



The University of Texas at Austin  
**Center for Identity**

# Identifying Real-World Credible Experts in the Financial Domain to Avoid Fake News

*Teng-Chieh Huang  
Razieh Nokhbeh Zaeem  
K. Suzanne Barber*

March 2020

*UTCID Report #2001*

# Identifying Real-World Credible Experts in the Financial Domain to Avoid Fake News

Establishing a solid mechanism for finding credible and trustworthy people in online social networks is an important first step to avoid useless, misleading or even malicious information. Social network users can hide their intention or fabricate their virtual personality to gain trust of others. There is a body of existing work studying trustworthiness of social media users and finding credible sources in specific target domains. However, most of the related work lack the connection between the credibility in the real-world and credibility on the Internet, which makes the formation of social media credibility and trustworthiness incomplete. In this paper, working in the financial domain, we identify attributes that can distinguish credible users on the Internet who are indeed trustworthy experts in the real-world. To ensure objectivity, we gather the list of credible financial experts from real-world financial authorities. By analyzing the distribution of attributes of social media users using the random forest classifier, we can find which attributes are related to real-world expertise, and which attributes have higher potential of being forged by malicious users.

CCS Concepts: • **Networks** → **Social media networks**; • **Information systems** → *Relevance assessment*.

Additional Key Words and Phrases: Source credibility, trust, financial market, fake news

## ACM Reference Format:

. 2019. Identifying Real-World Credible Experts in the Financial Domain to Avoid Fake News. In *Woodstock '18: ACM Symposium on Neural Gaze Detection, June 03–05, 2018, Woodstock, NY*. ACM, New York, NY, USA, 16 pages. <https://doi.org/10.1145/1122445.1122456>

## 1 INTRODUCTION

For most of the human history, trust among people has been judged through the physical world such as in-person communication, letters, or phone calls. The invention of the Internet and the emergence of social network, however, has thoroughly altered the way people communicate. In the blink of an eye, everyone can reach almost every corner of the world through online social networks. The quantity of sources a person can contact has grown enormously, while some of these sources may not deserve one's trust and might deceive him/her. This research seeks to answer an ever-pressing question for those who rely on social media for information – Who can I trust? Our contribution is to focus on the connection between trustworthiness and credibility of a user in the real-world and that user's credibility on the Internet.

By connecting the traditional definition of trust (e.g. reputation, expertise toward certain professional fields and experience related to the target domain) to many personal attributes retrieved from social networks (such as user information and social interconnection) we could establish a dependable way to recognize trustworthy users. This work aims to thoroughly investigate which user attributes could be retrieved as a differentiation tool for the degree of credibility.

To automatically identify social media users who are in fact credible experts in the real-world, we study social media attributes of a set of actual real-world experts. We focus on the financial market as the target domain. To begin with, we referenced some of the most trustworthy and prestigious journalists or media sources. In this

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

*Woodstock '18, June 03–05, 2018, Woodstock, NY*

© 2019 Association for Computing Machinery.

ACM ISBN 978-1-4503-9999-9/18/06...\$15.00

<https://doi.org/10.1145/1122445.1122456>

research, the financial experts are suggested by articles of these trusty sources. Combined with typical users and stock-related users, we are able to compare the differences between each group. We hope using actual financial experts in the real-world and their Twitter feed and the study of their trust filter scores (i.e., measures that evaluate different aspects of user credibility) can help us explain why some trust filters perform better than others. We also aim to understand which trust filters could be more suitable for specific target domains, which in our study is the financial market. We therefore analyze user attributes extracted from social networks and connect them with trust filters. Furthermore, by analyzing the distribution of user attributes among various groups, we achieve a better estimation of the importance of different attributes.

This paper makes the following contribution. We select testing groups based on objective and neutral third-party sources and hence provide an unbiased comparison of the user attributes. Instead of relying on self-defined or machine learning methods to identify credibility, this research provides an alternative way of understanding the differences between the experts and others.

The remaining of this paper is organized as follows. Section II reviews related literature including prediction by analyzing social media data, the history of the trust filter, and the definition of trust and credibility. Section III specifies how we retrieve data and select attributes. Section IV presents our results based on different user attributes versus types of social media users. Finally, section V concludes the paper and suggests some possible avenues for future work.

## 2 RELATED WORK

Given the sheer amount of research performed on social media data, in this section, we seek to review the most closely related work in order to place our work and highlight its differences and contributions.

### 2.1 The Credibility of Social Media Data

Any research or practical application that retrieves data from social media sources has to face a major concern: some social media data might be unreliable, untrustworthy or even malicious. Some of the most significant properties of social media (from the fast spread of information and the low cost and effort of broadcasting news to the ease of hiding one's identity) now become a nest for fake news. Some researchers [1, 24] have studied the spread of fake news on social media during election seasons. Many [6, 26] have applied data mining and machine learning techniques to detect fake news on social media. Some [15] have performed rumor detection on social media based on various data mining algorithms such as decision trees, random forest and SVM. More and more, researchers consider the credibility problem in social media as a major concern [1, 9, 22, 23, 26].

### 2.2 Social Media User Credibility

A body of work [9, 11, 19] exists that particularly looks at the credibility of Twitter *users*. For example, Castillo et al. [5] used a cascade of machine learning models (classifiers) to first find newsworthy and then identify credible tweets. In another work [4], they discovered that credible news are propagated through authors who write a large number of posts, originate at a single or few users in the network, and have many re-tweets. Since surely fake news can spread in a similar manner, they also point out that tweets which do not include URLs tend to be related to non-credible news, those with negative sentiment tend to be more credible, and those with question marks or smiling emoticons are more likely to spread non-credible information. Gupta et al. [10] found that only %17 of the dataset they considered contained credible situational awareness information. They used regression analysis to identify two sets of relevant features, namely content-based features (e.g., the number of unique characters or emoticons in a tweet) and user-based features (e.g., the number of followers or length of username). They suggested the use of those features as credibility scores. Canini et al. [3] performed a similar credibility ranking. Our work takes into account the features suggested by similar papers [4, 12, 21].

### 2.3 Application of Classifiers in Detecting Credible Users

Few researchers [14, 20, 27] have applied machine learning classifiers, such as random forest, to detect credible news sources on social media. Our contribution lies in our way of enlisting credible users and in the observation of the value distribution for features and not in the application of machine learning classifiers.

### 2.4 Trust Filters

The term trust filters was first introduced in [17], where the original idea was to filter and rank trustworthy social media users. There were six trust filters: Authority, Experience, Expertise, Identity, Proximity and Reputation. All six filters stem from the traditional way of estimating one’s trustworthiness. After that, an Internet bio-surveillance application [29] was developed to investigate and detect possible outbreak of diseases in early stages.

Even though there are fundamental differences between user trust and credibility, we use trust to mean credibility in this work. User trust is the subjective expectation of a user of the other; credibility is an objective description of the level of trustworthiness that one possess. Some surveys [2, 25] have covered conceptual differences between trust and credibility. A previous work [13] merges the objective credibility and subjective trust to make them interchangeable in terms of meaning. Furthermore, it applied this concept as a reference to a specific target domain of stock markets. By doing so, it achieved a more accurate prediction of stock prices based on Twitter sentiment analysis by leveraging trust filters.

