

Disambiguating Technique for Protecting Data in Social Network

P. Amudha

School of Engineering, Avinashilingam Institute for Home Science and Higher Education for Women, Coimbatore, Tamil Nadu

S.Sivakumari

School of Engineering, Avinashilingam Institute for Home Science and Higher Education for Women, Coimbatore, Tamil Nadu

ABSTRACT

A social network is a virtual environment powered by web technologies that enables users to publish and share all kinds of information and services with a global audience. Perhaps, social media is the vital area of the Internet, but, being open and social generates legitimate concerns about confidentiality and security. Private data is very valuable and privacy-preservation techniques provide access to sensitive contents by revealing more or less information to the users based on their credentials. This work focuses on the preservation of the output's utility and, the detection and sanitization of terms that may cause disclosure of sensitive data due to semantic correlation.

KEYWORDS-Social network, social media, privacy-preservation, sanitization

1. INTRODUCTION

A social network is made up of a set of social players and other social communications between players. Recently social networking sites such as Face book, Twitter, LinkedIn etc., have gained large popularity. Participating users of these sites form online social network, which provides sharing, organizing and finding contents and contacts. The relation between privacy and use of social network sites is very close and delicate [2]. Private information is very valuable when the information of many people gathered on social network sites. The popularity of online social network application increases serious problems about the security and privacy of their users. Privacy related with online social networking is influenced by the level of credentials of the data provided to its beneficiaries and its users. Even social networking sites that do not explicitly render the identities of their users may deliver adequate information to identify the profile of owner[1]. The commonly made mistakes that can expose an account in social network are: clicking on enticing Ads, connecting with strangers, using third party apps, exposing too much information, failing to utilize security settings, not logging out unknowingly. Growth of online social networks and publishing data in social network has led to the threat of leakage of private information of individuals. So there is need to protect owners profile and sensitive information[3].

The main aim of this paper is to prevent the risk of leakage of confidential information of individuals due to the development of online social networks and publication of social network. The rest of the paper is organized as follows: Section 2 provides literature review, section 3 presents dataset description, section 4 describes the methodology adopted, section 5 provides experimental results and discussions and finally conclusion is given in section 6.

4. METHODOLOGY

4.1 PRIVACY-PRESERVING GENERALIZATIONS

The privacy requirements of the user are attained when setting up the system into the user's environment. In this approach, the type of interactions between users of the social network contains three different privacy levels:

1. User1(FULL ACCESS): Can see all the posts
2. User2(PARTIAL ACCESS): Can see the generalized posts
3. User3(NO ACCESS): Can see any of the posts

The linguistic analysis is done on the input message to be published by the user and extracts potentially sensitive terms with regard to the thresholds of each privacy level. Replacing a sensitive term with a generalization (e.g., cancer -> disease), we are lowering the amount of information disclosed to a certain type of readers while retaining an amount of its semantics/utility.

Assuming privacy requirements with n levels $\{L_0, \dots, L_{n-1}\}$ and their corresponding n thresholds $\{T_{L_0}, \dots, T_{L_{n-1}}\}$, for each term t in a certain message m to be published do:

- The system will not do anything and t will be published as is, so that readers in the level L_0 (or higher) can use it.
- The system acquires the most informative generalization $g_0(t)$ from the databases in use such that $T_{L_0} \geq IC(g_0(t))$.
- The system acquires the most informative generalization $g_i(t)$ from the databases in use such that $T_{L_i} \geq IC(g_i(t))$.
- The system acquires the most informative generalization $g_0(t)$ from the databases in use such that $T_{L_0} \geq IC(g_0(t))$.

The sensitive words are identified by using linguistic analysis and then by using sanitization the words are generalized and then visible to user who is permitted with the partial access and then for no access. The terms or generalizations categorized in the lowermost privacy level (L_0) will get available in the social network.

4.2 PROTECTING SENSITIVE DATA

The sets $\{S_1, \dots, S_{n-1}\}$ contain the terms or generalizations (i.e., sanitized versions) that are accessible for the users belonging to each privacy level. The terms in $\{S_1, \dots, S_{n-1}\}$ are encrypted that each reader will get only the S_i linked to the privacy level L_i . The set of encrypted elements will be used by authorized readers to obtain the right cryptographic keys and get the corresponding protected terms. It is to avoid availability issues or problems inherent to distributed solutions; the system stores the protected information in the servers of the social network.

