

Toxicity Analysis in Twitter using Machine Learning Models

Álvaro Mazcuñán Herreros and Miquel Marín Colomé
Escuela Técnica Superior de Ingeniería Informática, ETSINF
Universitat Politècnica de València, Spain
{almazher, mimacol}@inf.upv.es

Abstract. In this paper we will explain briefly how the “Dembo” team approached the problem of detecting toxic language. This project covers two tasks. The first is to classify whether a text (set of tweets) is toxic or not and the second one is to classify that text into different levels of toxicity. Two models for those sub-tasks have been submitted to the competition. The first one is the hybrid stacking model in which SVM, Decision Tree, Random Forest and MLP models have been used with a logistic regression with the function of being metalearner. The second is the BETO model, which is a variant of the BERT model. All this is summarised in a table of results.

Keywords: Text-Classification, Natural Language Processing, Bag of Words, TF-IDF, Stemming, Lemmatization, Stacking, Support Vector Machines, Multi-layer Perceptron, Logistic Regression, Decision Tree, Random Forest, Transformers, BERT, BETO, Sklearn, Twitter

1. Introduction

The challenge of dealing with hate speech is ancient, but the scale, personalisation and speed of today's hate speech poses a uniquely modern quandary. While there is no precise definition of hate speech, it is generally speech that is intended not only to insult or ridicule, but to cause lasting by attacking something that is particularly important to the victim. Hate speech is widespread in online forums and social media. Some previous work has been carried out on the subject of Hate Speech. [1] [2]

Once the problem of detecting toxic language has been introduced, it should be said that in the different models made, except for the BERT model, the sklearn library has been used. Two datasets were available for this task. The first was the train dataset, with a total of 3958 tweets. On the other hand, in order to validate the quality of the model, we had the test dataset, which consisted of 891 tweets. Overall, this dataset had a wide variety of variables such as: constructiveness, positive / negative stance, stereotype, sarcasm, aggressiveness among others. However, for this competition, only the comment variable was used, i.e., the variable containing the different tweets from various Spanish newspapers such as ABC, elDiario.es, El Mundo, etc.

The aim of this competition was twofold. On the one hand, it was required to label the tweets in the test set with 0 or 1, i.e., whether these tweets were non-toxic or toxic, respectively. For this purpose, the variable toxicity was available in the training set to be able to evaluate the corresponding machine learning models. In addition, in the second subtask, the objective was a little more complicated, because in this case we were asked to label the same set of tests but, in this case, adding different levels of toxicity:

- 0 → Not toxic
- 1 → Mildly toxic
- 2 → Toxic
- 3 → Very toxic

As discussed above for the first subtask, in order to evaluate our models, the `toxicity_level` variable was available in the set of 3958 tweets.

2. System

2.1. Preprocessing

Before performing any separation of the data in order to train the corresponding models, it was decided to carry out a small amount of data cleaning/preprocessing. To do this, the first thing that was done was to remove from the messages those characters that were emojis, hashtags (#), URLs and other special characters. Once these parts of the tweets had been eliminated, we began **tokenizing** the text, i.e. separating the tweet into words and, with this, obtaining a list of these words.

With this list of words, the next step was to eliminate **stopwords**, i.e. words that have no meaning in themselves. This group of words usually consists of articles, pronouns, prepositions, adverbs and some verbs in particular. The next step is to apply the **stemming** and **lemmatization** algorithms. The former works by stemming the end or the beginning of the word, taking into account a list of common prefixes and suffixes that can be found in an inflected word. It has to be said that this approach can be successful in some occasions but not always.

Below is an example of how the stemming algorithm works:

studi**es** → -es → studi
study**ing** → -ing → study

On the other hand, the lemmatization algorithm takes into account the morphological analysis of words. An example is shown as in the previous case:

studies → third person, present tense of the verb study → study
studying → gerund of the verb study → study

Once all the pre-processing part of the tweets has been done, we can move on to the text representation part that has been used for this contest. However, before going into this, a **training and test partition** has to be carried out. Specifically, in the training dataset (the one containing 3958 messages) a training and validation partition will be made (around 10-20% depending on the algorithm used).

Once this is done, the remaining 891 tweets from the other set can be used to classify the messages. In addition to performing the corresponding partitions, it has to be observed whether the classes are balanced or not, as this can lead to problems when evaluating the subsequent models. In the variable containing the binary classification, **toxicity**, 2316 tweets can be observed with class 0, i.e. non-toxic, and the rest (1147) with class 1, meaning that they are toxic. According to the criteria considered by the team, it was decided not to carry out any class balancing task. However, the situation varies in the **toxicity_level** variable. This variable contains 2317 tweets with class 0, 808 with class 1, 269 with class 2 and, finally, 69 with the most toxic class. In this case, it was decided to carry out a balancing task while training the models. In the conclusions, some proposals for future work will be mentioned and one of them will be the balancing task before training the models. (the solution adopted will be briefly explained in the part on model evaluation).

