



**HAL**  
open science

## Online Testing of User Profile Resilience Against Inference Attacks in Social Networks

Younes Abid, Abdessamad Imine, Michael Rusinowitch

► **To cite this version:**

Younes Abid, Abdessamad Imine, Michael Rusinowitch. Online Testing of User Profile Resilience Against Inference Attacks in Social Networks. ADBIS 2018 - First International Workshop on Advances on Big Data Management, Analytics, Data Privacy and Security, BigDataMAPS 2018, Sep 2018, Budapest, Hungary. hal-01939277

**HAL Id: hal-01939277**

**<https://hal.science/hal-01939277>**

Submitted on 12 Dec 2018

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Online testing of user profile resilience against inference attacks in social networks

Younes Abid, Abdessamad Imine and Michaël Rusinowitch

<sup>1</sup> Lorraine University, Cnrs, Inria, 54000 Nancy, France  
firstname.lastname@loria.fr

**Abstract.** To increase awareness about privacy threats, we have designed a tool, SONSAI, for Facebook users to audit their own profiles. SONSAI predicts values of sensitive attributes by machine learning and identifies user public attributes that have guided the learning algorithm towards these sensitive attribute values. Here, we present new aspects of the system such as the automatic combination of link disclosure attacks and attribute prediction. We explain how we defined sensitive subjects from a survey. We also show how the extended tool is fully interfaced with Facebook along different scenarios. In each case a dataset was built from real profiles collected in the user neighbourhood network. The whole analysis process is performed online, mostly automatically and with good accuracy. It is 0.79 in AUC when inferring the political orientation.

**Keywords:** Online Social Network (OSN), Inference Attacks, Privacy, Link Disclosure

## 1 Introduction

Personal information if revealed may have serious consequences on social network users. This information can be exploited to carry out personalized spam attacks [7], identity theft attacks [4], cloning attacks [12], Sybil attacks [11], etc. They might cause serious damages to companies such as degradation of reputation, copyright infringement, loss of intellectual property, etc.

Social networks provide several solutions in order to safeguard the privacy of users. However, their main deficiencies are related to complicated, non-uniform, periodically updated and unintelligible privacy policies, long and ambiguous user charters, and non-ergonomic privacy management interfaces. Although most social networks offer similar services (creating profiles, pages and groups, establishing links and interactions), their visibility management and the definition of links (symmetrical, non-symmetrical) are different. These design differences may be confusing for users of multiple social networks that are careless with checking each network settings. Moreover the default parameters promotes public dissemination but increases the risk of sensitive information leakage.

Most importantly, social networks do not provide protection against inference of implicit information. Derived by correlating different public attributes or different public profiles, as in collaborative recommendation, this information is actually the main profit source of social networks' business model as

they can be exploited for targeted advertising. Therefore, knowledge accumulated in social networks about users goes beyond what is published and can be a threat to their privacy. In [13], Winter Mason has mined the cultural similarities between American Facebook users and their political view (Democrats or Republicans). He sampled Facebook users' profiles who liked the campaign pages of some Democrat or Republican politicians. Then, he collected their lists of liked pages. Finally, he statistically identified the page types that are most disproportionately liked by the supporters of one political view versus the other. The results of his work show that politics is highly correlated to musicians, landmarks, authors, books and TV shows.

*Contributions.* In order to combat privacy leakage, it is important to define which personal information is sensitive. Some researchers consider that all the unpublished values of attributes (masked or not specified) by a given user are sensitive for him [15,16]. While others select a few attributes and consider them to be sensitive such as sexual orientation [10], political affiliation [10,6] and age [14]. It is also possible to rely on the definition of sensitive information given by law. However, social networks are evolving faster than legislation. For instance, health data were not considered sensitive by the French law of January 6, 1978<sup>1</sup> related to computers, files and freedoms. It was considered sensitive much later. It is also possible to rely on a definition of sensitive subject given by social media themselves. But can we trust social networks in defining what is sensitive or not as they make most of their profit using personal information for targeted advertising?

