# Media Fakes: Finding Social Media Accounts of General Officers

CS4991 Capstone Report, 2024

Micah High
Computer Science
The University of Virginia
School of Engineering and Applied Science
Charlottesville, Virginia USA
mlh2rtk@virginia.edu

## ABSTRACT

Impersonating celebrity and political officials has been an issue since the inception of social media platforms. At the Army Cyber Institute, I created an algorithm that pulled Twitter and Facebook data related to the names of General Officers and found accounts impersonating them. We utilized the Twitter API and BeautifulSoup along with Selenium for webscraping to pull the data from these sites. With this, we were able to parse the data from both social media platforms and find major properties within each, like name, hometown, occupation, tweet data, image metadata, etc. We created a barebones scoring algorithm that was inaccurate for these datasets. For future work, we would need to find a scoring scheme to incorporate a neural network to determine if the account is real or fake, in order to identify these fake accounts and report them.

## 1. INTRODUCTION

In 2021, between April and June, Facebook reported removing 1.7 billion fake accounts from its site. In this same year, there were hundreds of accounts for one Naval Petty Officer, Mike Spency, whose face and name was used to swindle many people. There were cases where people sent $15,000 to these accounts impersonating Mike Spency.

This issue of digital impersonation is present on all social media platforms, not just Facebook. Twitter dealt with nearly 20 million accounts that were deemed fake and removed. These accounts try to trick people who know the individuals portrayed in these accounts. A newscaster in North Carolina who found a fake account of himself describes what these nefarious actors are after: "[They try to] get personal information…, credit card numbers, and any other data that can be used to steal money"(Sbraccia, 2021).

## 2. RELATED WORKS

The work on algorithms trying to detect algorithms is deep and vast. Twitter has been used as a primary dataset due to its API being able to pull accounts and derive posts and information from public accounts. The methods used to try to detect fake accounts include account-based fake detection, tweet-based detection, graph-based detection, and hybrid-based detection (Oumaima, et. al., 2020). These methods are based on the aspects that the API would pull and use as a scoring scheme base.

Along with creating these schemas for scoring and detecting fake accounts, some AI component has been incorporated into many models. For example, Ekosputra, et. al. (2021) created many AI models to detect fake accounts on Instagram, including linear regression, Bernoulli Naïve Bayes, and random forest. These methods have varying results, but have an accuracy above 0.89. These methods, combined with the techniques

for finding fake Twitter accounts, have the potential to get rid of many imposter profiles; However, the processes described in the studies above, did not take into consideration whether the accounts considered "fake" were trying to impersonate a particular individual. The metrics used only considered if the account was simply not a real person.

## 3. PROJECT DESIGN

This project had two large parts. The first was retrieval of the accounts for the respective social media platforms. This required pulling the name of the accounts, along with any information one can find on their Twitter or Facebook homepage. The second step was the scoring algorithm that decided whether or not the account is real or fake. The initial team on the project created a simple binary scoring system that will be explained later.

### 3.1 Data Acquisition: Twitter

In creating the algorithm, our project lead decided to have us proceed making the Twitter portion before the Facebook portion, due to the data acquisition process being easier. The data we needed to complete the project were profiles associated with the names of the general officers, tweets from that account, and any other public information on that account like occupation, hometown, high school, etc. To acquire this information we utilized the Twitter API to pull searches against the names of the officers on the Twitter database. The Twitter API has since been closed down significantly by throttling the amount of tweets one can now access, limiting the number of searches, and cutting down the amount of information one can pool. When this project was launched, the number of tweets we could pool was in the thousands, due to a license we used.

With this license and API, we looped through the list of names, and ran the names through a library that simplified the Twitter search functions in Python. We acquired all this information on accounts related to each name and stored them in a JSON file, which is a file with key value pairings. In this file was the key of "username" being linked with a specific Twitter account. From here, we looped over this file to pull information from each account individually, and ran the data from each account through the scoring scheme we created.

### 3.2 Scoring Algorithm: Twitter

We pulled tweets associated with each account, occupations, name, and any other general info that could be used to impersonate the general officers. These rules worked on a point system, where if the account has the feature that is present in some fake accounts, points will be added to the score. The rules are as follows: If we already have the real account of the general officer linked to that name, flagged immediately. If there is a username unrelated to the name, +40 points. Spelling typos in the bio: +25 points. Non-military-related links in bio: +15 points. There are many more rules, but this shows how the algorithm works on a basic level. We created functions for each of these point checks, which were easy to implement, except for how different the username is to the name of the individual. The difficult part was deciding on an algorithm to use to check how similar the name of the high profile individual and the names on the Facebook account are to the username. Since our time on the project was limited, our project had us skip this check altogether. The issue with implementing a lot of these checks was in the way the internship was setup. There were two groups of two for this project, and one group got there a week before the other and began the work. The first group that came had no experience with python or coding in general, so slowly teaching them in conjunction with completing our own work on the project slowed down the development process. That

all changed when they left, which is when our project lead had us start on the Facebook portion.

