CrossMark

# Influential users in Twitter: detection and evolution analysis

**Giambattista Amati[1] · Simone Angelini[1] · Giorgio Gambosi[2] · Gianluca Rossi[2] ·
Paola Vocca[3]** (ID)

## Abstract

In this paper, we study how to detect the most influential users in the microblogging
social network platform Twitter and their evolution over time. To this aim, we consider the
*Dynamic Retweet Graph (DRG)* proposed in Amati et al. (2016) and partially analyzed in
Amati et al. (IADIS Int J Comput Sci Inform Syst, 11(2) 2016), Amati et al. (2016). The
model of the evolution of the Twitter social network is based here on the retweet relation-
ship. In a DRGs, the last time a tweet has been retweeted we delete all the edges representing
this tweet. In this way we model the decay of tweet life in the social platform. To detect the
influential users, we consider the central nodes in the network with respect to the following
centrality measures: *degree, closeness, betweenness and PageRank-centrality*. These mea-
sures have been widely studied in the static case and we analyze them on the sequence of
DRG temporal graphs with special regard to the distribution of the 75% most central nodes.
We derive the following results: (a) in all cases, applying the closeness measure results into
many nodes with high centrality, so it is useless to detect influential users; (b) for all other
measures, almost all nodes have null or very low centrality and (c) the number of vertices
with significant centrality are often the same; (d) the above observations hold also for the
cumulative retweet graph and, (e) central nodes in the sequence of DRG temporal graphs
have high centrality in cumulative graph.

**Keywords** Graph analysis · Social media · Twitter graph · Retweet graph ·
Graph dynamics · Centrality

## 1 Introduction

One of the fundamental and most studied features in a social network is the detection of
central nodes, which can usually be considered as the *most important* nodes [7, 8, 13].

✉ Paola Vocca
vocca@unitus.it

Extended author information available on the last page of the article.

Centrality is widely-used for measuring the relative importance of nodes within a graph and it has many applications: in social networks to determine the most influential or well-connected people; in the Web graph to rank pages in a search; in a terrorist network, to detect agents that are critical for facilitating the transmission of information; for the dissemination of information in P2P Networks, Decentralized Online Social Networks and Friend-to-Friend Networks [11].

There is a plethora of centrality definitions: degree centrality [18], closeness centrality [5], graph centrality [15], stress centrality [19], betweenness centrality [12], each one of them useful to detect specific properties and with significantly different computational costs. Here we consider four of them: the *degree*, *closeness*, *betweenness*, and *PageRank*-centrality.

Degree centrality, i.e. the degree $d_v$ of a vertex $v$, is the simplest measure of centrality: it just takes into account how many direct, "one hop" connections each node has to other nodes of the network, hence it can be applied to detect popular individuals, agents who are likely to hold most information or individuals who can quickly connect with the wider network. The degree centrality of a node is very cheap to compute but, being a purely local notion, it is often unable to recognize the relevance of certain nodes. In this paper with the term *degree* we refer to *in-degree* centrality.

One of the most popular measures, even if computationally expensive for large graphs, is betweenness-centrality. It helps to detect nodes which act as "bridges" between other nodes in a network. It does this by computing the shortest path for each pairs of nodes and counting, for each node, the number of such paths which include it. Betweenness centrality is suitable for finding vertices who influence flows (such as information flow) in the network.

A third measure considered below is the closeness-centrality, which assigns to each node a score that is proportional to the reciprocal of the sum of all distances between the node and all other nodes. This definition of centrality is useful for quickly finding the agents who are in good position to influence the entire network. However, in a highly connected network most nodes often have a similar score.

Finally, PageRank-centrality was introduced in [9] and it recursively quantifies a "value", the PageRank, of a node based on: (i) the number of links it receives, (ii) the link propensity of the linkers (that is, the number of outgoing links of each in-going node), and (iii) the centrality of the linkers, that is their PageRank.

To study how influential users evolve over the time, we analyze the distribution of the centrality measures on an temporal evolutionary model of the Twitter network, the *Dynamic Retweet Graph (DRG)* proposed in [3] and partially analyzed in [2, 4].

This model has two major features: (i) it is based on the retweet graph, since that allows to better represent relationships among users and the information flow in Twitter [16, 17] and, (ii) once a tweet has been retweeted for the last time all the edges representing retweets of this tweet are deleted, to model the loss of relevance of the tweet content.

