

Black Lives Matter Sentiment Analysis

Francesco Di Salvo

Politecnico di Torino

s282417

s282414@studenti.polito.it

Abstract—In this report we will discuss a possible approach for a Twitter sentiment analysis. In particular this analysis regards the Black Lives Matter movement. This solution lays the foundation of Natural Language Processing by using some of the most common techniques, such as tokenization, stemming, stopwords elimination and n-grams implementation.

I. PROBLEM OVERVIEW

This sentiment analysis is based on a collection of 80.000 tweets collected from 2020-05-31 to 2020-06-08 about the Black Lives Matter movement. This movement is a decentralized political and social movement advocating for non-violent civil disobedience in protest against incidents of police brutality and all racially motivated violence against black people [1]. It was born in 2013 but it (sadly) becomes well known all over the world after the murder of George Floyd on 2020-05-25.

These tweets are already classified as:

- 1, if they shows a positive sentiment towards the movement
- 0, if they shows a negative sentiment towards the movement

Each tweet is described by 28 features and it is possible to see them on the Tweet Object Documentation [2]. The class label was added for the special purpose of the analysis. For this specific approach we focused just on the label and on the relative text, contained on the *full_text* feature.

This problem is well balanced, in fact we have 80.000 tweets, 40.069 classified as *0* and the remaining 39.931 as *1*. There was a strong activity during this week, in fact we have found in this sample a number of retweets equal to 981.031.288. As you can see in Figure 1, the daily activity increased until 2020-06-03 where it reached its peak and it strongly decreased until the last day to our disposal.

We have also analyzed the most frequent words inside both categories and you can see it in Figure 2. In particular they are both "coherent" with what we could expect from them. In fact for the "positive" tweets we can see words as "petitions", "everyone", "together", while for post against the movement there are words as "white", "division", "behind" and so on.

II. PROPOSED APPROACH

A. Data preprocessing

Natural Language Processing problems have a canonical preprocessing structure, in fact in order to analyze a collection of document we should consider several consecutive steps. In

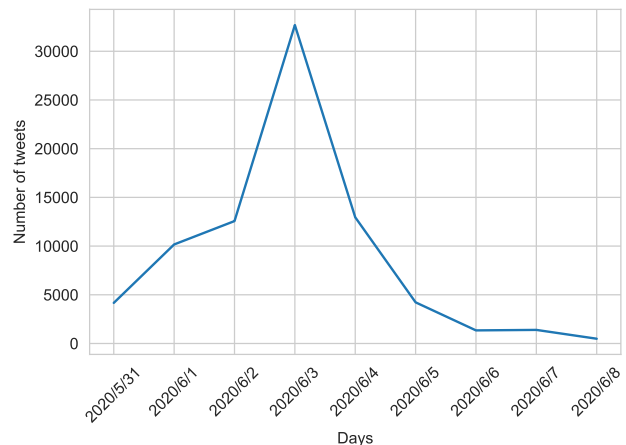


Fig. 1. Tweets per day



Fig. 2. Most frequent words on tweets supporting the cause

particular when we deal with huge documents, we can actually *split* them in shorter pieces in order to better analyze their content. Since this is not the case, we skipped to the second step, that is the *tokenization*: here we split each content into tokens (a.k.a words). These tokens will represent the "core" of our analysis, because some of these tokens will be strictly related to our class.

In particular we used the *Natural Language Tool Kit tokenizer*, that allows to easily convert a sentence like

be useful to focus on new models and their response to the different "combinations" of preprocessing.

REFERENCES

- [1] Wikipedia "Black Lives Matter" Phil. Trans. Roy. Soc. London, vol. A247, pp. 529–551, April 1955.
- [2] Giulio Angiani, Laura Ferrari, Tomaso Fontanini, Paolo Fornacciari, Eleonora Lotti, Federico Magliani, and Stefano Manicardi "A Comparison between Preprocessing Techniques for Sentiment Analysis in Twitter"