

TEXT MINING TWITTER: COLOMBIA 2018 ELECTIONS

JUAN SEBASTIAN GARCIA RODRIGUEZ
TRAVIS DUNLOP

Abstract

Between May and June of 2018 the people of Colombia will vote for their next president. As with any modern election, people are using Twitter, the social media platform, to support candidates they like, discredit the others, and debate who should win. Twitter provides a massive open forum to create dialogue across the country. Some groups take advantage of this by creating *social bots*, which automatically post political tweets to in an attempt to sway voters [7]. In our project, we use the text mining skills gained in this course to assess the influence of these bots. We hope to answer two questions: what percentage of users tweeting about the election are bots and what is the sentiment of the tweets?

INTRODUCTION

In order to answer these two questions we split the project into two sections: bot detection and sentiment analysis. For bot detection we leverage an existing algorithm and try both supervised and semi-supervised learning to extrapolate it's findings. With sentiment analysis, we hand-labeled a portion of the tweets and use those to help classify the rest into positive, negative, or neutral.

DATA SET

Since February we collected tweets that were related to the political situation and in specific about the Colombian elections for Senate and President. We collected 3.1 million of tweets from 3th of March to the 7th of May. We got information from 326,966 users and we have identified the following candidates: Gustavo Petro, Humberto de la Calle, Sergio Fajardo, Ivan Duque and German Vargas Lleras. Also, we collected tweets that were talking about politicians who are indirectly involved in the elections such as Alvaro Uribe and Claudia Lopez. We have this data set stored in a SQL data base and we used Python to process the data.

BOT DETECTION

Detecting social bots is a complicated phenomenon, because they usually aim to not be found. This results in a game of cat-and-mouse - researchers attempting more sophisticated classification strategies, and bot-makers emulating increasingly human-like behavior. One technique developed to identify bots is to use a *honeypot*. In this case, the honeypot is a set of researcher-created Twitter users who tweet mostly nonsense. The users that interact with them are likely to be bots - exploiting their desire to engage with many users. Of course, some real users just happen to message or follow one of these bots. In order to parse out which are the real, the researchers use unsupervised learning to cluster the data into groups. They find that some of the clusters seem more real than others. This particular strategy was developed

and implemented by Lee et. al [1]. Once this group of bots are identified, features are extracted to compare them with other users. The researchers use things like frequency of tweets, number of followers, time between posting, ratio of tweets to retweets, even the length of the username as features in the model. Another group from the Indiana University Network Science Institute has made this model available as a API [2]. Unfortunately, the API has limits on access and so, we label some users and then use a model to extrapolate our results to label other users.

Specifically

SENTIMENT ANALYSIS

In political times people tend to express in social media more fiercely their feelings about the current situation in a country. This time Colombians are divided between people who support the peace deal and people who support a reformulation of it. Also there is a lot of uncertainty about what's next for the country since the peace deal is being implemented and Colombia is still living with crime and war. To approach such negative feelings we implemented topic modeling to extract weekly topics in our data set and in parallel we train a model with a semi-supervised technique for learning the sentiment in those tweets. Respectively, we used a Latent Dirichlet Allocation(LDA) implementation from the package gensim and LabelSpreading from sklearn.

For getting the weekly topics we processed the tweets from each user as an entire document and then we gathered this as a weekly collection of documents. Since we noticed that the resulting topics from the model were more difused when we provided each tweet as a document. For the topics modeling we processed only the raw data for the week between April 16 to April 22 which was the one we trained initially the model. Additionally, during the training process we noticed that some words (such as, "gustavo petro", "colombia") that were highly frequent in each document, were present in all the resulting topics so we decided to add them as stopwords. Also, for the sentiment analysis we used a randomized subset of the data for labeling and we manually labeled 1800 tweets. With this data, we ran a supervised classification algorithm and a semisupervised classification algorithm.

Results

[width=]topics_{week}16.svg

CONCLUSION

In the end, we have shown some evidence that

REFERENCES

- [1] Twitter API . <https://developer.twitter.com/en/docs>. Accessed: 2018-07-19.
- [2] Clayton A. Davis, Onur Varol, Emilio Ferrara, Alessandro Flammini, and Filippo Menczer. Botornot: A system to evaluate social bots. *CoRR*, abs/1602.00975, 2016.
- [3] Kyumin Lee, Brian Eoff, and James Caverlee. Seven months with the devils: A long-term study of content polluters on twitter. In *ICWSM*, 2011.

- [4] HANNES MUELLER and CHRISTOPHER RAUH. Reading between the lines: Prediction of political violence using newspaper text. *American Political Science Review*, 112(2):358375, 2018.
- [5] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [6] Radim Řehůřek and Petr Sojka. Software Framework for Topic Modelling with Large Corpora. In *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*, pages 45–50, Valletta, Malta, May 2010. ELRA. <http://is.muni.cz/publication/884893/en>.
- [7] Jon Swaine. Twitter admits far more russian bots posted on election than it had disclosed. *The Guardian*, Jan 2018.