

Designing for Different Levels of Social Inference Risk

Sara Motahari, Sotirios Ziavras, Quentin Jones
New Jersey Institute of Technology
{sg262, quentin.jones, ziavras}@njit.edu

1. INTRODUCTION

Changes in the technological environment are creating numerous new and unaddressed risks to user privacy. Today's social computing applications such as Facebook enable users to exchange messages, reveal aspects of their profile, and even find profile-based matches. Location-based applications leverage location, mobility, or proximity information to support navigation, recommendations, match making, etc. The resulting use and sharing of such personal information raise many serious privacy concerns. Previous efforts to protect users' privacy have made considerable advances in terms of computer and network security, user control mechanisms [1], ethical considerations, and privacy policies [3]. However, the collaborative and pervasive nature of new mobile and social computing applications can give users the ability to leverage background knowledge about the social environment/context to make unwanted inferences.

The term *inference* as used in the privacy literature is the process of deducing unrevealed information as a consequence of being presented with authorized information. Inference is mostly known as a security threat to databases and sometimes as a privacy risk in data mining [2]. Although the inference problem as a threat to database confidentiality is discussed in many studies, Ubiquitous Social Computing (USC) raises new classes of inferences which we call *social inferences*. Social inferences are unwanted inferences that result from the use of social computing applications by the inferrer and are about user information associated with these applications such as identity, location, activities, social relations, and profile information. Threats to user privacy in mobile social computing systems have been placed into seven categories in [4]. We focus on the two categories that relate to social inferences:

1. *Instantaneous Social Inferences* (e.g. my cell phone shows that I have a romantic match, Bob, who is nearby and I can only see two people with a similar cell phone around me. One of them must be Bob, thus increasing my chance of identifying him).
2. *Historical Social Inferences* through persistent user observation (e.g. two nicknames are repeatedly shown on the first floor of the gym where the gym assistant normally sits. One of them must be the gym assistant).

As we collect more information about a user, such as his/her contextual situation, our uncertainty about other aspects such as his/her identity may be reduced, thus increasing our probability of correctly guessing these aspects. This uncertainty is measured by *information entropy* in information theory. In previous studies, we framed the social inference problem and explained the relation between social inferences and information entropy [5, 6]. We explained how to calculate the information entropy based on our modeling of the inferrer's background knowledge. We then suggested how to set entropy thresholds for each user based on users' privacy settings. If entropy is less than its threshold, there is high risk that the inferrer will infer the user's private information an appropriate action needs to be taken by the system to address the risk. In this note, we assess the risk and investigate the

applicability of different approaches to user information flow control in different classes of social computing applications. The goal of this research is to understand different design implications for initially-anonymous communication in Computer-Mediated Communication (A-CMC) and proximity-based applications.

2. METHOD

The first step to understand the design implications to understand the extent of the social inference risk to users of A-CMC and proximity based applications. We did this through 1) a laboratory experimentation, 2) a mobile phone field study, and 3) large scale simulations of the risk populated with empirical data obtained from the two user studies.

In our laboratory experiment, subjects participated in a study consisting of three phases: 1) online personal profile entry; 2) an experiment involving subjects chatting with an unknown online partner; and 3) a post chat survey about the subject's ability to guess their chat partner's identity. Five hundred and thirty students entered a personal profile, 304 participated in the chat session of which 292 subjects completed all three study components. A detailed presentation of Study 1 can be found in [6].

In the mobile phone field study, 175 students participated in a study of two phases: 1) an online pre-study survey and 2) carrying our Window Mobile phones for three weeks, and answering questionnaires on the phone that used Context-Aware Experience Sampling Method (CA-ESM). Prior to the experiment, a location estimation system and the 'Nearby' application were installed on the phones. The Nearby application shows the nicknames of the users in the vicinity of the phone user on campus. In the pop-up questionnaires, subjects were mainly asked questions about the nicknames they saw on their Nearby application, what they could guess about the identity of the nickname owner, and how they could map them to people in their vicinity. A detailed presentation of Study 2 can be found in [5].

As a final step we used simulations to investigate the problem and appropriate actions on a larger scale for various situations. The simulation models were populated with parameters derived from our user study providing a good approximation to real world deployments. To simulate the risk in CMC, we first simulated personal profiles. Parameters, such as the diversity of profile items, their statistical distribution, etc. were derived from the 532 user profiles obtained from Study 1. Additional information such as the number of courses, statistical distribution of the number of students in a class, and enrollment statistics were obtained from university admission statistics. We then simulated online interactions of the users. The probability of revealing profile items and users' desired degree of anonymity were derived from the user experiment. To simulate the risk in proximity-based applications, population density and its probability distribution were obtained from Study 2.

3. RESULTS

Study 1 showed that Identity inferences are frequently made in A-CMC and Study 2 showed that the risk of identity inference in a proximity-based application was lower than the risk in A-CMC. Simulations confirmed the experimental results on a larger scale. Figure 1 shows the probability that a user’s identity entropy is lower than its threshold. Entropy threshold was calculated based on users’ desired degree of anonymity (we call U a user’s *desired degree of anonymity* if he/she wishes to be indistinguishable from $U-1$ other users). The y-axis shows the percentage of users for whom entropy was less than the threshold. The x-axis was chosen to represent the population size. The depicted curves show this probability for desired degrees of anonymity of 2, 3, and 5. In user study 1, 80.8% of the users who wanted to stay anonymous desired a degree of anonymity of 2. As the figure shows, while in a small school the risk can be very high, in a campus of 10,000 students, it is still about 50% in online chats between students.

