

AperTO - Archivio Istituzionale Open Access dell'Università di Torino

## A centrality-based measure of user privacy in online social networks

### **This is the author's manuscript**

*Original Citation:*

*Availability:*

This version is available <http://hdl.handle.net/2318/1616511> since 2017-05-10T11:20:06Z

*Publisher:*

IEEE

*Published version:*

DOI:10.1109/ASONAM.2016.7752439

*Terms of use:*

Open Access

Anyone can freely access the full text of works made available as "Open Access". Works made available under a Creative Commons license can be used according to the terms and conditions of said license. Use of all other works requires consent of the right holder (author or publisher) if not exempted from copyright protection by the applicable law.

(Article begins on next page)

This is the author's final version of the contribution published as:

R.G. Pensa, G. Di Blasi. A centrality-based measure of user privacy in online social networks, in: Proceedings of 2016 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM), IEEE, 2016, 978-1-5090-2846-7, pp: 1438-1439.

The publisher's version is available at:

<http://xplore.staging.ieee.org/ielx7/7736513/7752180/07752439.pdf?arnumber=7752439>

When citing, please refer to the published version.

Link to this full text:

<http://hdl.handle.net/2318/1616511>

# A Centrality-based Measure of User Privacy in Online Social Networks

Ruggero G. Pensa and Gianpiero Di Blasi

Department of Computer Science, University of Torino, Italy

Email: ruggero.pensa@unito.it

**Abstract**—The risks due to a global and unaware diffusion of our personal data cannot be overlooked when more than two billion people are estimated to be registered in at least one of the most popular online social networks. As a consequence, privacy has become a primary concern among social network analysts and Web/data scientists. Some studies propose to “measure” users’ profile privacy according to their privacy settings but do not consider the topological properties of the social network adequately. In this paper, we address this limitation and define a centrality-based privacy score to measure the objective user privacy risk according to the network properties. We analyze the effectiveness of our measures on a large network of real Facebook users.

## I. INTRODUCTION

More than two billions people are estimated to be registered in at least one of the most popular online social network platforms. In view of these numbers, the risks due to a global and unaware diffusion of our sensitive personal data cannot be overlooked. Even though social networking sites notify their users about the risks of disclosing private information, most people are not aware of the dangers due to the indiscriminate disclosure of their personal data. Moreover, despite the fact that all social media provide some advanced tools for controlling the privacy settings of the user profile, such tools are not user-friendly and they are barely utilized, in practice. Some studies try to go beyond these limitations by “measuring” users’ profile privacy according to their privacy settings [1], [2], however these privacy measures do not consider the topological properties of the social network adequately.

Our assumption is that the actual privacy leakage risk of users is crucially affected by the properties of the social network they belong to. To explain this, let us consider two users  $u_1$  and  $u_2$  sharing the same attitude to their own privacy protection. User  $u_1$  is mostly surrounded by friends that care about their own (and others’) privacy, while  $u_2$  is principally connected to friends that do not care that much about their privacy leakage. According to these hypotheses, user  $u_2$  should be more exposed to privacy leakage than  $u_1$ . These considerations lead to the intuition that privacy risk in a social network may be modeled similarly as page authority in a hyperlink graph of web pages. In fact, it is a well-known fact that more authoritative websites are likely to receive more links from other authoritative websites. Our hypothesis is that we may transpose the concept of “importance” of a web-page into the concept of “privacy risk” of users in a social network as follows: the more an individual is surrounded by friends

that are careless about their privacy, the less the individual her/himself is likely to be protected from privacy leakage.

With the final goal of enhancing users’ privacy awareness in online social networks, in this paper we propose a new centrality-based privacy score based on Pagerank [3], one of the most popular algorithms to rank web pages based on their importance (or authority). We show the effectiveness of our privacy measure on a large network of real Facebook users.

