

Unraveling an Old Cloak: k -anonymity for Location Privacy

Reza Shokri[†], Carmela Troncoso[‡], Claudia Diaz[‡], Julien Freudiger[†], and Jean-Pierre Hubaux[†]

[‡] IBBT-K.U.Leuven, ESAT/COSIC,
Leuven-Heverlee, Belgium
firstname.lastname@esat.kuleuven.be

[†] LCA1, EPFL,
Lausanne, Switzerland
firstname.lastname@epfl.ch

ABSTRACT

There is a rich collection of literature that aims at protecting the privacy of users querying location-based services. One of the most popular location privacy techniques consists in cloaking users' locations such that k users appear as potential senders of a query, thus achieving k -anonymity. This paper analyzes the effectiveness of k -anonymity approaches for protecting location privacy in the presence of various types of adversaries. The unraveling of the scheme unfolds the inconsistency between its components, mainly the cloaking mechanism and the k -anonymity metric. We show that constructing cloaking regions based on the users' locations does not reliably relate to location privacy, and argue that this technique may even be detrimental to users' location privacy. The uncovered flaws imply that existing k -anonymity scheme is a tattered cloak for protecting location privacy.

Categories and Subject Descriptors

C.2.0 [Computer-Communication Networks]: General—Security and protection; K.4.1 [Computers and Society]: Public Policy Issues—Privacy

General Terms

Security, Measurement

1. INTRODUCTION

An increasing number of people own mobile devices with positioning capabilities, and use various location-based services (LBSs) to obtain all kinds of information about their surroundings. Privacy concerns have emerged because many of such services enable, by design, service providers to collect detailed location information about their users.

Protecting users' location privacy, while enabling them to still benefit from location-based services, is a challenging problem. The most popular technique for designing privacy-preserving SBSs consists in obfuscating the actual location from which a query is made by constructing cloaking regions

that contain the locations of k anonymous users. According to the k -anonymity metric, a user's level of location privacy directly depends on the number of other users that expose their location to the SBS using the same cloaking region and at the same time as the considered user does, while identity-wise they are indistinguishable from each other. This approach is an adaptation of the k -anonymity technique originally developed in the context of database privacy [14], to prevent the re-identification of anonymous people whose data were included in a published dataset. Taking for granted the *effectiveness* of the k -anonymity technique for protecting location privacy, a large body of literature has focused on maximizing its *efficiency* in a variety of system models.

In this paper, we provide a thorough security analysis of k -anonymity schemes for location privacy considering various adversaries, classified based on their knowledge. The unraveling of the scheme shows multiple incoherences. First, there is confusion about query anonymity and location privacy. We show that cloaking can help decouple a query and a user (query anonymity) but does not necessarily prevent the adversary from linking a location to a user (location privacy). Second, the absence of a clear adversary model in the security analysis affects the validity of privacy-preserving mechanisms. We show that given certain knowledge, the adversary can obtain the location of users hiding behind the cloaking regions. Finally, we show that k is *not* representative of the actual location-privacy of mobile users. In fact, the cloaking technique, constructed based on the k -anonymity metric, may even be counterproductive and give the illusion of a high location-privacy level, while the adversary is able to infer the users' locations.

The results of our analysis, that show the inconsistencies of the k -anonymity metric with respect to the users' actual location-privacy at microscopic level, in addition to previous works that show the inaccuracy of k -anonymity metric in reflecting the strength of mixing users' trajectories (i.e., location privacy at macroscopic level) [17], suggest that the k -anonymity scheme is inadequate for protecting location privacy. These negative results on such a popular scheme show that location privacy is not yet well understood and that more attention is needed to address this problem.

2. LOCATION PRIVACY

In location-based services, users share their location with a service provider in return for services. For example, many SBSs enable users to search for nearby points of interests (POIs). In such SBSs, users share their location with the service provider at the time they need the information. This

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

WPES'10, October 4, 2010, Chicago, Illinois, USA.

