

# Geographic Differential Privacy for Mobile Crowd Coverage Maximization

Leye Wang<sup>1</sup>, Gehua Qin<sup>2</sup>, Dingqi Yang<sup>3</sup>, Xiao Han<sup>4</sup>, Xiaojuan Ma<sup>1</sup>

<sup>1</sup>The Hong Kong University of Science and Technology, <sup>2</sup>Shanghai Jiao Tong University

<sup>3</sup>University of Fribourg, <sup>4</sup>Shanghai University of Finance and Economics

wly@cse.ust.hk, qingehua@gmail.com, dingqi@exascale.info, xiaohan@mail.shufe.edu.cn, mxj@cse.ust.hk

## Abstract

For real-world mobile applications such as location-based advertising and spatial crowdsourcing, a key to success is targeting mobile users that can maximally cover certain locations in a future period. To find an optimal group of users, existing methods often require information about users' mobility history, which may cause privacy breaches. In this paper, we propose a method to maximize mobile crowd's future location coverage under a guaranteed location privacy protection scheme. In our approach, users only need to upload one of their frequently visited locations, and more importantly, the uploaded location is obfuscated using a geographic differential privacy policy. We propose both analytic and practical solutions to this problem. Experiments on real user mobility datasets show that our method significantly outperforms the state-of-the-art geographic differential privacy methods by achieving a higher coverage under the same level of privacy protection.

## Introduction

Crowd coverage maximization is a classical problem in mobile computing: how to select  $m$  users from a candidate pool to maximize the probability of covering a set of target locations in a coming time period (e.g., one day or one week). This problem and its variants have a wide spectrum of applications in location-based advertising (Dhar and Varshney 2011), spatial crowdsourcing (Chen and Shahabi 2016; Zhang et al. 2014), urban computing (Zheng et al. 2014), etc. For example, it can help shop owners to offer electronic coupons to the set of mobile app users who may physically visit the region around the shop soon; it can also help crowdsourcing organizers to recruit the participants to cover the task area with the highest probability (Xiong et al. 2016).

One of the key steps in crowd coverage maximization is *mobility profiling*, i.e., predicting the probability of a user appearing at a certain location. A common practice is first dividing an area into fine-grained grids or sub-areas, and then counting the frequency of a user appearing in each grid based on trajectory history (Guo et al. 2017). One can use more sophisticated models like Poisson process to estimate users' occurrence distribution (Xiong et al. 2016). Existing mobility profiling methods often require access to

users' historical mobility traces, which may seriously compromise user privacy. For example, users' exposed location data may reveal sensitive information about their identities and social relationships (Cho, Myers, and Leskovec 2011; Rossi et al. 2015). Despite the importance of location privacy, as far as we know, there is little research effort combining location privacy, mobility profiling, and crowd coverage maximization up to date.

To fill this gap, this paper aims to explore how to protect the crowds' location privacy, while still optimizing their expected coverage of a set of locations. To achieve this goal, we propose a mobile crowd coverage maximization framework with a rigorous privacy protection scheme — *geographic differential privacy* (Andrés et al. 2013). A geographic differential privacy policy obfuscates a user's actual location to another with carefully designed probabilities, such that adversaries, regardless of their prior knowledge, can learn little about the user's true location after observing the obfuscated locations. However, with differential privacy protection, crowd coverage maximization can only be performed based on the obfuscated (inaccurate) locations, which leads to inevitable loss of the quality of the selected users. Therefore, we propose a method to generate the optimal location obfuscation policy which satisfies geographic differential privacy while minimizing such loss.

In summary, this paper has the following contributions:

(1) To the best of our knowledge, this is the first work studying the mobile crowd coverage maximization problem with location privacy protection.

(2) In our approach, users only need to upload one of their frequently visited locations, and more importantly, the uploaded location is obfuscated using the rigorous privacy policy — geographic differential privacy. We further formulate an optimization problem to obtain the optimal obfuscation policy that can maximize the expected future crowd coverage over a set of locations under a guaranteed level of differential privacy protection. As the optimization problem is non-convex, we first mathematically analyze the scenario when only one location needs to cover and then derive an optimal solution. Then, we extend this setting to the multi-location coverage scenario and propose a practical algorithm to obtain the optimal obfuscation policy.

(3) Experiments on real human mobility datasets verify that, by selecting the same number of users under the same

level of privacy protection, our method achieves a higher coverage than state-of-the-art differential privacy methods.

## Preliminaries

**Geographic differential privacy** (Andrés et al. 2013) introduces the idea of database differential privacy (Dwork 2008) into the location obfuscation context. Its key idea is: given an observed obfuscated location  $l^*$ , any two locations  $l_1$  and  $l_2$  have similar probabilities of being mapped to  $l^*$ . It is thus hard for an adversary to differentiate whether the user is at  $l_1$  or  $l_2$  by observing  $l^*$ .

**Definition 1** (Andrés et al. 2013). *Suppose the target area includes a set of locations  $\mathcal{L}$ , then an obfuscation policy  $P$  satisfies **geographic  $\epsilon$ -differential privacy**, iff.*

$$P(l^*|l_1) \leq e^{\epsilon d(l_1, l_2)} P(l^*|l_2) \quad \forall l_1, l_2, l^* \in \mathcal{L} \quad (1)$$

where  $P(l^*|l)$  is the probability of obfuscating  $l$  to  $l^*$ ,  $d(l_1, l_2)$  is the distance between  $l_1$  and  $l_2$ ,  $\epsilon$  is the privacy budget — the smaller  $\epsilon$ , the better privacy protection.

Note that the set of locations are usually constructed by dividing the target area into subregions, e.g., equal-size grids (Bordenabe, Chatzikokolakis, and Palamidessi 2014) or cell-tower regions (Xiong et al. 2016).

If  $P$  satisfies geographic differential privacy, it can be proven that for adversaries with *any* prior knowledge about users’ location distributions, their posterior knowledge after observing the obfuscated location can only be increased by a small constant factor (Andrés et al. 2013). Note that this protection is guaranteed even if the adversaries know  $P$ . Due to this rigorous protection effect, geographic differential privacy has seen many applications in location based services, spatial crowdsourcing, etc. (Bordenabe, Chatzikokolakis, and Palamidessi 2014; Wang et al. 2016; Wang et al. 2017).

**Mobility profiling** aims to estimate the probability of a user covering a certain location during a time period in the future. Specifically, a user  $u_i$ ’s mobility profile is denoted as  $M_i$ , and  $M_i(l_j)$ ,  $l_j \in \mathcal{L}$  means the estimated probability of  $u_i$  visiting  $l_j$  in a concerned future period (e.g., next week). Commonly used mobility profiling methods include frequency-based (Guo et al. 2017) and Poisson-based (Xiong et al. 2016) algorithms. We use the Poisson process to model user mobility given its better prediction performance in our experiments. More details can be found in the appendix.

