



KTH Electrical Engineering

DIVa: Decentralized Identity Validation for Social Networks

AMIRA SOLIMAN, LEILA BAHRI, BARBARA
CARMINATI, ELENA FERRARI, AND SARUNAS
GIRDZIJAUSKAS

Stockholm 2015

LCN/EES/KTH
Insubria University, Italy

DIVa: Decentralized Identity Validation for Social Networks

Amira Soliman^{*}, Leila Bahri^{**}, Barbara Carminati^{**}, Elena Ferrari^{**} and Sarunas Girdzijauskas^{*}

^{*} Royal Institute of Technology (KTH), Sweden
{aaeh, sarunasg}@kth.se

^{**} Insubria University, Italy
{leila.bahri, barbara.carminati, elena.ferrari}@uninsubria.it

April 16, 2015

Abstract

We suggested DIVa, a decentralized, unsupervised, and association rule mining based solution for the learning of fine-grained correlations between profile attributes in Online Social Networks. These correlations can be used for identity validation purposes as has been suggested in [1]. In this report, we provide the technical details and the security analysis proofs of the DIVa model.

1 Introduction

In this work, we have based on the results of [1] that showed that there exists correlations between profile schema attributes, which if identified, can be reliably used to estimate a profile's identity trustworthiness from its attribute values only. These results being promising in enhancing identity management in the realms of OSNs, one of the major inconveniences of the method in [1] is its reliance on community feedback for the learning of those profile attribute correlations. In fact, in that work, the authors suggest a two phase system. In the first one, they exploit a supervised crowd-based learning strategy to extract profile attribute correlations that could make sense from an identity validation perspective. They do this by gathering human feedback from a group of trusted users on a centralized profiles training dataset. After these correlations are identified, they are used in the second phase, while again engaging users' (raters') feedback, to estimate the identity trustworthiness of a target profile.

Given the proliferation of profiles in an OSN and their sizes in terms of number of profiles and of users, it might be not realistically scalable to rely on trusted users feedback for the learning of those attribute correlations. More importantly, it is also hard to identify the users that are trusted and that the system can reliably rely on for the accurate learning of attribute correlations. Moreover, the supervised learning assumes the existence of a central repository

of all profiles. This might be a limitation from a privacy perspective especially with the current shift towards considering decentralized solutions for online social computing.

To cope with these issues, we suggested DIVa that is an alternative to the learning phase of [1], and that ensures unsupervised, fully automated, and fully decentralized learning. In this technical report, we provide the details of the DIVa algorithm and we prove its security properties.

2 Background

Before detailing the DIVa model, we provide the needed background notations and definitions. We consider the social network as an undirected graph $G = (V, E)$, where V is the set of nodes and E is the set of edges. $e_{ij} \in E$ denotes a relationship between nodes v_i and $v_j \in V$. We denote with $S = \{A_1, A_2, \dots, A_m\}$, the profile schema adopted in the OSN. Given a node $v_i \in V$, p_i denotes the set of its profile values: $p_i = \{p_i.a_1, p_i.a_2, \dots, p_i.a_m\}$, where $p_i.a_k$ is the value provided by v_i for $A_k \in S$.

We denote by *Local Profile Collection*, the set of profiles of a node's friends. That is, given $v_i \in V$, and $DF_i = \{v_j \in V | e_{ij} \in E\}$ representing the set of v_i 's direct friends, $LPC_i = \{p_k | v_k \in DF_i\}$ denotes the collection of their profiles and is referred to as v_i 's local profile collection. Given LPC_i , the *Local Frequent Attributes* LFA_i is the set of attributes for which the values are highly repetitive in LPC_i . Formally we define:

Definition 2.1. Local Frequent Attributes. *Let $v_i \in V$ and LPC_i be its local profile collection. Let $A_k \in S$ be an attribute from the profile schema and let $P_k^\vartheta \subseteq LPC_i$ be the set of profiles in LPC_i having the same value for attribute A_k . That is, $P_k^\vartheta = \{p_m \in LPC_i | p_m.a_k = \vartheta\}$, where ϑ is a given value. Let $LFA_i \subseteq S$ be the set of attributes such that, $LFA_i = \{A_k \in S | \frac{|P_k^\vartheta|}{|LPC_i|} \geq \epsilon\}$, where ϵ is a global predefined threshold.*

For a given pair of attributes from LFA_i , its support is defined as:

Definition 2.2. Support of an attributes pair. *Let $v_i \in V$ be a node in the OSN. Let LPC_i be its local profile collection and let LFA_i be its local frequent attributes set. Let $BA = (A_j, A_h)$ be a pair of attributes from LFA_i . The support of BA defines the percentage of co-occurrence of the same paired values for the two attributes A_j and A_h to the total number of values in LPC_i :*

$$Support(BA) = \frac{\text{values-co-occurrence}(A_j, A_h)}{\text{all-values}(A_j, A_h, LPC_i)} \quad (1)$$

Where,

$\text{values-co-occurrence}(A_j, A_h) = |\{(p_e, p_m) \in LPC_i | p_e.a_j = p_m.a_j \wedge p_e.a_h = p_m.a_h\}|$.

and,

$\text{all-values}(A_j, A_h, LPC_i) = |\{\vartheta | \exists p \in LPC_i \text{ s.t., } p.a_j = \vartheta \vee p.a_h = \vartheta\}|$

Based on Definition 2.2, we define a Local Correlated Attribute Set as follows:

Definition 2.3. Local Correlated Attribute Set - LCAS. Let $v_i \in V$. Let $LFA_i \subseteq S$ be its local frequent attributes set. Let $BA = (A_j, A_h)$ be a pair of attributes from LFA_i . BA is a local correlated attribute set, denoted as LCAS, if: $Support(BA) \geq \beta$, where β is a global predefined threshold.

The threshold β can be set to whatever value that is most representative within the settings of the target network. That is, its value depends on the underlying characteristics of the OSN and on its nature. For example, an OSN targeting professional networking and focusing on professional profile information only might require higher values for the β variable. This is because the values co-occurrence in such a well scoped OSN are expected to be higher than in general purpose OSN, for example. However, as a general guideline, we base the set up of the β threshold, as well as per the ϵ threshold as in Definition 2.1, based on the 20-to-80 rule, or what is commonly known in the statistics and economics literature as the Pareto rule [2]. The rule states that 80% of the outcomes come mostly from 20% of inputs only. This rule has been demonstrated by the Italian economist Pareto in his research and as such it was named after him. There are many economic conditions that demonstrate this rule, such as the distribution of wealth on earth that is roughly estimated as about 80% of the planet’s resources being owned and/or controlled by only 20% of the population. Besides, this rule holds also at different micro levels and is used in statistics and economics as the basis of number of theories and working real-life solutions. As such, we consider in our scenario that a correlation between attributes that is pronounced in at least 20% of the community population is reflecting an existing correlation between these attributes and is not only the result of chance or of some randomness in population distribution. That is, we set our thresholds, as will be mentioned in the experiments section (see Section 5) to the value 0.2.

3 DIVa Model

DIVa operates in three phases. First, each node performs a local learning to identify its set of Local Attribute Sets (LCAS). Second, a decentralized community detection step takes place to allow each node to know the community or the communities to which it belongs. Finally, all the nodes belonging to a community communicate their LCAS to the community’s leader node (i.e., diva node), that aggregates all the received messages to compute and disseminate the final community CAS. We provide the technical details of each of these steps in what follows.

3.1 LCAS Learning

The learning of the LCAS is carried out by each node, v_i , independently of the rest of the network and is performed following these steps:

- v_i collects the profiles of all its direct friends to compose its LPC_i and then computes its LFA_i as per Definition 2.1.
- v_i computes the support for each attribute pair from LFA_i as per Definition 2.2.

- v_i computes its *LCAS* list based on Definition 2.3. The support threshold is a globally known variable to all the nodes in the social network apriori to the start of the execution of DIVa

3.2 Community Detection

Community detection is a well established discipline in itself, across different subject areas such as biology, physics, social networks analysis, etc, that is concerned with finding tightly knit groups of nodes within a target graph. Commonly, community detection is defined as:

Definition 3.1. Community Detection. *A community detection Φ , also known as graph clustering, is a mapping*

$$\Phi : G \rightarrow G'_1 \times \dots \times G'_c \quad (2)$$

that partitions G into c non-empty, node-disjoint subgraphs $G'_1 \times \dots \times G'_c$ representing a set of communities or clusters. A widely used quality measure for community detection is the modularity Q of the clustering $\Phi(G)$ [3], which is a mapping

$$Q : \Phi(G) \rightarrow \mathbb{R} \quad (3)$$

that assigns a quality value $q \in [-0.5, 1]$ to the clustering $\Phi(G)$, as defined by

$$q := \sum_i (e_{ii} - b_i^2) \quad (4)$$

Where $b_i = \sum_j e_{ij}$, and e_{ij} is the fraction of edges in community i for which the target node of the edge lies in community j . The higher the quality value q is, the better the detected community is. One possible definition for Φ is to maximize Q over all clustering $\Phi(G)$ [3], which was shown to be an NP-hard problem [4]

Community detection for social networks have been considerably studied and the literature offers many centralized and decentralized solutions [3, 5, 6, 7]. For DIVa, a compliant solution for community detection should answer the decentralization requirement by which every node can only know of and contact its direct neighbors. Both the algorithm Louvain [8], and the work by Rahimian et al. [9] seem to fit this criteria, as in both of them every node starts as a community by itself using its node ID as its community ID. Then, every node chooses to quit its current community and join one of its neighbour's if this brings some modularity gains.

As in [8] the calculation of the modularity gain assumes having global knowledge of all the edges in the graph, which is not applicable for DIVa, we opt for [9] that suggests the idea of dominant ID by which every node changes its community ID to the dominant one in its neighbours. Figure 1 depicts the adopted decentralized community detection approach with a toy example.

When the community detection algorithm converges, every node in the network is aware of the communities to which it belongs and is also aware of the path to the leader of each its communities. A community leader is the node with the dominant ID in it. We refer to leader nodes as *diva nodes*. However, recall that in DOSNs nodes are aware only of their direct friends. Therefore, nodes

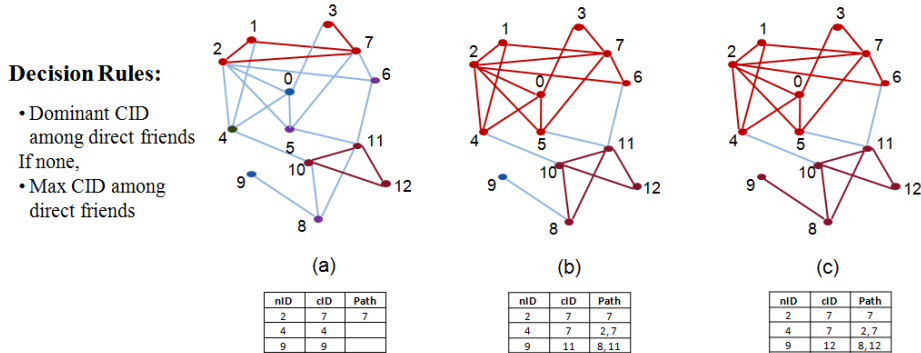


Figure 1: Visualization of community detection steps, (a) shows that every node starts by having community ID as the maximum ID among its direct friends. Next, in the following steps (b), and (c), nodes change their community ID either to dominant ID as for node 4 in step (b) or to maximum ID as for node 9. Afterwards, step (c) shows resulted communities after the algorithm converges, as no further changes are performed by the nodes.

must be aware of the path to reach their diva nodes by a hop-by-hop routing over their social ties. Accordingly, during the execution of the community detection part, nodes maintain paths leading to the diva of each community they belong to. Path construction is straightforward. First, a node checks if the diva is a direct friend. If it is not, the node creates the path by assembling the node IDs of intermediate nodes leading to the diva. Figure 2 depicts how nodes reach their communities diva nodes by following the constructed path towards them.

However, since communities in social networks are majorly overlapping, and considering that DIVa aims to learn correlations that are more representative of a node’s environment, we opt for a soft clustering approach. That is, a node can belong to more than one community and hence have more than one dominant community ID. More precisely, a node keeps track of the top dominant community IDs that it learns about when the community detection algorithm converges. As shown in Figure 2, DIVa instances have organized themselves into two overlapping communities. Further, node with the largest ID (i.e., node with darker shade) has been identified as diva node of the detected community.