The temporal model we consider coincides with other temporal models in the growing phase [6, 14], that is a new vertex is added when a new user first acts, either sending a tweet or retweeting an existing one, and a new directed edge $(a, b)$ is inserted when an user $a$ retweets for the first time a tweet of $b$, if an edge already exists then a timestamp is added to it. Conversely, as the decreasing stage happens when a tweet is no more retweeted, all the vertices and the edges, not involved in other retweeting processes, are deleted at once. As shown in previous experimentations [2, 4], this evolutionary model better captures the information flow in Twitter. DRGs seem to better represent the double nature of the Twitter platform: social network and news media [16, 17].

For what concerns the use of centrality measures to assess influential or authoritative users, Kwak et al. [16] compared three measures of influence: in-degree centrality, PageRank centrality in the following/follower network and the number of retweets on Twitter. In Cha et al. [10] three different measures of influence are compared: in-degree centrality, the number of retweets and mentions on Twitter. The results indicate that users with high in-degree were not necessarily influential.

In this paper, we study the evolution of the most influential users in the microblogging social network platform Twitter with respect to four centrality measures (betweenness, degree, closeness, and PageRank) and we analyze their behavior on the DRG evolutionary model of the retweet social networks proposed in [3].

We consider two different kind of data sets, first introduced in [1] and updated and refined in [3]: the *event driven* retweet graphs based on the events *Black Friday 2015* and the *World Series 2015* and the *Italian Sampling* that is the *firehose* retweet graph, filtered by language (i.e. Italian) from the whole Twitter stream.

The four centrality measures are analyzed in three different frameworks: (i) with respect to the sequence of DRG temporal graphs; (ii) with respect to the static cumulative graph, that is the graph that contains all the nodes and edges and (iii) with respect to the kind of networks considered, that is *event driven* or the firehose.

We derive that the DRGs allow to detect the most authoritative users, since:

1. in all cases the closeness centrality provides too many central nodes, hence it is useless to detect influential users;
2. with regard the other measures, almost all nodes have null or very low centrality;
3. vertices with centrality values above 75% of the maximum is a small set and they are often repeated in the three centrality measures;
4. the above observations hold also for the static graphs (the cumulative DRG);
5. central nodes in the sequence of DRG temporal graphs have high centrality in static graphs.

## 2 DRG temporal graphs

In this paper we will use a definition of Dynamic Retweet Graph (DRG) slightly different from the one in [4].

A DRG graph $G = (V, E, \ell)$ is defined as follows: nodes in $V$ are Twitter accounts and a direct edge $e \in E$ represents an interaction (a retweet) between two accounts. In particular, there is a directed edge from an account $a$ towards an account $b$, if $a$ has retweeted at least one tweet of $b$, that can be itself already a retweet. Observe that user $a$ may retweet several tweets of $b$. Information on such retweets is provided, for each edge $e = (a, b)$, by a list $l(e)$ of pairs $(i, t)$ where $i$ is the id of a tweet and $t$ is the time when a retweeted $i$ from $b$. Each list $l(e)$ is ordered in non-decreasing order with respect to the timestamp $t$.

For each tweet $i$, we define the *date of death* of $i$ (in short, $\mathtt{dod}(i)$) as the timestamp of the last retweet of $i$, that is

$$\mathtt{dod}(i) = \max_{e \in E}\{t : (i, t) \in \ell(e)\}.$$

Consequently, we define the *expiration date* of an edge $e$ (in short, $\mathtt{ed}(e)$) as the timestamp from which all tweets associated to $e$ will be dead. Formally,

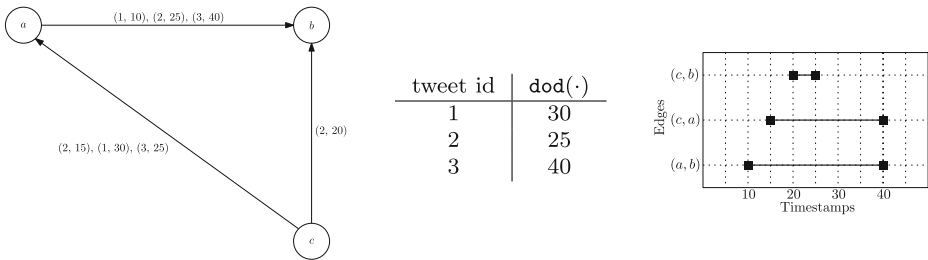$$\mathtt{ed}(e) = \max\{\mathtt{dod}(i) : (i, t) \in \ell(e)\}.$$