## Framework Overview

We present an overview of our privacy framework in Figure 1. The key idea of our framework is that users should expose their location information as little as possible, while we can still select a proper set of users for optimizing their coverage on certain target locations in the future.

The two main players in our framework are a server platform and its mobile client users. As we want users to expose their actual location information as little as possible, user mobility profiling runs locally on individuals’ smart devices. That means, the clients’ mobility profiles are only

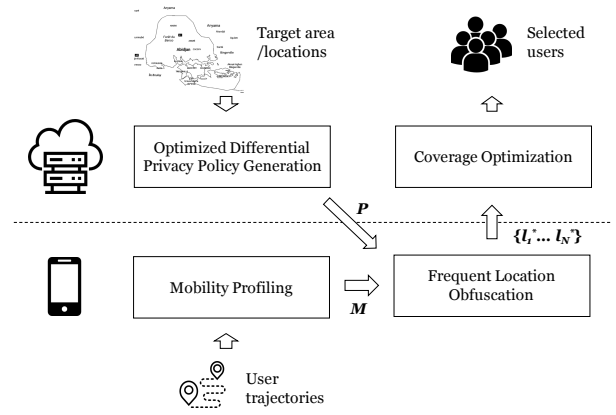


Figure 1: Framework overview.

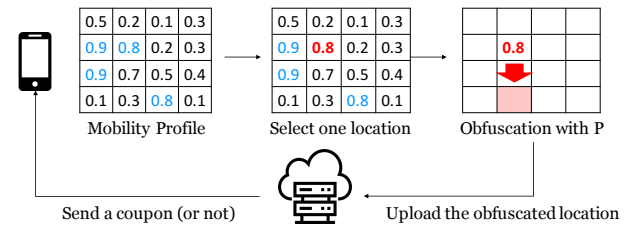


Figure 2: A running example of our framework.

known to themselves. As shown in the literature, only uploading frequent locations with high profiling probabilities (e.g.,  $> 80\%$ ) to the server can already help achieve a good future crowd coverage (Guo et al. 2017). To limit the potential location leakage, our framework only requires users to upload *one* of their frequent locations. Moreover, this frequent location is obfuscated by the geographic differential privacy policy  $P$  before being sent to the platform. The policy  $P$  is generated by the server based on which target locations need to be covered. Finally, according to the uploaded obfuscated frequent locations  $\{l_1^* \dots l_N^*\}$  (suppose  $N$  users), the platform aims to select a set of users to maximize the expected coverage of intended locations in the coming period.

A running example is shown in Figure 2, where the table represents a user’s mobility profile in a 2D spatial area split into uniform grids. Suppose that a location-based advertising platform needs to decide whether to send a Starbucks coupon to a user. The platform expects that a user receiving the coupon is a frequent visitor to the regions where Starbucks stores are located, so that the user will probably go to the stores. To achieve this goal, first, a user client computes its owner’s mobility profile locally. Second, from the set of locations whose probabilities are larger than a threshold (e.g.,  $80\%$ ), the user client randomly selects one location  $l_u$  to be uploaded to the server.<sup>1</sup> Third, according to the privacy policy  $P$  received from the server, the user client randomly obfuscates  $l_u$  to  $l_u^*$  and then sends  $l_u^*$  to the server.

<sup>1</sup>If there is no location with the probability larger than the threshold, then the user does not upload any location.

Finally, the server will decide whether to send the coupon to the user or not according to the uploaded  $l_u^*$ . In this case, a user's location privacy is preserved as the uploaded frequent location is rigorously obfuscated with differential privacy.

Location obfuscation would inevitably introduce certain loss of quality in selecting users for coverage optimization, as users' uploaded frequent locations contain deliberate noises. Hence, how the server generates the privacy policy  $P$  is the key challenge of our framework, which aims to minimize the loss of quality caused by privacy protection.

### Optimal Privacy Policy

In this section, we illustrate our solution that guarantees geographic  $\epsilon$ -differential privacy while minimizing the loss of quality in mobile crowd coverage optimization.

#### Single Location Coverage Problem (SLCP)

As the first step, we analyze the scenario where only one location needs to be covered. In location-based advertising, this reflects the scenario that the advertising only involves one specific site (e.g., a newly opened restaurant). In spatial crowdsourcing, this means that the task is only associated with one location (e.g., taking the photo of Statue of Liberty). Suppose the target location to cover as  $l_t$  and a user submits her/his obfuscated frequent location as  $l^*$ , then the probability of her/his frequent location being actually  $l_t$  is:

$$prob(l_t|l^*) = \frac{\pi(l_t)P(l^*|l_t)}{\sum_{l \in \mathcal{L}} \pi(l)P(l^*|l)} \quad (2)$$

where  $\pi$  is the overall distribution of all the users' frequent locations. Here we suppose that we can foreknow  $\pi$ , and later we will elaborate how to estimate it. Note that the denominator can be seen as the overall probability of a user reporting her/his frequent location as  $l^*$ .

Suppose we select a user reporting  $\hat{l}^*$  to cover the target location  $l_t$  in the coming time period, apparently we would like to maximize Eq. 2 so that the future probability of the user covering  $l_t$  is maximized. With this idea, we have the following optimization process to get the optimal privacy policy  $\hat{P}$ . Particularly, given  $l_t$  to cover, we aim to

$$\max_{\hat{l}^*, \hat{P}} \frac{\pi(l_t)\hat{P}(\hat{l}^*|l_t)}{\sum_{l \in \mathcal{L}} \pi(l)\hat{P}(\hat{l}^*|l)} \quad (3)$$

$$s.t. \quad \hat{P}(l^*|l_1) \leq e^{\epsilon d(l_1, l_2)} \hat{P}(l^*|l_2) \quad \forall l_1, l_2, l^* \in \mathcal{L} \quad (4)$$

$$\hat{P}(l^*|l) > 0 \quad \forall l, l^* \in \mathcal{L} \quad (5)$$

$$\sum_{l^* \in \mathcal{L}} \hat{P}(l^*|l) = 1 \quad \forall l \in \mathcal{L} \quad (6)$$

Eq. 4 is the constraint of geographic differential privacy; Eq. 5 and 6 are probability restrictions. By solving the above optimization problem, we can get the optimal privacy policy  $\hat{P}$ , as well as the user selection strategy, i.e., selecting the users reporting  $\hat{l}^*$  for future coverage maximization.

However, even given  $\hat{l}^*$ , Eq. 3 cannot be converted to a convex optimization problem with existing solutions (Boyd and Vandenberghe 2004). To overcome this difficulty, we

then analyze the relationship between the constraints and the objective function, and then deduce an optimal solution analytically.

#### An Analytic Solution to SLCP

Our analytic deduction includes three steps. First, we verify that the selection of  $\hat{l}^*$  will not affect the optimal objective value of Eq. 3. Second, we prove that Eq. 3 cannot exceed a certain upper bound. Finally, we show that this upper bound can be achieved by constructing a feasible solution of  $\hat{P}$ .

