

Methods of Profile Cloning Detection in Online Social Networks

Michał Zabielski¹, Rafał Kasprzyk¹, Zbigniew Tarapata^{1,a}, Krzysztof Szkółka¹

¹ *Institute of Computer and Information Systems, Military University of Technology, Warsaw, Poland*

Abstract. With the arrival of online social networks, the importance of privacy on the Internet has increased dramatically. Thus, it is important to develop mechanisms that will prevent our hidden personal data from unauthorized access and use. In this paper an attempt was made to present a concept of profile cloning detection in Online Social Networks (OSN) using Graph and Networks Theory. By analysing structural similarity of network and value of attributes of user personal profile, we will be able to search for attackers which steal our identity.

1 Introduction and Motivation

Problem of the privacy violation on the Internet exists from the time, when first mechanisms to store user personal data was developed [1]. This issue is an important part of discussion field on the Online Social Networks (OSN) domain, which are depending mostly on the Web.

The results of analyses indicate a wide range of techniques, that allows people to take an unauthorized access to user personal data. In general, we can classify them into two groups: local and global. Local mechanisms of privacy violation in OSNs are focused on analysing person and their direct neighbourhood i.e. friends in this network. This allows us, for example, to reveal hidden attribute values of user social profile. An example of this method is described in [2].

Using global techniques, by contrast, is a different approach to deal with that issue. They are based on the entire structure of the acquired social network. This allows for using additional groups of algorithms, starting from machine learning techniques to the ones based on the similarity of graphs and networks. As a result of such actions, we can gain data not only directly from the user's profile and other people within the network, but also from the similarities of network structures across multiple Online Social Networks. This method is capable of detecting phenomena, which are difficult to be examined with the use of local methods, such as cloning profiles. An example of using global techniques for finding evidence of identity theft was presented in [3].

A profile of the user in the Online Social Network can be represented as a set of features and links between other profiles that describe the person in that network. Depending on the type of OSN, a set of attributes and relationships that make up the profile may be different. This proves there is no profile of a unique social network

for a person on the Internet today, which gives the possibility of impersonating that particular person in different OSNs.

The complexity of the considerations about violating privacy in Online Social Networks intensifies the phenomenon of social network aggregation. It is a situation in which a person can have totally different profiles in different social networks. This potentially increases the chance of cloning a profile or impersonating another user, hindering the possibility for reliable comparisons of users' profiles in Online Social Networks. These and other problems led to creation of mechanisms for detecting violations of privacy.

As there is no clear definition of privacy, it can be variously interpreted. For the purposes of this study, privacy can be understood as an opportunity for individuals or group of individuals to maintain their personal data as well as habits and behaviours which are undisclosed to the public [4].

But in order to design solutions mentioned before, it is necessary to understand the phenomena with which we are dealing. In addition, it is complicated by the fact that violating privacy is not guided by one of the predetermined procedure. This process consists of a number of factors: from the structural description of the individual with the use of a social network profile to the dynamics of building relationships with others, and thus formation of links between the nodes. This proves that there are at least a few methods of privacy violation. One of them is the cloning the user profile.

The method of cloning a user profile in a social network involves impersonating a victim by creating the best possible copies of their social profile. In the simplest scenario, it includes building relationships and sharing profile data in the same manner as in the victim profile from another social network. As a result, an attacker is able to steal the identity of the victim in another social network.

^a Corresponding author: zbigniew.tarapata@wat.edu.pl

The paper is organized as follows. In section 2 we present short background for profile cloning detection problem being considered. Section 3 contains proposition for revealing hidden attributes of social profile. In section 4 we describe node similarity measure function. Section 5 contains description of graph structural similarity. In section 6 we present the concept of profile cloning detection using methods described in sections 3, 4 and 5.

2 Background

Nowadays there are some methods that allows us to detect a profile cloning event [3,5]. Most of them are based on Machine Learning algorithms [6,7]. This, in most cases, results in analysing persons in social network without taking into account aspects of relations between considered person and his friends.