## 3 METHODOLOGY

This section explains how we extract social media data, decide which attributes to be included and how many experimental groups of users to be used. The social media we utilize is Twitter. Given its property of short posts, we think Twitter would be a great fit for the fast-changing financial market. Moreover, we could access more posts and users with our available computational power and data size, which would be beneficial to this work when analyzing the distribution of social media users. We derive a set of attributes for each user and use the value of those attributes to judge the trustworthiness of users. Here an attribute refers to a user, text or social connection information of a single social media user. We comprehensively reviewed previous work [4, 12, 21] and selected a set of attributes for this research. We analyze the distribution of those attribute values for users and apply data mining techniques, such as random forest to evaluate trustworthiness and the best attributes to quantify it. The final goal is understanding how important any attributes is when establishing an individual’s trustworthiness on social media.

### 3.1 Attribute Selection

We chose 11 attributes (shown as bold in Table 1). These attributes are based on previous research [4, 12, 21] and use tweets content, Twitter user information, or social network structure. While previous research uses these attributes for a whole host of goals, we utilize them to verify information credibility. We chose attributes that are neutral and suitable for universal purpose and no specialized retrieval technique is required to calculate the attribute value for a user. These attributes can be directly applied to various different target domains such as politics or entertainment. Even in the case of *stock\_related\_tweet* we can simply switch the keywords to fit various target domains.

### 3.2 User Group Categorization

We retrieve data from a publicly available collection of tweets for the spritzer version from Internet Archive, which is a non-profit digital library. The spritzer version contains approximately a 1% sample of Twitter public posts. The sampled data is examined and shown to preserve enough information for the research and application

Table 1. Attributes Extracted from Tweets.

Attributes	Definition
<i>n_tweet</i>	Number of tweets collected in dataset
<i>stock_related_tweet</i>	Number of tweets that contain a stock symbol
<i>statuses_count</i>	Total number of posted tweets in user history
<i>followers_count</i>	Number of followers of a user
<i>friends_count</i>	Number of friends (following users) of a user
<i>avg_len_tweet</i>	Average tweet length in characters per tweet
<i>avg_n_word_tweet</i>	Average number of words per tweet
<i>avg_hashtag</i>	Average number of hashtag symbols per tweet
<i>avg_tweet_URL</i>	Average number of tweets that contain a URL
<i>avg_tweet_question</i>	Average number of tweets containing “?”
<i>avg_tweet_exclamation</i>	Average number of tweets containing “!”

based on the tweet or content statistics [28]. Our data sampling time period was from November 2015 to April 2016.

We consider three groups of users:

- (1) Typical users, who represent a random sampling of all Twitter users. This group stands for the baseline among all groups. Typical users serve as the control group and represent how the social media community should look like, so we can distinguish if there exists distinct distribution of attributes in the other groups. Therefore, there should be no condition of selecting users except the sampling time.
- (2) Stock-related users are those who post at least one tweet with a reference to a stock market symbol during the sampling time period. With this group, we aim to represent the users who might be interested in the stock market, while they may not be experts. Therefore, this group can help us differentiate between true financial experts and ordinary people with financial interest.
- (3) Financial expert users are retrieved from well-known online sources (the Business Insider’s article titled “The 129 finance people you have to follow on Twitter” [18], an article from CommodityHq.com titled “100 Insightful Futures Traders Worth Following on Twitter” [8], which is also quoted by NASDAQ.com, and Forbes’ “Must-Follow Twitter Feeds On Markets And Economics” [7]). These lists are compiled by their authors or by Wall Street analysts and journalists, who are traditionally considered financial authorities. Business Insider’s article asked Wall Street analysts and Business Insider journalists to list their must-follow tweeters. CommodityHq.com selected experts based on various criteria: activeness, knowledge, excellent investors, experience, expertise on a particular sector or asset and whether the account is directly run by the experts. For the article in Forbes, the list is suggested by the author himself. Therefore, the expert list is based on multiple aspects and is not just limited to a single source. Since our Twitter sampling time period is in 2015 and 2016, we only referenced articles posted between 2015 and 2016. Originally, there were 257 people/groups mentioned in these articles. However, we were unable to find some accounts, since they might have been deactivated by Twitter or their owner. After manually examining all the accounts, we kept 228 of them.

We sampled data from November 2015 to April 2016. However, given the sheer amount of data, this time period was not suitable for the general users group as it would result in a huge data set we could not effectively examine (682,271,432 tweets posted by 63,873,729 users, which contributed to 2.1 TB amount of data). Hence, the time interval we selected for typical users was only one day, November 2, 2015. For November 2 only, there are still

4,811,258 tweets posted by 2,655,283 users. The number of posts and users from November 2 is approximately of the same magnitude of the other two groups. A caveat is that some attributes would no longer be applicable if we have different time periods for different user types. Take `n_tweet` for example. In order to still use such attributes, we multiply the value of such attributes for typical users from November 2 by 183 (the total number of days from November 2015 to April 2016), in order to roughly put two groups of users with different time intervals into comparison. Nevertheless, the number of tweets a user posted in a single day does not always stay the same; on the contrary, it usually varies a lot. We will cover more about this problem in Section 4.

## 4 RESULTS

This section presents and discusses our results. To better compare data with different sizes of expert groups, we have normalized the charts. The Y-axis shows the percentage of number of users instead of the absolute number throughout this section. Because the time frame for typical users is only one day, normalizing the time axis is also necessary for time-sensitive attributes such as `n_tweet` and `stock_related_tweet`.

### 4.1 Attribute Distribution Analysis

We found from Fig. 1(right) and Fig. 2(right) that the distribution of typical users and stock-related users is almost the same. Experts shows a relatively higher proportion of prolific writers in terms of `n_tweet` and stock-related tweets.

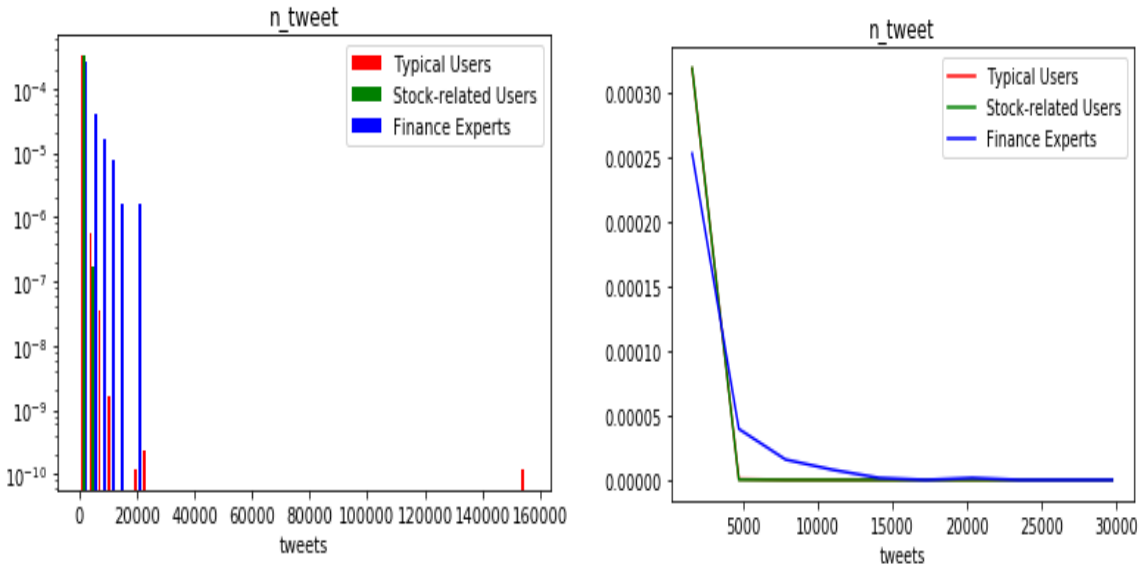


Fig. 1. (Left) Histogram of `n_tweet` for three user groups: (1) Typical Users. (2) Stock-related Users. (3) Financial Experts. (Right) Line chart shows the distribution of tweets under 30,000 using the same datasets. Typical users and stock-related users are almost identical and cannot be differentiated.