Hence, we have first conducted a questionnaire survey to define sensitive subjects based on the behaviour of french Net-surfers. This method has the advantage of being fast, objective, accurate and up-datable. The most sensitive subjects according to Facebook french users who have participated in our survey are *Religion, Money, Politics, Dating, Shopping* and *Health*.

Then, we present SONSAI, our application to help Facebook users to protect their privacy. To that end, SONSAI tests a user profile against privacy attacks and tracks their origin. This approach allows one to delimit the perimeter of threats and to design effective countermeasures. Concretely, SONSAI performs online inference attacks on the world largest social network, Facebook. The attacks have been tested by several real volunteer profiles. SONSAI allows users to identify public attributes that are correlated with sensitive attributes and therefore to prevent these attacks by modifying these public attributes. In the context of *General Data Protection Regulation* (GDPR)<sup>2</sup>, promulgated recently by the European Union, to stress on users' control over their personal data, our tool may contribute to increase user awareness as for the risks related to personal data processing.

*Outline.* In Sec. 2 we discuss the problems. In Sec. 3 we recall some related works. In Sec. 4 we present the result of a survey conducted to identify sensitive

<sup>1</sup> Loi n° 78-17 du 6 janvier 1978 relative à l'informatique, aux fichiers et aux libertés

<sup>2</sup> <https://www.eugdpr.org/>

subjects. In Sec. 5 we overview the architecture and component functionalities of SONSAI, a tool for users to test their profile against inference attacks. In Sec. 6 we describe experiments on real data. Finally we discuss accuracy of SONSAI attacks in Sec. 7, countermeasures in Sec. 8, and conclude in Sec. 9.

## 2 Discussion

Let us discuss briefly the problems that had to be solved in order to design SONSAI. In order to effectively track privacy leakage, it is important to combine link prediction and attribute prediction. For instance, an adversary can perform link prediction attacks in order to disclose the local network of his target (friends and group members). Then, he can perform more accurate attribute prediction attacks with extra information provided by the discovered local network. Online attribute prediction attacks encompasses two steps: (i) data collection and (ii) data analysis. Data collection must be fast, selective, passive and unnoticed. Since social networks are highly dynamic and contain a huge amount of data, random collection may result in useless data. On the other hand, massive collection is time wasting. A fast and selective sampling algorithm must be used in order to guide the collector toward most relevant data and speed up the process. Moreover, the adversary must limit his interaction with his target. He must perform his attack in a passive way in order to go unnoticed by the target to him. The adversary should only send legal requests to collect data and should not exceed some threshold to remain unnoticed by the social network. Data analysis should be fast, accurate and deal with sparsity. We recall that the system (collection and analysis) is meant to help users safeguard their privacy against real attacks. Hence, data analysis must not exceed a few minutes in order to rapidly put the hand on the origins of threats and quickly put countermeasures into action. Analysis results should be accurate in order to reduce false positive alerts and cover all threats. As the collector only samples a few data from an ocean of them, the analyser should deal with the fact that collected data may be sparse and incomplete.

## 3 Related works

In [5] the authors propose to combat attribute prediction attacks in social network by creating new links between users in order to reduce the difference of the distribution of attribute values in the user local network and in the global one. In [10], the authors propose a content&link-based classifier that outperforms both content-based classifiers and link-based classifiers when predicting the political views and the sexual orientation of Facebook users. In addition, they explore the effectiveness of sanitization techniques to prevent such attacks. In contrast to [5], sanitization solutions in [10] consist of removing contents and links. However, selecting the right contents and links to remove or to add without altering the utility of the social network is a challenging task. In our present work we aim to help users to identify the critical attributes that are correlated to sensitive subjects in their communities and that lead to the undesirable inference. It is then up to users to intervene by adding links in order to alter the accuracy of inference