**Fig. 1** On the left side, an example of a DRG retweet graph. Edges are labeled by pairs with the id of the tweet and the timestamp of the retweet. The center table shows the date of death of all tweets in the graph. On the right side, for each edge of $G$ is represented its creation and expiration date

On the contrary, the *creation date* of an edge $e = (a, b)$ (in short, $\mathtt{cd}(e)$) is the timestamp $b$ retweets $a$ for the first time, formally:

$$\mathtt{cd}(e) = \min\{t : (i, t) \in \ell(e)\}.$$

We define a *DRG temporal graph* at time $t$ the subgraph $G_t = (V_t, E_t)$ of the DRG $G$ at time $t$ as follows: $E_t$ contains any edge $e$ such that $\mathtt{cd}(e) \leq t \leq \mathtt{ed}(e)$; $V_t$ is the set of nodes induced by $E_t$.

For example if $G$ is the retweet graph represented in the left part of Fig. 1, $G_{30}$ contains edges $(a, b)$ and $(c, a)$ and the induced vertices since $(c, b)$ expires at timestamp 25. For all $20 \leq t \leq 25$, $G_t$ contains all edges of $G$.

## 3 Data sets

For the experiments we use the dataset of [3] that consists in two different classes of retweet graphs: the event driven retweet graph, filtered by topics about specific events (i.e. the Black Friday 2015 and the World Series 2015) and the Italian Sampling retweet graph, filtered by the Italian language from the whole Twitter stream. To obtain the Italian Twitter Sampling we use a list of the most used Italian stop words and the Twitter native selection function for languages. In Table 1, we show the dimensions of the three graphs.

In Fig. 2 we present the evolution of the dimensions of the three datasets over the period of observation. Note that the event-driven datasets (World Series and Black Friday) show a rapid growth close to the events, and then a slow decline. Differently, the Italian Sampling have a smooth and stable behavior, ignoring the border effects.

**Table 1** Dimensions of the dataset

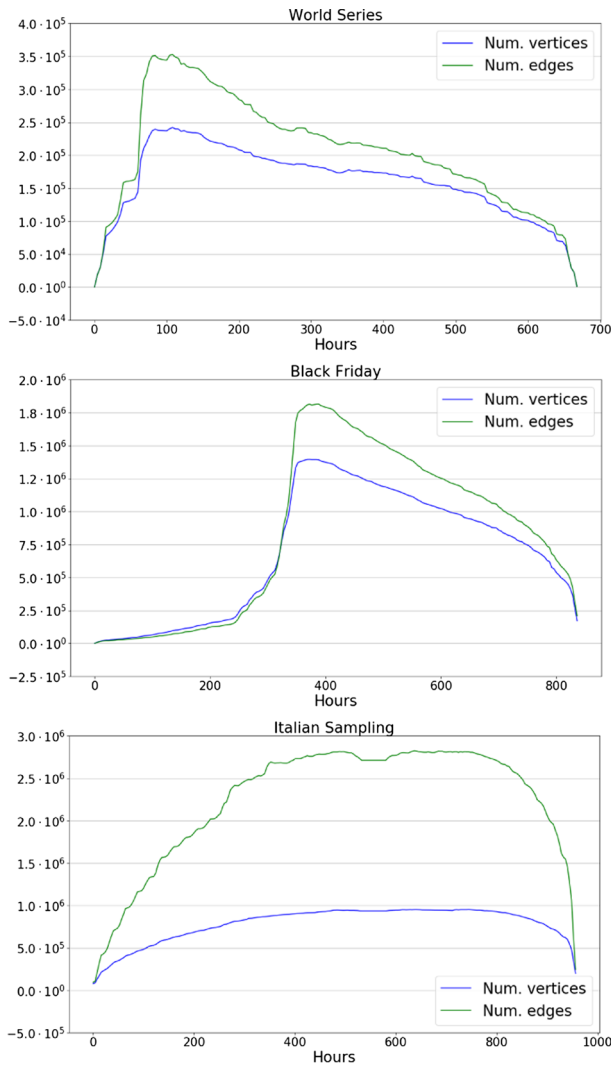|  | World Series | Black Friday | Italian Sampling |
|---|---|---|---|
| Vertices | $4.74 \cdot 10^5$ | $2.7 \cdot 10^6$ | $2.541739 \cdot 10^6$ |
| Edges | $8.40 \cdot 10^5$ | $3.8 \cdot 10^6$ | $1.3708317 \cdot 10^7$ |
| Tweets/edges | 2.3 | 2.603 | 5.45 |
| Tweets/vertices | 4 | 3.66 | 29.4 |