**Lemma 1.** For any two locations  $l_1^*, l_2^* \in \mathcal{L}$ , the optimal objective values of Eq. 3 are the same if we set  $\hat{l}^* = l_1^*$  or  $l_2^*$ .

*Proof.* For  $\hat{l}^* = l_1^*$  or  $l_2^*$ , we can always find a pair of  $P_1, P_2$ , where  $P_1(l_1^*|l) = P_2(l_2^*|l)$ ,  $P_1(l_2^*|l) = P_2(l_1^*|l)$ , and  $P_1(l^*|l) = P_2(l^*|l)$  for other  $l^*$ ;  $P_1$  and  $P_2$  lead to the same objective value. A detailed proof is in the appendix.  $\square$

*Remark.* Lemma 1 demonstrates that we can use any location  $l \in \mathcal{L}$  as the obfuscated location  $\hat{l}^*$  for user selection without impacting the achievable optimal coverage utility.

**Lemma 2.** The optimal value of Eq. 3 cannot exceed

$$\frac{\pi(l_t)}{\sum_{l \in \mathcal{L}} \pi(l)e^{-\epsilon d(l, l_t)}} \quad (7)$$

and this value can only be achieved if

$$\hat{P}(\hat{l}^*|l) \propto e^{-\epsilon d(l, l_t)}, \quad \forall l \in \mathcal{L} \quad (8)$$

*Proof.* With geographic differential privacy constraints,

$$\frac{\pi(l_t)\hat{P}(\hat{l}^*|l_t)}{\sum_{l \in \mathcal{L}} \pi(l)\hat{P}(\hat{l}^*|l)} = \frac{\pi(l_t)}{\sum_{l \in \mathcal{L}} \pi(l) \frac{\hat{P}(\hat{l}^*|l)}{\hat{P}(\hat{l}^*|l_t)}} \quad (9)$$

$$\leq \frac{\pi(l_t)}{\sum_{l \in \mathcal{L}} \pi(l) \frac{e^{-\epsilon d(l, l_t)} \hat{P}(\hat{l}^*|l_t)}{\hat{P}(\hat{l}^*|l_t)}} = \frac{\pi(l_t)}{\sum_{l \in \mathcal{L}} \pi(l)e^{-\epsilon d(l, l_t)}} \quad (10)$$

$\square$

*Remark.* Lemma 2 points out an upper bound of the optimal objective value and the condition (Eq. 8) that  $\hat{P}$  must satisfy for getting the upper bound value. However, whether we can find a feasible  $\hat{P}$  satisfying Eq. 8, as well as Eq. 4-6 is still unknown. Next, we prove that this  $\hat{P}$  exists.

**Lemma 3.** If  $\hat{P}(\hat{l}^*|l)$  satisfies Eq. 8, then

$$\hat{P}(\hat{l}^*|l_1) \leq e^{\epsilon d(l_1, l_2)} \hat{P}(\hat{l}^*|l_2) \quad \forall l_1, l_2 \in \mathcal{L} \quad (11)$$

*Proof.* Considering that  $d$  is a distance metric, then

$$\frac{\hat{P}(\hat{l}^*|l_1)}{\hat{P}(\hat{l}^*|l_2)} = e^{\epsilon(d(l_2, l_t) - d(l_1, l_t))} \leq e^{\epsilon d(l_1, l_2)} \quad (12)$$

$\square$

*Remark.* Lemma 3 proves that when Eq. 8 stands, Eq. 4 of  $l^* = \hat{l}^*$  must also hold for any  $l_1, l_2$ .

**Theorem 1.** Given any  $\hat{l}^*$ , we can get a feasible  $\hat{P}$ ,

$$\hat{P}(\hat{l}^*|l) = \theta e^{-\epsilon d(l, l_t)}, \quad \forall l \in \mathcal{L} \quad (13)$$

$$\hat{P}(l^*|l) = \frac{1 - \theta e^{-\epsilon d(l, l_t)}}{|\mathcal{L}| - 1}, \quad \forall l^*, l \in \mathcal{L} \text{ and } l^* \neq \hat{l}^* \quad (14)$$

which can achieve the upper bound Eq. 7. Here,  $\theta$  can be any positive constant value smaller than or equal to a threshold  $\tau$ , where

$$\tau = \min_{l_1, l_2 \in \mathcal{L}} \frac{e^{\epsilon d(l_1, l_2)} - 1}{e^{-\epsilon(d(l_2, l_t) - d(l_1, l_2))} - e^{-\epsilon d(l_1, l_t)}} \quad (15)$$

The proof is in the appendix.

Note that while Theorem 1 gets an optimal solution, in reality, there may not be enough users who report  $\hat{l}^*$  for selection (if  $\tau$  is too small and the total user number is limited). Later we will propose a practical solution overcoming this shortcoming, when addressing the multi-location scenario.

### Multi-Location Coverage Problem (MLCP)

A more complicated setting for mobile crowd coverage problem includes a set of locations that need to be covered. Real-life examples include delivering coupons of chain stores to users who will probably visit any of them in the next time period. Denote the set of locations to cover as

$$\mathbb{L} = \{l_t^1, l_t^2, \dots, l_t^z\} \subset \mathcal{L} \quad (16)$$

then the probability of a user's actual frequent location belonging to  $\mathbb{L}$  is:

$$\sum_{l_t \in \mathbb{L}} \text{prob}(l_t | l^*) = \frac{\sum_{l_t \in \mathbb{L}} \pi(l_t) P(l^* | l_t)}{\sum_{l \in \mathcal{L}} \pi(l) P(l^* | l)} \quad (17)$$

Then, we can maximize Eq. 17 with the constraints Eq. 4-6 to get the optimal privacy policy  $\hat{P}$ , and the obfuscated location  $\hat{l}^*$  for future crowd coverage maximization.

$$\max_{\hat{l}^*, \hat{P}} \frac{\sum_{l_t \in \mathbb{L}} \pi(l_t) \hat{P}(\hat{l}^* | l_t)}{\sum_{l \in \mathcal{L}} \pi(l) \hat{P}(\hat{l}^* | l)} \quad (18)$$

$$\text{s.t. Eq. 4-6} \quad (19)$$

Similar to the single location coverage problem, we can prove the following lemmas.

**Lemma 4.** For any  $l_1^*, l_2^* \in \mathcal{L}$ , the optimal objective values of Eq. 18 are the same if we set  $\hat{l}^* = l_1^*$  or  $l_2^*$ .

Lemma 4 is a straightforward extension of Lemma 1 to the multiple location coverage scenario.