Copyright 2010 ACM 978-1-4503-0096-4/10/10 ...\$10.00.

is done by sending *queries* that include a user pseudonym, her location (as the search domain of the query), and the body of the query (e.g., what type of POI she is looking for). The pseudonym of a user can be of different types (explicitly given, such as her application username, or implicitly inferable from the content of her packets such as her IP address), each of which reveals different information about the user’s real name. Permanent pseudonyms (that do not change over time) make the user’s queries linkable to each other, hence, eventually enabling re-identification.

In this model, the LBS provider is able to link users with their visited locations, and thus is capable of inferring sensitive private information. Malicious or incompetent service providers are thus a threat to users’ location privacy. The adversary may have prior information about the users’ permanent pseudonyms, the space in which users move, or their mobility patterns, along with publicly available information such as their homes and work places. This knowledge can help the adversary to infer more information about users’ locations from their queries and to perform attacks that leads to absence/presence disclosure of the users’ locations [16].

Following the terminology introduced in [16], location privacy is defined in two levels: *microscopic* and *macroscopic*. Microscopic location privacy is defined as the users’ location privacy on a small scale, i.e., corresponding to a single query, and reflects how accurately the adversary can infer the users’ locations after observing their individual queries, given his a priori knowledge. Macroscopic location privacy represents users’ privacy on a large scale, e.g., given multiple (possibly correlated) queries from users as they move.

Multiple privacy-preserving mechanisms have been proposed so far. Users may *hide* their current locations from the adversary by abstaining from sending their queries to the LBS for a short amount of time and *anonymize* their queries by removing their real names and changing their pseudonyms (e.g., mix zones [1]); they may *obfuscate* the queries by decreasing the accuracy or precision of their location/time (e.g., location perturbation [7], time perturbation [11]); or they may add some *dummy* queries that are indistinguishable from the real queries [3].

At both microscopic and macroscopic levels, many of the protection mechanisms are based on anonymization and obfuscation methods. Among them, *k*-anonymity scheme is by far the most employed protection scheme for location privacy, mostly due to its simplicity [4, 6, 12, 15, 19, 20, 21]. Taking for granted the effectiveness of the *k*-anonymity technique in preserving location privacy, researchers have mainly focused on adapting it to variants of the basic system model and on improving its efficiency (i.e., minimizing the cost of obfuscation on the system utility, while guaranteeing a *k*-anonymity level for the users).

In this paper, we focus on analyzing *k*-anonymity for location privacy at the microscopic level, as it has already been proven ineffective at the macroscopic level [8, 17].

3. K-ANONYMITY

In this section, we introduce the original concept of *k*-anonymity and its extension to the field of location privacy.

3.1 The Concept of *k*-anonymity

The concept of *k*-anonymity was originally proposed by Samarati and Sweeney in the field of database privacy [13, 14, 18]. Databases are typically populated with person-

specific data entries such as names, birth date, and gender. Many situations call for the release of these data. For example, a medical database may need to be shared or made public in order to study the incidence of diseases. When releasing data, the privacy of the individuals who provided it should be protected: database entries should not be linkable to individuals. The mere removal of *explicit identifiers*, such as individuals’ names, is insufficient because individuals can be re-identified by linking their distinctive attributes (e.g., date of birth) to publicly available information. These subsets of attributes are called *quasi-identifiers* because they facilitate the indirect re-identification of individuals.

To overcome this problem, the approach of *k*-anonymity suggests the suppression and generalization (obfuscation) of quasi-identifiers to make an individual’s data entry indistinguishable from others. By definition [13, 14, 18], a database provides *k*-anonymity if explicit identifiers are removed from the database and, additionally, the quasi-identifiers of each individual in the database cannot be distinguished from those of at least $k - 1$ other individuals.

In essence, the concept of *k*-anonymity relies on a simple protection mechanism: obfuscation. It then measures the provided privacy with a single parameter *k*. The value *k* determines the privacy protection in place: the larger the *k* is, the higher the privacy protection is. Thus, it is this tight coupling of the privacy preserving mechanism and the metric that builds the *k*-anonymity scheme.