**Fig. 2** Number of vertices (blue) and number of edges (green) of: World Series, Black Friday, and Italian Sampling, as functions of hours

## 4 Experimentation

For each graph $G$ of our dataset, we consider the sequence of DRG temporal graphs $(G_{t_i})_{i \geq 0}$ where $t_{i+1} - t_i$ is 4 hours. For each $G_t$ we compute the four centrality values (betweenness, closeness, degree, and PageRank centrality) of each vertex of the graph.

Given the centrality measure $c$, the *relative centrality value* with respect to $c$ of a vertex $u$ is the ratio between $c(u)$ and the maximum value of $c(\cdot)$.

**Preliminary considerations** Regarding the closeness centrality, we find that the relative centrality values are above 0.9 for almost one third of the nodes (see Fig. 3a). As a
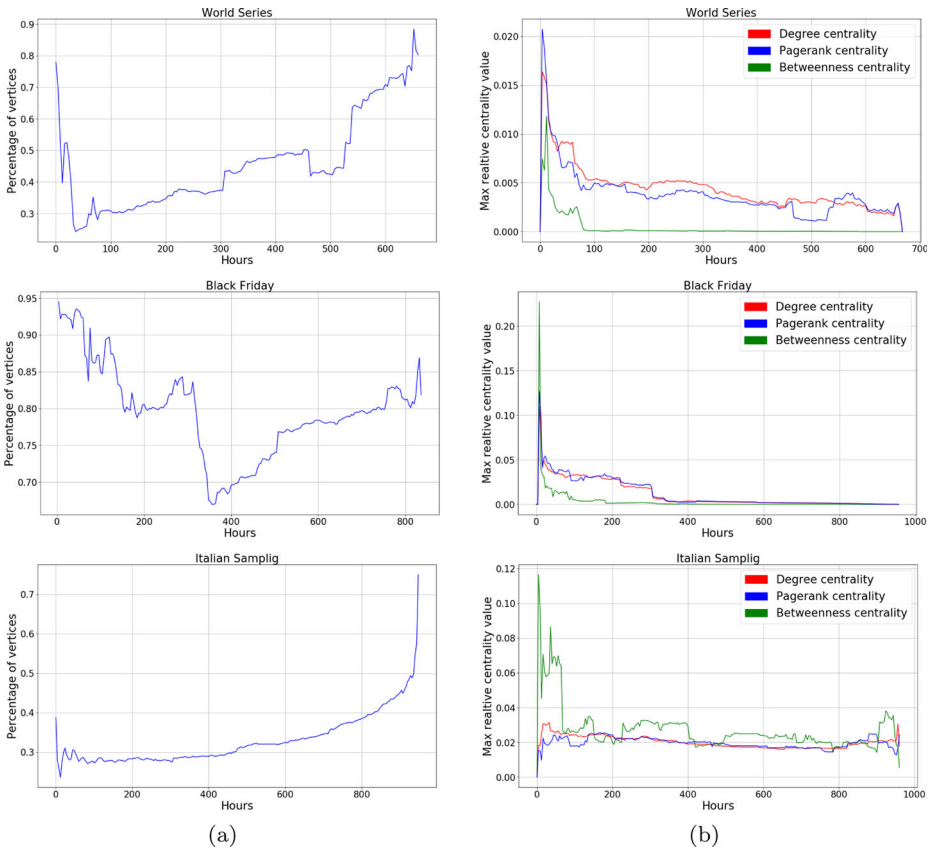
**Fig. 3** **a** Trend over time of percentage of nodes with relative closeness centrality greater than 0.9. **b** The maximum relative degree, PageRank and betweenness centrality values of the 99.9-th percentile over time

consequence the closeness centrality is quite useless to determine the more influential nodes in the graph.

