studenti.polito.it

Sevde Bekdik
Politecnico Di Torino
s298317
s298317@studenti.polito.it

Abstract—In this article, we propose a method to predict the sentiment of the tweets with a given textual data on Twitter. The proposed approach consists of a machine learning pipelines involving the use of Linear Regression, Linear Support Vector Classifier and Naive Bayes along with using Term Frequency-Inverse Document Frequency (TF-IDF). F1 score and accuracy are used in evaluating the performance of the classifier.

I. PROBLEM OVERVIEW

Nowadays, through the online communities, large volumes of sentiment rich data is generated in different formats. Sentiment analysis or also known as 'opinion mining' is the process of extracting subjective information which can help in determining the polarity of the data. This project aims to build a classification model to predict the sentiment of tweets on Twitter either positive or negative. The data set is split into two portions:

- 1) "Development set", collection of 224597 tweets in tabular format and characterized by:
 - *ids*: a numerical identifier of the tweet;
 - *date*: the publication date;
 - *flag*: the query used to collect the tweet;
 - *user*: the username of the original poster;
 - *text*: the text of the tweet;
 - *sentiment*: the sentiment of the tweet which is positive if equals 1 and negative if equals to 0.
- 2) "Evaluation set", containing 74872 tweets without the *sentiment* feature.

We have already labeled data and our goal is to classify the sentiment of the tweets represented in the '*text*' column of the Evaluation set, by using the proposed models. The models are trained and validated on the Development set.

First of all, we imported the data set and libraries required and did some preliminary exploration of the development data in order to identify possible problems. As we can see from the Fig. 1, the data set is balanced and no **duplicate** or **missing** values were found in it. Next, we performed data preprocessing in order to clean and convert textual data to numeric data, and after that, we split our data into a training set which is used for training the model, and a testing set in order to evaluate the model. Once data is divided, we used three different algorithms to train our model; Linear Regression, SVM, and Naive Bayes. Lastly, after training and testing the model, we used f1-score and accuracy in order to evaluate the performance.

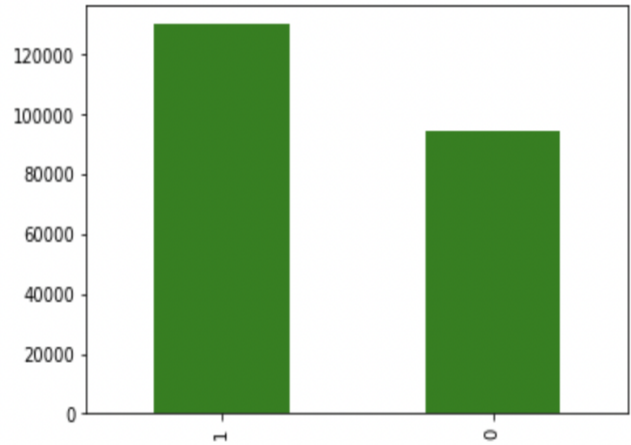


Fig. 1. Distribution of the target variable in the Development Set

II. PROPOSED APPROACH

A. Preprocessing

The first step of our work is concatenating both the development and evaluation set in order to be able to preprocess both sets at once. In addition, we drop the "user", "flag", "date" and "id" columns since we believe that they do not have any effect on what sentiments users are expressing in a certain tweet. Our outcome is based solely on the "text" column, thus we are only using these two columns in the further processing. The raw text can contain a lot of things which are not important for our analysis, such as punctuation, numbers, stopwords etc. Also, one word can have different forms due to grammatical reasons and it is important to normalize the word to its root form. To tackle these issues, we applied the following steps to clean the data:

- Converting all the letters to lower
- Removing all URLs
- Replacing emoticons with their sentiments
- Removing usernames
- Removing any extra empty spaces
- Removing any numbers
- Replacing any three or more consecutive letters with two
- Lemmatization (converting words to their root form)

We also explored additional preprocessing techniques, such as removing stopwords and stemming, but in our case this particular set of preprocessing steps has shown to be the most

[illegible]

Word	Frequency (approx.)
you	5900
to	5800
the	5400
it	5200
my	3800
word	3700
and	3400
that	3200
is	3100
for	3100
me	3000

[illegible]

Model	Parameter	Parameter values
Logistic Regression	C penalty	{ 0.01, 0.1, 1, 2, 5, 10, 100, 1000} {'l1', 'l2'}
LinearSVC	C	{0.01, 0.1, 1, 10}
Naive Bayes	alpha	{0.01, 0.1, 0.5, 1.0, 10.0}

TABLE I
HYPERPARAMETERS CONSIDERED

The results of the search for each model are shown below:

Model	Parameter	Calculated value
Logistic Regression	C	2
	penalty	l2
LinearSVC	C	0.1
Naive Bayes	alpha	0.5

TABLE II
HYPERPARAMETERS TUNING WITH RANDOMIZED SEARCH

III. RESULTS

By analyzing and comparing all the results after evaluating the model, we are now able to determine the best algorithm in order to make a prediction.

Considering the accuracy of the model, it is observed that the best performing algorithm is Logistic Regression, while SVM follows it and outperforms Naive Bayes as can be seen in Fig. 5.

The f1-score for both positive and negative sentiments:

- For 0: Logistic Regression (f1-score = 0.74) > SVM (f1-score = 0.73) > Naive Bayes f1-score = 0.70)
- For 1: Logistic Regression (f1-score = 0.82) > SVM (f1-score = 0.81) > Naive Bayes (f1-score = 0.80)

Therefore, Logistic Regression is concluded as the best algorithm for the prediction.