3.3 CAS aggregation

The diva nodes are the ones responsible over collecting their community’s LCASes to generate from them the community’s set of CASes. This is because all the nodes in a given community know the path to the diva node. Algorithm 1 presents the process by which a diva node aggregates the supports of all the LCASes it received from the nodes in its community to form the community CASes. Basically, a diva calculates the aggregate support for each distinct LCAS it receives by summing its different supports across all the nodes in which it appeared (lines 1-5). After summing up the supports of equal LCAS from the list of all received LCAS, $LIST_L$, these occurrences of the same LCAS are removed from the list using the method *remove-repeated()* (line 4). Finally, the

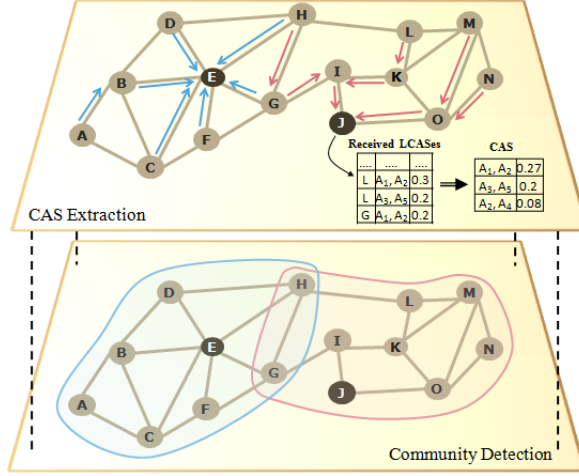


Figure 2: Identifying two overlapping communities and communicating with diva nodes to perform community-level aggregation.

final LCAS support is computed by dividing it by the total number of nodes that communicated the diva node (line 6-7). Then, if the aggregate support of an LCAS is higher or equal to a minimum support (sup_{lowest})¹, it is considered a CAS (line 9).

Algorithm 1 Aggregate Community CAS by diva node

Require: list of nodes contacted: $participants$, list of received LCAS: $LIST_L$,
minimum support: sup_{lowest}

Ensure: community CAS list, $LIST_{CAS}$

- 1: **for all** $LCAS \in LIST_L$ **do**
 - 2: $suppport(LCAS) \leftarrow LCAS.c$
 - 3: $suppport(LCAS) \leftarrow \sum_{e \in LIST_L | e=LCAS} e.c$
 - 4: $remove - repeated(LCAS, LIST_L)$
 - 5: **end for**
 - 6: **for all** $LCAS \in LIST_L$ **do**
 - 7: $suppport(LCAS) \leftarrow \frac{suppport(LCAS)}{participants.size}$
 - 8: **if** $suppport(LCAS) \geq sup_{lowest}$ **then**
 - 9: $insert(LIST_{CAS}, LCAS)$
 - 10: **end for**
-

4 Security Analysis

Assuming a malicious adversary model, DIVa enjoys two security properties. The first property guarantees that malicious adversaries cannot introduce fake CAS sets to a community unless they make the majority in that community. This is expressed in the following theorem:

¹The value of sup_{lowest} is set, as common in ARM [10], based on the achieved supports in the community.

Theorem 4.1. Let $\bar{C} \in G$, $\bar{C} = (\bar{C}.V, \bar{C}.E)$, be a community of size n ($|\bar{C}.V| = n$). Let sup_{lowest} be the lowest support by which a CAS is accepted in \bar{C} . For a new CAS, CAS_{new} to appear in \bar{C} , it must be inserted a group of fake nodes C_f that successfully join \bar{C} and that show profile information confirming CAS_{new} such that:

$$z = |C_f| \geq \frac{sup_{lowest}}{(1-sup_{lowest})} * n.$$

Proof. 1: Consider C_f of size z ($|C_f| = z$) is carrying a correlation $CAS_f = \{A, B\}$, that is unknown to the nodes in \bar{C} . Assume all C_f successfully joins \bar{C} . Therefore $\bar{C}.V = \bar{C}.V \cup C_f$ and $|\bar{C}.V| = n + z$. That is, the aggregate support of CAS_f in \bar{C} would be: $support(CAS_f) = \frac{values-co-occurrence(A,B)}{n+z}$. Since all nodes in C_f carry the correlation in CAS_f that is unknown to \bar{C} initial n nodes, the support for CAS_f will be: $support(CAS_f) = \frac{z}{n+z}$. According to the proposed method, for CAS_f to be recognized as a CAS in \bar{C} ($support(CAS_f) \geq sup_{lowest}$), this inequality shall hold: $z \geq \frac{sup_{lowest}}{(1-sup_{lowest})} * n$. \square

By Theorem 4.1, we can clearly see that the required adversary effort to introduce a fake CAS to a community is directly proportional to the size of the community and to the lowest support by which it accepts a new CAS. That is, malicious nodes can introduce new CAS to a community only if their number is big enough compared to legitimate nodes. Also, it is clear that the higher the support threshold is, the higher the percentage of fake nodes to legitimate ones is needed for an attack to succeed.