Conversely, the other centrality measures (degree, betweenness, and PageRank) reveal an opposite behavior: excluding the first and last timestamp, 99.9% of vertices always have centrality values below the 20% of the maximum. Figure 3b shows the evolution over time of the three centrality values below which the 99.9% of all values fall (99.9-th percentile). Observe that, from Fig. 3b it results that the highest values are at the very beginning of time sequences, when there is still much instability. After that, values fall below 0.05.

**Analysis of temporal graphs** We say that a node is *central* (with respect to a centrality measure) if its centrality value is at least 75% of the maximum. From the previous

| | World Series | Black Friday | Italian Sampling |
|---|---|---|---|
| Betweennes | 15 | 44 | 31 |
| Degree | 4 | 11 | 10 |
| PageRank | 12 | 16 | 11 |

**Table 2** Number of nodes that have been central at least once for dataset and centrality measure

observations it follows that if we restrict ourselves to the betweenness, degree and PageRank measures, the number central nodes is so small that we can study them one by one. Let $G$ be a DRG, $c$ be a centrality and $t$ be a timestamp, we define $A_{G,c,t}$ as the set of central node of $G_t$ with respect to $c$. Table 2 shows the number of central nodes for each dataset and centrality measure.

In Fig. 4, are shown the sets $A_{G,c,t}$ for the three datasets. The $x$-axis represent the time and in the $y$-axis are reported the vertex ids. In the same plot are collected the informations regards the three centrality measures each of which is represented by a color: green for the betweenness; red for the degree; and blue for the PageRank centrality. An horizontal segment in correspondence to node $u$ that intersects timestamp $t$ means that $u \in A_{G,c,t}$ where $c$ is the centrality measure associated to the segment color.

Nodes that are central for more than one centrality measure are grouped together in the lower portion of the plots. We have observed that there are nodes central only with respect the betweenness centrality measure and for a very short period. These nodes do not add any useful information therefore, in order to make the picture clearer, we have grouped together and represented them by a pseudo-node denoted as more. In the World Series the more node is the union of the segments corresponding to 6 nodes; in the Black Friday 29; and in the Italian Sampling 21. From the above analysis we get the following observations:

–   For all datasets, the degree centrality always produces a total number of central nodes lower than the other measures. Conversely, betweenness centrality is the one that produces more.
–   For all datasets and all the centrality measures, there are nodes that are central for long periods: this trend is more prominent for degree and PageRank centrality.
–   A significant overlap between the central vertices with respect the three measures. For example vertex 572 in Italian Sampling is central for most of the time also for all the three measures (see the third plot of Fig. 4).
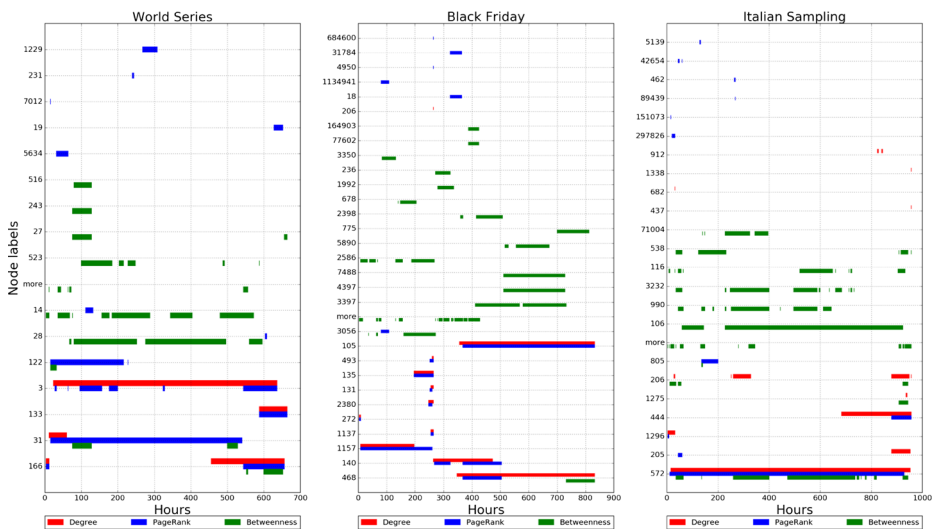