2.2. Text representation

Having considered the issue of class balancing, we now move on to explain the text **representation techniques** used. The first of these is the **bag of words** [3]. It is a way of representing the vocabulary that we will use in our models and consists of creating a matrix in which each column is a token and the number of times each token appears in each sentence is counted.

Below is a small example of how this technique works:

We have 3 sentences which are as follows:

This movie is very scary and long
 This movie is not scary and is slow
 This movie is spooky and good

With these 3 sentences we get a vocabulary with all unique words as follows: "This", "movie", "is", "very", "scary", "and", "long", "not", "slow", "spooky", "good".

Applying bag-of-words we would be left with the following matrix:

| | 1 This | 2 movie | 3 is | 4 very | 5 scary | 6 and | 7 long | 8 not | 9 slow | 10 spooky | 11 good | Length of the review(in words) |
|-------------|-----------|------------|---------|-----------|------------|----------|-----------|----------|-----------|--------------|------------|--------------------------------------|
| Review 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 7 |
| Review 2 | 1 | 1 | 2 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 8 |
| Review 3 | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 6 |

Table 1 - Matrix example to represent BOW

The problem with this technique is that it only considers unigrams. For this reason, the **n-grams** technique was also used, as word order could be considered in this way. It was decided to use between 2 and 3-grams (bigrams-trigrams). The procedure would be the same as before but taking into account the latter approach.

The third and last text representation technique used was **Term-Frequency - Inverse Document Frequency** (TF-IDF) [4]. This technique consists of measuring how important a word is within a text, i.e. each word will have an associated weight, depending on its importance. It should be noted that this technique was used for the Stacking model and was applied to the initial text. The previous techniques (BOW and Ngrams) were used for individual models such as Support Vector Machines, Logistic Regression, among others.

2.3. Method

Having explained the techniques used to represent tweets, we now move on to mention the different Machine Learning **models** that were used: Support Vector Machine (SVM), Decision Tree, Logistic Regression, Multi-layer Perceptron (MLP), Random Forest, Stacking and BETO.

Since the Stacking and BETO models have been submitted as runs in the competition, the strategy used in both models will be briefly explained.

It must be said that the **Stacking** model [5] is a combination of some of the previous models, specifically we used Support Vector Machine (SVM), Decision Tree, Random Forest and Multi-layer Perceptron (MLP) as base models and, as a meta learner model, logistic regression. In addition, in some of the previous methods, in order to obtain the ideal parameters for each of them, fine tuning was used, specifically, employing the **Grid search** technique [6].

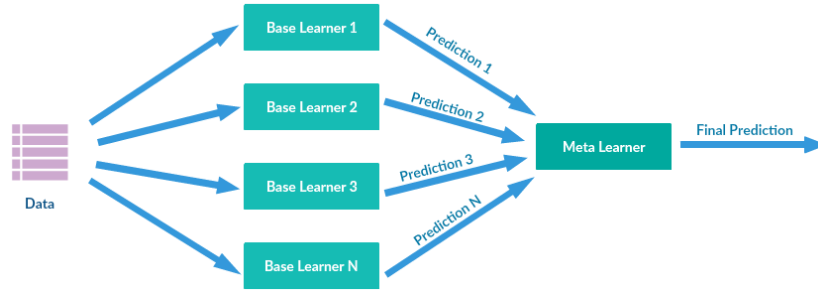


Fig. 1. General structure of the Stacking system

As previously mentioned, the issue of imbalance is a problem to be taken into account when obtaining the labels for the different types of toxicity (toxicity_level). Therefore, for all the previous models, except for the BETO model, an extra parameter called *stratify* of the sklearn library was used at the time of training these models, which allowed the classes to be balanced while performing the training task.

Finally, the **BERT** [7] technique was also used, specifically the `dccuchile/bert-base-spanish-wwm-uncased` model of Hugging Face for tweets in Spanish (BETO)¹. BERT is a Transformer that uses an attention mechanism that learns the contextual relationships between words in a certain text (in this case tweets). Moreover, a Transformer comprises two structures: an encoder that reads the text input and a decoder that produces a task prediction. Since the goal of BERT is to generate a language model, only the encoder mechanism is needed. In this model, a maximum tweet length of 200 characters and a batch size of 16 were used for training.