Fig. 3 shows the status count, i.e., the number of Tweets a user has ever posted. It shows a similar trend as Fig. 1 and Fig. 2: there is no meaningful difference between typical users and stock related users. One notable thing is that stock-related users surpass typical users for tweets over 160,000. However, a bigger percentage of financial experts have posted between 50,000 and 100,000 times.

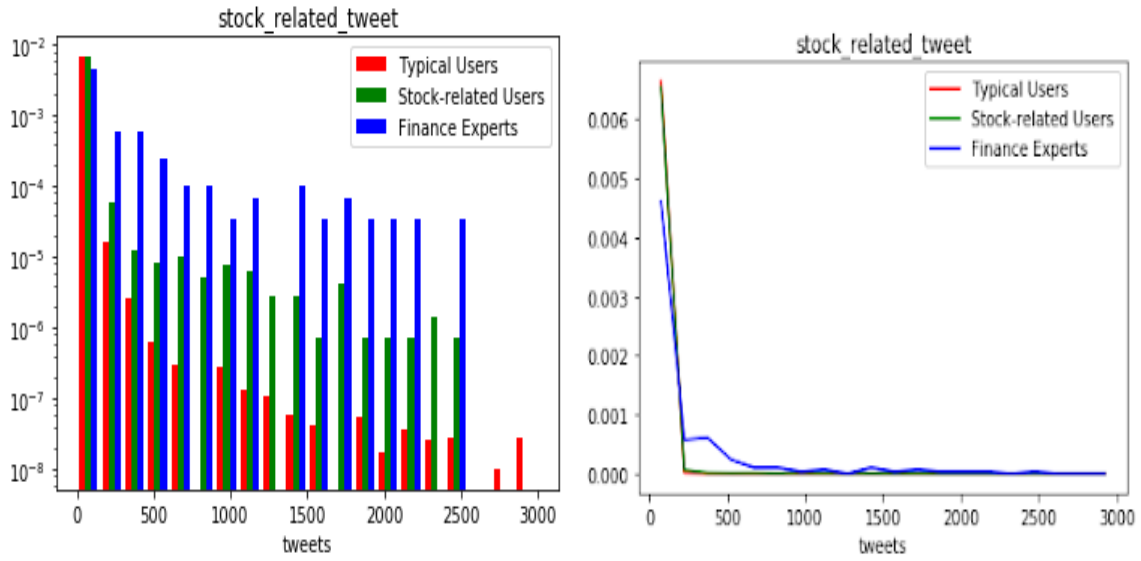


Fig. 2. (Left)Histogram of stock\_related\_tweet. (Right) Line chart shows the distribution of tweets under 3,000 using the same datasets. Typical users and stock-related users are almost identical and cannot be differentiated.

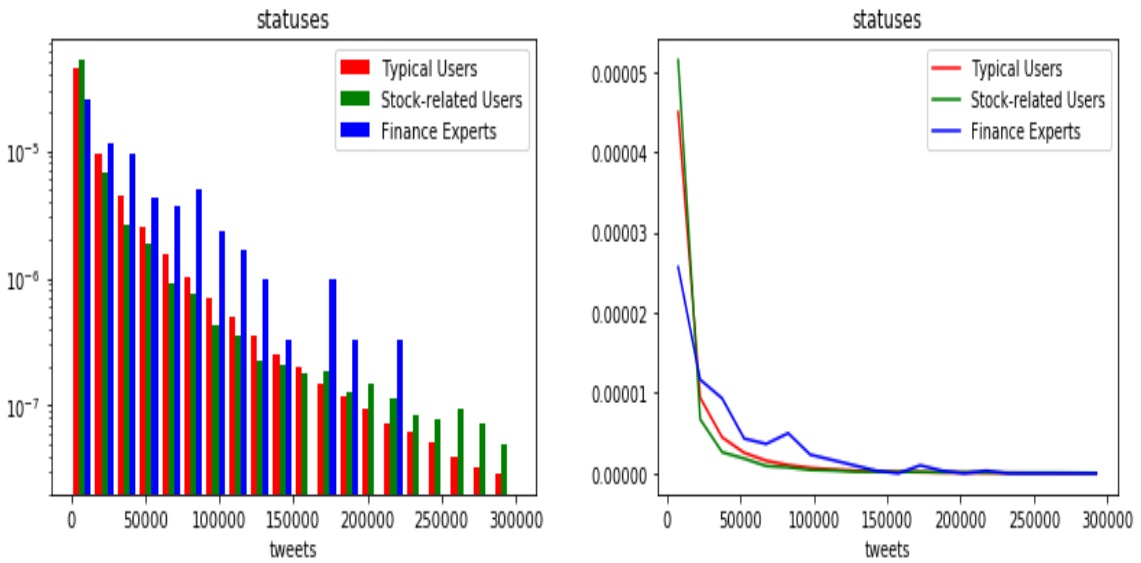


Fig. 3. (Left)Histogram of statuses. (Right)Line chart shows the distribution of tweets under 300,000 using the same datasets.

Fig. 4 and Fig. 5 are followers count and following (friends) count, respectively. The experts in followers count have much higher proportion of large follower users, which surpasses the other two groups by a lot. This trend

sustains while not so overwhelming in the following count. Compared to the experts, the difference between typical users and stock-related users is not so obvious.