## II. COMPUTING PRIVACY SCORES

We consider a set of  $n$  users  $U = \{u_1, \dots, u_n\}$  participating in a social network, here represented as a directed graph  $G(V, E)$ , where  $V$  is a set of  $n$  vertices  $\{v_1, \dots, v_n\}$  such that each vertex  $v_i \in V$  is the counterpart of user  $u_i \in U$  and  $E$  is a set of directed edges  $E = \{(v_i, v_j)\}$ . Given a pair of users  $u_i, u_j \in U$ ,  $(v_i, v_j) \in E$  iff there exists a link from  $u_i$  to  $u_j$  (e.g., users  $u_j$  is in the friend list/circle of  $u_i$  or  $u_i$  follows  $u_j$ ). Without loss of generality, we assume that the link between two users is always reciprocal (if there is a link from  $u_i$  to  $u_j$  then there is also a link from  $u_j$  to  $u_i$ ). Hence, the social network here is represented as an undirected graph  $G(V, E)$ , where  $E$  is such that if  $(v_i, v_j) \in E$ , then  $(v_j, v_i) \in E$ . Finally, each user is characterized by an *intrinsic privacy risk*  $\rho_p(u_i)$ , which is defined as the user propensity to privacy leakage. The assumption is that some users are more prone to disclose their personal data than others. This propensity is reflected in the way users configure their privacy settings. Assuming that users’ activity in a social network is known, measuring their intrinsic privacy risk is not trivial. In this work we rely on the privacy score (denoted *P-Score*) defined by Liu and Terzi [1]. It is based on a mathematical model leveraging the item response theory (a well known theory in psychometrics).

By definition, the *intrinsic privacy risk*  $\rho_p(u_i)$  does not consider the topology of the social network. However, the actual privacy leakage risk of users is crucially affected by the properties of the social network they belong to: two users sharing the same attitude to their own privacy protection are not necessarily subject to the same risk. If a user is mostly surrounded by friends that do not care that much about their privacy leakage, then she should be more exposed to privacy leakage than a user who is principally connected to friends that care about their own (and others’) privacy. This consideration leads to the intuition that privacy risk in a social network may be modeled similarly as page authority in a hyperlink graph of web pages. Hence, we transpose the concept of “importance”

of a web-page into the concept of “privacy risk” of users in a social network as follows: the more an individual is surrounded by friends that are careless about their privacy, the less the individual her/himself is likely to be protected from privacy leakage. One of the most popular algorithms to rank web pages based on their centrality (or authority) is PAGERANK [3]. In particular, our setting is similar to the definition of *personalized PAGERANK* [4], used to create a personalized view of the relative importance of the nodes. We can now introduce our centrality-based privacy score (called *CP-Score*), defined by the following distribution:

$$\mathbf{P}^\rho = d\mathbf{A}\mathbf{P}^\rho + \frac{(1-d)}{\sum_{k=1}^n \rho_p(u_k)} \boldsymbol{\rho} \quad (1)$$

where  $\mathbf{P}^\rho = [p^\rho(v_1), \dots, p^\rho(v_n)]^\top$  is the PAGERANK vector ( $p^\rho(v_i)$  being the PAGERANK associated to vertex  $v_i$ ),  $d = [0, 1]$  is the damping factor (the  $1-d$  quantity is also known as restart probability),  $\boldsymbol{\rho} = [\rho_p(u_1), \dots, \rho_p(u_n)]^\top$ , and  $\mathbf{A}$  is a  $n \times n$  matrix such that each element  $a_{ij} = 1/\text{deg}(v_i)$  ( $\text{deg}(v_i)$  being the degree of  $v_i$ ) if  $(v_i, v_j) \in E$  ( $a_{ij} = 0$  otherwise).

Equation 1 provides a set of values that can not be directly interpreted as a privacy score, since they are not in the same scale. Hence, the following re-scaling operation is required to compute the correct values of the privacy score:

$$\rho'_p(u_i) = \rho_{min} + (\rho_{max} - \rho_{min}) \cdot \frac{p^\rho(v_i) - p_{min}^\rho}{p_{max}^\rho - p_{min}^\rho} \quad (2)$$

where  $\rho_p$  denotes the intrinsic privacy score value,  $p^\rho(v_i)$  is the centrality-based privacy score value for node  $v_i$ , and  $\rho'_p$  denotes the recomputed privacy score value. Moreover,  $\rho_{min} = \min_j \{\rho_p(u_j)\}$ ,  $\rho_{max} = \max_j \{\rho_p(u_j)\}$ ,  $p_{min}^\rho = \min_j \{p^\rho(v_j)\}$  and  $p_{max}^\rho = \max_j \{p^\rho(v_j)\}$ .

