

Protecting Location Privacy against Inference Attacks

Kazuhiro Minami
National Institute of Informatics
Tokyo, Japan
minami@nii.ac.jp

Nikita Borisov
University of Illinois at Urbana-Champaign
Urbana, IL 61801-2307
nikita@illinois.edu

ABSTRACT

GPS-enabled mobile devices are a quickly growing market and users are starting to share their location information with each other through services such as Google Latitude. Location information, however, is very privacy-sensitive, since it can be used to infer activities, preferences, relationships, and other personal information, and thus access to it must be carefully protected. The situation is complicated by the possibility of inferring a users' location information from previous (or even future) movements. We argue that such inference means that traditional access control models that make a binary decision on whether a piece of information is released or not are not sufficient, and new policies must be designed that ensure that private information is not revealed either directly or through inference. We provide a formal definition of location privacy that incorporates an adversary's ability to predict location and discuss possible implementation of access control mechanisms that satisfy this definition. To support our reasoning, we analyze a preliminary data set to evaluate the accuracy of location prediction.

Categories and Subject Descriptors: C.2.4 [Distributed Systems]: Distributed applications; K.6.5 [Management of Computing and Information Systems]: Security and Protection

General Terms: Security

Keywords: Location privacy, access control, the Markov model

1. INTRODUCTION

Soon the vast majority of mobile devices will be equipped with some form of localization capability; already, most smart phones include a GPS receiver. This has led to the rise of location-based services on a number of mobile platforms, including Symbian, iPhone, and Android. Novel applications, such as Google Latitude [5], have opened up the possibilities of sharing location information with other users [4, 5, 7].

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.

This theme has been picked up in social networks as well; e.g., Twitter recently announced support for embedding a location in each post [14].

Location sharing raises significant privacy concerns [1], since a location, such as a bar or a hospital, can be used to infer a user's personal activities. Therefore, location-sharing services (LSSs) have introduced an access control mechanism that allows the user to specify what location data may be shared with whom. For example, Google Latitude allows a user to authorize others' access to his or her location; it also allows a user to enter a decoy location manually. Glympse [4] specifies a time duration during which location information is shared. These interfaces provide coarse-grained controls. Researchers in pervasive computing have proposed more fine-grained access control schemes [6, 8, 12] that make use of context information such as location, time of day, and so on. These rules both better represent the users' actual sharing desires and at least partially automate the decisions to provide seamless integration of location sharing into people's daily lives.

One additional danger of sharing location information, however, is that it can lead to inference of previous or past locations. For example, a person traveling along a trajectory is likely to remain along that path. Things get significantly more complex as more background data is introduced. For example, walking and driving paths follow a predictable pattern, following streets and sidewalks; furthermore, each person exhibits more specific patterns in their activities. For example, Figure 1 shows two potential walking paths leading to a hospital and a library. Given background knowledge, it is possible to infer that a user traveling towards the intersection (black circles) is likely to visit one of these two places. A user who turns left at the intersection (white circles) may then be assumed to be going to the hospital. Therefore, if the user wishes to hide visits to the hospital, it is important to stop revealing his or her location earlier as well.

We, therefore, propose to develop a new access-control scheme that prevents such inference attacks. Our basic approach is to model an adversary as a location predictor that predicts future movements of a target user from his previous movements with certain probabilities. Intuitively, our access control scheme discloses a user's location information only if an unauthorized user cannot predict that the user moves to some private location with a sufficiently high probability. To model outside knowledge, we conservatively assume that the adversary has access to a complete history of the previous movements of the target user. We base our predictions on higher-order Markov models, which have shown to be

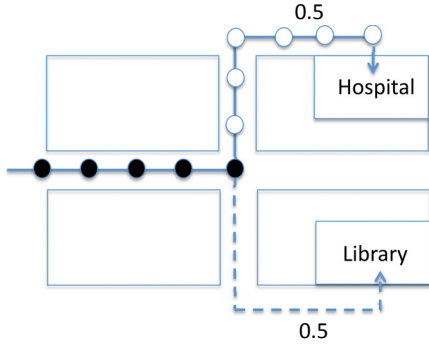


Figure 1: Example safe disclosure of location information. The solid line represents an actual path of a user visiting a hospital. We assume that the hospital is a private place and the library is a public place. A safe LSS would disclose location points denoted by black nodes.

very good at predicting user location movements [13]. Our preliminary results with actual GPS traces of a single user show that we can predict the user’s next movement with the accuracy of 60% using a first-order Markov model, and we can improve the accuracy by 10% by considering multiple previous movements with a higher-order Markov model.

The rest of the paper is organized as follows. Section 2 introduces our system model for LSSs in this paper, and Section 3 describes a location predictor based on the Markov model. We present our preliminary results of experiments in Section 4. We cover related work in Section 5, discuss key issues for future research in Section 6, and finally conclude in Section 7.