Figure 2 shows the probability that a user is at the risk of instantaneous identity inference in a proximity-based application. The y-axis shows the percentage of users whose identity entropy was lower than its threshold. The x-axis represents the desired degree of anonymity. Each curve depicts the risk for a different mean of nearby population density. We see that assuming mass usage, the risk of identity inference is about 7% for a desired degree of 3, and 20% for a desired degree of 5.

4. DESIGN IMPLICATIONS

The user studies and simulations show that social inferences are common and pose a serious threat to user privacy. User experiments showed that even when users are in complete control of the information they reveal, they are not able to maintain their desired degree of anonymity [6], as they are unable to correctly judge inference risks. Further, the nature of these risks is quite different for online A-CMC profile exchanges as compared to proximity information. Collectively, it is possible to derive a number of design implications from these findings.

First, we need to provide users with the means to set their desired level of anonymity since we observed that users of USC applications can have a wide range of anonymity preferences [6].

Second, for A-CMC systems user interfaces should be built to improve user inference-risk judgments through techniques such as risk visualizations or warning messages (perhaps applied to customized introduction tools). This implication is derived from the results that indicate that identity inference risks are quite

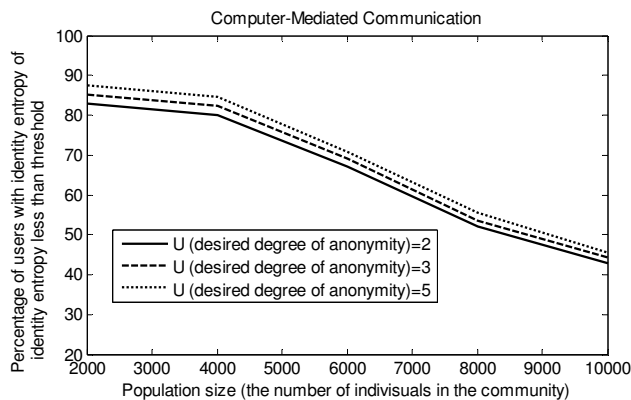


Figure 1. Risk of Identity Inference in A-CMC

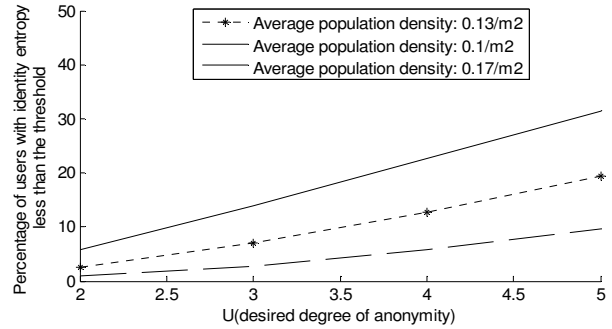


Figure 2: Risk of Identity Inference for proximity-based applications

common in A-CMC, which means that automatic control of information exchange in such applications can degrade system usability or be frustrating for the user. Furthermore, such applications are designed for users to consciously exchange information and users may be willing to compromise their anonymity settings to have more meaningful and productive communication. The systems can show users how uniquely they have specified themselves so far, or send a warning message when revealing a piece of information would enable their partner to invade their desired degree of anonymity.

The final design implication is derived from the finding that the prevalence of situations with identity inference risks in location-based services is lower than in A-CMC. In location-aware mobile applications, inference protection systems should modify information exchange, for example, by lowering the granularity of revealed information, rejecting a query, or blocking information exchange. In most cases lowering the information granularity, such as revealing the location at floor precision instead of room precision, or showing an anonymous name instead of a nickname can address the inference risk. This technique should not overly interfere with information exchange, or overly burden the user with privacy management actions.

Acknowledgments- This research is partially supported by the National Science Foundation Grant NSF IIS DST 0534520 and NSF CNS 0454081. The opinions expressed are those of the authors and may not reflect those of the NSF.

5. REFERENCES

- [1] Cornwell, J., et. al., User-controllable security and privacy for pervasive computing, in *Proceedings of the IEEE Workshop on Mobile Computing Systems & Applications*, 2007.
- [2] Farkas, C. and Jajodia, S. The inference problem: a survey. *SIGKDD Explorer Newsletter*, 4 (2). 6-11, 2002.
- [3] Langheinrich, M., Privacy by Design – Principles of Privacy-Aware Ubiquitous Systems. *UbiComp 2001*, 273-291.
- [4] Motahari, S., Manikopoulos, C., Hiltz, R. and Jones, Q., Seven privacy worries in ubiquitous social computing. *Proceedings of the 3rd symposium on Usable privacy and security* (2007), 171-172.
- [5] Motahari, S., Zivras, S., Naaman, M., Ismail, M. and Jones, Q. Social Inference Risk Modeling in Mobile and Social Applications, *The IEEE International Conference on Privacy, Security, Risk and Trust (PASSAT 2009)*, 2009.
- [6] Motahari, S., Zivras, S., Schular, R. and Jones, Q. Identity Inference as a Privacy Risk in Computer-Mediated Communication. *IEEE Hawaii International Conference on System Sciences (HICSS-42)*, 2008.