**Lemma 5.** The optimal value of Eq. 18 cannot exceed

$$\left(1 + \sum_{l \notin \mathbb{L}} \sum_{l_t \in \mathbb{L}} \frac{\pi(l)}{\pi(l_t)} e^{-\epsilon d(l, l_t)}\right)^{-1} \quad (20)$$

and this value can be achieved only if

$$e^{\epsilon d(l, l_t)} P(\hat{l}^* | l) = P(\hat{l}^* | l_t), \quad \forall l_t \in \mathbb{L}, \forall l \notin \mathbb{L} \quad (21)$$

The detailed proof is in the appendix.

Although Lemma 5 seems to be an extension of Lemma 2 for the multi-location scenario, they have a significant difference that the optimal value Eq. 20 may not always be feasible, i.e., Eq. 21 may not stand. Take a toy example of  $\mathbb{L}$  containing two locations, it means that, for any  $l \notin \mathbb{L}$

$$e^{\epsilon d(l, l_t^1)} P(\hat{l}^* | l) = P(\hat{l}^* | l_t^1) \quad (22)$$

$$e^{\epsilon d(l, l_t^2)} P(\hat{l}^* | l) = P(\hat{l}^* | l_t^2) \quad (23)$$

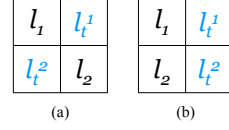


Figure 3: Toy examples with two locations to cover.

Then, for any two locations  $l_1, l_2 \notin \mathbb{L}$ , let  $l = l_1$  or  $l_2$ , then

$$\frac{P(\hat{l}^* | l_t^1)}{P(\hat{l}^* | l_t^2)} = \frac{e^{\epsilon d(l_1, l_t^1)}}{e^{\epsilon d(l_1, l_t^2)}} = \frac{e^{\epsilon d(l_2, l_t^1)}}{e^{\epsilon d(l_2, l_t^2)}} \quad (24)$$

$$\Rightarrow d(l_1, l_t^1) - d(l_1, l_t^2) = d(l_2, l_t^1) - d(l_2, l_t^2) \quad (25)$$

Hence, if Eq. 20 is feasible, Eq. 25 must hold. Figure 3 shows two examples, in one of which Eq. 25 stands (Figure 3a) and the other does not (Figure 3b, considering the Euclidean distance). This shows that whether Eq. 20 can be achieved depends on the distribution of the target locations.

### A Practical Solution to MLCP

While we cannot always obtain the upper bound value of Eq. 20 for the multi-location coverage problem, here we propose a practical solution which can work in real scenarios.

Revisiting the objective function of the multi-location coverage problem, Eq. 18, we can see that the main difficulty in solving the optimization problem is that the denominator includes  $\hat{P}$  in it. To address this issue, we propose to add one more constraint to the optimization process by setting the denominator to a constant value,

$$\sum_{l \in \mathcal{L}} \pi(l) \hat{P}(\hat{l}^* | l) = \beta \quad (26)$$

where  $\beta$  is a constant between 0 and 1; we will later elaborate how to set  $\beta$ . With Eq. 26, the objective function is,

$$\max_{\hat{l}^*, \hat{P}} \frac{\sum_{l_t \in \mathbb{L}} \pi(l_t) \hat{P}(\hat{l}^* | l_t)}{\beta} \quad (27)$$

Lemma 4 has shown that we can set  $\hat{l}^*$  to any  $l \in \mathcal{L}$  without affecting the optimal objective value. Since Eq. 4-6 are all linear constraints, we can then use state-of-the-art linear programming tools (e.g., Mosek and Gurobi) to solve the optimization problem to get the optimal privacy policy  $\hat{P}$ .

**Setting  $\beta$  with Binomial Distribution.** We then discuss how to set  $\beta$  in real-life scenarios. First, we prove that if we want to get the objective value as high as possible, we should set  $\beta$  as small as possible.

**Theorem 2.** Given  $\hat{l}^*$ , suppose  $v_1, v_2$  are the two optimal objective values of Eq. 27 when we set  $\beta$  to  $b_1, b_2$ , respectively, and  $b_1 < b_2$ , then  $v_1 \geq v_2$ .

*Proof.* We denote the optimal  $\hat{P}$  when setting  $\beta$  to  $b_1, b_2$  as  $\hat{P}_1, \hat{P}_2$ , respectively. Then, we construct a new solution of  $\hat{P}'_1$  when  $\beta = b_1$  as follows:

$$P'_1(\hat{l}^* | l) = \theta \hat{P}_2(\hat{l}^* | l) \quad \forall l \in \mathcal{L}$$

$$P'_1(l' | l) = \hat{P}_2(l' | l) + (1 - \theta) \hat{P}_1(l' | l) \quad \forall l \in \mathcal{L}, l' \neq \hat{l}^*$$

---

**Algorithm 1:** Optimal policy for multi-location coverage.

**Input** :  $\pi$ : overall user spatial distribution.  
 $\epsilon$ : differential privacy budget.  
 $\mathcal{L}$ : whole set of locations.  
 $L$ : set of target locations to cover.  
 $N$ : total number of users.  
 $\alpha$ : number of users to select.  
 $\rho$ : probability threshold for user selection.  
**Output**:  $\hat{P}$ : optimal differential privacy policy.  
 $\hat{l}^*$ : the obfuscated location to select users.

- 1  $\hat{l}^* \leftarrow l_1$  (or any other  $l \in \mathcal{L}$ );
- 2  $\beta \leftarrow$  the minimum value that can ensure  $Pr(X \geq \alpha) \geq \rho$  for the Binomial distribution  $B(X, N, \beta)$ ;
- 3 Solve the linear program to get optimal  $\hat{P}$ :

$$\begin{aligned} & \max_{\hat{P}} \frac{\sum_{l_t \in \mathcal{L}} \pi(l_t) \hat{P}(\hat{l}^* | l_t)}{\beta} \\ \text{s.t. } & \hat{P}(l^* | l_1) \leq e^{\epsilon d(l_1, l_2)} \hat{P}(l^* | l_2) \quad \forall l_1, l_2, l^* \in \mathcal{L} \\ & \hat{P}(l^* | l) > 0 \quad \forall l, l^* \in \mathcal{L} \\ & \sum_{l^* \in \mathcal{L}} \hat{P}(l^* | l) = 1 \quad \forall l \in \mathcal{L} \\ & \sum_{l \in \mathcal{L}} \pi(l) \hat{P}(\hat{l}^* | l) = \beta \end{aligned}$$

4 **return**  $\hat{l}^*, \hat{P}$ ;

---

where  $\theta = b_1/b_2$ . All the constraints of Eq. 4-6 still stand for  $P'_1$ . As the optimal objective value is  $v_1$  when  $\beta = b_1$ ,

$$v_1 \geq \frac{\sum_{l_t \in \mathcal{L}} \pi(l_t) P'_1(\hat{l}^* | l_t)}{b_1} = \frac{\sum_{l_t \in \mathcal{L}} \pi(l_t) \hat{P}_2(\hat{l}^* | l_t)}{b_2} = v_2$$

□

*Theorem 2* is very important for our practical solution, because it tells us that to get the optimal solution, we only need to solve the linear program *once* by setting  $\beta$  to the smallest value that we can accept, rather than enumerating all the possible  $\beta$ . On the other hand,  $\beta$  can be seen as the overall probability that a user will report her/his frequent location as  $\hat{l}^*$ . Since we need to select users from such users, we cannot set  $\beta$  to a too small value, which will lead to very few people reporting their locations as  $\hat{l}^*$ . Therefore, we propose a method to set  $\beta$ , with a guarantee that the platform can find  $\alpha$  users with a probability of  $\rho$  (e.g., 95%) as follows.