due to data disagreement or by deleting links in order to disrupt inferences by lack of data. In [6], the authors design a classifier to predict the political alignment of Twitter users based on the tweet contents and the re-tweet network. They show that such classifiers widely outperform classifiers that are based only on content. Our first tests [3] show that some attributes such as political alignment are correlated to the network structure (e.g., friendship links). However, this is not the case for other ones such as gender and relationship status (e.g., married, ...). In [18], the authors introduce an information re-association attack in order to predict the values of sensitive attributes of users. This attack consists in combining web search with information extraction and data mining techniques. This study shows that the attack is more successful when including information about the target university networks. In addition, it shows that Facebook graduated users from top schools are more vulnerable under this attack than random users. In our work, we quantify the correlation between attributes. SONSAT generates inference rules that depends on the behaviours of users in the target neighborhood. For instance, SONSAT can automatically decide whether the university networks around the target is an important factor for inference success or not. In [15] the authors show that an adversary can infer sensitive attribute values of a target based only on the target local network (1-hop friendship network) and the public attribute within it. The proposed predictor takes into consideration the network structure by quantifying the importance of friendship relations. Then, it measures the power of each attribute value according to the importance of the target friend that publishes it. In [9], the authors extend the attribute-augmented social network model that is introduced in [17]. In the initial model, attribute values are represented by nodes. The users that publishes a particular value of an attribute are linked to its representing node in the model. The extended model adds negative links between users and their hidden attribute values and mutex links between mutually exclusive values of the same attribute such as male and female. This model is used with both supervised and unsupervised methods to predict links between users as well as links between users and attribute values. In [16], the authors design a classifier to predict the missing attribute values of a Google+ user. The classifier only takes into consideration users that are one hop distant from the target. In addition to attributes, the designed classifier exploit the direction of links (follower or followings), the type of links (acquaintance, family, friend ...) and the tie-strength of the links. In [8], the authors extend the attribute-augmented social network model [17] by adding behaviours nodes to the framework. Then, they design a vote distribution attack to predict attribute values. They show that by taking into account social friendship, attributes and behaviours, the accuracy of attacks is considerably increased.

In our work, we analyse the local friendship network of the target (direct friends). When the target hides friends we first perform link disclosure attack with certainty to disclose his local network. This combination of link disclosure (with certainty) and attribute inference coordinated within the same system (fully interfaced with Facebook) seems to be unique. Each attribute is represented by a bipartite graph where edges connect users to the attribute values

they like. The system also relies on our previous works: specific graph comparison techniques to measure attribute correlations [3], a clustering algorithm to group similar sensitive attribute values (for instance similar politicians) and a shallow neural network to infer semantic proximity between public values and sensitive ones [2].

## 4 Definition of sensitive subjects

We have conducted a questionnaire survey to define sensitive subjects according to the behaviours of french net-surfer. We have analysed the behaviours of 232 users of social media aged between 20 and 78 years that live in 21 different French regions. We have classified the subjects discussed on social media according to four criteria: rate of discussion on social networks, rate of discussion on forums and websites, rate of anonymous publication and avoided subjects. Based on those criteria, we have proposed a definition of sensitive subjects.

Category	Discussion on social networks (in %)	Discussion on forums and websites (in %)	Anonymous publications (in %)	Avoided subject (in %)
Money	0.94	54.42	25	10.14
Religion	5.63	26.05	14.28	33.33
Shopping	1.88	66.05	9.09	0
Dating	5.16	24.19	0	21.74
Health	17.37	66.05	9.09	5.8
Politics	25.82	54.42	25	50.72

Table 1: Statistics related to the sensitive subjects.

**Sensitive subjects.** Let  $V$  be the set of avoided subjects,  $N$  the set of subjects whose rate of anonymous publications on forums and websites is above average and  $D$  the set of subjects whose discussion rate on the forums/websites **or**<sup>3</sup> social networks are below the threshold of the mean of all discussions minus the standard deviation on that media. A discussed subject on social media is **sensitive**, if and only if, it belong to at least two sets from the defined sets ( $V$ ,  $N$  and  $D$ ). The most sensitive subjects according to french Facebook users that participate in our survey are *Religion*, *Money*, *Politics*, *Dating Shopping* and *Health*. Table 1 details the statistics related to the defined sensitive subjects with regard to the analysed criteria. Additionally, the analysis of the participants behaviours results in the following privacy attack vector statistics:

- 52.05% use the same e-mails and 65.75% use the same user-names on different social networks.
- 90% have the same friends over different networks.
- 72.16% do not cleanly delete their profiles when leaving social networks.
- 15.96% publish photos without asking the consent of people appearing in these photos.