**Fig. 4** Temporal evolution of $A_{G,c,t}$ for the three datasets with respect to the betweenness (in green), degree (in red), and PageRank (in blue) centrality measure

**Table 3** Percentage of vertices in the cumulative DRGs whose relative centrality value is at most 0.01

|                 | Betweenness | Degree  | PageRank |
|-----------------|-------------|---------|----------|
| World Series    | 99.93%      | 99.95%  | 99.93%   |
| Black Friday    | 99.97%      | 99.96%  | 99.97%   |
| Italian Sampling| 99.93%      | 99.78%  | 99.84%   |

**Comparison with the static cumulative DRGs** The *static cumulative DRG* or *cumulative DRG* is obtained from the DRG ignoring the temporal data: it consists of a direct graph where the presence of edge $(a, b)$ means that user $a$ has retweeted user $b$ at least once.

This last analysis involves the centrality measures of the static cumulative DRGs $G$ representing the three datasets. Similarly to DRGs temporal graphs, a large portion of vertices, varying form 28% (for World Series) to 50% (for Black Friday), have closeness centrality above 90% of the maximum, hence, we discard it.

On the contrary, for betweenness, degree, and PageRank centralities almost all the nodes have a value below 1% of the maximum. Table 3 shows, for each dataset and for each measure, the percentage of vertices whose relative centrality value is at most 0.01.

Our goal is to compare the centrality measures in the cumulative DRGs with the ones in the temporal DRGs. Given a dataset and a centrality measure $c$, for each node $u$, we have:

- a single centrality value $c(u)$, for the cumulative DRGs;
- a sequence of centrality values $c_0(u), c_1(u), \ldots$, for the temporal DRGs, one value for timestamp.

In order to make the two data sets comparable, we aggregate the sequence of centrality values $c_0(u), c_1(u), \ldots$ into a single value $s(u)$ given by the sum of all $c_i(u)$. Finally, we can compare the sequence of centrality values of the cumulative DRGs with the sequence of the $s(\cdot)$ values of the temporal DRGs. Table 4 shows the Pearson correlation coefficients between these observations. It turns out that there is a strong correlation in the case of degree and PageRank centralities. Instead, for betweenness centrality the correlation coefficient varies considerably.

From an in-depth analysis we discover that also for the betweenness centrality there is a strong relationship between nodes that are central in both the cumulative DRGs and the temporal DRGs. Table 5 shows the relative betweenness centrality value in the Black Friday cumulative DRG of all central nodes in the Black Friday temporal DRGs with respect to the same centrality measure.

It is interesting to note that 31 of the 44 listed nodes belong to the 0.03% $(=100 - 99.97,$ see Table 3) of vertices whose relative betweenness centrality is at least 0.01. That is a large majority of nodes that are central in temporal graphs are also central in the whole graph. Such behavior is even more pronounced in the World Series and Italian Sampling datasets. Table 6 and 7 show the analogue of Table 5 for the World Series and Italian Sampling datasets. In the World Series (resp. Italian Sampling) case there is only 1 out of 15 (resp.

**Table 4** Pearson correlation between the centrality values in the cumulative DRGs and the aggregated centrality values in the temporal DRGs

|                  | Betweenness | Degree | PageRank |
|------------------|-------------|--------|----------|
| World Series     | 0.59        | 0.89   | 0.91     |
| Black Friday     | 0.16        | 0.84   | 0.97     |
| Italian Sampling | 0.72        | 0.75   | 0.88     |

**Table 5** Relative betweenness centrality in the cumulative Black Friday dataset of nodes that are central in the temporal graphs

| Vertex id | Relative centrality |
| --- | --- |
| 236 | 0.53 |
| 2398 | 0.39 |
| 5890 | 0.16 |
| 7780 | 0.14 |
| 12605 | 0.12 |
| 17414 | 0.11 |
| 16426 | 0.10 |
| 3397, 2607 | 0.09 |
| 9451 | 0.07 |
| 2586, 16417, 4397 | 0.06 |
| 7488, 3056 | 0.05 |
| 7082, 5806, 2542, 146487 | 0.04 |
| 3350, 678, 56750, 56759, 4946, 6118 | 0.03 |
| 37990, 37982, 56760, 164903 | 0.02 |
| 77602, 682 | 0.01 |
| 8714, 9450, 4846, 24191, 22726, 16411, 25530, 118159, 1992, 468, 775, 6077, 170197 | < 0.01 |

4 out of 31) central nodes in the temporal graph with a relative centrality in the cumulative graph less than 0.01.