The *k*-anonymity model may fail, in some cases, to guarantee the privacy to the level that its metric promises. This is because making the quasi-identifiers of a user identical to those of $k - 1$ other users does not reflect how and to what extent her sensitive information is hidden from the adversary, e.g., all the *k* users might have cancer (considering disease as sensitive information). Additional properties such as *l*-diversity [10] and *t*-closeness [9] complement *k*-anonymity by considering how the users’ sensitive information is different and remote from that of other users with whom she shares the same obfuscated quasi-identifiers.

3.2 *k*-anonymity for Location Privacy

In the context of location privacy, the *k*-anonymity metric was initially adapted to measure microscopic location-privacy by Gruteser and Grunwald [6]. In this model, each query sent to the LBS (including the user’s pseudonym, her position and the query time) is equivalent to one entry in a database, and the location-time information in the query serves as the quasi-identifier. In order to protect a user’s location privacy using *k*-anonymity, each of her queries must be indistinguishable from that of at least $k - 1$ other users. To this end, the pseudonyms of these *k* users are removed from their queries, and the location-time pair in their queries is obfuscated to the same location-area and time-window, large enough to contain the users’ actual locations.

The *k*-anonymity scheme for location privacy has become very popular, mainly due to its simplicity. A large body of research has focused on increasing the efficiency of *k*-anonymity schemes and reducing their cost of query obfuscation [4, 12, 19, 20], extending the obfuscation method to protect traces [2] (i.e., location privacy at the macroscopic level), or adapting the architecture presented in [6] to different scenarios [15, 21].

All of these systems can be represented by the initial model introduced in [6]. We present this model in Fig. 1:

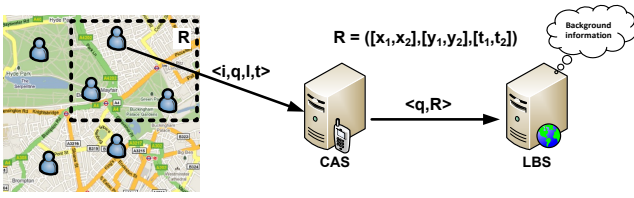


Figure 1: Basic k -anonymity model

There is a set of users who access a LBS through a trusted Central Anonymity Server (CAS). Users send their LBS queries $\langle i, q, l, t \rangle$ to the CAS, where i is the identity of the user, q is her query, l is her precise location (expressed as a point with coordinates (x, y) in a 2-dimensional space), and t is the time at which the query is generated. In order to protect users' privacy, the CAS removes the identity i of the users. Furthermore, it obfuscates the location $l = (x, y)$ and the time t at which the queries were generated. For this, it constructs a cloaking region $\mathcal{R} = ([x_1, x_2], [y_1, y_2], [t_1, t_2])$ such that there are at least k users ($k = 3$ in the figure) in \mathcal{R} whose location $l = (x, y)$ at time t satisfies that $x_1 \leq x \leq x_2$, $y_1 \leq y \leq y_2$, and $t_1 \leq t \leq t_2$.

We note that, decentralized approaches [15] can also be represented by our model, by considering that the cloaking region \mathcal{R} is computed by a set of entities (e.g., users themselves) in a distributed manner (i.e., they collectively play the role of the CAS). Moreover, our model can accommodate both the systems in which users have to continuously report their location to the CAS [12], in order to build optimal regions, and the systems that rely only on user-triggered discrete queries for that purpose [4].

The k -anonymity location obfuscation technique aims at achieving two properties: *query anonymity* and *location privacy*. Achieving query anonymity implies that it is not possible for the adversary to link the identity i of a user to her query q , based on the location information (cloaking region) associated with the query. Location privacy is achieved when it is not possible for the adversary to learn the location l of a user i at time t , using queries he receives from the users and his a priori knowledge.

We consider an adversary that controls the LBS and, in addition to the received queries, has access to some background information. For example, one of the threats considered in [6] is "restricted space identification". In this threat scenario the adversary knows that a given location corresponds (exclusively) to a user address, meaning that a query coming from that precise location would be linked to the user who resides at that address. Another considered threat in [6] is that the adversary (in addition to controlling the LBS) may deploy antennas in the vicinity of the users and thus knows that a given user is in location l at time t .

