

Twitter Sentiment Analysis Classification

Giuseppe Concialdi

Politecnico di Torino

Student id: s294666

giuseppe.concialdi@studenti.polito.it

Christian Montecchiani

Politecnico di Torino

Student id: s303681

christian.montecchiani@studenti.polito.it

Abstract—

I. PROBLEM OVERVIEW

The objective of the competition was to build a model that is able to classify whether a tweet contains positive or negative sentiments. The dataset provided is arranged as follow:

- A **development** set composed by 224,994 records of tweets. Each sample has six different features, including the *sentiment* attribute that is the target of the classification.
- A **evaluation** set consisting of 74,999 samples. Its dimension is one-third of the development set and it does not feature the target variable.

*****move**The dataset is quite large, so the time required to manage the operations on the data should be taken into account because it could be not trivial. It is essential to retrieve some meaningful information of the data by exploiting its features. In particular, every sample is characterized by:

- *ids*: the unique identifier of the tweet. It is represented by a progressive integer number that is related to the timestamp of the tweet. The lowest value of the *ids* attribute is 1,467,811,193 while the highest is 2,329,205,038. It is uncertain if the dataset includes only a subset of the actual number of posted tweets or if the increase of the values follows some pattern. In fact, by digging in the past twitter documentation, we found out that the tweet id is generated with a snowflake schema, invented by Twitter for the generation of sequential ids for their tweets.
- *date*: the timestamp of each tweet. The date is encoded as a string not in the ISO 8601 standard [1] but with the format:

weekday month day hour:min:sec tz year

From this schema, the information can be easily extracted and exploited to train the model. The dates in the dataset range from April 6 to June 25 of 2009. Knowing the temporal ranges will help during the preprocessing phase.

- *flag*: a string whose significance is unsure. It is present in the whole dataset with the unique value of "NO_QUERY". Given its absence of meaning, this feature will be removed without a second thought, but it is possible that it would report the query used to retrieve the tweets when they were extracted using some Twitter APIs.

- *user*: the username of the creator of the tweet. Even though there are almost 225 thousand tweets, there are only 10,647 different usernames. Therefore, on average, the users are very active and they have probably posted several tweets in this period.
- *text*: the text of the tweet. This is the core part of the analysis, it embodies a lot of insights that can be extracted and analyzed to retrieve the overall polarity of the tweet's sentiments. In 2009 the maximum length of a tweet was 140 characters [2]. However 1 shows that some tweets exceed this threshold, so there are likely issues with the encoding of the text extracted.
- *sentiment*: the target variable of the classification. It assumes two possible integer values: 0 and 1 that represent respectively negative and positive sentiments. The dataset is fairly unbalanced, as shown in figure 2, there are more positive sentiments than negative ones.

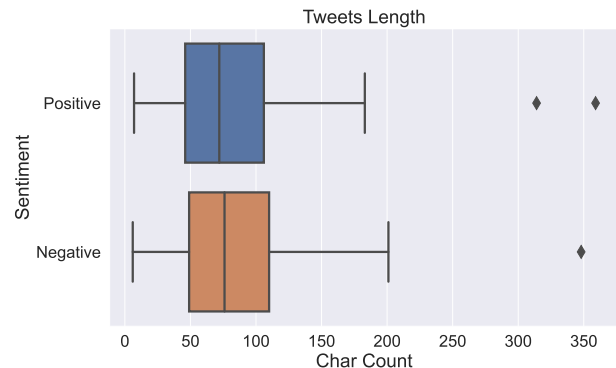


Fig. 1: Boxplot showing the number of chacters of the tweets per sentiment

In the II we will assess what are the most relevant features to carry on the classification and how we will extract the relevant information that lies within each feature.

II. PROPOSED APPROACH

*****Move in overview maybe with a table**The dataset does not feature any missing value, but it contains redundant information that is useless for the analysis. We decided to extract the time-related information from the *date* attribute. In this manner, we added new features to the dataset and also some aggregated measures like the time of the day (morning,

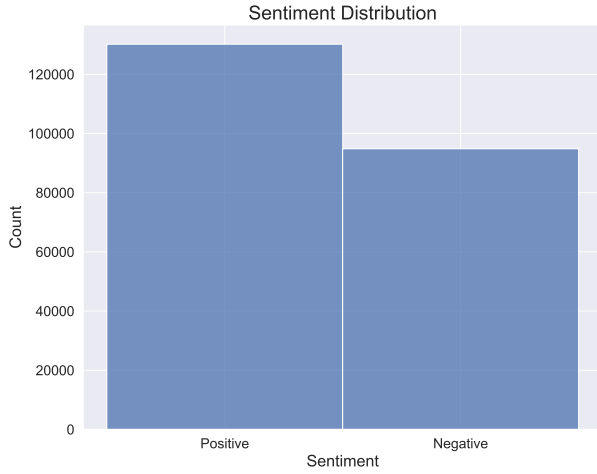


Fig. 2: Distribution of the tweets' sentiments in the training dataset

afternoon, night) and the time of the time (workday, weekend) during which the tweet was posted. At this stage was still unknown whether all these attributes extracted from the *date* variable would have been useful or not. That is because the *ids* feature is not useless, but it already encodes the timestamp of the tweet. Nonetheless, we decided to leave those features and figure out later if their relative importance would matter.