The second security property of DIVa is related to the resilience of its valid CAS sets once they are identified. The following theorem formalizes this property:

Theorem 4.2. Let $\bar{C} \in G$, $\bar{C} = (\bar{C}.V, \bar{C}.E)$, be a community of size n ($|\bar{C}.V| = n$). Let sup_{lowest} be the lowest support by which a CAS is accepted in \bar{C} . For a valid CAS, CAS_{valid} with support S_v , to disappear from \bar{C} , it must be inserted in \bar{C} a group of fake nodes, C_f , that does not have profile information confirming CAS_{valid} such that:

$$z = |C_f| > \frac{S_v * n}{sup_{lowest}} - n.$$

Proof. 2: Let $CAS_v = \{A, B\}$ be a valid CAS in \bar{C} with aggregate support S_v : $S_v = \frac{m}{n} \geq sup_{lowest}$, where $m = values - co - occurrence(A, B)$. Let C_f of size z ($|C_f| = z$) be not carrying the correlation between attributes A and B. Assume all the nodes in C_f successfully join \bar{C} . Therefore $\bar{C}.V = \bar{C}.V \cup C_f$ and $|\bar{C}.V| = n + z$. That is, the aggregate support of CAS_v in \bar{C} becomes: $S_{v1} = \frac{values-co-occurrence(A,B)}{n+z}$. Since all nodes in C_f do not carry the correlation in CAS_v , $values - co - occurrence(A, B)$ is still equal to m ; therefore $S_{v1} = \frac{m}{n+z}$. For CAS_v to no more be a valid CAS, its new support shall be: $S_{v1} < sup_{lowest}$. That is, $\frac{m}{n+z} < sup_{lowest}$. From where $m < sup_{lowest} * (n + z)$. Dividing the inequality by n ($n \in \mathbb{N}^+$ and $n > 0$), we get: $\frac{m}{n} < \frac{n * sup_{lowest} + z * sup_{lowest}}{n}$. Therefore, dividing the inequality by the positive number sup_{lowest} gives, $\frac{S_v}{sup_{lowest}} < \frac{n+z}{n}$. From that, $z > \frac{S_v * n}{sup_{lowest}} - n$. \square

By Theorem 4.2, a valid CAS is as vulnerable as is the percentage of nodes in the community that carry it. That is, the lower the support of a CAS is, the smaller the number of nodes in the community carrying it. Hence, the easier it becomes for an adversary to lessen it. However, it is also interesting to note that a weak CAS (with weak support) should initially be not of much importance to the identity validation. That is to say, this result is crucial to be considered when devising the validation phase by adding weights to the different CAS available in a community when the nodes want to use them to evaluate the trustworthiness of new potential contacts, for example.

5 Related Work

Personal identity, its formation process, and its components have been the subject of scientific discussion and research work across multiple scientific disciplines such as sociology [11], psychology [12], criminology [13], etc. With the growth of the Internet as a world wide virtual platform that connects data, devices, people, etc, new dimensions for humans' interactions have seen the light of day. Online human to human interactions developed from basic open chat rooms connecting virtual personas to nowadays popular and widespread OSNs with more sophisticated communication and data exchange forms. Within these emerging online socializing realms, identity has had its place as a pole of attraction for researchers from different disciplines. From a computer science perspective, resolving identities in the sense of differentiating between real and fake ones has been the main research concern related to identity. As a result, we find many pieces of work studying and formalizing online identities patterns with the objective of classifying them as good or bad.

This gave birth to classifications for bad identities such as sybil (a fake identity operated, along with many other sybils, by one same physical entity)[14, 15], clone (an identity created by a malicious entity based on information collected about another honest entity)[16], compromised (an honest identity but taken control of by a malicious entity)[17], etc. Therefore, we find works such as SybilyGuard [14] and SybilLimit [15] that study OSN topological properties to detect sybil identities. We find [16], a framework for the detection of clone identities based on attribute and friends' network similarities, or [17] where the authors address identity theft across multiple social networks. These works, with others on the same line, share the common goal of detecting malicious nodes classified under formalized identity attack trends. However, identity concerns on OSNs go beyond binary classification. For example, some 'good' identities are created with the aim of fooling a category of users, such as child abuse over social networks [18][19].