Suppose totally  $N$  users report their frequent locations, then we can estimate the number of users who will report their obfuscated frequent locations as  $\hat{l}^*$  with the Binomial probability  $Pr(X = m) = B(m, N, \beta)$ . Then, the probability that we can find at least  $\alpha$  users is that,

$$Pr(X \geq \alpha) = \sum_{m=\alpha}^N B(m, N, \beta) \quad (28)$$

And thus we would like to set  $\beta$  to the smallest value that ensures  $Pr(X \geq \alpha) \geq \rho$ .

We describe the pseudo-code of our practical solution for the private multi-location coverage problem in Algorithm 1. Note that since covering one location is a special case of covering multiple locations, Algorithm 1 can also solve the single location coverage problem, without the need to assume that we will always have enough users reporting  $\hat{l}^*$ .

---

**Algorithm 2:** User selection with dynamic estimating  $\pi$ .

**Input** :  $k$ : number of user groups to split.  
other inputs like Algorithm 1, except that  $\pi$  is unknown.

**Output**:  $U^*$ : selected users.

- 1  $\pi \leftarrow$  uniform distribution (or other proper initial distribution);
  - 2  $U_1, U_2, \dots, U_k \leftarrow N$  users are split into  $k$  groups, each with  $N/k$  users;
  - 3 **for**  $i = 1, 2, \dots, k$  **do**
  - 4      $\hat{l}^*, \hat{P} \leftarrow$  run Algorithm 1 with  $\pi$ ;
  - 5     **foreach**  $u \in U_i$  **do**
  - 6         /\*  $u$  downloads  $\hat{P}$  to the mobile client \*/
  - 7          $l_u \leftarrow$  a randomly selected frequent location;
  - 8          $l_u^* \leftarrow$  obfuscating  $l_u$  by  $\hat{P}$ ;
  - 9         /\*  $u$  uploads  $l_u^*$  to the server \*/
  - 10          $\pi'_u(l) \leftarrow \frac{\pi(l) \hat{P}(l_u^* | l)}{\sum_{l' \in \mathcal{L}} \pi(l') \hat{P}(l_u^* | l')}, \forall l \in \mathcal{L}$ ;
  - 11     **end**
  - 12      $\pi \leftarrow$  the mean value of  $\pi'_u$  over  $u \in U_i$ ;
  - 13 **end**
  - 14  $U^* \leftarrow \emptyset$ ;
  - 15 **for**  $j = k, k-1, \dots, 1$  **do**
  - 16     **foreach**  $u \in U_j$  **do**
  - 17         **if**  $u$ 's obfuscated location is  $\hat{l}^*$  **then**
  - 18              $U^* \leftarrow U^* \cup \{u\}$ ;
  - 19             **if**  $|U^*| == \alpha$  **then**
  - 20                 **return**  $U^*$ ;
  - 21             **end**
  - 22     **end**
  - 23 **end**
  - 24 **return**  $U^*$ ;
- 

**Estimating Overall Location Distribution  $\pi$ .** Previously, we assume that we have known the overall frequent location distribution  $\pi$ . This may be possible when we have other sources to infer  $\pi$ , e.g., mobile call logs (Blondel et al. 2012). However, if we do not have such data, other methods are required to estimate  $\pi$  along with user selection. We thus propose a Bayes rule based method to do user selection and  $\pi$  estimation simultaneously, as shown in Algorithm 2.

Our basic idea is using users' uploaded obfuscated locations to refine  $\pi$ . Note that our mechanism requires that each user uploads the obfuscated location only once to ensure differential privacy protection (Andrés et al. 2013). Hence, to preserve differential privacy, we split all the users into  $k$  groups, get users' obfuscated locations group by group, and iteratively refine  $\pi$  with the obfuscated locations from previous user groups. The key update formula of  $\pi$  is the Bayes rule in line 8. In such a way, the estimated  $\pi$  gradually reaches the actual  $\pi$  after iterative refinements. As  $\pi$  generally becomes more and more accurate, the final user selection is biased to the users in the groups who upload locations later (line 13-22). The number of groups  $k$  balances the trade-off between algorithm running efficiency and solution quality — larger  $k$  updates  $\pi$  more frequently, but costs more time as it involves  $k$  iterations of running Algorithm 1.

Note that in real implementation, users who do not have any frequent locations can still upload 'NULL' to the server. Then, we can estimate the percentage of users who can report locations from previous user groups. This can help us to set an appropriate  $\beta$  used in the optimization so as to finally find  $\alpha$  users with a probability of  $\rho$ .

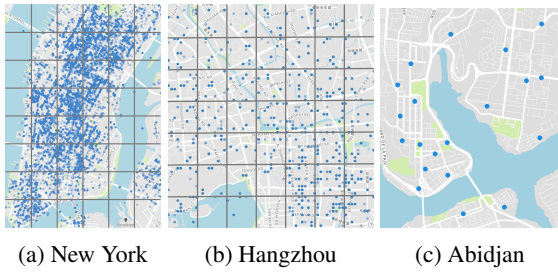


Figure 4: Experiment areas. Points in NY and Hangzhou are user locations, and points in Abidjan are cell towers.

## Experiments

In this section, we conduct empirical studies on three real user mobility datasets. We use Algorithm 2 for both single and multi-location coverage scenarios given its practicality (no need to foreknow  $\pi$ ).

### Baselines

- *Laplace*. The state-of-the-art method to achieve geographic differential privacy is based on the Laplace distribution (Andrés et al. 2013).
- *NO*. We use the No-Obfuscation (NO) policy, i.e., the users upload one of their real frequent locations to the server, to show an upper bound of the coverage.
- *Random*. We use the random user selection to serve as the lower bound of the coverage that can be achieved.