<sup>3</sup> Or indicates an inclusive disjunction.

- 8.45 % add strangers to their friend lists only because they have common friends.
- In a test, 6.10 % are not able to recognize a person added randomly to their friend lists.

## 5 Social Networks Sensitive Attribute Inference

In this section we detail the *SOcial Networks Sensitive Attribute Inference* system (SONSAI). Algorithm 1 details the flow of tasks performed by SONSAI in order to detect privacy threats. SONSAI first collects the 1-hop friendship network around the user  $u$  and their attributes (lines 1-6). If the user  $u$  hides his friend list then SONSAI performs link disclosure attacks as detailed in [1] in order to disclose with certainty some of his friends. The attacker model is passive as the attack does not require interaction with users. The preparation of the link disclosure attack consists of sampling users that have high probability to be friend with  $u$  and that publish their friend list. To that end, SONSAI explores the group network at distance 2 from  $u$ . After that, it performs friendship and mutual friend disclosure attacks by taking advantages of queries provided by the social networks APIs. A friendship attack consists of disclosing the links between  $u$  and the members of the explored groups. A mutual friend attack consists of disclosing the links between  $u$  and the friends of the members of the explored groups. The results of our previous study [1] show that about half of Facebook users are exposed to the danger of friendship disclosure through their membership to groups of less than 50 members .

In Section 4 we have identified the most sensitive subjects for french users. Assume a user  $u$  wants to check with SONSAI whether a sensitive subject information can be inferred from his profile. He first selects through a Combobox an attribute correlated to these sensitive subjects (line 7). For instance he can select politicians or political parties for *Politics* subject. Correlation of attributes is quantified by comparing their bipartite graph representation (see [3]).

The sensitive attribute is selected from a displayed list of attributes discovered in the user local social network. To simulate an inference attack on the selected sensitive attribute, the user is asked to provide two pieces of information: (i) *top\_n*, the percentage of attributes to be selected for learning and (ii) whether the sensitive attribute values have to be clustered by similarity to reduce the search space (lines 9-11).

SONSAI uses random walks (line 12) and Word2Vec algorithm (line 13) in order to infer the sensitive values of the target based on his preferred values for the selected attributes for learning. The results of the inference attack help the users understand the source of information leakage. SONSAI ranks the list of the selected attribute for learning according to their correlation to the sensitive attributes (line 14).

Finally, SONSAI interacts with the users in order to check the accuracy of the inference (i.e., whether the inferred values are correct) and assess the risk of privacy leakage (line 15). Verdicts returned by the SONSAI depends solely on the collected data around the user. This verdict is given as a score quantifying the risk of inferring correct values for a given sensitive attribute. The score is

obtained by comparing the ranking of sensitive values to the values that are really liked by user  $u$ . We use *Area Under the Curve* to compute this score. The risk of disclosing values of a sensitive attribute is considered to be high when the score is higher than 65%, moderate when the score is between 50% and 65% and low when the score is less than 50%.

**Data:**  $target$  ▷ target profile (user input)  
 $top\_n$  ▷ learning attributes ratio to be selected (user input)  
 $cluster\_b$  ▷ boolean for clustering option (user input)  
 $sensitive\_att$  ▷ sensitive attribute (user input)  
**Result:**  $correlated\_attributes$  ▷ attribute ranking  
 $ranked\_values$  ▷ sensitive value ranking  
 $risk\_level$  ▷ sensitive value disclosing risk

```

1 if  $masked\_friend\_list(target)$  then
2 |  $friend\_list \leftarrow friend\_disclosure\_attack(target)$ ;
3 else
4 |  $friend\_list \leftarrow get\_public\_friend\_list(target)$ ;
5 end
6  $attributes \leftarrow crawl(friend\_list)$ ; ▷ attributes are stored as bipartite graphs
7  $correlated\_atts \leftarrow graph\_comparison(sensitive\_att, attributes)$ ;
8  $selected\_atts \leftarrow select\_atts(top\_n, correlated\_atts)$ ;
9 if  $cluster\_b$  then
10 |  $cluster\_values(sensitive\_att)$ ;
11 end
12  $walks\_text\_files \leftarrow random\_walk(sensitive\_att, selected\_data)$ ;
13  $embeddings \leftarrow word2vec(walks\_text\_files)$ ;
14  $ranked\_values \leftarrow ranks(embeddings, target, sensitive\_att.values)$ ;
15  $risk\_levels \leftarrow compare\_ranks(ranked\_values, user\_real\_values)$ ;