Another way of thinking about aspect of profile cloning is designing a set of methods how to prepare this kind of cyberattack. It is explained by the fact that having a knowledge about how hackers works gives us an opportunity to detect susceptibilities of their method, which eventually implicates developing a solution to stop them. This approach is well described in [8].

Using the results from presented solutions, we developed method that take into account not only static social profile part, such attributes values, but also a dynamical part, which consists of relation between person and other people in OSN. Implemented methods of profile cloning attacks give us basics to analyse them and take into consideration in our methods of detecting them.

3 Revealing hidden attribute values of social profile

For detecting a profile cloning it is important to have a complete information about user personal profile in Online Social Network.

Most often the reason of hidden attribute value are mechanisms, provided by Online Social Networks, which give us an opportunity to show values of attributes only for user-defined group of people in network. In most cases this approach is implemented by Access Control Lists (ACL). Despite their usability in protecting privacy aspect, this situation complicates our analysis in sense of profile cloning detection.

To provide a situation, when we know all about attribute values for each personal profile in particular Online Social Network, we provide a method, based on network model and k-nearest neighbours' algorithm, which allows us to estimate a hidden values of attributes for appropriate personal profile. A general procedure for this show Figure 1.

First of all, we have to provide analysing network, where nodes represents persons and edges describe a relation between two people. Moreover, a set of functions on nodes, which represents an attributes of social profile and one function on edge, which defines force of relation has to be defined. In that network we choose a node, for which we want to reveal hidden attribute values – we denote it as w_s .

Next, we get the nearest neighbours of node w_s – that is an set of direct friends of considered person – and calculate for them an similarity measure between person w_s and his neighbors. The similarity function can be any function that gives a higher value for nodes that are similar in sense of attribute values. In our approach, we use function described by equation:

$$s(w_s, w) = d(w_s, w) * \sum_{a_i \in A} p_{a_i}(w_s, w) \quad (1)$$

where:

$d(w_s, w) \in (0,1]$ – force of relation between person w_s and w . It can be measured, for example, as a number of messages sent to person w by w_s divided by numbers of whole messages sent in OSN by node w_s ,

$$p_{a_i}(w_s, w) = \begin{cases} 1 & \text{if } a_i(w) = a_i(w_s) \wedge a_i(w) \text{ known} \\ 0 & \text{if } a_i(w) \neq a_i(w_s) \vee a_i(w) \text{ unknown} \end{cases} \quad (2)$$

$a_i(w)$ – value of the i -th attribute of personal profile of person w ;

A – set of personal profile attributes.

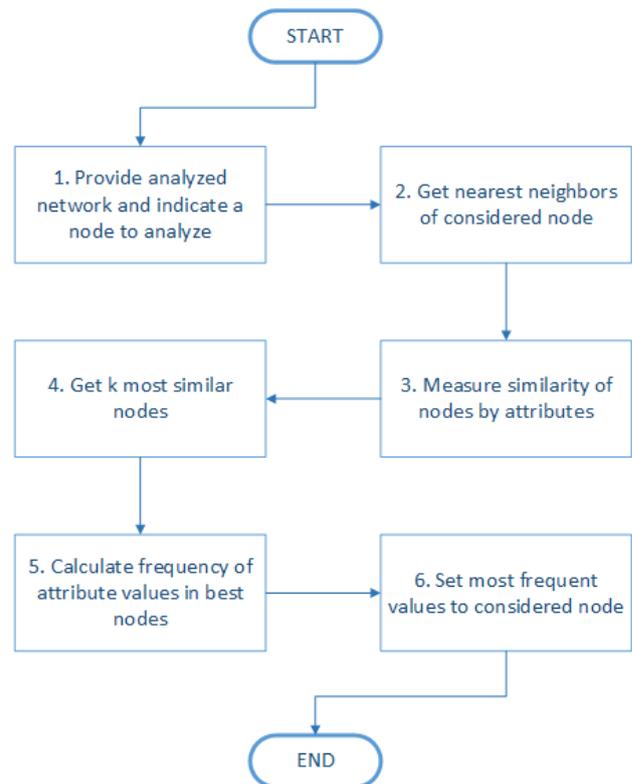


Figure 1. Revealing hidden attributes of social profile procedure.