Although the *date* attribute could hide some useful insights about the sentiment of the tweet, the major knowledge lies within the *text* attribute. The information extraction from this feature can be performed in different ways and with different processes. We decided to tackle the problem from different points of view and we managed to embed all this data into our model. Firstly, we cleaned up and fix some problems present in the tweets' text, then we exploited the text performing:

- A **sentiment intensity analysis** [] on the text, obtaining new features on the polarization of the sentiments.
- A **tf-df** [] of the text in order to get the words that mainly influence the sentiments of the tweets.
- A **word embedding** [] approach with the FastText [] library to retrieve the morphological relationships between the words in a sentence.

Figure 3 shows a summary of our approach. The last technique resulted in a very powerful tool, that is able, on its own, to perform the classification of the sentiment of the evaluation set with an f1-score [?] higher than 0.8. This was our baseline for the development of the model, and we accomplished higher performance by integrating the likelihoods of the FastText supervised learning classification into our dataset.

Finally, we wanted to add the *user* attribute into the equation. We thought that due to the relative low cardinality of this attribute (only 10,000 different usernames) the tweets posted by the same author may reflect its personality, thus it is probable that a sentiment is predominant to the other

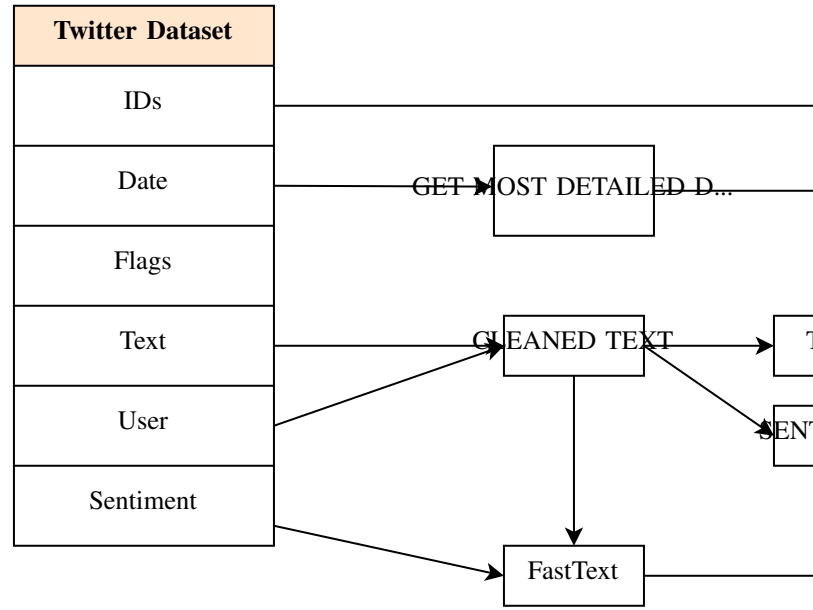


Fig. 3

one. We did not want to encode the information of each user because of the subsequent dimensionality increase of the dataset with an approach like the One-Hot-Encoding [3]. *****to rephrase** So we chose to incorporate the user author of the post at the beginning of the text of the tweet. The text already contains mentioned user whose name is preceded by an at-sign (@). In this way, the word embedding classification will retrieve information about the author of the post like if he was mentioned within it.

Sections II-A and II-B will dig into the details of the three different text mining techniques and some comments about the word embedding can be found in the Discussion section IV of the paper.

A. Preprocessing

B. Model selection

C. Hyperparameters tuning

III. RESULTS

Here you will present your results (models & configurations selected, performance achieved)

IV. DISCUSSION

Any relevant discussion goes here.

REFERENCES

- [1] ISO 8601-1:2019, *Part 1: Basic Rules: Date and time – Representations for information interchange*. ISO, Geneva, Switzerland.
- [2] K. Gligorić, A. Anderson, and R. West, "Adoption of twitter's new length limit: Is 280 the new 140?," 2020.
- [3] D. Harris and S. Harris, *Digital Design and Computer Architecture*. Engineering professional collection, Elsevier Science, 2013.