```

**Algorithm 1:** SONSAI crawling and analysis steps.

## 6 Inference scenarios

In the following we detail two scenarios of inferring the pages of politicians liked by two real Facebook users  $u_1$  and  $u_2$ . Table 2 gives a sample of values of the pages of *Musicians/Bands*, *News/Media Websites* and *Communities* that are liked by  $u_1$  and  $u_2$ . Target  $u_1$  is a french user. He publishes his friend list on Facebook. Target  $u_2$  is a canadian user that hides his friend list. SONSAI performs then link disclosure attack in order to first disclose the friends of  $u_2$ . Then, it crawls the his friendship network and collects the values of attributes liked by his direct friends.

After the analysis, a first table that summarizes the list of the top correlated attributes to the sensitive one is displayed. The correlations are quantified in percentage. A second table presenting sorted values of the sensitive attribute according to the probability they are liked by the user is displayed too. Values in the second table are grouped in clusters in a way that maximizes the similarity



Targets	Musicians/Bands	News/Media Websites	Communities
$u_1$	Clean Bandit Dillon Francis Monsieur Monsieur Max Vangeli DJ Fresh Cazzette Charlotte de Witte RL Grime Bassjackers Tritonal	Spi0n BuzzFilGeek My Little Paris confidentielles.com Boiler Room MY Secret NY Street FX Motorsport & Graphics Hitek Le Figaro PIX GEEKS	Pour le retrait du timbre Femen Bordel De Droit Soigner Dans la Dignité ADDM - Respect it Enjoy it Entourage - Réseau Civique Soutien au bijoutier de Nice Valls Dégage Banamak Fab Bike Pour la démission de Hollande
$u_2$	Justice The Prodigy Queen Crystal Castles Le husky Daniela Andrade Boys Noize Gunther Heroik Les Poignards	TED InfopresseJobs Too Close To Call Radio-Canada Information Isarta - Emplois NowThis Infos Insolites Faits et Causes Progrès Villeray - Parc-Extension	Keep Calm & Be Real Es-tu game? Nos casseroles contre la loi spéciale The Voyage North Arcade MTL Nous sommes les 68% Larping.org Pierre Céré Commodore 64 Astuces de Mac Gyver

Table 2: Some attribute values of targets  $u_1$  and  $u_2$ .

between values inside a cluster and minimizes it between values from different clusters. The size of clusters is controlled so that any of them never gets twice larger than any other one. The similarity inside a cluster is the mean of all similarity indexes of its elements. It is computed by Equation 1.

$$similarity(c) = 2 \frac{\sum_{(v,v') \in c \times c} similarity\_index(v,v')}{|c|(|c|-1)} \quad (1)$$

$|c|$  is the number of values in Cluster  $c$ . The similarity index computes the similarity between two values. In order to select an adequate similarity index to be implemented in SONSAI, we have first tested four well-known indexes. These indexes are defined in Table 3.  $\Gamma(v)$  is the set of users that like value  $v$ .

Jaccard	$\frac{ \Gamma(v) \cap \Gamma(v') }{ \Gamma(v) \cup \Gamma(v') }$
Adamic Adar	$\sum_{z \in \Gamma(v) \cap \Gamma(v')} \frac{1}{\log  \Gamma(z) }$
Common Neighbours	$ \Gamma(v) \cap \Gamma(v') $
Preferential Attachment	$\frac{ \Gamma(v)  \times  \Gamma(v') }{ \Gamma(v)  +  \Gamma(v') }$

Table 3: Similarity indexes

SONSAI implements the algorithm detailed in [2] to define clusters. The clustering algorithm is greedy. It defines one cluster at a time. At each step, it adds the most similar value from the set of non clustered values to the last created cluster. When size conditions are met and when adding any non clustered values will only decrease the average similarity of the cluster values, a new cluster is initialized that contains the most liked non clustered value. Overall, the algorithm computes the  $n(n-1)/2$  similarities between all couples of values.