When a similarity is calculated for every neighbour, we choose k nodes (the best nodes) with the highest similarity measure to consider their attributes values. Using a parameter k gives us confidence that calculated value of hidden attributes for node w_s will not be dominated by only one in Online Social Network. After choosing appropriate nodes, we calculate a frequency of appearing every value from the best nodes that are connected with hidden attributes. The most frequent value is taking as an answer. If we have more options to

choose a value, we take a value with better similarity measure. We can repeat this procedure for every attribute that value is hidden.

The whole algorithm of revealing hidden values of attribute can be used for every node in a network, giving us an opportunity to estimate value for every attribute in every personal profile in OSN.

There are some important issues, that are worth discussing for best use of this model. Firstly, because this method is partially derived from k -neighbours' algorithm, it is good to choose a parameter k that is odd. This approach ensures that chance to get a value using only a frequency of appearing value in best nodes is high, which gives us better estimation of attribute hidden value. From our experiments, it is also good to choose parameter k from range [3,9], depend on nearest neighbour amount of considered node. The more neighbours' person has the higher value of parameter k should be. Secondly, it is good to know which kind of attribute values we can get. A similarity measure defined in (1) is good for discrete values, but it can give worst results for continuous values. If we have knowledge about a value range for attributes, we can use mix of similarity measures, which probably gives us a better estimation for continuous values of attributes.

4 Node similarity

After revealing all values of attributes in user social profile, we have to provide a measure, that helps us to detect an appropriate person from pattern network in our analysed network. To achieve this, we have to use a function which fulfil following conditions:

- it allows to simply compare attribute values between persons;
- it should be easy to interpret, which means that by using them we will be able to decide which personal profile, in sense of attribute values, is most similar to considered person's profile without any additional transformation of measuring function values;
- it should use a characteristic which are easy to obtain from network model.

According to the presented assumptions, for our purposes, we use a node similarity function, which is described by (3):

$$ID(w_s, w) = \sum_{i=1}^m I_i p_{a_i}(w_s, w) \in [0,1] \quad (3)$$

where:

m – number of attributes in personal profile;

$I_i \in [0,1]$ – the importance of the i -th attribute $\sum_{j=1}^m I_j = 1$;

$p_{a_i}(w_s, w)$ – comparing function defined by (2).

Function $ID(w_s, w)$ meet all the conditions defined for node similarity measure function. By using only simple comparing of attribute values, it is reduced to compare vectors of attributes that creates user personal profile. Interpretation of (3) is easy, because $ID(w_s, w)$ can take values from range [0,1]. Finally, using only attribute value and significance, we provide a method, which base only on characteristics described directly on network.

Having a measure of similarity between nodes we can obtain a global measure of node similarity between two networks (pattern and analysed networks) S and S' , which gives us a better understanding of similarity of whole network. To do this we can formulate and solve (using for example Hungarian algorithm) optimal assignment problem to find the best allocation matrix $X = [x_{ij}]_{|W_G| \times |W_{G'}|}$ of nodes from network S to nodes from network S' :

