

An Analysis of Political Social Networks and Bot Detection Based on Twitter User Data

Berat Biçer, Vahid Namakshenas

Social media platforms, since their inception, have grown significantly in terms of both their popularity and wide spread usage, and their ability to impact the socio-economic, political, and cultural discourse of countries around the globe. Such an influence undoubtedly attracts agents with nefarious intent, ranging from individuals to political activists and state actors, who seek to obtain a platform to further their political or other related agendas. In particular, recent years suggest social media platforms host a growing number of such agents, which became evident during the 2016 U.S. Presidential Elections via the suspected interference of the Russian state and political consulting firms such as Cambridge Analytica. All these factors necessitate a unified effort to combat such actors seeking to alter the public's perception through various means.

This project, in particular, seeks to study how the political discourse on Twitter is influenced by or is otherwise subjected to bot activity and whether it's possible to identify them through computational methods. Specifically, we're interested in a small political network based on the broad bot-detection corpus called [TwiBot-20](#). Our network is a small subset of the original TwiBot-20 corpus. Subjects are prominent figures from contemporary U.S. politics and their followers, which enables us to construct a social network based on followings/followers relationships. Moreover, a certain number of tweets of these accounts also exist within the dataset. Thus, we are able to study the network dynamics and its real-world applications as well as the tweets of these accounts.

For bot detection, we will primarily develop a machine learning-based pipeline that follows the common NLP recipe which starts by tokenization and is followed by a NN or similar ML-based classifier for bot detection based on an analysis of the tweets of the accounts. However, due to the size of the dataset, this may be impractical and prone to overfitting, which could lead to erroneous conclusions such as overfocusing on artifacts like excessive usage of emojis or Unicode characters, etc. Therefore, we propose complementing this method by studying the network itself. For example, a node which has a high degree is unlikely to be a bot account due to its popularity in the political context. However, such a node would likely be followed by many bot accounts seeking to expand their influence. Studying the actors based on the characteristics of their social circle may complement our predictions, which can be done by including but not limited to clustering, visualization, and overall metric-based high-level analysis.

In terms of the expected challenges, we suspect filtering out the original corpus may be time-consuming and forming a social network based on this connectivity information may be nontrivial. We also expect difficulties in training the machine learning model which could require a pretraining step beforehand, due to limited supply of tweets available. Lastly, we expect some disconnectivity amongst users: Despite, in truth, two nodes may be connected, due to limited connections within the dataset, these two nodes may seem disconnected.