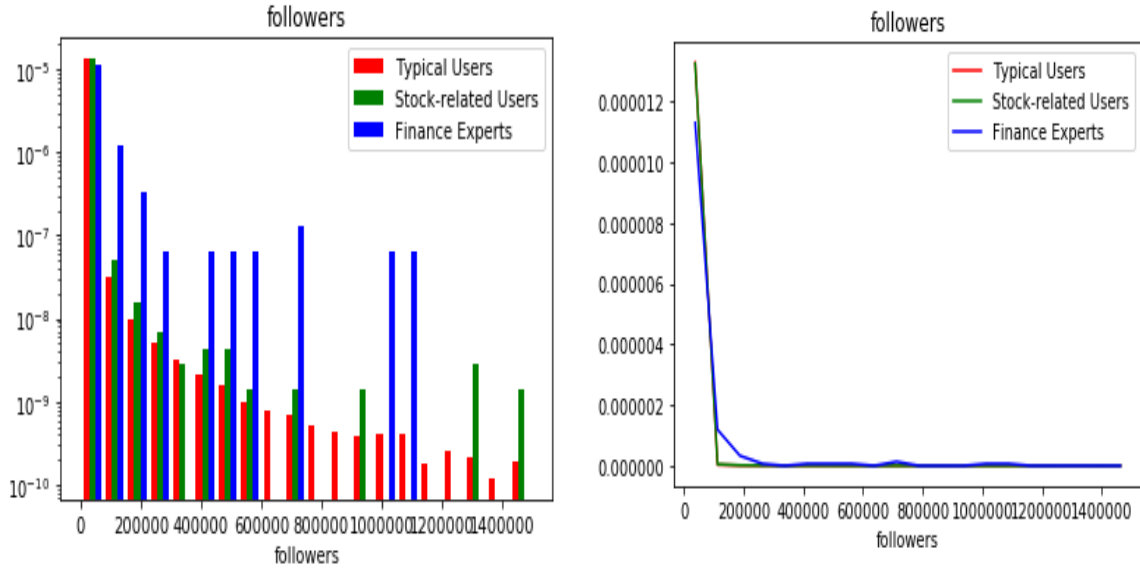


Fig. 4. (Left)Histogram of followers. (Right)Line chart shows the distribution of followers under 1,400,000 using the same datasets.

The distribution of `avg_len_tweet` for stock-related users and experts is similar (Fig. 6), which both have peak at around 110 characters. The typical users, however, have a very different distribution for their average tweet length. Consequently, the average tweet length is an attribute that can distinguish between typical users and experts or people with an interest in the stock market. Similarly, the distribution of `avg_n_word_tweet` in Fig. 7 peaks around 15 words for stock-related users and experts.

For `avg_hashtag` in Fig. 8, there are some peaks for typical users, while the percentage of users with a given average number of hashtags per tweet decreases steadily for stock-related users and experts. Our best guess is there exists many users who only post one or two tweets and make the integer number in the distribution distinct.

As for `avg_tweet_URL` in Fig. 9, stock-related users and experts have a greater proportion of users between 0.1 to 0.9 compared to typical users.

Fig. 10 and Fig. 11 show the number of users with a given average number of exclamation point and question mark per tweet. Stock-related users and experts have similar distributions for exclamation point, and both of their distributions are mostly higher than the one for typical users. The distribution of `avg_tweet_question` in experts is higher than the distribution of stock-related users for values larger than 0, while both of them are much higher than the distribution of typical users. We can see that typical users have a peak when `avg_tweet_URL`, `avg_tweet_exclamation` and `avg_tweet_question` equal 1 ( Fig. 9, Fig. 10 and Fig. 11). This happens because of the higher ratio of typical users composing merely one tweet in our entire sampling. If the one tweet for such users contains a URL, an exclamation point or a question mark, the average number will equal one.



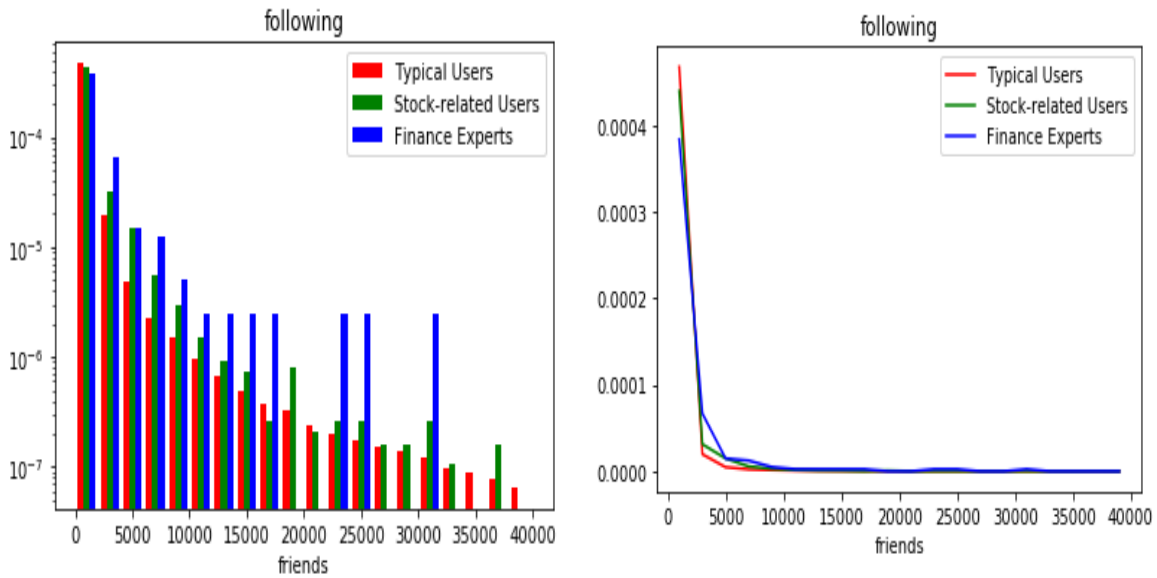


Fig. 5. (Left)Histogram of following users. (Right)Line chart shows the distribution of followings under 40,000 using the same datasets.

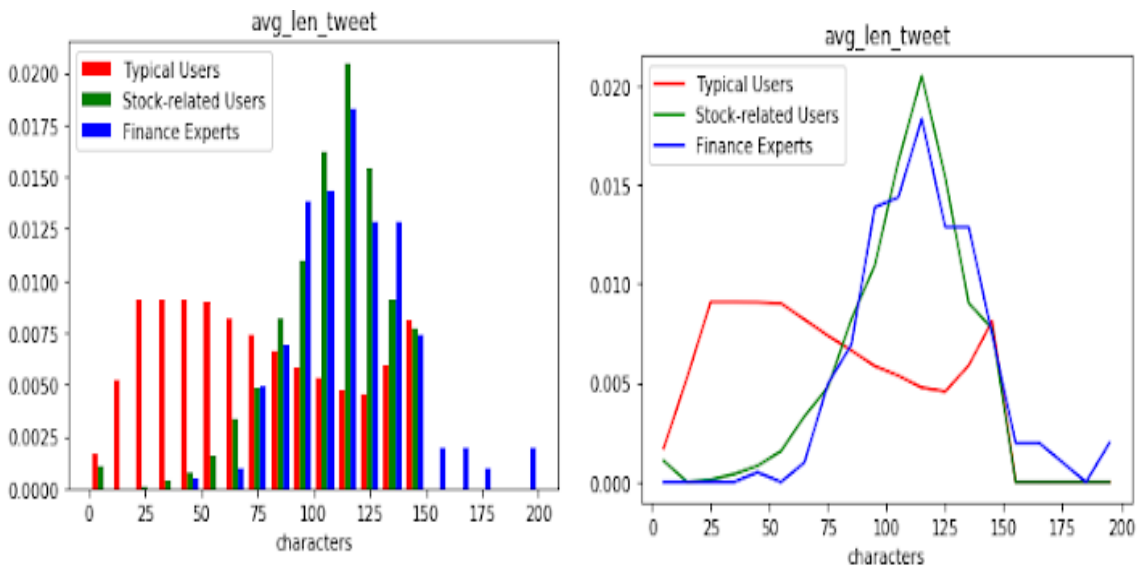


Fig. 6. (Left)Histogram of average length of characters per tweet. (Right)Line chart shows the distribution of characters under 200 using the same datasets.