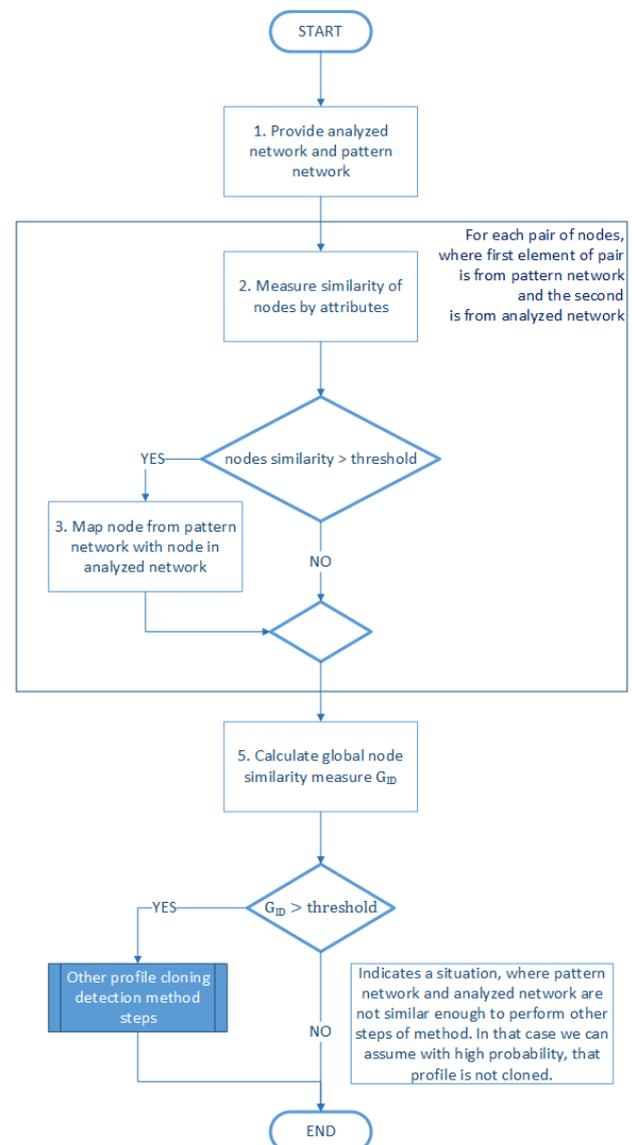


Figure 2. Algorithm of mapping persons from pattern network to analysed network in profile cloning detection method.

$$G_{ID} = \sum_{i=1}^{|W_G|} \sum_{j=1}^{|W_{G'}|} ID(i, j) \cdot x_{ij} \rightarrow \max \quad (4)$$

with constraints:

$$\sum_{i=1}^{|W_G|} x_{ij} \leq 1, j = \overline{1, |W_{G'}|} \quad (5)$$

$$\sum_{j=1}^{|W_{G'}|} x_{ij} \leq 1, i = \overline{1, |W_G|} \quad (6)$$

$$\bigwedge_{i \in \{1, \dots, |W_G|\}} \bigwedge_{j \in \{1, \dots, |W_{G'}|\}} x_{ij} \in \{0,1\}$$

where:

W_G – set of nodes in pattern network S ;

$W_{G'}$ – set of nodes in analysed network S' .

Thanks that, we will be able to assess a “percentage” of similarity between considered networks and their persons.

It is important to remember that solutions provided in (3) and (4) are only some kind of proposals how to calculate a similarity at the vertices and network level, which we use in our specific implementation. For profile cloning detection method, which will be described in detail in chapter 6, we can define any other methods, which fulfils the conditions described earlier. Some of them are defined in details in [9]. In general, a procedure which is used in our profile cloning method is illustrated in Figure 2.

5 Graph structural similarity

There are a lot of measures that we can use to assess similarity of graphs nowadays. In this section, we present most popular of it. Other methods and models presented in this paper are independent of the algorithm that is used for measuring graph similarity.

We propose general classification of graph similarity measures as shown in Figure 3.

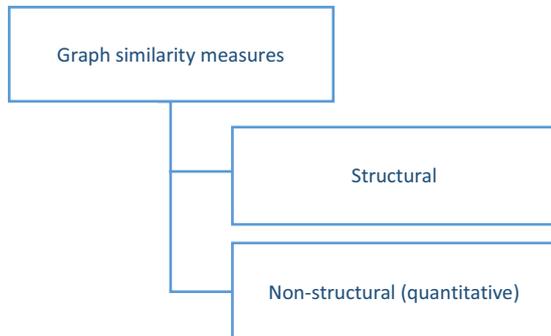


Figure 3. Taxonomy of similarity measures.