5. EXPERIMENTAL RESULTS AND DISCUSSIONS

Input is syntactically analyzed and protected messages in social networks and their percentage of information preservation shows that User has published a message that most of the members are affected by cancer. Therefore, according to the permission, L_2 readers can know everything, thus $T_{L_2} = \infty$. L_1 readers can only know the generalized term which, according to the generalizations of CANCER in the knowledge base's in use would correspond to "disease" $T_{L_0} = IC("disease") = -\log_2(126E6/17E9) = 7,1$. finally, any external entity, which in this case would be only the social network operator that corresponds to L_0 , will not even know the type, that is $T_{L_0} = 0$. This does not allow any information accessed by the user. Reader L_0 only know encrypted message.

Terms and generalizations published in the social network and stored in the S_i corresponding to each privacy level L_i ; $\$$ states that the generalization for S_i is the same as such already published. The last row shows the size (in bytes) of each S_i to be stored in attached images. Table 1 shows, “In NASA many people having disease problem”, in this tweet sensitive word cancer is generalized to disease and it is stored in images to publish to other users.

Table 1. Size of Each S_i to be stored in attached images

PUBLISHED	S_1	S_2
Many	\$	Many
People	\$	People
Having	\$	Having
Disease	Disease	Cancer
Problem	\$	Problem
NASA	\$	NASA
Bytes to store	15 bytes	59 bytes

Ontology based disambiguating technique is adopted which gives highest score for each word and the different levels of access provided for the user is shown in Fig 3, 4 and 5. Fig 3 shows that when the user is provided with full access, they can view full tweets. Fig 4 shows that the user provided with partial access can view generalized tweets by using ontology technique wherein Fig 5 indicates that user permitted with no access can view only encrypted message.

Also a comparison of sanitization and ontology method in terms of Precision, Recall, F-Measure is done and is shown in Fig 6. Precision value of sanitization method is 80% and of Ontology method is 83%. Recall value of sanitization is 70% and ontology is 77% and an F-Measure value of ontology is 79.8% and sanitization is 74.6%.