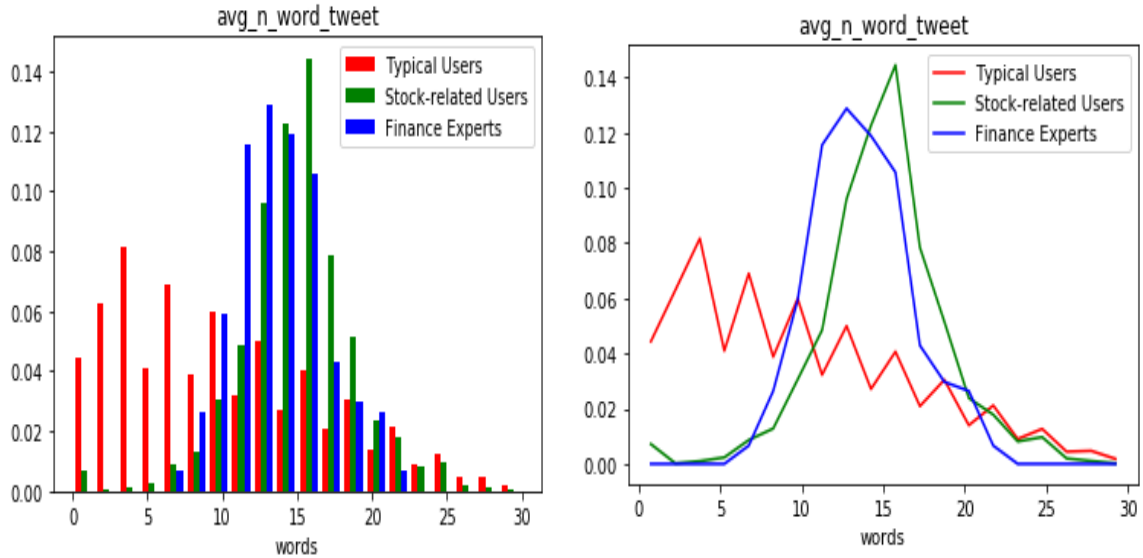


Fig. 7. (Left)Histogram of average length of words per tweet. (Right)Line chart shows the distribution of words under 30 using the same datasets.

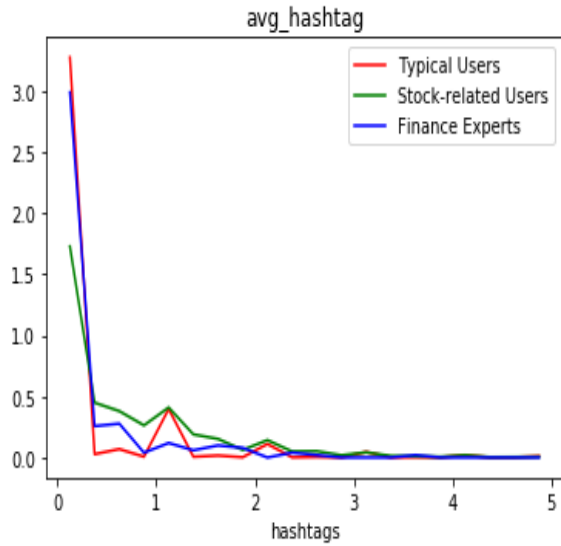


Fig. 8. Line chart of average number of hashtags per tweet.

#### 4.2 Correlation Between Attribute Pairs

Based on the above results, in order to better understand the correlation between attributes and the testing groups, we chose some attribute pairs to investigate the distribution of the Twitter users. Those attribute

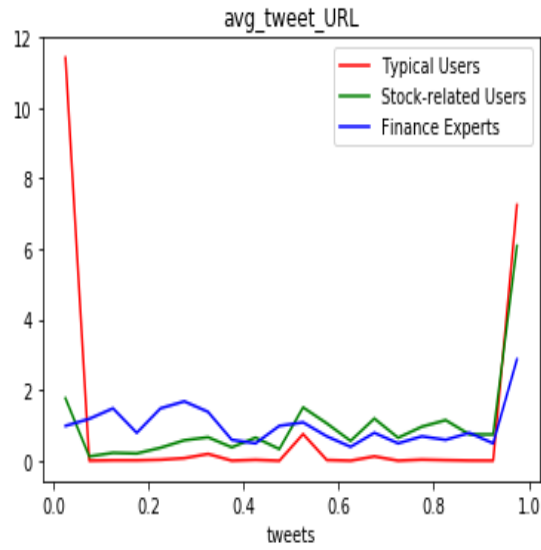


Fig. 9. Line chart of average number of URLs per tweet.

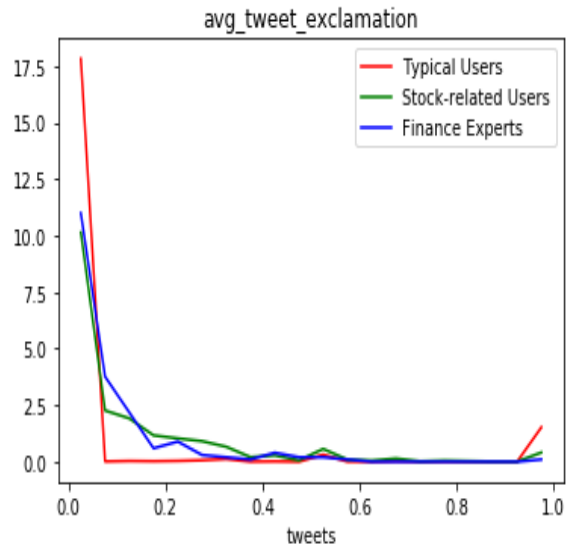


Fig. 10. Line chart of average number of ! per tweet.

pairs are 1) followers\_count vs. friends\_count 2) avg\_len\_tweet vs. avg\_n\_word\_retweet and 3) n\_tweet vs. stock\_related\_tweet.

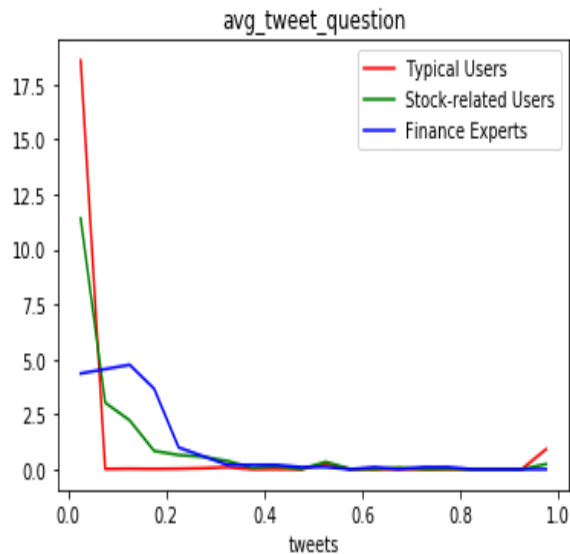


Fig. 11. Line chart of average number of ? per tweet.