In short, structural similarity measures are used for graph and quantitative similarity measures for network (weighted graph) analysis. In the former only graph structure (topology) is examined. On the other hand, with quantitative similarity measures, we can consider node and edge weights (attributes). Graph isomorphism [12][13], maximum common subgraph [12], minimum common supergraph [12], edit distance [12][13], topological measures [12][13] and iterative methods [10][12][13] are examples of structural graph similarity measures. Non-structural graph similarity measures, in most cases, are based on iterative structural measures. Examples of these are presented in [5][11].

For the purpose of this research we propose graph structural similarity measure based on node neighbour matching. It is classified as one of iterative measures. The idea of it is as follows: two nodes are considered similar if they neighbours are similar [8][9][14].

Let structural similarity measure over the nodes of two networks S and S' be represented as [9][13]:

$$P_{node}(S, S') = [s_{ij}]_{|W_G| \times |W_{G'}|} = \lim_{k \rightarrow +\infty} Z_{2k} \tag{7}$$

where:

$$Z_{k+1} = \frac{BZ_k A^T + A^T Z_k B}{\|BZ_k A^T + A^T Z_k B\|}, k \geq 0,$$

A and B - transition matrices of S and S' , $Z_0=1$ (matrix with all elements equal 1); A^T - matrix A transposition;

$\|y\| = \sqrt{\sum_{i=1}^{|W_G|} \sum_{j=1}^{|W_{G'}|} y_{ij}^2}$. Element z_{ij} of the matrix Z describes similarity score between the i -th node of the S' and the j -th node of the S . The greater value of z_{ij} the greater similarity between the i -th node of the S' and the j -th node of the S .

Having matrix $P_{node}(S, S')$, optimal assignment problem can be formulated and solved (using for example Hungarian algorithm) to find the best allocation matrix $X = [x_{ij}]_{|S| \times |S'|}$ of nodes from graphs S and S' :

$$P(S, S') = \sum_{i=1}^{|W_G|} \sum_{j=1}^{|W_{G'}|} s_{ij} \cdot x_{ij} \rightarrow \max \tag{8}$$

with constraints:

$$\sum_{i=1}^{|W_G|} x_{ij} \leq 1, j = \overline{1, |W_{G'}|} \tag{9}$$

$$\sum_{j=1}^{|W_{G'}|} x_{ij} \leq 1, i = \overline{1, |W_G|} \tag{10}$$

$$\bigwedge_{i \in \{1, \dots, |W_G|\}} \bigwedge_{j \in \{1, \dots, |W_{G'}|\}} x_{ij} \in \{0, 1\}$$

The $P(S, S')$ is the value of structural similarity measure of graphs S and S' .

It is worth notice that analysed nodes could be part of the same or different graph. In the proposed method only shared nearest neighbours (first level neighbours) are taken into consideration. It can be easily extended to include neighbours of nearest neighbours (second level neighbours) and further. As an another extension we could propose to include node neighbours weights (attributes) when similarity of nodes is calculated.

6 The concept of profile cloning detection

For better understanding of provided approach, we introduce a following notation:

S – pattern network;

S' – analysed network;

$t_s \in [0, 1]$ – threshold parameter for structural similarity of networks;

$t_{ID} \in [0, 1]$ – threshold parameter for nodes global similarity;

w_s – person in analysed network that we are going to check if his identity was stole nor not;

$P(S, S')$ – structural similarity of network S and S' ;

i^{max} – maximum amount of iterations that algorithm should take;

ε - minimal allowed increase of structural similarity, that continues iteration of algorithm;
 $t_s^{old} \in [0,1]$ – structural similarity value between network S and S' from previous iteration.

With this in mind, we provide an approach to detect profile cloning, by using a method illustrated on a diagram in Figure 4.

