

# Twitter Trending Topic Classification

MES: Maria Ludovica Costagliola - Emanuele De Santis - Serena Ferracci

**Abstract**— The Web Information Retrieval course allowed us to develop the described project. It is about classification of topics that are trending on twitter in a given moment. We have dealt with (i) text based classification using single tweets (ii) network based classification exploiting the graph structure of Twitter.

The used dataset is composed by tweets divided by trending topics. We classify the trending topics into 8 categories: *Event, Health, Movie, Music, Politics, Science, Society* and *Sport*.

The first approach gives good result in terms of precision and recall. Instead, the second approach does not give conclusive results because it requires too many resources to be implemented on a single machine.



## 1 INTRODUCTION

We chose this project because Twitter is one of the most popular social network. It is used every day by millions of users expressing their opinions about several fields. When an user looks for something, the first thing Twitter displays to him is the list of trending topics of the moment. Often the user can not know what the topic is about, so it is as to manually search for tweets belonging to that trending topic to better understand it.

It is interesting for the user to have a way to know what the Trending Topic is about without further searches. Our work tries to replicate the results obtained by [3]. The general categories used in this project are 8, named: *Event, Health, Movie, Music, Politics, Science, Society* and *Sport*.

## 2 RELATED WORK

The paper we refer to is the work done by Lee et al. for tweets classification [3]. They have used 18 general categories to classify each trending topic. They address the problem following two different approaches. The first one is a supervised learning technique, named Multinomial Naïve Bayes. Instead, the second one is network based and it uses Personalized PageRank to compute the top-k influencers and then it computes the intersection between the influencers to classify the new trending topic.

All trending topics they used were downloaded from *what the trend*.

## 3 DATASET

Since it was not possible to get a suitable dataset for our purpose, we needed to build it manually using Twitter APIs [1] through Tweepy [2].

In `tweets_retrieve.py` we first retrieved trending topics from USA using USA WOEID (Where On Earth Identifier, that is 23424977 for USA). Then, for each trending topic, we queried Twitter to get all tweets belonging to the given topic. We pay attention to retrieve the entire text of the tweet; in case of retweeted statuses we consider only the retweeted text. Twitter imposes some limitation on GET requests, in particular we could retrieve 180 tweets every 15 minutes. So, we have to handle exceptions fired due to the reached limit. In order to submit a request we have to set up Tweepy authorization handler with our consumer key and our access token, we take from our twitter developer account. As result, we get a file for each tweet where its name is the ID of the author and the body is the text of his tweet. This file is placed in a folder whose name is composed by the name of the trending topic. At the end, we collect about 20.000 tweets belonging to 139 trending topics.

After this first phase, we had to retrieve the link structure for the users that wrote the tweets collected in the previous phase. In order to do so, we developed two python files, named `followers_retrieve` and `friends_retrieve`. We used Tweepy to retrieve all followers and friends for each author ID being aware of all exceptions. Since the limitation, in this case, is 15 queries every 15 minutes, we had to find a way to speed up the process. To do so we decided to parallelize the computation using more than one key (precisely 6 keys), contrarily with respect to the previous retrieving phase, and working in data separation. We associated each key to a subset of IDs to be processed, this mapping was done using a simple hash function. We saved the list of followers and the list of friends for an user  $x$  in the files named  $x$

- followers.txt and x - friends.txt. These two files are saved in the same folder that contains the tweet tweeted by x.

We faced three main problems to retrieve the dataset:

- some users changed their privacy settings between the first and the second phase of our retrieving. For this reason we couldn't get neither their followers nor their friends.
- a subset of users have a plethora of followers, of the order of millions of users. It took a lot of time to complete the retrieving of their followers, sometimes almost an entire day at full capacity. It also happened that the execution stopped due to connection problems or other unpredictable events, so it was not possible to retrieve the entire list of IDs.
- there wasn't a retrievable ground-truth for the trending topics we collected.

The last problem was the only one that we could manage to solve.

We tried to automatically get a classification using some pretrained machine learning algorithms, but when we manually checked the classification done we saw that almost all the predictions was wrong.

The solution, that we came up with, was to label every trending topic by hand. We took one directory at a time, we read some tweets to understand what the topic was about, we assigned to the directory a single category among the 8 defined.

## 4 CLASSIFICATION USING NAIVE BAYES

Multinomial Naive Bayes is a very useful Machine Learning algorithm when we are dealing with text classification. Naive Bayes is based on the assumption that the probability that a document (composed of terms  $(t_1, \dots, t_T)$ ) belongs to a class  $c_k$ , namely  $P(c_k | t_1, \dots, t_T)$ , can be written as  $P(c_k) \cdot \prod_{i=1}^T P(t_i)$ . The probability that a term is in a class is conditionally independent from the probability that another term is in the same class.

So Naive Bayes classifier makes a prediction  $\hat{y}$  using the following formula

$$\hat{y} = \arg \max_{k \in \{1, \dots, K\}} P(c_k) \cdot \prod_{i=1}^T P(t_i | c_k)$$

We rearranged the data in the dataset in order to have a folder that contains 8 subfolders (one for each category we considered). Each one contains all the tweets belonging to a trending topic that belongs to the

category of the subfolder. This operation is done by training\_set.py python script.

We developed also naive\_bayes.py. It takes as dataset the folder yet mentioned.

sklearn allows to make an automatic division of the dataset, splitting it into train set and test set (25% for the test set). Applied a vectorization combined with stemming\_tokenization, we trained the Multinomial Naive Bayes classifier. We, then, evaluate it using the test set and we get the following results:

## 5 NETWORK BASED CLASSIFICATION

### REFERENCES

- [1] Standard search api - twitter developers. <https://developer.twitter.com/en/docs/tweets/search/api-reference/get-search-tweets>.
- [2] Tweepy documentation. <http://docs.tweepy.org/en/v3.5.0/>.
- [3] K. Lee, D. Palsetia, R. Narayanan, M. M. A. Patwary, A. Agrawal, and A. Choudhary. Twitter trending topic classification. *2011 IEEE 11th International Conference on Data Mining Workshops*, 2011.