### Datasets

- *FS* dataset (Yang et al. 2016) contains 1083 Foursquare users’ check-ins in New York, USA across near one year. We set the time period to a *weekly* granularity, that is, the selected users are expected to visit the target locations in the next week. The studied area (Figure 10a) is split into 1km\*1km grids. Among the 45 weeks of user mobility data, we use the last five weeks as the test time period, and first 40 weeks for mobility profiling.
- *CMCC* dataset contains 1315 users’ GPS trajectories in Hangzhou, China, for one month from one mobile operator. The time period is set to a *daily* granularity. The studied area (Figure 10b) is split into 1km\*1km grids. We use the first 18 weekdays for mobility profiling and the remaining four weekdays for testing.
- *D4D* dataset (Blondel et al. 2012) includes 5378 users’ two-week mobile phone call logs with cell tower locations in Abidjan, Côte d’Ivoire. The time period is set to a *daily* granularity. The studied area (Figure 10c) is split into cell-tower-based regions (Xiong et al. 2016; Wang et al. 2017). We use the first nine weekdays for mobility profiling and the last one weekday for testing.

Table 1 summarizes the experimental parameters. Note that the default differential privacy budget  $\epsilon$  is set to  $\ln(4)$  as suggested by the original paper (Andrés et al. 2013).

## Results on FS

**Single Location Coverage.** We first evaluate the scenario where only one location (grid) needs to be covered. Our

Table 1: Experimental parameters.

Notation	Values	Description
$\epsilon$	$\ln(2), \ln(4), \ln(6), \ln(8)$	differential privacy level
$\delta$	0.5, 0.6, 0.7, 0.8	threshold for frequent locations
$N$	1083 (FS), 1315 (CMCC) 5378 (D4D)	total number of users
$\alpha$	$5\% \cdot N$	number of selected users
$\rho$	95%	probability for user selection
$k$	6	number of user groups

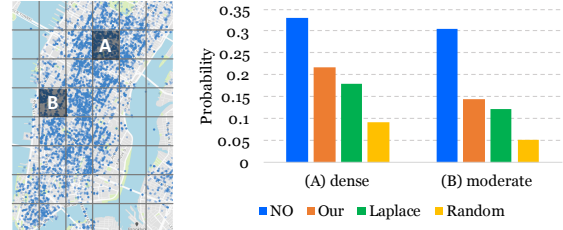


Figure 5: Experiment of single location coverage on two different populated locations on FS ( $\epsilon = \ln(4)$ ,  $\delta = 0.7$ ).

evaluation metric is the probability that a selected user will actually appear at the target location in the next week.

Figure 5 shows the results on two target locations with different population sizes when  $\epsilon = \ln(4)$  and  $\delta = 0.7$ . In both target locations, our proposed method can achieve a larger coverage probability (up to 5% improvement) than the Laplace mechanism. Compared to the no-obfuscation method, the coverage probability of our method drops from 32.9% to 21.7% for the densely populated target location. For the less densely populated one, the drop is bigger (from 30.5% to 14.5%). A possible explanation is that when the target location is densely populated, even if our mechanism mis-selects a user whose frequent location is not the target one, s/he still may go to the target location by chance.

Figure 6a illustrates how the coverage probability changes when we vary the privacy budget  $\epsilon$  for the densely populated target location. As a trade-off between privacy and coverage, when  $\epsilon$  increases (i.e., lower level of privacy), we can get a higher coverage probability. More specifically, the improvement of our method over Laplace is more significant for a lower  $\epsilon$ , i.e., higher privacy protection guarantee.

Figure 6b shows the change of coverage probability when the threshold of frequent locations  $\delta$  varies. The coverage probabilities of all the methods rise with the increase of  $\delta$ . While a higher  $\delta$  benefits coverage probability, the number of users who can upload their (obfuscated) frequent locations (i.e., candidates for selection) is smaller, because only users with at least one location profiling probability larger than  $\delta$  will upload frequent locations. Based on experiment results, setting  $\delta$  to around 0.7-0.8 is appropriate for our method, as the coverage probability is relatively satisfactory while a large portion of users can be involved.

**Multi-Location Coverage.** We evaluate the scenario where multiple target locations exist. We randomly select 2, 4, 6 and 8 locations as the targets. Figure 7 shows the actual



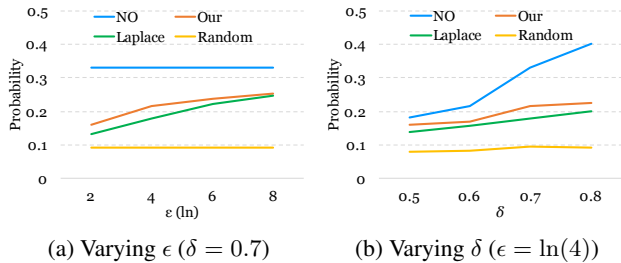


Figure 6: Single location coverage results on FS.

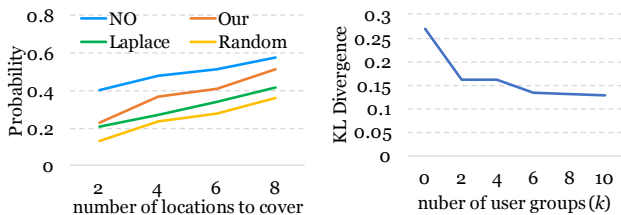


Figure 7: Multi-location coverage results on FS. Figure 8: Results of estimating  $\pi$  on FS.

coverage probability that we can get, i.e., the probabilities of selected users covering any one of the target locations in the coming week. The results show that our proposed method consistently outperforms Laplace under the same level of privacy protection. Moreover, with an increasing number of the target locations, we find that the performance gap between our method and no-obfuscation becomes smaller. This indicates that, when there are more locations to cover, using our mechanism is more profitable, as the performance loss incurred by the geographic differential privacy protection becomes smaller.

**Estimation of  $\pi$ .** We evaluate whether our proposed Bayes rule based method can estimate  $\pi$  correctly. We use KL divergence (Kullback and Leibler 1951) to quantify the similarity between the estimated  $\pi$  and the actual  $\pi$ . The smaller KL divergence is, the more similar they are. Figure 8 shows the change of KL divergence with  $k$  (the number of user groups), and  $\pi$  is initialized to a uniform distribution. In Figure 8,  $k = 0$  represents the KL divergence between the uniform and the actual distribution. When  $k$  is small, we have fewer iterations to update  $\pi$ , leading to a larger KL divergence. In our experiment,  $k = 6$  is a good setting, as KL divergence achieves a relatively low value, while the algorithm can complete execution within a reasonable time.

**Runtime Efficiency.** We use Gurobi 7.5 (Gurobi 2014) as the linear programming solver engine to run Algorithm 1 for getting the optimal policy  $\hat{P}$ . It takes about 450 seconds on a commodity laptop with i5-5200U (2.2 GHz), 8G memory. We split all the users to six groups, meaning that Algorithm 1 is executed six times, which sums up to about 45 minutes. As the optimal privacy policy generation can be an offline process, such runtime efficiency is totally acceptable for real applications. Note that this running time is not affected by the number of users, so our method can serve mobile applications with a large number of users.