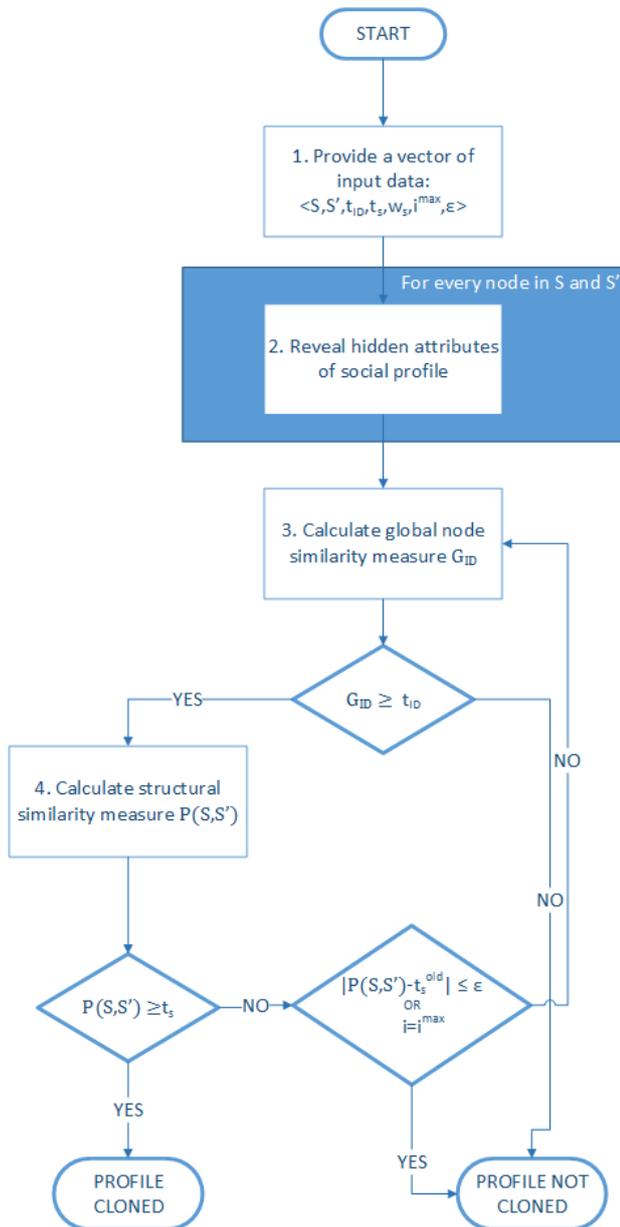


Figure 4. Profile cloning detection algorithm.

The general procedure is as follows. After providing an input data, we are using a method to reveal hidden values of attributes for every node in analysed network. A mechanism how to do it was described in chapter 3. Next, we calculate a global node similarity function to check whether analysing network is similar enough to pattern network. If this similarity is low, then we can assume that in analysed network there are not enough information to detect a profile cloning or the structure of this network, in sense of social profile, is too different so that the possibility of cloning a profile is low.

In situation when the G_{ID} value is greater than defined threshold, we can consider another steps of method. In that case, we calculate a structural similarity measure between pattern network and structural network. If this value is high enough, then we can assume, that the personal profile of person w_s is nearly the same in analysed network than in pattern network. This situation could happen in two scenarios. First is when the considered person create an account in analysed OSN and started to make a relations with his friends. Another case is when someone creates an account in analysed network and tries to recreate a social profile of person w_s . Second situation appears when a profile cloning has place, so we can say that we detect profile cloning.

Other part of algorithm consists of conditions to stop the method when a profile cloning didn't have a place. For this purpose we use two approaches. First is to set an maximum for amount of iteration. This gives us an opportunity to stop executing of method, but will also cause a risk, that algorithm will stop before detect a profile cloning. Another option used in our solution is to check whether an increasing of similarity measures are significant enough to continue searching. In that case, we assume that if we cannot obtain appropriate similarity of networks, then there is not any possibility that user's personal profile, in sense of attributes and relations between him and other people in OSN, was cloned. Because there is a risk that a second approach will not converge enough to stop executing an algorithm, a composition of first and second solution is used.