### III. EXPERIMENTAL RESULTS

In this section we report and discuss the results of the experiments that we conducted on a Facebook graph generated leveraging an online experiment that enabled us to collect the ego-networks of 185 volunteers<sup>1</sup>. The social network consisting of all participants plus their friends is an undirected graph with 75,193 nodes and 1,377,672 edges, an average degree of 36.644 and a clustering coefficient of 0.613. The participants had also to indicate to which people (no one, close friends, friends except acquaintances, all friends, friends of friends, everyone on Facebook) they were willing to allow the access to five topics with different levels of sensitivity. From December 2015 to February 2016, 101 out of 185 participants answered all questions of the survey.

We conducted our experiments as follows. First we compute the intrinsic privacy score of each node using two different strategies: in a first set of experiments, the intrinsic risk for the nodes corresponding to the participants in our survey is computed according to the privacy score (P-Score) obtained by processing their answers [1]. For all other 75,193 – 101 nodes the intrinsic privacy risk is uniformly set equal to the

<sup>1</sup><http://kdd.di.unito.it/privacyawareness/>

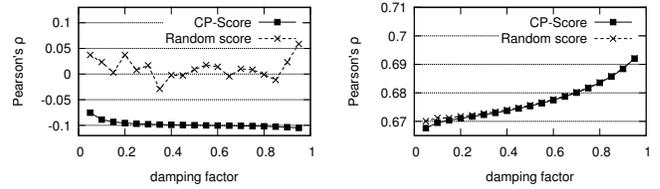


Fig. 1: Pearson’s correlation values of scores w.r.t. original scores (left) and eigenvector centrality (right).

mean of the P-Score’s. In the second set of experiments, the intrinsic privacy score of the participants is drawn from a Gaussian distribution having the same mean and variance than those observed for the P-Score of the 101 participants in the network, while the intrinsic privacy score of all other nodes is set to the mean of the P-Score’s. Then, for each experimental setting we compute the centrality-based privacy score (CP-Score) using the power-iteration method [5]. We repeat the experiments by varying the values of the damping factor in the range  $[0.05, 0.95]$ . We then measure the Pearson’s correlation between the intrinsic privacy risk and the CP-Score on the set of 101 participants. To achieve significant results, we run each experiment 30 times. The results are reported in Figure 1 (left). The CP-Score exhibits a slightly negative correlation w.r.t. the P-Score. This result probably means that the privacy score defined in [1] is not always a good estimate of the objective privacy risk of the user. It is worth noting that the random score always fluctuates around the zero correlation value. Figure 1 (right) also shows that the centrality-based privacy scores computed on the 101 participants are positively correlated with their eigenvector centrality, despite their intrinsic privacy score. In this case, the CP-Score and random score exhibit a similar behavior when the damping factor  $d$  is greater than 0.2. These preliminary results seem to confirm our initial claim: to measure the objective privacy risk, any privacy metric should be contextualized within the social graph by considering the influence of the network on each user.

As future work, we will use simulated data to better analyze the behavior of our centrality-based privacy measure.

### ACKNOWLEDGEMENT

The work presented in this paper is co-funded by Fondazione CRT (grant number 2015-1638). The authors wish to thank all the volunteers who participated in the survey.

### REFERENCES

- [1] K. Liu and E. Terzi, “A framework for computing the privacy scores of users in online social networks,” *TKDD*, vol. 5, no. 1, p. 6, 2010.
- [2] Y. Wang, R. K. Nepali, and J. Nikolai, “Social network privacy measurement and simulation,” in *Proceedings of ICNC 2014*. IEEE, 2014, pp. 802–806.
- [3] S. Brin and L. Page, “The anatomy of a large-scale hypertextual web search engine,” *Computer Networks*, vol. 30, no. 1-7, pp. 107–117, 1998.
- [4] G. Jeh and J. Widom, “Scaling personalized web search,” in *Proceedings of WWW 2003*. ACM, 2003, pp. 271–279.
- [5] G. H. Golub and H. A. van der Vorst, “Eigenvalue computation in the 20th century,” *Journal of Computational and Applied Mathematics*, vol. 123, no. 1–2, pp. 35–65, 2000.