



# A Feature Based Approach on Behavior Analysis of the Users on Twitter: A Case Study of AusOpen Tennis Championship

Niloufar Shoeibi<sup>1</sup>(✉), Alberto Martín Mateos<sup>1</sup>, Alberto Rivas Camacho<sup>1</sup>, and Juan M. Corchado<sup>1,2</sup>

<sup>1</sup> Bisite Research Group, University of Salamanca, Salamanca, Spain  
{Niloufar.shoeibi, alberto\_martin}@usal.es

<sup>2</sup> Air Institute, IoT Digital Innovation Hub, Salamanca, Spain

**Abstract.** Due to the advancement of technology, and the promotion of smart-phones, using social media got more and more popular. Nowadays, it has become an undeniable part of people's lives. So, they will create a flow of information by the content they share every single moment. Analyzing this information helps us to have a better understanding of users, their needs, their tendencies and classify them into different groups based on their behavior. These behaviors are various and due to some extracted features, it is possible to categorize the users into different categories. In this paper, we are going to focus on Twitter users and the AusOpen Tennis championship event as a case study. We define the attributions describing each class and then extract data and identify features that are more correlated to each type of user and then label user type based on the reasoning model. The results contain 4 groups of users; Verified accounts, Influencers, Regular profiles, and Fake profiles.

**Keywords:** Social media analytics · Behavior analysis · User behavior mining · Feature extraction · Twitter · Verified · Influencers · Regular and fakes

## 1 Introduction to Social Media

Social networks have consolidated as a source of communication and transmission of the information at a global level over the last few years. An infinite number of topics can be dealt with, so there is a huge amount of information spreading by users.

In some cases, the objective may be to detect certain profiles, for example, with behaviors that induce unethical thinking or activities, such as sexist ideologies [1]. Other times, the aim may be to detect relevant or current issues in real-time [2]. Because of these mentioned reasons and many more, the importance of behavior analysis on social networks is undeniable.

## 1.1 Behavior Analysis

Each creature has a behavior towards the environment and others. By analysis of the behavior, it is possible to discover the patterns of thinking in different situations. However, social media can be considered as an environment for human beings so they can express themselves through their interactions. The analysis of the data extracted from the user's behavior is an important block in this type of case study knowing well the environment in which the events occur to extract information and knowledge [3]. Besides, different studies try to identify possible types of profiles in the networks: bots, fakes, and so. [4] is an example, but since there is no robust concept about what characteristics each profile has, it is complicated to reach a unique solution.

## 1.2 Data Extraction on Social Media

Each social network has different casuistry, and these different types of profiles that have been discussed behave differently. Thus, it is important to make a good analysis; for instance, finding similar users, topics, or arguments in different social networks [4]. The challenge is unstructured data extracted from different platforms [5]. Therefore, data mining techniques will play a decisive role.

In this paper, we are going to focus on Twitter, which is a platform that allows two-way communication in which any user can interact with another quickly and easily, so we don't have the problem of different sources. Besides, messages can be spread through "tweets", related to any daily aspect and in which only users are identified as verified or not verified. To go one step further, each account must be analyzed in more detail. Therefore, in this research, we are going to detect different types of profiles according to their behavior on the Twitter social network, with the hope that in the future, we will be able to address this study in a more detailed way and detect possible profiles with different purposes.

This paper organized as follows: Sect. 2, related work. Section 3, pertains to data and proposed method which describes the data extraction and feature selection and the technique that has been used for labelling the users. In Sect. 4, the results will be more explained. In Sect. 5, the conclusion of the results and future work has been presented. And finally, Sect. 6 is indicated to the references.

## 2 Related Work

Most of the research in this area is related to taking advantage of social network data, either through machine learning algorithms (supervised and unsupervised), deep learning, graph theory, etc.

Rashidi et al. studied opportunities and challenges for exploring the capacity of modelling travel behavior. They used the data extracted from social networks to obtain information based on features like trips, their purpose, mode of transport, duration, etc. through surveys. However, the processing time is very slow [6].

Large organizations try to influence choices in social networks, which is going to cause a lack of freedom of expression. In this article, Subrahmanian et al. developed an algorithm for detecting bots based on tweet parameters, profiles, and environment [7].

Erşahin et al. create a twitter fake account detection based on supervised discretization techniques by the reason for the increase of the exposure of incorrect information through fakes profiles has increased [8].

A method that tries to group different profiles according to influential words with the complexity of semantics was developed by Sundararaman et al. [9].

One study has the main objective to characterize the behavior of cancer patients. In this article, Crannell et al. match different types of cancer-patient with several sentiments [10].