Figure 1 gives the similarity variations with respect to the sizes of clusters. The algorithm clusters 4 589 politician pages. The x-axis represents the minimal sizes of clusters. The y-axis represents the similarities inside clusters computed using the corresponding similarity index. We notice that Adamic Adar, Common Neighbours and Preferential Attachment based clustering method are very

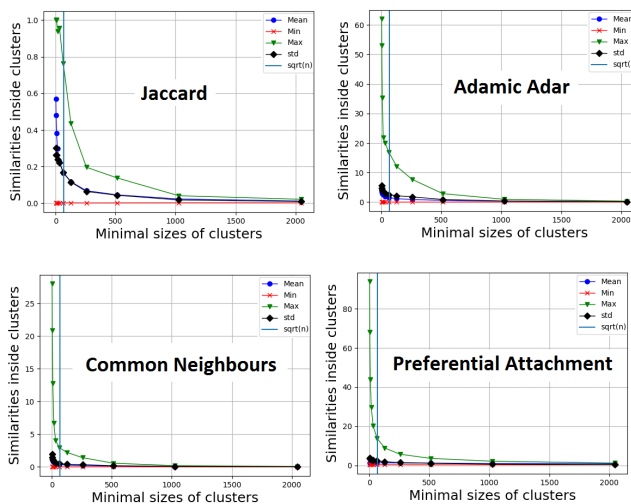


Fig. 1: Variations of similarities with respect to the size of clusters.

sensitive to the size of clusters. In fact, when the size of cluster increases the similarities dramatically decreases. On the other hand, Jaccard based clustering method maintains high similarities inside clusters even when the size of cluster is  $\sqrt{n}$ , with  $n$  is the total number of values. Four volunteers have manually checked the similarities of clusters based on their political backgrounds. They confirmed that Jaccard based clustering method generates better clusters than Adamic Adar, Common Neighbours and Preferential Attachment based clustering method. Hence, SONSAI uses Jaccard index to measure the similarity between sensitive attribute values. Moreover, the mean Jaccard index,  $MJ$ , between the values of a cluster  $C$  can be interpreted as follows: If a user likes a particular value from the cluster  $C$  then he tends to like  $MJ \times 100\%$  of values of the rest of values from the same cluster  $C$ .

SONSAI allows users to open any cluster displayed on the table and click on any value to open its corresponding Facebook page. For each cluster, the user can specify the number of his real liked values inside. The algorithm measures then the accuracy of the ranking using the Area Under the Curve (AUC). Table 4 summarizes the politicians that are liked by  $u_1$  and  $u_2$ . Target  $u_1$  is a french user politically right-oriented. The inference accuracy for the politician pages he like is 0.72. Target  $u_2$  is a canadian user of left political orientation. The inference accuracy for the politician pages he like is 0.97.

## 7 Accuracy of SONSAI inferences

We have crawled the friendship network of 100 Facebook profiles of users that live in North-East France. For each crawled profiles we have collected the list of liked pages, the list of friends, the gender and the relationship status. The dataset contains 1 926 different types of pages, 1 022 847 different pages and 15

Targets	$u_1$	$u_2$
Liked Politicians	Marine Le Pen Jean-François Copé Laurent Wauquiez Bruno Le Maire Jean-Marie Le Pen Nicolas Dupont-Aignan Xavier Bertrand Nathalie Kosciusko Morizet Francois Fillon Marion Maréchal-Le Pen	Simon Marcell Martine Ouellet Amir Khadir Jack Layton Jocelyn Beaudoin Alain Therrien Justin Trudeau Bernard Drainville Robert Aubi Alexandre Boulerice
Accuracy of inference in AUC	0.72	0.97

Table 4: Politicians liked by targets  $u_1$  and  $u_2$ .