7 Conclusions

Presented approach shows that detecting a profile cloning using analytical methods is possible. This allows us to develop a solutions, which gives us opportunity to automate the process of discovering identity stealing in Online Social Networks. Moreover, a revealing hidden attributes values model, which is a part of detecting profile cloning procedure, can be also useful in analysing networks with lack of information about users and gives us also opportunity to analyse a procedure used by attackers to detect hidden attribute values. If we know how hackers are working to steal people identity, we can define a vulnerabilities of their action, which helps us to design methods to prevent from an unauthorized access to our attribute values.

Despite the effectiveness of considered approach, it is still place to extend presented method. For example, it will be good to design specialized structural similarity measures dedicated directly for Online Social Networks, which take into account a social profile characteristics and attributes. In our future work, we also plan to add ability to predict creation of relation between people in considered OSN. An example of predictor for making relations between people in Online Social Networks is described in details in [16]. This allows us to forecast an incident of profile cloning before this process will be completed by attacker. Thanks that, we will be able to prevent hacker from stealing user identity, which will definitely increase a value of implemented solution.

References

1. L. Faith Cranor, Internet privacy, *Communications of the ACM*, **42**, p. 28-38, (1999)
2. M. Mo, D. Wang, B. Li, D. Hong, I. King, Exploit of Online Social Networks with Semi-Supervised Learning, *The 2010 International Joint Conference on Neural Networks (IJCNN)*, 1-8, (2010)
3. M. Khayyambashi, F. Rizi, An approach for detecting profile cloning in online social networks, *e-Commerce in Developing Countries: With Focus on e-Security (ECDC)*, (2013)
4. H. Jeff Smith, *Managing Privacy: Information Technology and Corporate America*, UNC Press Books, (1994)
5. G. Kontaxis, I. Polakis, S. Ioannidis and E. P. Markatos, Detecting social network profile cloning, *IEEE International Conference on Pervasive Computing and Communications Workshops (PERCOM Workshops)*, Seattle, WA, 295-300, (2011)
6. X. Zhu, Z. Ghahramani, J. Lafferty, Semi-Supervised Learning Using Gaussian Fields and Harmonic Functions, *Proceedings of the Twentieth International Conference on Machine Learning*, 912-919, (2003)
7. A. Blum, Sh.Chawla, *Learning from Labeled and Unlabeled Data using Graph Mincuts*, Computer Science Department, paper 163, (2001)
8. L. Bilge, T. Strufe, D. Balzarotti, E. Kirida, All your contacts are belong to us: automated identity theft attacks on social networks, *Proceedings of the 18th international conference on World wide web (WWW '09)*, ACM, New York, NY, USA, 551-560, (2009)
9. Z. Tarapata, R. Kasprzyk, An application of multicriteria weighted graph similarity method to social networks analyzing, *Proceedings of the 2009 International Conference on Advances in Social Network Analysis and Mining*, Athens (Greece), IEEE Computer Society, 366-368, (2009)
10. V. Blondel, A. Gajardo, M. Heymans, P. Senellart, P. Van Dooren, A Measure Of Similarity Between Graph Vertices: Applications To Synonym Extraction And Web Searching, *SIAM Review*, **46** (4), 647-666, (2004)
11. M. Nikolić, Measuring similarity of graph nodes by neighbor matching, *Intelligent Data Analysis*, **16** (6), 865-878, (2012)
12. G. Jeh, J. Widom, SimRank: a measure of structural-context similarity, *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*, Edmonton, 538-543, (2002)
13. Z. Tarapata, Multicriteria weighted graphs similarity and its application for decision situation pattern matching problem, *Proceedings of the 13th IEEE/IFAC International Conference on Methods and Models in Automation and Robotics*, 1149-1155, (2007)
14. L. Zager, *Graph similarity and matching*, PhD thesis, MIT, (2005)
15. C. Bartosiak, R. Kasprzyk, Z. Tarapata, Application of Graphs and Networks Similarity Measures for Analyzing Complex Networks, *Biuletyn Instytutu Systemów Informatycznych*, **7**, 1-7, (2011)
16. D. Liben-Nowell, J. Kleinberg, The Link-Prediction Problem for Social Networks, *Journal of the American Society for Information Science and Technology*, **58** (7), 1019-1031, (2007)