4. EVALUATING K-ANONYMITY

This section evaluates the k -anonymity scheme with respect to the properties stated previously: *query anonymity* (i.e., concealing the link between user i and her query q) and *location privacy* (i.e., concealing the link between user i and her location l at time t). The evaluation is twofold: first, we analyze the consistency of the k -anonymity metric over time and space, and second, its coherence with the users' location privacy given adversaries with different knowledge.

4.1 Consistency

Consider that the k users in the cloaking region of a given user are situated next to each other in a small place (e.g., a bar). In this case, the adversary learns the actual location of the users (e.g., all of them are at the bar). In contrast, consider a large cloaking region that encompasses the same number of users. In this case, users' location privacy is better protected because the adversary has more uncertainty about their exact locations. Hence, the number of users in the obfuscated region is not a consistent metric for location privacy. The independence of the value of k and the accuracy of a user's location estimation by the adversary, implies that k is irrelevant to their actual location privacy.

4.2 Adversarial Knowledge

A security analysis must define the knowledge of the adversary by considering the information the adversary has access to. We consider three types of background information the adversary could have and examine which properties are provided in each of the cases. The first and the third type capture extreme scenarios: the worst case scenario (i.e., the adversary knows everything), and the ideal case scenario (i.e., the adversary knows nothing). The second type is generic and captures most realistic adversary models.

4.2.1 Real-Time Location Information

We first consider a scenario in which the adversary has access to real-time information on the location of users. In this case, as also mentioned in [6], the adversary could eavesdrop on the communications between users and CAS, and localize them (e.g., using multiple directional antennas) thus being able to obtain the location from which users send queries.

Upon receiving a query $\langle q, R \rangle$ in which k users are present, the LBS may not distinguish which of the k users is the sender of query q . Thus, the originator of the query is *at best* k -anonymous. However, as the adversary knows the users' exact location they have *no* location privacy.

4.2.2 Statistical Information

Let us now consider an adversary who is unaware of the real-time position of users, but who has access to statistical information about their mobility patterns. For example, the adversary may have access to publicly available information on users' homes and work places [5] and knows that, with a high probability, users will be at home during the night, and at their work places during office hours. As the adversary does not have access to the *actual* location of users, to perform an attack he only relies on the queries $\langle q, R \rangle$ forwarded by the CAS and the available background information.

In this case, the success of the adversary in pinpointing users' actual locations in the obfuscated regions depends *only* on his statistical background information. Thus, we argue that computing cloaking regions based on actual locations does not necessarily improve users' location privacy and hence is neither efficient nor effective. Let us consider a neighborhood as the one shown in Fig. 2(a), and assume that the adversary knows that with a high probability all users are at home (for instance, late in the evening). When user A sends a query $\langle q, R \rangle$ to the LBS, it is unaware of the current location of users, and can only use the available statistical information to infer who/where is the sender of q . Thus, user A is 4-anonymous independently of whether or not B , C , and D are *currently* using the system, or even



Figure 2: Statistical background information: efficiency (a) and additional information (b)

present at their home locations. Therefore, there is no need for the CAS to execute complex algorithms (e.g., [4, 12]) to compute or select the region \mathcal{R} on-the-fly. Instead, regions \mathcal{R} can be pre-computed (taking into account the background information available to the adversary) and later be selected by users uniquely based on their own location, regardless of whether or not $k-1$ other users are currently in the vicinity.

Further, computing cloaking regions based on the users' current locations can be counterproductive. In order to optimize the accuracy of LBS, previous proposals aim at minimizing the area of the region \mathcal{R} in the query. When \mathcal{R} is computed according to the current location of users, this minimization allows the adversary to make inferences about their current position (as there must be at least k users in the region). Consider the example in Fig. 2(b) in which only users A , B , E , and F are active (i.e., using the system), and assume that the adversary has the same information as in the previous case. When user A sends a query, the CAS forwards $\langle q, R' \rangle$ to the LBS. Upon receiving this information, the adversary learns that A , B , E , and F are currently in their home locations, and that C and D are either inactive or absent. This is because had C and D been active in the system, the minimal region sent to the LBS would have been R (as in Fig. 2(a)). Thus, the only configuration that results in R' is that A , B , E , and F are active and at home. In this case, although the query q is still 4-anonymous, A , B , E , and F have no location privacy, due to the information revealed to the adversary by the cloaking region itself.