Studying identity related attacks is unquestionably an important thread of work, but there is also a parallel need for empowering users themselves to evaluate the trustworthiness and the validity of the online identities they interact with. The literature provides us with works such as [20] where it is suggested to evaluate an identity on a given network based on feedback of her connections on another one. [21] suggests people to people recommendations for friendships' acceptance by relying on collaborative filtering techniques. In [22], users are suggested to be identified from their typing patterns; whereas chatting patterns are exploited for users' identification in [23]. More recently, [24] suggests identi-

fyng users across networks based on geo-location and time-stamp information attached to their posts and on their writing styles. All these pieces of work still do not provide users with a framework to evaluate, by themselves, their perceived trustworthiness of their new online contacts. At this level comes [1] to suggest using community feedback to assign trustworthiness levels to identities on a social network. More precisely, identities in [1] are validated based on community validations of homogeneity between values of some defined correlated profile attributes. However, [1] relies on a central repository of all the profiles of the OSN, on the existence of a group of trusted users for the learning of the correlated profile attributes, and on the responsiveness of the OSN community to evaluate available target identities.

In this work, we adopt the same idea leveraged on in [1], but based on more realistic and commonly observed assumptions. Moreover, our solution is fully automated, fully decentralized, and proves efficiency and effectiveness with real OSN data. To the best of our knowledge, this work is a first in addressing identity validation based on fully unsupervised and fully decentralized learning from profile information only.

6 Conclusion

We have suggested a decentralized association rule mining based approach for the learning of correlations between profile attributes within OSN communities. These correlations can be exploited, as has been shown in [1] to allow users estimate identity trustworthiness values for their potential contacts on the OSN.

We have evaluated DIVa, extensively, using two real profile datasets from two OSN giants: Facebook and Google+. The evaluation results show that DIVa scales well, succeeds in learning profile correlations that can be exploited for identity validation, and outperforms a centralized learning in unveiling fine-grained and community-specific correlations that are easily discriminated in a centralized learning approach but that are also better descriptive for given communities.

Acknowledgment

This work is under the umbrella of the iSocial EU Marie Curie ITN project (FP7-PEOPLE-2012-ITN). The authors also thank Naeimeh Laleh for her help with Facebook data sanitizing.

References

- [1] L. Bahri, B. Carminati, and E. Ferrari, “Community-based identity validation on online social networks,” in *Distributed Computing Systems (ICDCS), 2014 IEEE 34th International Conference on*. IEEE, 2014, pp. 21–30.
- [2] J. M. Herbener, “The pareto rule and welfare economics,” *The Review of Austrian Economics*, vol. 10, no. 1, pp. 79–106, 1997.