Model	f1-score
Logistic Regression	0.79
SVM	0.78
Naive Bayes	0.76

Fig. 5. Accuracy of Model Selection

The confusion matrices of the algorithms with the ratios between predicted and actual values, for both outcomes, are given in Fig. 6, Fig. 7 and Fig. 8.

After fitting the models with the tuned parameters, we had a slight improvement in the performance of all the algorithms but Logistic Regression was still outperforming the other two.

Model	Before tuning	After tuning
Logistic Regression	0.7871	0.7873
LinearSVC	0.7790	0.7860
Naive Bayes	0.7642	0.7646

TABLE III
ACCURACY IMPROVEMENT WITH HYPERPARAMETERS TUNING

IV. DISCUSSION

During our work on this project, we tried different types of preprocessing and what came as a surprise to us was that removing the stop words resulted with lower performance. We attempted this by using both pre-defined list from the nltk

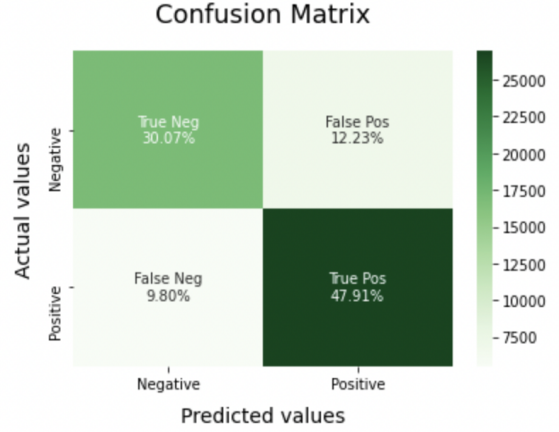


Fig. 6. Confusion Matrix for Linear SVC

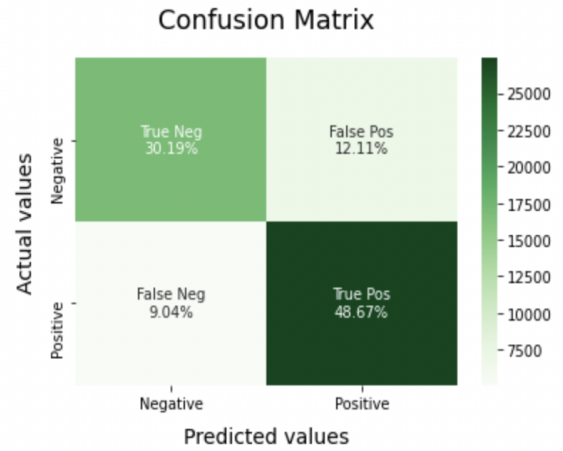


Fig. 7. Confusion Matrix for Logistic Regression

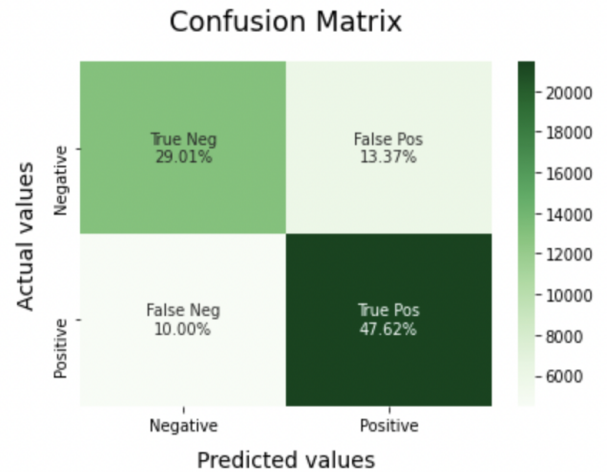


Fig. 8. Confusion Matrix for Naive Bayes

library and a custom list of stopwords (without negation words which are actually present in the nltk library). For the sake of the performance, we decided to keep all the words in the final model. Also, we observed that considering the sentiment of the emoticons in our analysis has shown to considerably improve the overall sentiment score on both positive and negative opinions. Due to large size of the dataset and our limited resources we were unable to perform GridSearch and only used a small subset of random parameters for tuning. We believe that there are better arrangements of parameters that could result in improved performance. After all, we can see that all three models performed with very small differences in the accuracy. Based on this we can conclude that preprocessing is a key factor for effective sentiment analysis.

REFERENCES

- [1] Sampooram, K. P., Hemalatha B., Kaaviyasri, M., Jaya Priyanka, K., "Sentiment analysis using Machine learning techniques on python", *Conf. Series: Materials Science and Engineering*, doi:10.1088/1757-899X/1084/1/012028 .
- [2] Agarwal, A., Xie, B., Vovsha, I., Rambow, O., Passonneau, R., "Sentiment analysis of Twitter data", *2003 IEEE MTT-S Int. Microwave Symp. Dig.*, vol. 3, pp. 32-33, 2020.
- [3] Saif, H., Fernandez, M., He, Y., Alani, H., O., "On Stopwords, Filtering and Data Sparsity for Sentiment Analysis of Twitter", *Knowledge Media Institute, The Open University, UK*
- [4] Scikit-learn: "Machine Learning in Python", *Pedregosa et al., JMLR 12* pp. 2825-2830, 2011
- [5] B. Liu, "Sentiment analysis and subjectivity". *Handbook of natural language processing*, 2010.