NLP-based word embedding grouping method for publishing health surveillance was published by Dai et al. This method is tested versus other bags of words methods [11]. Lastly, Kaneko et al. presented a method based on using keyword burst and image clustering instead of only text analysis for event photo mining [12].

### 3 Data Extraction from Twitter API

The analysis of the data is an important block in this type of case study, knowing well the environment in which the events occur to extract information and knowledge.

Tables 1 and 2 show the list of variables considered from the tweets of each of the users and about the information of their accounts.

**Table 1.** Features obtained about the tweets of each Twitter profile.

Features (tweets)	Definition
Text	Tweet text
Favorites	Number of favorites that have a tweet
Retweet	Number of retweets that have a tweet
Created at	Date of the publication of a tweet

**Table 2.** Features obtained about each Twitter profile.

Features (user profile)	Definition
Name	Name the user, as they have defined it
Screen name	Name of the twitter account, that it is unique
Listed count	Number of public lists that a user is a member
Biography profile	Biography profile text
Followings	Number of followings that have an account
Followers	Number of followers that have an account
Favourites count	Number of favorite tweets that have an account
Statuses count	Number of tweets (RT + own tweets) that have an account
Creation_at	Date of the creation of an account

From all the default variables allowed by the twitter API (Tables 1 and 2) and the possible characteristics that can identify different types of user profiles, the variables of Table 3 have been added to the study and, consequently, those of Table 4 that are directly related.

**Table 3.** New variables from the data analysis (metadata).

New features	Definition
Number of own Tweets	Number of tweets published by a user in the last week
Number of Retweets	Number of tweets retweeted by a user in the last week
Number of own Tweets	Number of tweets (published and retweeted) by a user in the last week
Favorited Tweets count	Number of favorites that have a tweet published
Retweet & Tweets count	Number of retweets that have a tweet published
Mentions	Number of mentions inside tweets published by a user in the last week
Tweets URL	Number of tweets with a URL in it text inside tweets published by a user in the last week
Time between Tweets	Minutes between tweet publications
Twitter years	Number of years that have an account

**Table 4.** Rates from new variables (metadata).

Rates	Definition
Tweets per Retweets rate	The ratio of tweets published per tweets retweeted in the last week
Tweets year ratio	The ratio of tweets (published and retweeted) per year
Time between Tweets	Minutes between tweets (mean) in the last week
Followings per Follower	The ratio of followings per followers' rate

## 4 Feature Extraction Using Graph Theory and Analysis

In this paper, we proposed analyzing the content of the tweet to find the *mentioned users* or the *retweet* as a means of defining the relationship between the users. After extracting 5000 number of tweets about our case study subject, “AusOpen”, which is the Australian tennis championship, we explore the content, and if the tweet is a *retweet*, our proposed method extracts the screen names of the profiles whose tweets have been retweeted by others. If not, we will extract the profiles that have been tagged in the tweet. Then define the graph of relations from the source, which is the user who is doing retweets to the target which his/her content has been retweeted. After building the graph, we do the graph analysis and measure new features extracted from the graph and defined relationship. These features are represented in Table 5.

**Table 5.** Features Extracted from Graph Analysis [13, 14].

Rates	Definition
Eccentricity	The maximum shortest distance of one node from others. The less Eccentricity, the more influencing the power of the node
Clustering Coefficient Centrality	Which nodes in a network are tending to be in the same cluster based on the degree of the nodes. $cc = \frac{n}{t}$
Closeness Centrality	Indicates how close a node is to other nodes in a network by capturing the average distance based on one vertex to another. $cl = \frac{1}{\sum_{v \neq u} d(u,v)}$
Betweenness Centrality	Shows how influential is the node. The more the value of betweenness centrality is, the more important that node would be to the shortest paths through the network. So, if that node is removed, so many connections will be lost. $b = \sum_{s \neq v \neq t} \frac{\delta_{st}(u)}{\delta_{st}}$
Harmonic Closeness Centrality	This measure is so similar to closeness centrality, but it can be used in networks that are not connected to each other. It means that, when two nodes are not connected, the distance will be infinity, and Harmonic Closeness is able to handle infinity just by replacing the average distance between the nodes with the harmonic mean
In-Degree Centrality	This centrality indicates the importance via the number of edges entering the node
Out-Degree Centrality	This centrality indicates the importance via the number of edges going out of the node
Degree Centrality	This measures how many connections a node has. In other words, it's the summation of the In-Degree and Out-Degree of the node and shows how important a node is, by the number of connections. $Deg(v) = InDeg(v) + OutDeg(v)$

n: No of connection between neighbors of a particular node  
t: Total number of possible connections among all the neighbors of the node  
d(u, v): the geodesic distance between u, v.  
s: source, t: destination  
st: number of shortest paths between (s, t)  
st(u): number of shortest paths between (s, t) that pass-through u.