Finally, if we consider the degree and PageRank centrality measures the above described behavior is even more evident: all central nodes (except for node 3) in the temporal graphs belong to the $\approx 0.2\%$ of nodes with relative centrality in the cumulative graph higher than 0.01. The three exceptions are related the PageRank measure for Black Friday (two nodes) and Italian Sampling (one node).

# 5 Discussion and conclusions

In this paper we have studied the evolution of four centrality measures (betweenness, degree, closeness, and PageRank) on the DRG temporal retweet graphs based on three datasets:

**Table 6** Relative betweenness centrality in the cumulative World Series dataset of nodes that are central in the temporal graphs

| Vertex id | Relative centrality | Vertex id | Relative centrality |
| --- | --- | --- | --- |
| 299 | 1.00 | 243 | 0.19 |
| 31 | 0.69 | 340 | 0.18 |
| 27 | 0.67 | 126 | 0.10 |
| 122 | 0.62 | 516 | 0.07 |
| 14 | 0.49 | 521 | 0.05 |
| 46 | 0.25 | 66050 | 0.03 |
| 28 | 0.23 | 166 | < 0.01 |
| 523 | 0.20 | | |

**Table 7** Relative betweenness centrality in the cumulative Italian Sampling dataset of nodes that are central in the temporal graphs

| Vertex id | Relative centrality |
|---|---|
| 106 | 1.00 |
| 3232 | 0.45 |
| 572 | 0.44 |
| 206,4853 | 0.32 |
| 990 | 0.30 |
| 3306 | 0.27 |
| 2567 | 0.22 |
| 653 | 0.21 |
| 372 | 0.16 |
| 116 | 0.15 |
| 1125 | 0.12 |
| 538 | 0.11 |
| 1275 | 0.10 |
| 5960, 71004, 493, 1851, 1511 | 0.08 |
| 645 | 0.07 |
| 209 | 0.06 |
| 6039 | 0.05 |
| 805, 8998, 5741 | 0.04 |
| 1849, 34521 | 0.01 |
| 22854, 41134, 273383, 52488 | < 0.01 |

Black Friday, World Series, and Italian Sampling. Our main results can be summarized as follows: (i) too many nodes are central with respect closeness centrality, hence this measure is useless to detect influential users; (ii) for the other measures, the number of nodes with very low centrality is very high and the sets of central nodes (with centrality values above 75% of the maximum) are very small and quite similar in the three measures; (iii) similar results hold also for the static cumulative graphs where the sets of nodes with relevant centrality contain central nodes in the sequence of DRG temporal graphs.

As pointed out in [4], the DRG temporal graphs derived from our datasets are quite sparse: this could explain the small number of central nodes respect to the three centrality measures.

According to the above analysis the approach based on the DRG temporal graph and the centrality measures represent a promising approach for detecting influencers in the microblogging Twitter platform. However, this method must be further refined in order to exclude from the influencer the nodes that are central for too little time.

# References

1. Amati G, Angelini S, Bianchi M, Costantini L, Marcone G (2014) A scalable approach to near real-time sentiment analysis on social networks. In: CEUR-WS.org, editor DART 2014 Information Filtering and

Retrieval, Proceedings of the 8th international workshop on information filtering and retrieval co-located with XIII AI*IA symposium on artificial intelligence (AI*IA 2014), vol 1314, pp 12–23

2. Amati G, Angelini S, Capri F, Gambosi G, Rossi G, Vocca P (2016) Modelling the temporal evolution of the retweet graph. IADIS Int J Comput Sci Inform Syst, 11(2):19–30

3. Amati G, Angelini S, Capri F, Gambosi G, Rossi G, Vocca P (2016) Twitter temporal evolution analysis: comparing event and topic driven retweet graphs. In: BIGDACI 2016 - Proceedings of the international conference on big data analytics, data mining and computational intelligence, Volume 1, Funchal, Madeira, Portugal, July 2–4, 2016