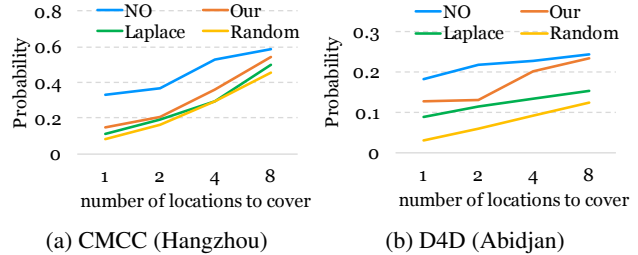


Figure 9: Experiment results on CMCC and D4D with different number of locations to cover ( $\epsilon = \ln(4)$ ,  $\delta = 0.8$ ).

## Results on CMCC and D4D

To test the robustness of our proposed method, we also conduct experiments on CMCC and D4D datasets. The results are shown in Figure 9a and 9b, where we randomly select 1, 2, 4, and 8 locations to cover. The results verify that our proposed method can always outperform the Laplace mechanism in attaining a higher coverage probability of the selected users. Moreover, the results show that when the number of target locations to cover increases to 8, our privacy mechanism almost achieves the same coverage probability as no-obfuscation, especially for the D4D dataset. This further emphasizes the practicability of our mechanism, as user privacy is gained with a nearly negligible quality loss. Note that the achieved coverage probability on D4D is smaller than FS or CMCC in general, because the phone call locations on D4D are intrinsically more difficult to predict. Please refer to the appendix for detailed mobility prediction results.

## Related Work

Selecting a set of users who can cover a set of locations in the near future is a very important problem for real applications like spatial crowdsourcing (Chen and Shahabi 2016; Zhang et al. 2014) and location-based advertising (Dhar and Varshney 2011). In most of previous research works, users' moving histories are known and hence their mobility patterns can be effectively modeled for predicting their future locations (Xiong et al. 2016; Guo et al. 2017; Yang et al. 2015).

As user privacy is becoming more and more important nowadays, some pioneering works have started to model users' mobility or activity patterns based on privacy-preserving data. Geo-indistinguishability mechanisms are proposed for location-based query systems where users can submit their differentially obfuscated locations (Andrés et al. 2013; Bordenabe, Chatzikokolakis, and Palamidessi 2014). PrivCheck is designed to enable personalized location-based advertising or recommendation with obfuscated user check-ins, so that users' sensitive information (e.g., age and gender) cannot be inferred by adversaries (Yang et al. 2016). In spatial crowdsourcing data acquisition, recent works also incorporate privacy mechanisms to protect participants' precise locations (Wang et al. 2016; Wang et al. 2017; To, Ghinita, and Shahabi 2014; Vergara-Laurens, Mendez, and Labrador 2014; Pournajaf et al. 2014). While these studies have various applica-

tions, they usually focus on obfuscating users' *current* locations. As far as we know, little previous work has studied the privacy-preserving future crowd coverage maximization problem based on users' obfuscated *historical* mobility profiles which we specifically focus on in this paper.

## Conclusion

In this paper, we study the crowd coverage maximization problem under the privacy protection on user locations. The key idea is to select users who will probably visit certain locations in near future with their differentially obfuscated locations. To maximize the quality (coverage probability) of selected users under such a privacy protection scheme, an optimization problem is formulated to obtain the optimal privacy policy. We mathematically analyze the problem, and then propose a practical algorithm to obtain the optimal privacy policy. Experiments on various real user mobility datasets have verified the effectiveness of our privacy mechanism. As future work, we plan to study the problem when a user can upload multiple obfuscated frequent locations.

## Acknowledgment

This research is partially supported by NSFC Grant no. 71601106, State Language Commission of China Key Program Grant no. ZDI135-18, Hong Kong ITF Grant no. ITS/391/15FX, the European Research Council (ERC) under the European Unions Horizon 2020 research and innovation program (grant agreement 683253/GraphInt).

## References

- [Andrés et al. 2013] Andrés, M. E.; Bordenabe, N. E.; Chatzikokolakis, K.; and Palamidessi, C. 2013. Geo-indistinguishability: Differential privacy for location-based systems. In *Proc. CCS*, 901–914.
- [Blondel et al. 2012] Blondel, V. D.; Esch, M.; Chan, C.; Clérot, F.; Deville, P.; Huens, E.; Morlot, F.; Smoreda, Z.; and Ziemlicki, C. 2012. Data for development: the d4d challenge on mobile phone data. *arXiv preprint arXiv:1210.0137*.
- [Bordenabe, Chatzikokolakis, and Palamidessi 2014] Bordenabe, N. E.; Chatzikokolakis, K.; and Palamidessi, C. 2014. Optimal geo-indistinguishable mechanisms for location privacy. In *Proc. CCS*, 251–262.
- [Boyd and Vandenberghe 2004] Boyd, S., and Vandenberghe, L. 2004. *Convex optimization*. Cambridge university press.
- [Chen and Shahabi 2016] Chen, L., and Shahabi, C. 2016. Spatial crowdsourcing: Challenges and opportunities. *IEEE Data Eng. Bull.* 39(4):14–25.
- [Cho, Myers, and Leskovec 2011] Cho, E.; Myers, S. A.; and Leskovec, J. 2011. Friendship and mobility: user movement in location-based social networks. In *Proc. KDD*, 1082–1090.
- [Dhar and Varshney 2011] Dhar, S., and Varshney, U. 2011. Challenges and business models for mobile location-based services and advertising. *Communications of the ACM* 54(5):121–128.
- [Dwork 2008] Dwork, C. 2008. Differential privacy: A survey of results. In *International Conference on Theory and Applications of Models of Computation*, 1–19.
- [Fawcett 2006] Fawcett, T. 2006. An introduction to roc analysis. *Pattern recognition letters* 27(8):861–874.
- [Guo et al. 2017] Guo, B.; Liu, Y.; Wu, W.; Yu, Z.; and Han, Q. 2017. Activecrowd: A framework for optimized multi-task allocation in mobile crowdsensing systems. *IEEE Transactions on Human-Machine Systems* 47(3):392–403.
- [Gurobi 2014] Gurobi. 2014. Inc.,gurobi optimizer reference manual, 2014. URL: <http://www.gurobi.com>.
- [Kullback and Leibler 1951] Kullback, S., and Leibler, R. A. 1951. On information and sufficiency. *The annals of mathematical statistics* 22(1):79–86.
- [Pournajaf et al. 2014] Pournajaf, L.; Xiong, L.; Sunderam, V.; and Goryczka, S. 2014. Spatial task assignment for crowd sensing with cloaked locations. In *Proc. MDM*, volume 1, 73–82.
- [Rossi et al. 2015] Rossi, L.; Williams, M. J.; Stich, C.; and Musolesi, M. 2015. Privacy and the city: User identification and location semantics in location-based social networks. In *Proc. ICWSM*, 387–396.
- [To, Ghinita, and Shahabi 2014] To, H.; Ghinita, G.; and Shahabi, C. 2014. A framework for protecting worker location privacy in spatial crowdsourcing. *Proc. of the VLDB Endowment* 7(10):919–930.
- [Vergara-Laurens, Mendez, and Labrador 2014] Vergara-Laurens, I. J.; Mendez, D.; and Labrador, M. A. 2014. Privacy, quality of information, and energy consumption in participatory sensing systems. In *Proc. PerCom*, 199–207.
- [Wang et al. 2016] Wang, L.; Zhang, D.; Yang, D.; Lim, B. Y.; and Ma, X. 2016. Differential location privacy for sparse mobile crowdsensing. In *Proc. ICDM*, 1257–1262.
- [Wang et al. 2017] Wang, L.; Yang, D.; Han, X.; Wang, T.; Zhang, D.; and Ma, X. 2017. Location privacy-preserving task allocation for mobile crowdsensing with differential geo-obfuscation. In *Proc. WWW*, 627–636.
- [Xiong et al. 2016] Xiong, H.; Zhang, D.; Chen, G.; Wang, L.; Gauthier, V.; and Barnes, L. E. 2016. icrowd: Near-optimal task allocation for piggyback crowdsensing. *IEEE Transactions on Mobile Computing* 15(8):2010–2022.
- [Yang et al. 2015] Yang, D.; Zhang, D.; Zheng, V. W.; and Yu, Z. 2015. Modeling user activity preference by leveraging user spatial temporal characteristics in lbsns. *IEEE Transactions on Systems, Man, and Cybernetics: Systems* 45(1):129–142.
- [Yang et al. 2016] Yang, D.; Zhang, D.; Qu, B.; and Cudré-Mauroux, P. 2016. Privcheck: privacy-preserving check-in data publishing for personalized location based services. In *Proc. UbiComp*, 545–556.
- [Zhang et al. 2014] Zhang, D.; Wang, L.; Xiong, H.; and Guo, B. 2014. 4w1h in mobile crowd sensing. *IEEE Communications Magazine* 52(8):42–48.
- [Zheng et al. 2014] Zheng, Y.; Capra, L.; Wolfson, O.; and Yang, H. 2014. Urban computing: concepts, methodologies,