## 5 Different User Groups Base on Behaviors

Based on the behavior of users, we have defined four different categories; “Verified Accounts”, “Influencer Users”, “Regular Users” and “Fake Accounts”.

- **“Verified Accounts”**, which are the ones that have been verified by Twitter, and they have the verified blue mark beside their names on twitter panel. We call these users “Verified” people who are politicians, famous artists, TV shows, special sports events, etc. It is possible to detect this class by checking the “Verified” term in the JSON file extracted from the tweet object.
- **“Influencer Users”** are the users that have a great influence on other users, and others believe in the content they are sharing feelings, thoughts, or expressions. However, they are not verified by twitter. It isn’t very easy to detect this category of users because they are regular users, but sometimes, they have the same behavior as the famous accounts, like verified accounts. To detect this category, in addition to the features appearing in the JSON file extracted from Twitter API, we need to go deeper and analyze the behavior using graph analysis, extracting more complicated features and analyze the content. Betweenness Centrality, In-Degree Centrality, and many more. An Influencer account should have too many retweets, so we can say that In-Degree is a feature that defines how important that specific node is. Then more value of in-Degree shows the more times the tweet of that profile has been retweeted. And after applying this filter, so many of the nodes have been dropped because they didn’t have retweets.
- **“Regular Users”** are the ones that express their thoughts, opinion, and feelings, and they are not aiming to make an influence on other people’s thoughts and opinions. They are not having too many followers and followings, and their interaction ratio is not high.
- **“Fake Accounts”** which are the accounts that usually have an irregular behavior with their content, like fake news, spam, incoherent tweets, etc. So, they usually aim to put false influence or to change statistics in society. This type of profile can be defined with several different features, and it is hard to tell them apart from regular profiles and make a strong definition of it. For example, some of them can retweet so much, others can have a default image profiles, and others can have both features. From this study, we created new variables that can help as a good filter to identify these profiles, like the number of tweets (published and retweeted) per day, based on the features that appear in the JSON file [15, 16].

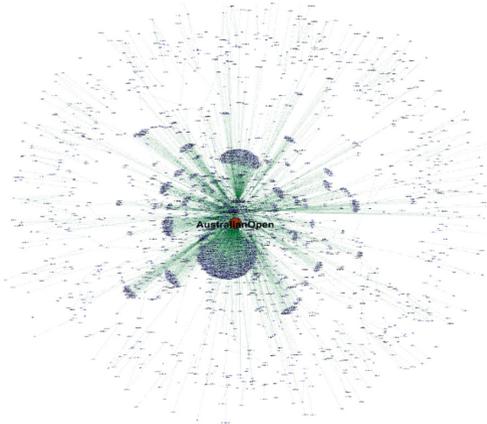
In the Table 6, we can see a summary of the characteristics of each of the profiles:

**Table 6.** Different categories of users based on their behavior.

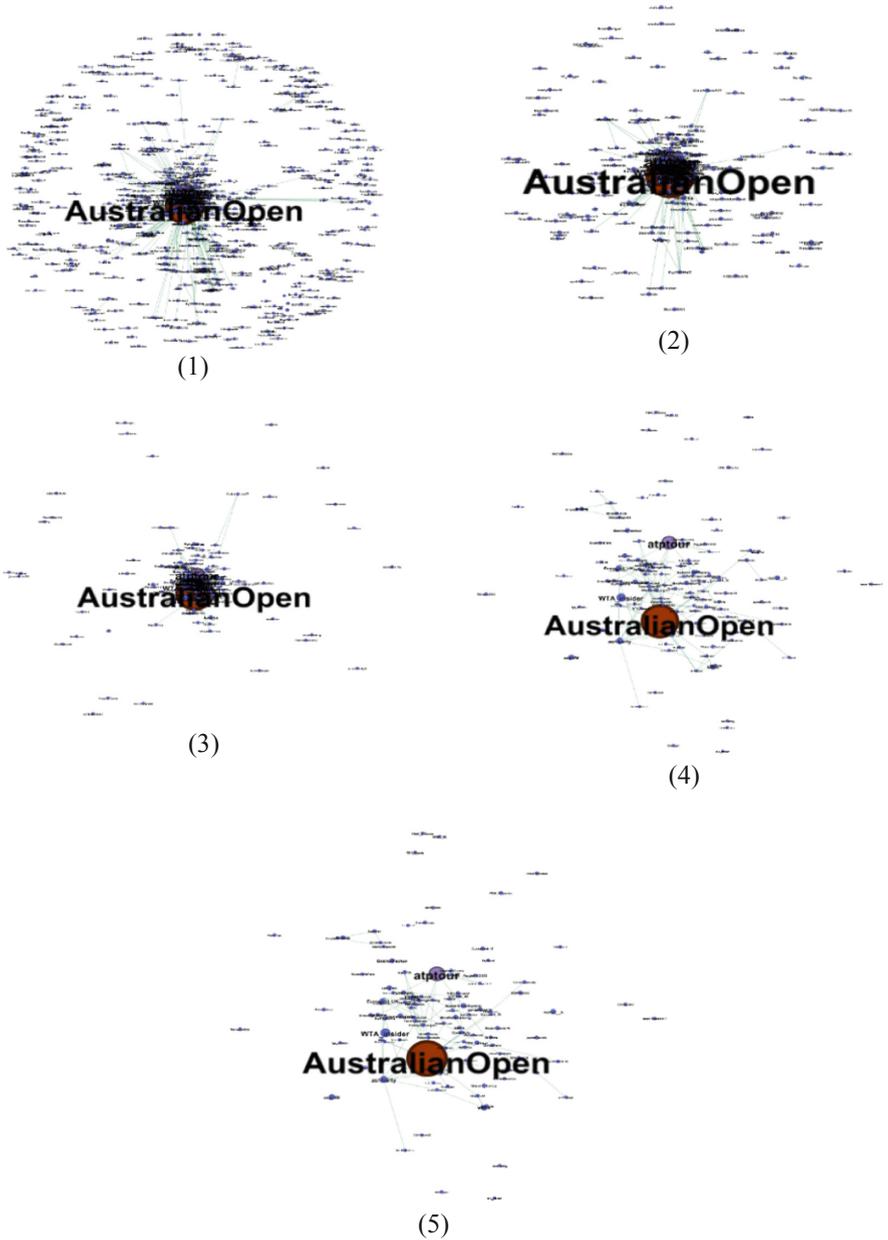
Profile type	Characteristics	Related features
Verified	- Celebrities - Famous Sport Players - Politics	- <i>Verified</i> = True
Influencers	People who are not verified by twitter but their content influence other thoughts	- A high number of their tweets - A high number of followers - The low time between tweets - The high number of interactions with their own tweets - High in-degree centrality
Regulars	People who are no verified by twitter, publish a few contents, favorites, with a balanced number of followers and followings, not in large numbers	
Fakes	People who are not verified by twitter, but their contents are fake news, spam, incoherent tweets, etc.	- The high number of Retweets - The low time between tweets - Default image profile - No biography - Numbers on its account name - The small number of followers - 2001 followings - Tweets duplicated - Self-Loop - High outdegree number - Low indegree number

## 6 Results

With the application of several filters based on the behavior analysis, we can label the data. We apply two different filters to our dataset based on features from the graph analysis and then, from the metadata. Graph analysis is leading us to find the filters

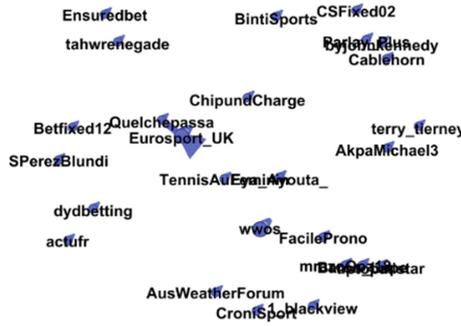


**Fig. 1.** The graph of relationship of people who tagged each other or retweeted other’s tweets in the concept of AusOpen Tennis Championship



**Fig. 2.** Graph revolution of finding *Influencers* and *Verified Profiles* after applying in-degree centrality filters in order to remove *Fakes* and *Regulars*.

helping us to have more accurate labeling. Figure 1 is representing the graph of the relationship between all the users who tweeted about the AusOpen tennis championship. As can be understood, there are nodes (users) with high in-degree centrality, which possibly are the verified profiles and influencers. There are nodes with less in-degree and higher out-degree or self-loops, which probably are Fakes. In Fig. 2. we demonstrated filtering the data based on the in-degree, and after filtering five times, the remained group of users are more probable to be verified profiles and influencers. Figure 3 was created by applying the self-loop filter on the whole network.