### 3.3    Data Acquisition: Facebook

Due to Facebook's Cambridge Analytica scandal, their API was almost completely shut down. Because of this, we needed a backdoor approach to obtain similar information to what we did with the Twitter accounts. We found out that Facebook's search bar utilized a GET query to obtain its results. This means it manipulates the URL, and the program behind Facebook looks up in their database for profiles close to the search query string.

Using this, we utilized a library in Python called Selenium, which functions as a programmable web browser. We had Selenium send the URL so we could programmatically add the name of the general officers where the text in the search bar would be. From there, we used a JavaScript inject to scroll down on the page to pull up more profiles and pulled the html of the site. We then used BeatifulSoup, a Python library that can pull elements from html. From this, we found the <div id> in the html that had the username for our search results and saved them to a file, then repeated the same process for all the general officers.

The process to find the information for all the profiles looked similar but was slightly more complicated. We tried the same iteration method to get the accounts, but our Facebook got banned, so we needed to create a new one to continue due to bot checkers Facebook had in place. Going forward, we had to have the program pause for five seconds every three profiles due to the account getting flagged and being blocked. This made the program run time skyrocket to nearly twelve hours to pull all the information.

Nonetheless, we proceeded with pulling the html from the profiles themselves. We found the <div id> for each of the different aspects of the profile, like high school, occupation, hometown, bio text, etc. We saved all this info into a JSON, where the structure would be the username as the main value, then each informational aspect, like name, school, etc. is nested within the username in the JSON. This is where the project was concluded due to the time constraints of the internship.

## 4. ANTICIPATED RESULTS

The goal for the end of the internship was to lay the groundwork for an AI algorithm to detect digital imposters on Twitter and Facebook. As shown, due to a lack of programming knowledge, the AI route was abandoned. This is why the scoring algorithm was based on boolean point system: Does the account meet these constraints? If so add points. What we did accomplish, though, is the ability to pull information based on names for both Facebook and Twitter, which can then be used to analyze via a Neural Network, or one of the other methods described in Related Works.

This information-gathering algorithm is being utilized by the project lead in his ongoing research on digital imposters on social media. Since impersonating high-level figures is too common, especially for individuals in the defense sector, his further research is designed to deter this trend.

## 5. CONCLUSION

Social media has given people the special opportunity to impersonate someone within the amount of time required to make an account and find photos. People do this to scam people out of money, spread misinformation, or to slander the individual. The project described above, while not being a particularly strong model to predict whether an account is real or fake, provides

researchers with a method to pull Twitter and Facebook accounts to conduct further research. Currently, researchers have been running tests on data sets of accounts from almost six years ago. From personal experience, imposter accounts now have gotten more sophisticated than a misspelled bio and slapping pictures from a Google search of someone's name.

## 6. FUTURE WORK

The next steps for researchers would be creating a model and testing to determine whether it can detect real accounts from the imposters. Currently, the datasets used are not all the same person, but an amalgamation of accounts where some are real and some are fake. They are not the same person, but just carry the label of real or fake. Some other aspects that can be expanded are more efficient retrieval methods. Currently the Facebook method takes over eight hours to gather all the accounts, so other methods should be looked into.

## REFERENCES

M. J. Ekosputra, A. Susanto, F. Haryanto and D. Suhartono, "Supervised Machine Learning Algorithms to Detect Instagram Fake Accounts," (2021). *4th International Seminar on Research of Information Technology and Intelligent Systems (ISRITI)*, Yogyakarta, Indonesia, pp. 396-400, doi: 10.1109/ISRITI54043.2021.9702833.key words: {Support vector machines;Seminars;Machine learning algorithms;Social networking (online);Multimedia Web sites;Artificial neural networks;Feature extraction;social media;Instagram;fake account;supervised machine learning;user's profile}

Sallah, A., Abdellaoui Alaoui, E. A., Agoujil, S., & Nayyar, A. (2022). Machine Learning Interpretability to Detect Fake Accounts in Instagram. *International Journal of Information Security and Privacy*, *1*, 1–25. https://doi.org/10.4018/ijisp.303665

Sbraccia, S. (2022, January 10). *Be aware of fake social media accounts: More than 1 billion were ousted in 2021 | CBS 17*. CBS17.Com; CBS17.com. https://www.cbs17.com/news/investigators/be-aware-of-fake-social-media-accounts-more-than-1-billion-were-ousted-in-2021/

CBS News. (2021, October 21). *How fake social media profiles are fueling scams and getting people "duped out of money"* CBS News - Breaking News, 24/7 Live Streaming News & Top Stories; https://www.cbsnews.com/news/fake-social-media-accounts/