012 different crawled Facebook profiles. It counts 4 589 different pages of politicians that are liked by 2 554 user profiles. We have performed several tests to evaluate the accuracy of inferences made by SONSAL. We have first generated a new auxiliary dataset from the original dataset by selecting 10% of the users that publish their liked pages of politicians. Then we have removed all their preferences concerning politicians. The experiments have consisted then in inferring back the deleted preferences by analysing the new auxiliary dataset. When the analyser selects the top 23 most correlated attributes to the attribute pages of politicians, the precision of inference is equal to 0.79 in AUC. In other words, in average, the inferred set of pages of politicians by the analyser is 79% similar to the ones that are really liked by the target. However, inference accuracy is only 41% when the 23 attributes are selected randomly.

## 8 Toward efficient countermeasures

SONSAI discloses friendship networks, quantifies the correlation between attributes and analyses the behaviours of users in order to infer the values of a target sensitive attribute. It helps users safeguard their privacy by identifying attributes that are correlated to the sensitive one. Users should then act on these correlated attributes to prevent these sensitive correlations. However, SONSAL is not fine-grained enough to identify which attribute values must be modified to hinder a sensitive value inference. First, to derive effective countermeasures one has to somewhat trade network utility for privacy by hiding some information. Second, a collaboration has to be established between users since their respective private data are interrelated. Figure 2 depicts an example of social network where music is correlated to politics. It is easy to correlate the Beatles music to Democrats and George Strait music to Republicans. Consequently, it is easy to infer that the target is more likely to like democrats as he likes the Beatles. One solution to hinder this inference is to delete the link between the target and the Beatles. However, the target will remain connected to Beatles through Lady Gaga that is in its turn correlated to Beatles. David and the target need to collaborate and delete their preference for Beatles in order to safeguard their privacy. It is obvious that any decision taken by David or the target can affect their mutual privacy. This shows that deriving a general solution is a challenging problem.

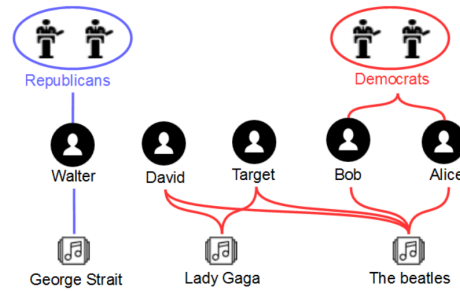


Fig. 2: Musics and politics.

## 9 Conclusion

We have presented a first prototype for self-auditing Facebook profile resilience against inference attacks. As it will probably take time before legal enforcement of privacy is fully implemented by social networks, SONSAI may contribute in the period to user awareness against privacy threats. E-reputation awareness raising has already motivated the recent development of tools to analyse how a person or a brand is perceived on social medias. “Soyez net sur le net”<sup>4</sup> and “mes datas et moi”<sup>5</sup> are such platforms oriented toward teenagers. These tools are based on explicit content processing. However many companies and their recruiters are nowadays equipped with AI systems that are also able to extract latent information by machine learning. Our proposed approach is to defeat or at least degrade the performance of these systems by letting basic users, and also entities or companies, anticipate the leakage of their information. Hence our results may contribute too to a new generation of E-reputation management applications.

**Acknowledgements.** This work is supported by MAIF Foundation<sup>6</sup>.

## References

1. Y. Abid, A. Imine, A. Napoli, C. Raïssi, and M. Rusinowitch. Online link disclosure strategies for social networks. In *Risks and Security of Internet and Systems - 11th International Conference, CRiSIS 2016, Roscoff, France, September 5-7, 2016, Revised Selected Papers*, pages 153–168, 2016.
2. Y. Abid, A. Imine, A. Napoli, C. Raïssi, and M. Rusinowitch. Two-phase preference disclosure in attributed social networks. In *Database and Expert Systems Applications - 28th International Conference, DEXA 2017, Lyon, France, August 28-31, 2017, Proceedings, Part I*, pages 249–263, 2017.
3. Y. Abid, A. Imine, and M. Rusinowitch. Sensitive attribute prediction for social networks users. In *2nd International workshop on Data Analytics solutions for Real-Life APplications March 26th, 2018 Vienna, Austria*, 2018.