**Fig. 3.** The group of probable Fake Profiles after applying self-loop filter

On the other hand, it is necessary to analyze the metadata based on the feature, which has been explained in Sect. 3. After doing the analysis, we detected 415 Fakes, 49 Influencers, and 2266 regular accounts, such shows in Table 7. It is an essential and useful approach to build a robust dataset in an almost chaotic environment. Besides, it can help us to discover new patterns or outliers in the different types of profiles and various subjects.

**Table 7.** Profile types classification

Profile type	Graph analysis	Metadata analysis	Total
Verified	179	0	179
Regulars	9	2257	2266
Fakes	26	389	415
Influencers	17	32	49
Others	2678	0	0
			2909

## 7 Conclusion and Future Work

In this paper, we defined the graph of relationship based on the interaction of users called “retweet”. This relationship, which is from the person (source) who retweeted another person’s tweet (target), has been demonstrated in a directed graph in which the nodes are users. The edges are which are arrows from the source to target, showing retweets. By using graph analysis, we generated more features then added them to the original ones that have been extracted from Twitter and each profile of the user. The results show that adding these new features helped the behavior analysis of the users.

In the future, we are looking forward to trying other social media like Facebook, Instagram, and so many more and making the various machine learning models for training them with the dataset generated by the method this paper proposed in this paper and then measure the type of the account.

## References

1. Waseem, Z., Hovy, D.: Hateful symbols or hateful people? Predictive features for hate speech detection on twitter. In: Proceedings of the NAACL Student Research Workshop, pp. 88–93, June 2016
2. Xie, W., Zhu, F., Jiang, J., Lim, E.P., Wang, K.: TopicSketch: real-time bursty topic detection from twitter. *IEEE Trans. Knowl. Data Eng.* **28**(8), 2216–2229 (2016)
3. Oh, C., Roumani, Y., Nwankpa, J.K., Hu, H.F.: Beyond likes and tweets: consumer engagement behavior and movie box office in social media. *Inf. Manag.* **54**(1), 25–37 (2017)
4. El, A., Azab, A.M.I., Mahmoud, M.A., Hefny, H.: Fake account detection in twitter based on minimum weighted feature set. *World Acad. Sci. Eng. Technol. Int. J. Comput. Inf. Eng.* **10**(1) (2016)
5. Injadat, M., Salo, F., Nassif, A.B.: Data mining techniques in social media: a survey. *Neurocomputing* **214**, 654–670 (2016)
6. Rashidi, T.H., Abbasi, A., Maghrebi, M., Hasan, S., Waller, T.S.: Exploring the capacity of social media data for modelling travel behaviour: opportunities and challenges. *Transp. Res. Part C: Emerg. Technol.* **75**, 197–211 (2017)
7. Subrahmanian, V.S., Azaria, A., Durst, S., Kagan, V., Galstyan, A., Lerman, K., Menczer, F.: The DARPA Twitter bot challenge. *Computer* **49**(6), 38–46 (2016)
8. Erşahin, B., Aktaş, Ö., Kılınc, D., Akyol, C.: Twitter fake account detection. In: 2017 International Conference on Computer Science and Engineering (UBMK), pp. 388–392. IEEE, October 2017
9. Sundararaman, D., Srinivasan, S.: Twigraph: discovering and visualizing influential words between Twitter profiles. In: International Conference on Social Informatics, pp. 329–346. Springer, Cham, September 2017
10. Crannell, W.C., Clark, E., Jones, C., James, T.A., Moore, J.: A pattern-matched Twitter analysis of US cancer-patient sentiments. *J. Surg. Res.* **206**(2), 536–542 (2016)
11. Dai, X., Bikdash, M., Meyer, B.: From social media to public health surveillance: word embedding based clustering method for twitter classification. In: SoutheastCon 2017, pp. 1–7. IEEE, March 2017
12. Kaneko, T., Yanai, K.: Event photo mining from twitter using keyword bursts and image clustering. *Neurocomputing* **172**, 143–158 (2016)
13. Perez, C., Germon, R.: Graph creation and analysis for linking actors: application to social data. In: Automating Open Source Intelligence, pp. 103–129. Syngress (2016)

14. Deverashetti, M., Pradhan, S.K.: Identification of topologies by using harmonic centrality in huge social networks. In: 2018 3rd International Conference on Communication and Electronics Systems (ICCES), pp. 443–448. IEEE, October 2018
15. Bovet, A., Makse, H.A.: Influence of fake news on Twitter during the 2016 US presidential election. *Nat. Commun.* **10**(1), 1–14 (2019)
16. Gurajala, S., White, J.S., Hudson, B., Matthews, J.N.: Fake Twitter accounts: profile characteristics obtained using an activity-based pattern detection approach. In: Proceedings of the 2015 International Conference on Social Media & Society, pp. 1–7, July 2015