Fig. 12 shows the distribution of each group in terms of followers count and friends count. There is a clear boundary in Fig. 12(left) because of two Twitter follower regulations: (1) The 5,000 friends limit: If a user has less than 5,000 followers, the maximum he/she can follow is 5,000 accounts. (2) The 10% limit: If a user follows more than 5,000 users, it would be limited to 10% more than the number of people that follow the user. For example, if 5,000 people follow user X, then X can follow a max of 5,500 people; if 10,000 people follow user Y, then Y can follow up to 11,000 people, and so on. In Fig. 12(right) a set of trendlines for each group is depicted. The slope of the trendline for financial experts is smaller than the other two groups. This new attribute could also be an identifiable characteristic for various users.

When comparing `avg_n_word_tweet` versus `avg_len_tweet` (Fig. 13), the slope of the trendline for finance experts is quite different than the other two which are steeper. This means that the experts tend to use less words for a given length of tweet, i.e., they use longer and more sophisticated words.

Comparing `n_tweet` and `stock_related_tweet` (Fig. 14), the slope of the trendline for finance experts is higher than typical users, but smaller than stock-related users. It is reasonable because the definition of a stock-related user is whoever contributed any stock-related tweet. Considering that, finance experts do tend to contribute much more stock-related tweets than regular users.

### 4.3 Random Forest Analysis

The random forest algorithm [16] has been broadly applied on classification and regression tasks. Random forest is a technique combining decision tree, bagging and random selection of features. One of the benefits of random forest is it can rank each variable according to its importance, which can help us understand how suitable the set of attributes is for the target domain. Random forest also can handle large quantities of data very well since it can be implemented with parallelism, and it can maintain fast prediction and training speed.

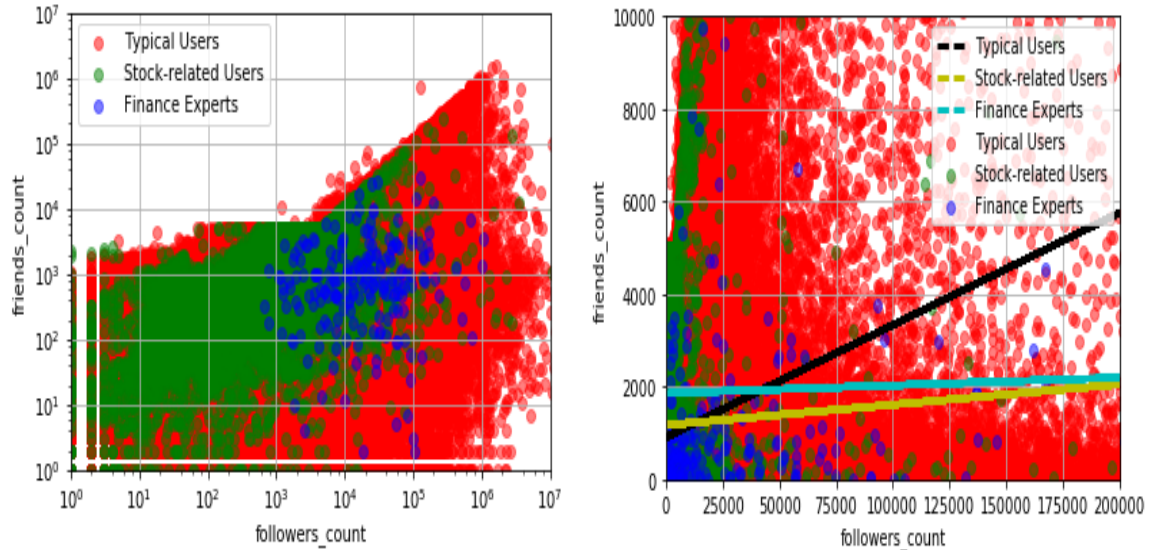


Fig. 12. (Left)The distribution of users in terms of friends\_count and followers\_count. Each dot represents one user and different colors mean different groups. (Right)The distribution set with smaller focus region. There are three trend lines representing the three user groups, respectively. Slope for each trendline: Typical - 0.0243, Stock-related - 0.0044, expert - 0.0017.

Table 2 shows the results using random forest<sup>1</sup> models to distinguish each group: typical users from experts in the second column, stock-related users from experts in the third column, and all three groups from each other in the last column. We consider only time-independent attributes because the sampling time lengths are different among the three groups. Therefore, there are only 9 attributes listed in Table 2. After testing different parameters from 5 to 500 (Fig. 15), we set the number of trees to 50 for the balance of speed and performance. Beside the accuracy and AUC (Area Under the Curve of Receiver Operating Characteristic curve), the feature importance is also shown in Table 2. AUC provides an aggregate measure of performance across all possible classification thresholds. It measures how well predictions are ranked, rather than their absolute values. This helps us understand which attributes are more influential during the comparison process. When comparing between two groups (typical users & experts, stock-related users & experts), the importance values have larger variation (range from 0.0460 to 0.2577 and from 0.0687 to 0.2285) and the top three most important attributes (0.2577, 0.1726, and 0.1351 in the second column, and 0.2285, 0.2027 and 0.1124 in the third column) are totally distinct among the two comparing sets. However, when taking all three groups into account, the importance values become more evenly distributed (ranging from 0.0720 to 0.1664) and there is no particularly overwhelming attribute. We also observe that even the least important attribute still cannot be ignored, which implies the entirety of the set of attributes we selected are mostly effective.

There are three attributes which frequently have lowest importance values: statuses, friends and avg\_n\_word\_tweet. The former two are easier to manipulate since users can post as many tweets as they wish, or follow as many friends as they want, except for the friends limit mentioned above. These kinds of attributes are more unidirectional and do not require feedback from others. If a user wants to improve his or her trustworthiness on

<sup>1</sup>We used the Scikit-learn library.

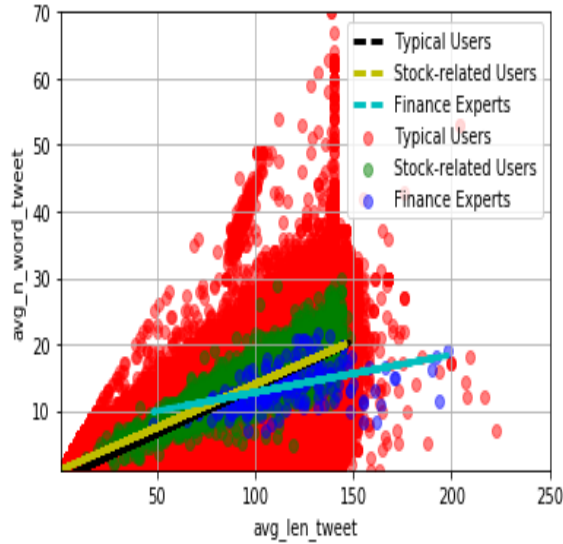


Fig. 13. The distribution of users in terms of avg\_len\_tweet and avg\_n\_word\_tweet. Each dot represents one user and different colors mean different groups. There are three trend lines representing the three user groups, respectively. Slope for each trendline: Typical - 0.14, Stock-related - 0.13, expert - 0.056.

Table 2. Feature importance of each attribute. There are three comparison groups: Typical users & Experts, stock\_related users & experts and all three compared together. Highest/lowest values of importance in each column are shown in bold/underlined respectively. Prediction accuracy and AUC are also show in the last two rows.