2. SYSTEM MODEL

Figure 2 shows our system model for LSSs. We assume that a Alice is interested in receiving Bob’s location movements. Bob, carrying a GPS-enabled mobile devices periodically sends location-timestamp pairs (loc_k, t_k) to the LSS for $k \in \mathcal{N}$. In our model, the LSS is completely trusted and receives all of the pairs:

$$L = \{(loc_k, t_k) \mid k \in \mathcal{N}\}.$$

Bob also defines an access-control policy to protect his location information, with the LSS implementing the policy. We represent access control policies by the function

$$acl : \mathcal{P} \times \mathcal{W} \rightarrow 2^{\mathcal{P}}$$

where \mathcal{P} is a set of all users and \mathcal{W} is a finite set of all locations. The function acl takes a user identity X and a location name l as inputs and outputs a set of users who are authorized to learn that “Bob is at location l .” In other words, the LSS releases Bob’s location movement (l_k, t_k) to principal X only if X belongs to set $acl(\text{Bob}, l_k)$, and thus user X receives a subset of events $L'(X) \subseteq L$

$$L'(X) = \{(loc_k, t_k) \mid X \in acl(\text{Bob}, loc_k)\}.$$

To simplify the presentation, we consider access policies that depend on location only, but it would be easy to incorporate other pieces of context, such as the timestamp t_k .

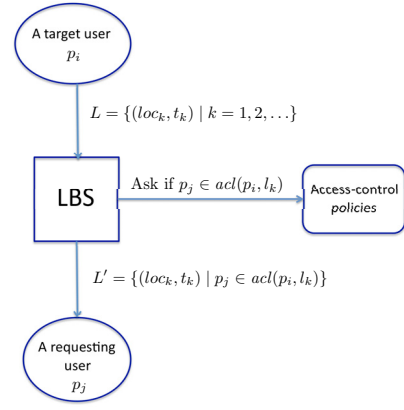


Figure 2: System model.

We next provide an informal definition of privacy that incorporates inference. In short, if Alice is not authorized to see when Bob is in a certain location, then she should not be able to infer this fact from other information released by the LSS.

DEFINITION 1 (PRESERVATION OF LOCATION PRIVACY.)
We say that a LSS preserves a user X ’s location privacy against another user Y if Y cannot infer X ’s movement (l, t) for any l such that $Y \notin acl(\text{Bob}, l)$ from the released location-timestamp pairs $L'(Y)$.

In next section, we describe how we model Alice’s inferences by a location predictor based on the Markov model, and give a more precise definition of privacy metrics based on probabilistic inference.

3. LOCATION PRIVACY AGAINST INFERENCE ATTACKS

We can represent Bob’s potential locations with random variables

$$X_1, X_2, X_3, \dots$$

where each X_i has a value drawn from the finite set of locations \mathcal{W} . For simplicity, we assume that location information is updated at regular intervals and dispense with the timestamp t_k . We will use Markov models to predict the location, which assumes that a location depends *only* on the last k states. So, for example, for a Markov chain of order 1,

$$Pr(X_{n+1}|X_1, \dots, X_n) = Pr(X_{n+1}|X_n)$$

We can represent this Markov chain as a $|\mathcal{W}| \times |\mathcal{W}|$ transition matrix, indexed by locations in \mathcal{W} . Each matrix entry represents the probability of moving from location l_i to l_j :

$$M_{i,j} = Pr(X_{n+1} = l_i | X_n = l_j)$$

for every pair of l_i and l_j in set \mathcal{W} . The probability of moving from location l_i to l_j in n time steps can be computed by multiplying the transition matrix M n times as follows:

$$Pr(X_{n+1} = l_i | X_1 = l_j) = M_{i,j}^n.$$

Since it is likely that we can improve the accuracy of location predictions by considering multiple previous movements, we also consider a location predictor based on a Markov model of a higher order. If we use the second-order Markov model, the transition matrix becomes a $|\mathcal{W}|^2 \times |\mathcal{W}|$ matrix, where:

$$M_{(j-1)*|\mathcal{W}|+k,i} = Pr(X_{n+1} = l_i | X_n = l_j, X_{n-1} = l_k)$$

To create the transition matrix, we compute the probabilities based on the past history of a user, as stored by the LSS. Note that we include in the history locations that are not part of $L'(X)$ to be conservative: even though these locations are not released through the LSS, the adversary may have some external knowledge of a user’s locations that can be factored in as well. Thus, we define location privacy against an unauthorized user (i.e., an adversary) with the knowledge of transition matrix M as follows:

DEFINITION 2 ((M, δ) -LOCATION PRIVACY.). *Given a transition matrix M corresponding to the first order Markov model, and a probability threshold $\delta > 0$, we say that a LSS preserves a user X ’s (M, δ) -location privacy with respect to a user Y if, whenever a pair (l_i, t) is released from the LSS to Y , for every l_j such that $Y \notin acl(X, l_j)$,*

$$M_{i,j}^n \leq \delta \text{ for all } n = 1, 2, \dots$$

Intuitively speaking, the above definition requires that user Y cannot predict that the target user X is at some private location l_j in some future time with probability p , which is greater than the threshold value δ . It is easy to generalize this definition to higher-order Markov models as well. Note that, although the definition involves examining points indefinitely in the future, a Markov chain will (under common conditions) converge to a stationary distribution, $\pi = \lim_{n \rightarrow \infty} M^n \vec{v}$ for any starting state vector \vec{v} , thus $\max_n M_{i,j}^n$ can be estimated accurately after a finite number of steps.

4. EXPERIMENTAL RESULTS

We conducted experiments with actual GPS traces to study how accurately we can predict future location movements using location predictors based on the Markov model. One of the authors collected GPS traces by carrying a GPS device over a period of 50 days. We recorded a GPS data point every three seconds, and collected 111K data points in total. The data captures location movements while driving as well as walking.

We consider GPS data whose data points reside within a rectangular region, which covers the campus of University of Illinois and its surrounding off-campus areas. The dimension of the region is 4.8 kilometers times 4.0 kilometers. We divide each coordinate into 40 units and define 1,600 rectangular location regions, each of which is about the size of a typical building in town.

We used half of the data to construct a state transition matrix M and used the other half to compute the accuracy of the predictions with matrix M . We implement the rows and columns of M as balanced binary trees to maintain a largely sparse matrix of a large size efficiently. When we construct M , we do not consider movements within the same location; we only consider movements to a different location. Figure 3 shows our experimental results. The X-axis shows how many

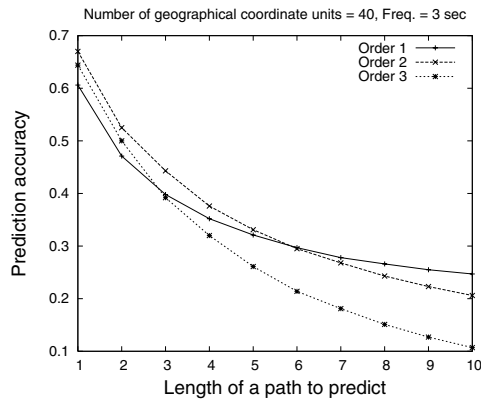


Figure 3: Accuracy of location predictions.

steps we predict ahead, and the Y-axis shows the accuracy of our predictions. We computed the accuracy of predicting every next location and took its average. When we predict a next location in a single time step with a 1-order Markov model, our predictions are about 60% accurate. However, as we try to predict a location reachable in greater number of steps, the prediction accuracy decreases. We compare the results of Markov models of three different orders. As we can see, when we predict locations reachable in a fewer number of time steps, we can improve the accuracy by 10% by using a higher-order Markov model, which considers multiple previous movements. Thus, we believe that the threat of an adversary with a location prediction is real considering these preliminary results.

When we predict a location multiple steps ahead, using a higher-order model is not useful. Since our current dataset is small, we plan to study this issue with a larger data set in the future. We also plan to represent historical data in a higher level format (e.g., a series of connected segments punctuated by turns), as proposed in [3, 11], to improve the accuracy of predictions as our next step.

5. RELATED WORK

Location privacy has been studied heavily in the context of the anonymization and obfuscation of location data (See [10] for a comprehensive survey). The focus of research in this area is to ensure that no anonymized and/or obfuscated data is associated with an individual. This inference problem concerning location privacy is different from ours since we consider the inference problem in access-control systems for LSSs, which release identifiable location data of mobile users.

6. DISCUSSION

Having studied our initial proposal, we identified several key questions for future research.

Background knowledge of an adversary. What sources of background information should be available to a prediction model of an adversary? We make the conservative assumption that an adversary knows all the previous movements of a target user. However, this may be too strong an assumption. For example, if Alice visits a doctor’s office on a regular basis, but wants to hide

it from Bob, it may not make sense to assume that Bob knows about her propensity to in fact be there. It may be possible to eliminate private locations from the training data for the predictor, but augment it with map data or collection of traces from other users to enable baseline predictions of trajectories and paths.

Privacy metrics. How should location privacy be measured?

Our privacy metrics in Section 3 considers the probability of visiting private locations. However, a location predictor could assign a non-trivial prior probability of a sensitive location event, making it nearly impossible to achieve (M, δ) -privacy for small δ . Conversely, it may be the case that Bob can predict $P(\text{“Alice is in a bar”})$ with a probability greater than δ , even though the probability of her being in any particular bar is below the threshold, technically satisfying the privacy constraint.