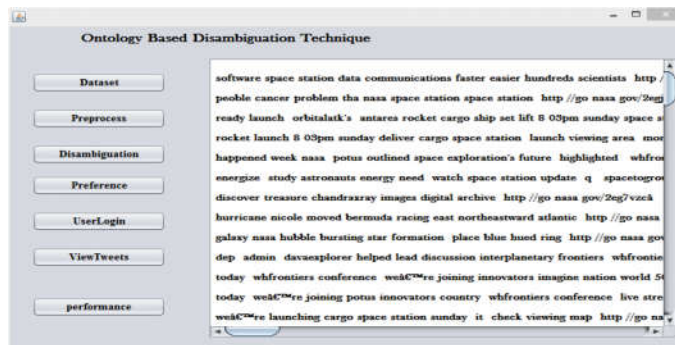


Fig 3: Disambiguating Technique for Full Access



Fig 4: Disambiguating Technique for Partial Access

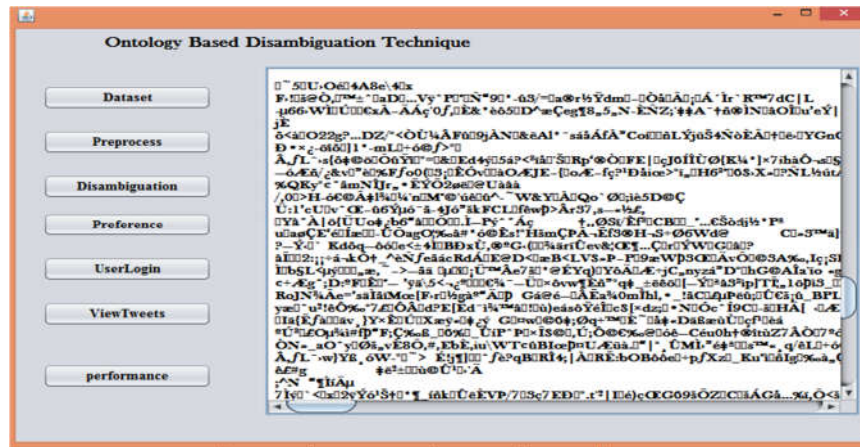


Fig 5: Disambiguating Technique for No Access

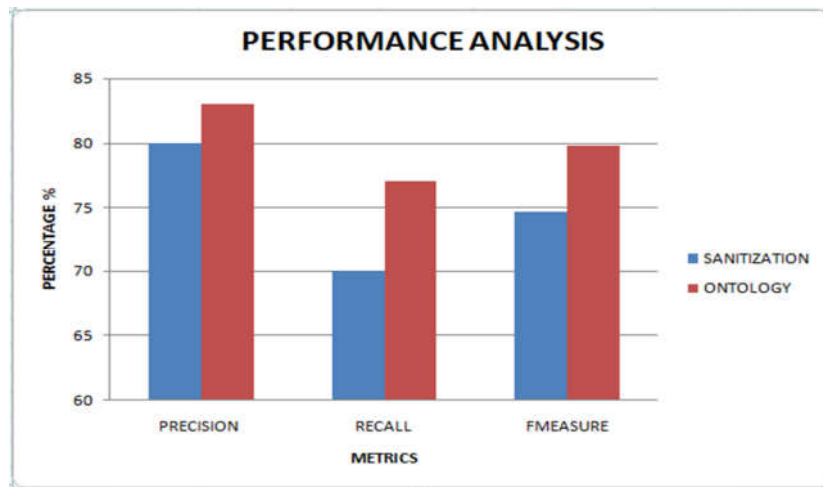


Fig 6: Overall Performance comparison

6. CONCLUSION

In the recent years, privacy of data in online social networks data has been of ultimate concern. In this proposed work, the preservation of the output’s utility and, the detection and sanitization of terms are focused it may cause disclosure of sensitive data due to semantic correlation. In addition to that ontology’s are used to retrieve the generalization. The experiments are conducted in terms of recall, precision and F-measure. From the experimental results it is proved that the proposed method has high precision, recall and F-measure than the existing methods. The scope of the privacy preserving in online social networks still to be explored.

REFERENCES

[1] Viejo, A., & Sánchez, D. 2014. Profiling Social Networks to Provide Useful and Privacy-Preserving Web Search. Journal of the Association for Information Science and Technology, 65, 2444-2458.
 [2] Viejo, A., Sánchez, D., & Castellà-Roca, J. 2012. Preventing automatic user profiling in Web 2.0 applications. Knowledge-Based Systems, 36, 191-205, 2012.
 [3] Sánchez, D., Batet, M., & Viejo, A. 2014. Utility-preserving privacy protection of textual healthcare documents. Journal of Biomedical Informatics, 52, 189-198.

- [4] Sánchez,D., Batet, M., & Viejo. 2014. A.Utility-preserving sanitization of semantically correlated terms in textual documents. In Proceedings of the 2nd international conference on Information Sciences, 279, .77–93.
- [5] Sánchez, D.,&Batet, M. 2016. C-sanitized: A privacy model for document reduction and sanitization. Journal of the Association for Information Science and Technology, 67, 148–163.
- [6] Anandan,B.,&Clifton,C. 2011. Significance of term relationships on anonymization. InIEEE/WIC/ACM International Joint Conference on Web Intelligence and Intelligent AgentTechnology - Workshops, 253–256.
- [7] Velásquez, J. D. 2013. Web mining and privacy concerns: Some important legal issues to be consider before applying any data and information extraction technique in web-based environments. Expert Systems with Applications, 40, 5228–5239.
- [8] Batet, M. ,Harispe, S. , Ranwez, S. , Sánchez, D. , &Ranwez, V. 2014. An information theoretic approach to improve semantic similarity assessments across multiple ontologies. Information Sciences, 283, 197–210.
- [9] Carminati, C. , Ferrari, E. , Heatherly, R. , Kantarcioglu, M. , &Thuraisingham, B. 2016. Semantic web-based social network access control. Computers & Security, 30, 108–115