	Typical Users & Expert	Stock_Related Users & Experts	All
<i>Statuses</i>	<u>0.0621</u>	<u>0.0687</u>	<u>0.0730</u>
<i>Followers</i>	0.0863	<b>0.2285</b>	0.1208
<i>Friends</i>	0.0823	<u>0.0689</u>	<u>0.0848</u>
<i>avg_len_tweet</i>	0.0869	<b>0.2027</b>	<b>0.1319</b>
<i>avg_n_word_tweet</i>	<u>0.0460</u>	0.0809	<u>0.0720</u>
<i>avg_hashtag</i>	<b>0.1726</b>	0.0785	<b>0.1409</b>
<i>avg_tweet_URL</i>	<u>0.0710</u>	<b>0.1124</b>	0.1071
<i>avg_tweet_question</i>	<b>0.1351</b>	<u>0.0780</u>	0.1029
<i>avg_tweet_exclamation</i>	<b>0.2577</b>	0.0814	<b>0.1664</b>
<b>Accuracy</b>	0.999969	0.994524	0.999945
<b>AUC</b>	0.789474	0.821429	0.631579

the social network, it is also intuitive that the user would post a lot or eagerly extend their social connection. As a result, statuses and friends become relatively non-ideal attributes while classifying user groups. As for avg\_n\_word\_tweet, it might be affected by avg\_len\_tweet which possesses similar meaning among users.

On the other hand, when it comes to attributes with highest importance values, it is surprising the top attribute, avg\_tweet\_exclamation, is related to writing style. Our guess is experts tend to refrain themselves from revealing

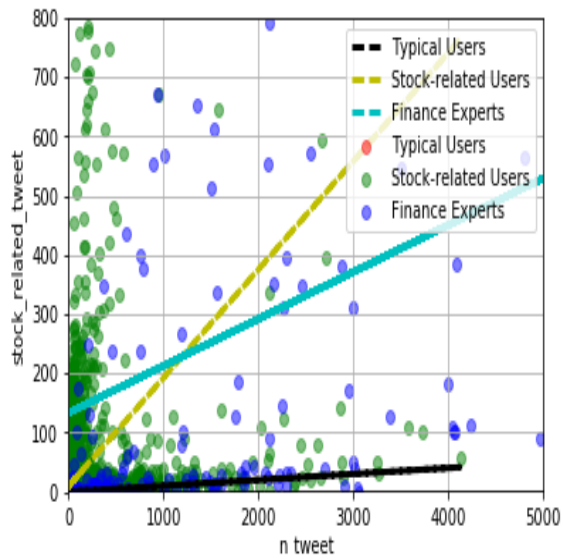


Fig. 14. The distribution of users in terms of `n_tweet` and `stock_related_tweet`. Each dot represents one user and different colors mean different groups. Because the quantity of `n_tweet` and `stock_related_tweet` in typical users are much smaller than stock-related users and experts, the red spots are hardly seen in this figure. However, we can still see there are three trend lines representing the three user groups, respectively. Slope for each trendline: Typical - 0.01, Stock-related - 0.183, expert - 0.079.

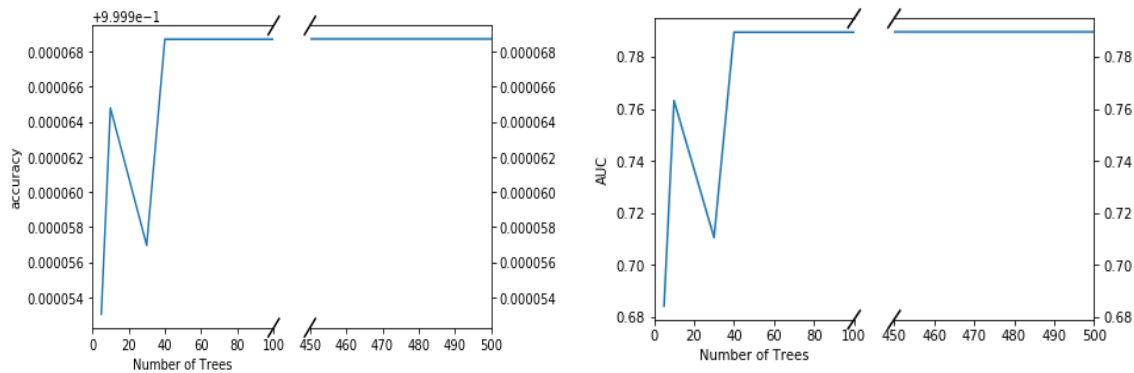


Fig. 15. (Left) Accuracy with different number of trees for random forest algorithm. (Right) AUC with different number of trees for random forest algorithm.

too much expression in the words. Other than that, we notice attributes relevant to text content are mostly within higher importance values. This implies even without much social network connection, we could still distinguish experts from others by simply investigating their writing style.

#### 4.4 Threats to Validity

**Internal Validity:** Because the typical users are selected randomly through the sampling mechanism of the spritzer database, there might exist experts in the typical user group. Our assumption was that the proportion of experts would be negligible in the typical users group. Therefore, the influence of the existence of experts in typical user distribution and in the application of random forest is limited and negligible. To verify this assumption, we split the typical users and reserve 100,000 of 2,655,385 users as “untouched data” for verification. Using the random forest prediction model, we found only 3 of 100,000 typical users are determined as experts. This very small amount of experts would not interfere with the final results. In addition, there are 252 detected experts among the 9,425 stock-related users, which contributes to 2.67% of stock-related users.

**External Validity:** The most important threat to external validity is our identification of experts. To address this threat we selected the list of experts from third party independent publications.

## 5 CONCLUSION

This research presented a distinctly novel view to differentiate credible and trustworthy users on social media. We referenced some of the most prestigious media listings of credible users instead of self-identifying them. We further utilized those listings to study the distribution of user attributes as well as for the training phase of classification methods. We focused on the financial market and 1% of all public tweets in a six-month period from Twitter as the target expertise domain and target social media, respectively. We categorized social media users in three different groups: typical, stock-related and expert users, based on their relevant expertise on the financial market. We considered a set of user attributes based on our extensive review of related work to analyze the attribute distribution of each group. The experimental results show that based on attribute value distribution, some attributes are better indicators of expertise, such as the average length of tweets, the average number of words per tweet, and the number of followers, which display different distribution patterns in experts/stock-related users than typical users. We further applied random forest classification to enhance verification of this observation. We found that the frequency of using exclamation points (!) and hashtags (#) are strong differentiators of experts from typical users and the number of followers and the average tweet length are best indicators for identifying stock-related users from experts using random forest classification. In addition, even among those users with high involvement in certain expertise domains (e.g., financial market in this research), our application of random forest classification can still identify true experts from merely interested users by considering multiple user attributes. Our work paves the way for analysing news source credibility to combat the effect of fake news on social media.