An alternate approach would be to use information theory to model privacy. For example, we can define random variables Y_1, Y_2, \dots to be a projection of the location variables X_1, X_2, \dots , where, whenever $X_i = l$ for some non-private location $Bob \in acl(\text{Alice}, l)$, $Y_i = \perp$, and otherwise $Y_i = X_i$. Then we could impose a bound on the mutual information between the the observations the LSS shares with Bob and $Y^{(i)}$, or the information gain regarding Alice’s position among the private locations. However, the result will depend on the prior knowledge that Bob has; in theory, one could attempt to ensure the bound holds for an arbitrary prior, but this may be difficult to verify computationally.

Information gain through a request denial. Simply not revealing Alice’s location may tell Bob something sensitive. Moreover, if Bob can learn the exact release policy used by the LSS, he can use that to update his predictor. For example, in Figure 1, suppose that the LSS always reveals Alice’s location on the way to the library, but never on the way to the hospital. Then, when Bob sees updates from Alice stop after she reaches the intersection, Bob can deduce that she is in fact going to the hospital. A potential approach to remedy this is to not reveal Alice’s location even on the way to the library, at least probabilistically, but this would degrade the utility of the LSS and require careful probabilistic modeling to ensure privacy is preserved. Another alternative is to provide fake location data that looks plausible [2, 9] rather than denying requests.

7. CONCLUSIONS

In this paper, we study an issue of inference attacks on GPS traces when we support mobile users’ privacy policies for LSSs. We define an adversary who has access to a mobile user’s previous location data as a location predictor based on the Markov model, and then gives location privacy metrics under the presence of such inference attacks. Our preliminary experimental results show that it is possible to predict a mobile user’s future location with high accuracy and thus we need an additional mechanism for enforcing the privacy policies of users by hiding routes towards private locations in a proper way.

Future work includes collecting a large dataset of GPS traces from many mobile users and evaluating different location prediction algorithms with that dataset. We also study other location privacy metrics considering an adversary with different external knowledge and inference methods.

Acknowledgments

This research is in part supported by grant from the Promotion program for Reducing global Environmental load through ICT innovation (PREDICT) of the Ministry of Internal Affairs and Communications in Japan.

8. REFERENCES

- [1] Denise Anthony, Tristan Henderson, and David Kotz. Privacy in location-aware computing environments. *IEEE Pervasive Computing*, 6(4):64–72, 2007.
- [2] Richard Chow and Philippe Golle. Faking contextual data for fun, profit, and privacy. In *Proceedings of the 8th ACM Workshop on Privacy in the Electronic Society (WPES)*, pages 105–108, New York, NY, USA, 2009. ACM.
- [3] Jon Froehlich and John Krumm. Route prediction from trip observations. In *Society of Automotive Engineers (SAE) 2008 World Congress*, April 2008.
- [4] Glympse. <http://www.glympse.com>.
- [5] Google latitude. <http://www.google.com/latitude>.
- [6] Urs Hengartner and Peter Steenkiste. Access control to people location information. *ACM Transactions on Information and System Security (TISSEC)*, 8(4):424–456, 2005.
- [7] Instamapper. <http://www.instamapper.com>.
- [8] Apu Kapadia, Tristan Henderson, Jeffrey J. Fielding, and David Kotz. Virtual Walls: Protecting Digital Privacy in Pervasive Environments. In *Proceedings of the Fifth International Conference on Pervasive Computing (Pervasive)*, volume 4480 of *LNCS*, pages 162–179. Springer-Verlag, May 2007.
- [9] John Krumm. Realistic Driving Trips For Location Privacy. In *Pervasive Computing*, number 5538 in Lecture Notes In Computer Science. Springer, 2009.
- [10] John Krumm. A survey of computational location privacy. *Personal and Ubiquitous Computing*, 13(6):391–399, 2009.
- [11] John Krumm and Eric Horvitz. Predestination: Inferring destinations from partial trajectories. In *Proceedings of the 8th International Conference on Ubiquitous Computing (UbiComp)*, pages 243–260, September 2006.
- [12] Ginger Myles, Adrian Friday, and Nigel Davies. Preserving privacy in environments with location-based applications. *IEEE Pervasive Computing*, 2(1):56–64, January-March 2003.
- [13] Libo Song, David Kotz, Ravi Jain, and Xiaoning He. Evaluating next cell predictors with extensive Wi-Fi mobility data. *IEEE Transactions on Mobile Computing*, 5(12):1633–1649, December 2006.
- [14] How to tweet with your location. <http://twitter.zendesk.com/entries/122236>.