<sup>4</sup> <http://www.ereputation.paris.fr/>

<sup>5</sup> <https://www.mesdatasetmoi.fr/>

<sup>6</sup> [www.fondation-maif.fr/](http://www.fondation-maif.fr/)

4. L. Bilge, T. Strufe, D. Balzarotti, and E. Kirda. All your contacts are belong to us: automated identity theft attacks on social networks. In *Proceedings of the 18th International Conference on World Wide Web, WWW 2009, Madrid, Spain, April 20-24, 2009*, pages 551–560, 2009.
5. S. Chester and G. Srivastava. Social network privacy for attribute disclosure attacks. In *International Conference on Advances in Social Networks Analysis and Mining, ASONAM 2011, Kaohsiung, Taiwan, 25-27 July 2011*, pages 445–449, 2011.
6. M. Conover, B. Gonçalves, J. Ratkiewicz, A. Flammini, and F. Menczer. Predicting the political alignment of twitter users. In *PASSAT/SocialCom 2011, Privacy, Security, Risk and Trust (PASSAT), 2011 IEEE Third International Conference on and 2011 IEEE Third International Conference on Social Computing (SocialCom), Boston, MA, USA, 9-11 Oct., 2011*, pages 192–199, 2011.
7. M. I. A. P. Garrett Brown, Travis Howe and K. Borders. Social networks and context-aware spam. In *ACM conference on computer supported collaborative*, 10, 2008.
8. N. Z. Gong and B. Liu. Attribute inference attacks in online social networks. *ACM Trans. Priv. Secur.*, 21(1):3:1–3:30, 2018.
9. N. Z. Gong, A. Talwalkar, L. W. Mackey, L. Huang, E. C. R. Shin, E. Stefanov, E. Shi, and D. Song. Joint link prediction and attribute inference using a social-attribute network. *ACM TIST*, 5(2):27:1–27:20, 2014.
10. R. Heatherly, M. Kantarcioglu, and B. M. Thuraisingham. Preventing private information inference attacks on social networks. *IEEE Trans. Knowl. Data Eng.*, 25(8):1849–1862, 2013.
11. I. Kayes and A. Iamnitchi. A survey on privacy and security in online social networks. *CoRR*, abs/1504.03342, 2015.
12. G. Kontaxis, I. Polakis, S. Ioannidis, and E. P. Markatos. Detecting social network profile cloning. In *Ninth Annual IEEE International Conference on Pervasive Computing and Communications, PerCom 2011, 21-25 March 2011, Seattle, WA, USA, Workshop Proceedings*, pages 295–300, 2011.
13. W. Mason. Politics and culture on facebook in the 2014 midterm elections, 2014.
14. B. Perozzi and S. Skiena. Exact age prediction in social networks. In *Proceedings of the 24th International Conference on World Wide Web Companion, WWW 2015, Florence, Italy, May 18-22, 2015 - Companion Volume*, pages 91–92, 2015.
15. E. Ryu, Y. Rong, J. Li, and A. Machanavaajhala. curso: protect yourself from curse of attribute inference: a social network privacy-analyzer. In *Proceedings of the 3rd ACM SIGMOD Workshop on Databases and Social Networks, DBSocial 2013, New York, NY, USA, June, 23, 2013*, pages 13–18, 2013.
16. B. S. Vidyalakshmi, R. K. Wong, and C. Chi. User attribute inference in directed social networks as a service. In *IEEE International Conference on Services Computing, SCC 2016, San Francisco, CA, USA, June 27 - July 2, 2016*, pages 9–16, 2016.
17. Z. Yin, M. Gupta, T. Weninger, and J. Han. A unified framework for link recommendation using random walks. In *International Conference on Advances in Social Networks Analysis and Mining, ASONAM 2010, Odense, Denmark, August 9-11, 2010*, pages 152–159, 2010.
18. L. Zhang and W. Zhang. An information extraction attack against on-line social networks. In *2012 International Conference on Social Informatics (SocialInformatics), Washington, D.C., USA, December 14-16, 2012*, pages 49–55, 2012.