4.2.3 No Information

Finally, we consider an adversary that does not (and will not) have any background information. Assuming that the only information available to the adversary is the queries forwarded by the CAS, both query anonymity and location privacy are achieved by simply removing the users' identities from their queries. Thus, in this case, constructing cloaking regions based on the k -anonymity technique does not provide any additional protection and it only reduces the performance of the system in terms of accuracy and computational load. Given the availability of public location information, this is an unrealistically weak adversary model, included here for the sake of completeness.

Acknowledgment

C. Diaz and C. Troncoso are funded by the Fund for Scientific Research in Flanders (FWO). This work was supported in part by the IAP Programme P6/26 BCrypt of the Belgian State.

5. REFERENCES

- [1] A. R. Beresford and F. Stajano. Mix zones: User privacy in location-aware services. In *PERCOMW*, 2004.
- [2] C. Bettini, X. S. Wang, and S. Jajodia. Protecting privacy against location-based personal identification. In *VLDB Workshop SDM*, 2005.
- [3] R. Chow and P. Golle. Faking contextual data for fun, profit, and privacy. In *WPES*, 2009.
- [4] B. Gedik and L. Liu. Protecting location privacy with personalized k -anonymity: Architecture and algorithms. *IEEE Trans. on Mobile Computing*, 2008.
- [5] P. Golle and K. Partide. On the anonymity of home/work location pairs. In *Pervasive*, 2009.
- [6] M. Gruteser and D. Grunwald. Anonymous usage of location-based services through spatial and temporal cloaking. In *ACM MobiSys*, 2003.
- [7] B. Hoh and M. Gruteser. Protecting location privacy through path confusion. In *SECURECOMM*, 2005.
- [8] B. Hoh, M. Gruteser, H. Xiong, and A. Alrabady. Preserving privacy in GPS traces via uncertainty-aware path cloaking. In *ACM CCS*, 2007.
- [9] N. Li and T. Li. t -closeness: Privacy beyond k -anonymity and l -diversity. In *ICDE07*, 2007.
- [10] A. Machanavajjhala, D. Kifer, J. Gehrke, and M. Venkatasubramanian. L -diversity: Privacy beyond k -anonymity. *ACM TKDD*, 2007.
- [11] J. Meyerowitz and R. Roy Choudhury. Hiding stars with fireworks: location privacy through camouflage. In *MobiCom*, 2009.
- [12] M. F. Mokbel, C.-Y. Chow, and W. G. Aref. The new casper: query processing for location services without compromising privacy. In *VLDB*, 2006.
- [13] P. Samarati. Protecting respondents' identities in microdata release. *IEEE Trans. on Knowl. and Data Eng.*, 2001.
- [14] P. Samarati and L. Sweeney. Generalizing data to provide anonymity when disclosing information (abstract). In *PODS*, 1998.
- [15] K. Sampigethaya, L. Huang, M. Li, R. Poovendran, K. Matsuura, , and K. Sezaki. Caravan: Providing location privacy for vanet. In *ESCAR*, 2005.
- [16] R. Shokri, J. Freudiger, and J.-P. Hubaux. A unified framework for location privacy. Technical Report EPFL-REPORT-148708, EPFL, Switzerland, 2010.
- [17] R. Shokri, J. Freudiger, M. Jadhwal, and J.-P. Hubaux. A distortion-based metric for location privacy. In *WPES*, 2009.
- [18] L. Sweeney. k -anonymity: a model for protecting privacy. *Int. J. Uncertain. Fuzziness Knowl.-Based Syst.*, 10(5), 2002.
- [19] K. W. Tan, Y. Lin, and K. Mouratidis. Spatial cloaking revisited: Distinguishing information leakage from anonymity. In *SSTD*, 2009.
- [20] T. Xu and Y. Cai. Feeling-based location privacy protection for location-based services. In *CCS*, 2009.
- [21] G. Zhong and U. Hengartner. A distributed k -anonymity protocol for location privacy. *PerCom*, 2009.