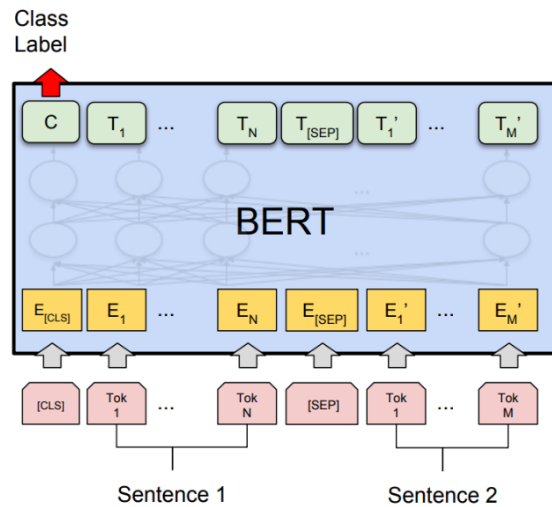


Fig. 2. BERT Model

¹ <https://huggingface.co/dccuchile/bert-base-spanish-wwm-uncased>

3. Results

Once a brief description of the models used has been given, we move on to the results. However, it should first be noted that in order to evaluate the quality of these models, different evaluation measures have been used. In the first subtask, which refers to the detection of whether a message is toxic or not, the **F1** score measure was used. However, in the second subtask, which refers to the detection of a text according to its level of toxicity, more measures are used: **CEM** (Closeness Evaluation Metric), which is used for ordinal ranking tasks, **RBP** (Rank Biased Precision) [8], which is suitable when we are retrieving highly toxic comments from large texts, **Pearson's coefficient** and finally the **accuracy**.

For the DETOXIS competition, due to the fact that only a maximum of 5 runs could be sent, we decided to send the Stacking and BETO models. The results for the first subtask were as follows:

| System | F1-Score |
|----------|----------|
| BETO | 0.4632 |
| Stacking | 0.3893 |

Table 2. Model performance on the testing set (toxicity task)

On the other hand, the results for the second subtask were the following:

| System | CEM | RBP | Pearson | Accuracy |
|----------|--------|--------|---------|----------|
| BETO | 0.6703 | 0.1037 | 0.2677 | 0.6936 |
| Stacking | 0.6258 | 0.0999 | 0.1529 | 0.7160 |

Table 3. Model performance on the testing set (toxicity levels task)

4. Conclusion and Future work

In the DETOXIS [9] **ranking** we have finished in the following positions: **8th** for the first subtask (toxicity) and **11th** for the second one (toxicity_level).

Throughout this project, different approaches could be adopted in order to obtain the best possible results in labelling tweets with their corresponding toxicity values. However, due to lack of time, some of the improvements we had in mind could not be implemented. Therefore, knowing this, the possible improvements of this project are as follows:

1.- Due to the fact that in BETO the accuracy is not entirely good, what could be done is the following: because that class is very unbalanced, around 2000 samples and, in this case, we could predict the level of toxicity of a tweet according to whether the model has previously predicted it as whether that tweet is toxic or not. Therefore, we could first make a prediction on the toxicity variable and, once we have obtained the predictions, we would obtain a new *new_predictions* column and with this we will only work with those tweets that the model has predicted as toxic.

2.- Perform balancing tasks before training the models.

3.- Use more variables such as sarcasm, aggressiveness, etc. and not just the information of the tweet comment itself.

References

1. Basile, V., Bosco, C., Fersini, E., Deborá, N., Patti, V., Pardo, F. M. Rangel, Rosso P. & Sanguinetti, M. (2019). Semeval-2019 task 5: Multilingual detection of hate speech against immigrants and women in Twitter. 13th International Workshop on Semantic Evaluation (pp. 54-63), Association for Computational Linguistics.
- Kumar, R., Ojha, A. K., Malmasi, S., & Zampieri, M. (2018). Benchmarking aggression identification in social media. In Proceedings of the First Workshop on Trolling, Aggression and Cyberbullying (TRAC-2018) (pp. 1-11)
2. Sai Saketh Alurul , Binny Mathew , Punyajoy Saha , and Animesh Mukherjee, «Deep Learning Models for Multilingual Hate», Indian Institute of Technology Kharagpur, 2020
3. Yin Zhang, Rong Jin, Zhi-Hua Zhou, «Understanding Bag-of-Words Model: A Statistical Framework», International Journal of Machine Learning and Cybernetics, 2010
4. Joon-Min Gil, Sang-Woon Kim, Research paper classification systems based on TF-IDF and LDA schemes, Human-centric Computing and Information Sciences volume, 2019
5. Alexandre Alves, Stacking machine learning classifiers to identify Higgs bosons at the LHC, Journal of Instrumentation, 2016
6. Siji George, B.Sumathi, Grid Search Tuning of Hyperparameters in Random Forest Classifier for Customer Feedback Sentiment Prediction, International Journal of Advanced Computer Science and Applications, 2020
7. Jacob Devlin, Ming-Wei Chang, Kenton Lee, Kristina Toutanova, «BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding, arXiv , 2018
8. Moffat, A., & Zobel, J. (2008). Rank-biased precision for measurement of retrieval effectiveness. ACM Transactions on Information Systems (TOIS), 27(1), 1-27.
9. Taulé, Mariona, Alejandro Ariza, Montserrat Nofre, Enrique Amigó, Paolo Rosso (2021). ‘Overview of the DETOXIS Task at IberLEF-2021: DEtection of TOXicity in comments In Spanish’, Procesamiento del Lenguaje Natural, Vol. 67.