- [3] M. E. Newman, “Modularity and community structure in networks,” *Proceedings of the National Academy of Sciences*, vol. 103, no. 23, pp. 8577–8582, 2006.
- [4] U. Brandes, D. Delcling, M. Gaertler, R. Gorke, M. Hoefer, Z. Nikoloski, and D. Wagner, “On modularity clustering,” *Knowledge and Data Engineering, IEEE Transactions on*, vol. 20, no. 2, pp. 172–188, 2008.
- [5] F. Radicchi, C. Castellano, F. Cecconi, V. Loreto, and D. Parisi, “Defining and identifying communities in networks,” *Proceedings of the National Academy of Sciences of the United States of America*, vol. 101, no. 9, pp. 2658–2663, 2004.
- [6] H. Meyerhenke, B. Monien, and T. Sauerwald, “A new diffusion-based multilevel algorithm for computing graph partitions of very high quality,” in *Parallel and Distributed Processing, 2008. IPDPS 2008. IEEE International Symposium on*. IEEE, 2008, pp. 1–13.
- [7] P. Sanders and C. Schulz, “Distributed evolutionary graph partitioning.” in *ALLENEX*. SIAM, 2012, pp. 16–29.
- [8] V. D. Blondel, J.-L. Guillaume, R. Lambiotte, and E. Lefebvre, “Fast unfolding of communities in large networks,” *Journal of Statistical Mechanics: Theory and Experiment*, vol. 2008, no. 10, p. P10008, 2008.
- [9] F. Rahimian, S. Girdzijauskas, and S. Haridi, “Parallel community detection for cross-document coreference,” in *Web Intelligence (WI) and Intelligent Agent Technologies (IAT), 2014 IEEE/WIC/ACM International Joint Conferences on*, vol. 2. IEEE, 2014, pp. 46–53.
- [10] R. Agrawal, T. Imieliński, and A. Swami, “Mining association rules between sets of items in large databases,” in *ACM SIGMOD Record*, vol. 22, no. 2. ACM, 1993, pp. 207–216.
- [11] J. E. Stets and P. J. Burke, “A sociological approach to self and identity,” *Handbook of self and identity*, pp. 128–152, 2003.
- [12] R. E. Spears, P. J. Oakes, N. E. Ellemers, and S. Haslam, *The social psychology of stereotyping and group life*. Blackwell Publishing, 1997.
- [13] M. J. Lynch, R. J. Michalowski, and W. B. Groves, *The new primer in radical criminology: Critical perspectives on crime, power, and identity*. Criminal Justice Press Monsey, NY, 2000.
- [14] H. Yu, M. Kaminsky, P. B. Gibbons, and A. Flaxman, “Sybilguard: defending against sybil attacks via social networks,” in *ACM SIGCOMM Computer Communication Review*, vol. 36, no. 4. ACM, 2006, pp. 267–278.
- [15] H. Yu, P. B. Gibbons, M. Kaminsky, and F. Xiao, “Sybillimit: A near-optimal social network defense against sybil attacks,” in *Security and Privacy, 2008. SP 2008. IEEE Symposium on*. IEEE, 2008, pp. 3–17.

- [16] L. Jin, H. Takabi, and J. B. Joshi, “Towards active detection of identity clone attacks on online social networks,” in *Proceedings of the first ACM conference on Data and application security and privacy*. ACM, 2011, pp. 27–38.
- [17] B.-Z. He, C.-M. Chen, Y.-P. Su, and H.-M. Sun, “A defence scheme against identity theft attack based on multiple social networks,” *Expert Systems with Applications*, vol. 41, no. 5, pp. 2345–2352, 2014.
- [18] C. Hope, “Facebook is a ‘major location for online child sexual grooming’, head of child protection agency says,” Oct. 2013, the Telegraph. [Online]. Available: <http://www.telegraph.co.uk/technology/facebook/10380631/Facebook-is-a-major-location-for-online-child-sexual-grooming-head-of-child-protection-agency-says.html>
- [19] M. Chorley, “How facebook and social networking sites are used by child abuse gangs to groom victims for ‘sex parties’,” Nov. 2012, mail Online. [Online]. Available: <http://www.dailymail.co.uk/news/article-2236208/How-Facebook-social-networking-sites-used-child-abuse-gangs-groom-victims-sex-parties.html>
- [20] M. Sirivianos, K. Kim, J. W. Gan, and X. Yang, “Assessing the veracity of identity assertions via osns,” in *Communication Systems and Networks (COMSNETS), 2012 Fourth International Conference on*. IEEE, 2012, pp. 1–10.
- [21] X. Cai, M. Bain, A. Krzywicki, W. Wobcke, Y. S. Kim, P. Compton, and A. Mahidadia, “Collaborative filtering for people to people recommendation in social networks,” in *AI 2010: Advances in Artificial Intelligence*. Springer, 2011, pp. 476–485.
- [22] P. Chairunnanda, N. Pham, and U. Hengartner, “Privacy: Gone with the typing! identifying web users by their typing patterns,” in *Privacy, security, risk and trust (passat), 2011 ieee third international conference on and 2011 ieee third international conference on social computing (social-com)*. IEEE, 2011, pp. 974–980.
- [23] G. Roffo, C. Segalin, A. Vinciarelli, V. Murino, and M. Cristani, “Reading between the turns: Statistical modeling for identity recognition and verification in chats,” in *Advanced Video and Signal Based Surveillance (AVSS), 2013 10th IEEE International Conference on*. IEEE, 2013, pp. 99–104.
- [24] O. Goga, H. Lei, S. H. K. Parthasarathi, G. Friedland, R. Sommer, and R. Teixeira, “Exploiting innocuous activity for correlating users across sites,” in *Proceedings of the 22nd international conference on World Wide Web*. International World Wide Web Conferences Steering Committee, 2013, pp. 447–458.