## REFERENCES

- [1] Hunt Allcott and Matthew Gentzkow. 2017. Social media and fake news in the 2016 election. *Journal of economic perspectives* 31, 2 (2017), 211–36.
- [2] Majed Alrubaian, Muhammad Al-Qurishi, Atif Alamri, Mabrook Al-Rakhami, Mohammad Mehedi Hassan, and Giancarlo Fortino. 2018. Credibility in Online Social Networks: A Survey. *IEEE Access* 7 (Dec 2018), 2828–2855. <https://doi.org/10.1109/access.2018.2886314>
- [3] Kevin R Canini, Bongwon Suh, and Peter L Pirolli. 2011. Finding credible information sources in social networks based on content and social structure. In *Privacy, Security, Risk and Trust (PASSAT) and 2011 IEEE Third International Conference on Social Computing (SocialCom), 2011 IEEE Third International Conference on*. IEEE, 1–8.
- [4] Carlos Castillo, Marcelo Mendoza, and Barbara Poblete. 2011. Information Credibility on Twitter. In *Proceedings of the 20th International Conference on World Wide Web (WWW '11)*. ACM, New York, NY, USA, 675–684. <https://doi.org/10.1145/1963405.1963500>
- [5] Carlos Castillo, Marcelo Mendoza, and Barbara Poblete. 2013. Predicting information credibility in time-sensitive social media. *Internet Research* 23, 5 (2013), 560–588.
- [6] Niall J Conroy, Victoria L Rubin, and Yimin Chen. 2015. Automatic deception detection: Methods for finding fake news. *Proceedings of the Association for Information Science and Technology* 52, 1 (2015), 1–4.
- [7] Simon Constable. 2015. “Must-Follow” Twitter Feeds On Markets And Economics. Retrieved November 7, 2019 from <https://www.forbes.com/sites/simonconstable/2015/08/14/must-follow-twitter-feeds-on-markets-and-economics/#46f90dc0b09e>



- [8] Jared Cummins. 2015. *100 Insightful Futures Traders Worth Following on Twitter*. Retrieved November 7, 2019 from <https://commodityhq.com/investor-resources/100-insightful-futures-traders-worth-following-on-twitter/>
- [9] Daniel Gayo-Avello. 2012. "I Wanted to Predict Elections with Twitter and all I got was this Lousy Paper"—A Balanced Survey on Election Prediction using Twitter Data. *arXiv preprint arXiv:1204.6441* (2012).
- [10] Aditi Gupta and Ponnurangam Kumaraguru. 2012. Credibility ranking of tweets during high impact events. In *Proceedings of the 1st workshop on privacy and security in online social media*. ACM, 2.
- [11] Aditi Gupta, Ponnurangam Kumaraguru, Carlos Castillo, and Patrick Meier. 2014. Tweetcred: Real-time credibility assessment of content on twitter. In *International Conference on Social Informatics*. Springer, 228–243.
- [12] Manish Gupta, Peixiang Zhao, and Jiawei Han. 2012. Evaluating event credibility on twitter. In *Proceedings of the 2012 SIAM International Conference on Data Mining*. SIAM, 153–164.
- [13] Teng-Chieh Huang, Razieh Nokhbeh Zaeem, and K. Suzanne Barber. 2019. It Is an Equal Failing to Trust Everybody and to Trust Nobody: Stock Price Prediction Using Trust Filters and Enhanced User Sentiment on Twitter. *ACM Trans. Internet Technol.* 19, 4, Article 48 (Sept. 2019), 20 pages. <https://doi.org/10.1145/3338855>
- [14] Rizki Kurniati and Dwi H Widyantoro. 2017. Identification of Twitter User Credibility Using Machine Learning. In *2017 5th International Conference on Instrumentation, Communications, Information Technology, and Biomedical Engineering (ICICI-BME)*. IEEE, 282–286.
- [15] Sejeong Kwon, Meeyoung Cha, Kyomin Jung, Wei Chen, and Yajun Wang. 2013. Prominent features of rumor propagation in online social media. In *2013 IEEE 13th International Conference on Data Mining*. IEEE, 1103–1108.
- [16] Andy Liaw, Matthew Wiener, et al. 2002. Classification and regression by randomForest. *R news* 2, 3 (2002), 18–22.
- [17] G. Lin, R. N. Zaeem, H. Sun, and K. S. Barber. 2017. Trust filter for disease surveillance: Identity. In *2017 Intelligent Systems Conference (IntelliSys)*. 1059–1066. <https://doi.org/10.1109/IntelliSys.2017.8324259>
- [18] Linette Lopez and Lucinda Shen. 2015. *The 129 finance people you have to follow on Twitter*. Retrieved November 7, 2019 from <https://www.businessinsider.com/117-finance-people-to-follow-on-twitter-2014-9>
- [19] Marcelo Mendoza, Barbara Poblete, and Carlos Castillo. 2010. Twitter Under Crisis: Can we trust what we RT?. In *Proceedings of the first workshop on social media analytics*. ACM, 71–79.
- [20] Srijith Ravikumar, Kartik Talamadupula, Raju Balakrishnan, and Subbarao Kambhampati. 2013. Raprop: Ranking tweets by exploiting the tweet/user/web ecosystem and inter-tweet agreement. In *Proceedings of the 22nd ACM international conference on Information & Knowledge Management*. ACM, 2345–2350.
- [21] Eduardo J. Ruiz, Vagelis Hristidis, Carlos Castillo, Aristides Gionis, and Alejandro Jaimes. 2012. Correlating Financial Time Series with Micro-blogging Activity. In *Proceedings of the Fifth ACM International Conference on Web Search and Data Mining (WSDM '12)*. ACM, New York, NY, USA, 513–522. <https://doi.org/10.1145/2124295.2124358>
- [22] Derek Ruths and Jürgen Pfeffer. 2014. Social media for large studies of behavior. *Science* 346, 6213 (2014), 1063–1064.
- [23] Harald Schoen, Daniel Gayo-Avello, Panagiotis Takis Metaxas, Eni Mustafaraj, Markus Strohmaier, and Peter Gloor. 2013. The power of prediction with social media. *Internet Research* 23, 5 (2013), 528–543.
- [24] Chengcheng Shao, Giovanni Luca Ciampaglia, Onur Varol, Alessandro Flammini, and Filippo Menczer. 2017. The spread of fake news by social bots. *arXiv preprint arXiv:1707.07592* (2017), 96–104.
- [25] Wanita Sherchan, Surya Nepal, and Cecile Paris. 2013. A Survey of Trust in Social Networks. *ACM Comput. Surv.* 45, 4, Article 47 (Aug. 2013), 33 pages. <https://doi.org/10.1145/2501654.2501661>
- [26] Kai Shu, Amy Sliva, Suhang Wang, Jiliang Tang, and Huan Liu. 2017. Fake news detection on social media: A data mining perspective. *ACM SIGKDD Explorations Newsletter* 19, 1 (2017), 22–36.
- [27] Kai Shu, Suhang Wang, and Huan Liu. 2019. Beyond news contents: The role of social context for fake news detection. In *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining*. ACM, 312–320.
- [28] Yazhe Wang, Jamie Callan, and Baihua Zheng. 2015. Should We Use the Sample? Analyzing Datasets Sampled from Twitter's Stream API. *ACM Trans. Web* 9, 3, Article 13 (June 2015), 23 pages. <https://doi.org/10.1145/2746366>
- [29] Razieh Nokhbeh Zaeem, David Liau, and K Suzanne Barber. 2018. Predicting Disease Outbreaks Using Social Media: Finding Trustworthy Users. In *Proceedings of the Future Technologies Conference*. Springer, 369–384. [https://doi.org/10.1007/978-3-030-02686-8\\_29](https://doi.org/10.1007/978-3-030-02686-8_29)