4. Amati G, Angelini S, Capri F, Gambosi G, Rossi G, Vocca P (2017) On the retweet decay of the evolutionary retweet graph. In: Smart objects and technologies for social good: second international conference, GOODTECHS 2016, Venice, Italy, November 30 – December 1, 2016, Proceedings. Springer International Publishing, Cham, pp 243–253

5. Bavelas A (1950) Communication patterns in task-oriented groups. J Acoustic Soc America 22(6):725–730

6. Bhattacharya D, Ram S (2012) Sharing news articles using 140 characters: a diffusion analysis on twitter, pp 966–971

7. Bonacich P (1987) Power and centrality: a family of measures. Am J Sociol 92(5):1170–1182

8. Borgatti. SP (2005) Centrality and network flow. Soc Netw 27(1):55–71

9. Brin S, Page L (1998) The anatomy of a large-scale hypertextual web search engine. Comput Netw ISDN Syst 30(1-7):107–117

10. Cha M, Haddadi H, Benevenuto F, Gummadi PK (2010) Measuring user influence in twitter the million follower fallacy. Icwsm 10(10-17):30

11. Conti M, De Salve A, Guidi B, Ricci L (2014) Epidemic diffusion of social updates in Dunbar-Based DOSN. In: Proceedings of parallel processing workshops: Euro-Par 2014 international workshops, pp 311–322

12. Freeman LC (1977) A set of measures of centrality based on betweenness. Sociometry:35–41

13. Freeman LC (1978) Centrality in social networks conceptual clarification. Soc Netw 1(3):215–239

14. Guille A, Hacid H, Favre C, Zighed DA (2013) Information diffusion in online social networks: a survey. SIGMOD Rec 42(2):17–28

15. Hage P, Harary F (1995) Eccentricity and centrality in networks. Soc Netw 17(1):57–63

16. Kwak H, Lee C, Park H, Moon S (2010) What is twitter, a social network or a news media? In: Proceedings of the 19th international conference on world wide web, WWW '10. ACM, New York, pp 591–600

17. Myers SA, Sharma A, Gupta P, Lin J (2014) Information network or social network?: the structure of the twitter follow graph. In: Proceedings of the 23rd international conference on world wide web, WWW '14 companion. ACM, New York, pp 493–498

18. Nieminen J (1974) On centrality in a graph. Scand J Psychol 15:322–336

19. Shimbel A (1953) Structural parameters of communication networks. Bullet Math Biophys 15(4):501–507

**Giambattista Amati** received a PhD in Computing Science from the University of Glasgow, founding the DFR (Divergence From Randomness) models of Information Retrieval and the open source search engine Terrier. Current interests include sentiment analysis and Big Data.

**Simone Angelini** is Big Data Architect and Analyst at Fondazione Ugo Bordoni. He received his Bachelor Diploma in Computer Science fronm the University of Rome "Tor Vergata". Current interests include software development for big data analysis.



**Giorgio Gambosi** Full professor at the University of Rome "Tor Vergata". His research interest include the design and analysis of algorithms and data structure with particular reference to their application to networks and distributed systems.



**Gianluca Rossi** Assistant Professor in computer science at the University of Rome "Tor Vergata". He received the Ph.D. degree in Mathematical Logic and Theoretical Computer Science at the Department of Mathematics of the University of Siena jointly with the Computer Science Department of the University of Florence. His research interests include the design and analysis of algorithms for graphs, networks and distributed systems.

**Paola Vocca** is an associate professor at the University of "La Tuscia", Viterbo. She received a PhD in Computer Science from the University of Rome "La Sapienza". Current interests include algorithms, data structure and large graphs analysis, microblog and wireless networks.

## Affiliations

**Giambattista Amati[1] · Simone Angelini[1] · Giorgio Gambosi[2] · Gianluca Rossi[2] · Paola Vocca[3]** (iD)

> Giambattista Amati
> gba@fub.it

> Simone Angelini
> sangelini@fub.it

> Giorgio Gambosi
> giorgio.gambosi@uniroma2.it

> Gianluca Rossi
> gianluca.rossi@uniroma2.it

[1]    Fondazione Ugo Bordoni, Rome, Italy
[2]    University of Rome "Tor Vergata", Rome, Italy
[3]    University of Tuscia, Viterbo, Italy