## Appendix

### Detailed Proof of Lemma 1

Suppose we have two different  $\hat{l}^*$ , i.e.,  $l_1^*, l_2^* \in \mathcal{L}$ , and get two different optimal objective values

$$\frac{\pi(l_t)\hat{P}_1(\hat{l}_1^*|l_t)}{\sum_{l \in \mathcal{L}} \pi(l)\hat{P}_1(\hat{l}_1^*|l)} < \frac{\pi(l_t)\hat{P}_2(\hat{l}_2^*|l_t)}{\sum_{l \in \mathcal{L}} \pi(l)\hat{P}_2(\hat{l}_2^*|l)}$$

We now construct a new solution of  $\hat{P}'_1$  when  $\hat{l}^* = l_1^*$  as follows:

$$\hat{P}'_1(l_1^*|l) = \hat{P}_2(l_2^*|l), \forall l \in \mathcal{L}$$

$$\hat{P}'_1(l_2^*|l) = \hat{P}_2(l_1^*|l), \forall l \in \mathcal{L}$$

$$\hat{P}'_1(l^*|l) = \hat{P}_2(l^*|l), l^* \neq l_1^*, l^* \neq l_2^*, \forall l \in \mathcal{L}$$

We can verify that all the constraints of the optimization still stand, and then  $\hat{P}'_1$  is a feasible solution when  $\hat{l}^* = l_1^*$ , and then

$$\frac{\pi(l_t)\hat{P}'_1(\hat{l}_1^*|l_t)}{\sum_{l \in \mathcal{L}} \pi(l)\hat{P}'_1(\hat{l}_1^*|l)} = \frac{\pi(l_t)\hat{P}_2(\hat{l}_2^*|l_t)}{\sum_{l \in \mathcal{L}} \pi(l)\hat{P}_2(\hat{l}_2^*|l)} > \frac{\pi(l_t)\hat{P}_1(\hat{l}_1^*|l_t)}{\sum_{l \in \mathcal{L}} \pi(l)\hat{P}_1(\hat{l}_1^*|l)}$$

This violates that  $\hat{P}_1$  is the optimal solution when  $\hat{l}^* = \hat{l}_1^*$ .  $\square$

### Detailed Proof of Theorem 1

With the following way to construct  $\hat{P}$ ,

$$\hat{P}(\hat{l}^*|l) = \theta e^{-\epsilon d(l, l_t)}, \quad \forall l \in \mathcal{L} \quad (29)$$

$$\hat{P}(l^*|l) = \frac{1 - \theta e^{-\epsilon d(l, l_t)}}{|\mathcal{L}| - 1}, \quad \forall l^*, l \in \mathcal{L} \text{ and } l^* \neq \hat{l}^* \quad (30)$$

Then, for any  $l \in \mathcal{L}$ ,

$$\sum_{l^* \in \mathcal{L}} \hat{P}(l^*|l) = \theta e^{-\epsilon d(l, l_t)} + (|\mathcal{L}| - 1) \frac{1 - \theta e^{-\epsilon d(l, l_t)}}{|\mathcal{L}| - 1} = 1$$

So the probability sum constraint stands. We then prove that differential privacy constraint also stands. Note that *Lemma 3* has proved that the differential privacy constraint holds if  $l^* = \hat{l}^*$ . Therefore, we only need to show that the differential privacy constraint also stands for  $l^* \neq \hat{l}^*$ . Next, we show how to select  $\theta$  to ensure that this is true for any  $l^* \neq \hat{l}^*$ ,

$$\frac{\hat{P}(l^*|l_1)}{\hat{P}(l^*|l_2)} = \frac{1 - \theta e^{-\epsilon d(l_1, l_t)}}{1 - \theta e^{-\epsilon d(l_2, l_t)}} \leq e^{\epsilon d(l_1, l_2)} \quad (31)$$

$$\Rightarrow \theta \leq \frac{e^{\epsilon d(l_1, l_2)} - 1}{e^{-\epsilon(d(l_2, l_t) - d(l_1, l_2))} - e^{-\epsilon d(l_1, l_t)}} \quad (32)$$

It is worth noting that both the numerator and denominator in the right side of Eq. 32 are larger than zero when  $\epsilon > 0$ . Hence, we can set  $\theta$  to any positive value smaller than or equal to

$$\min_{l_1, l_2 \in \mathcal{L}} \frac{e^{\epsilon d(l_1, l_2)} - 1}{e^{-\epsilon(d(l_2, l_t) - d(l_1, l_2))} - e^{-\epsilon d(l_1, l_t)}} \quad (33)$$

and then for any  $l^* \neq \hat{l}^*$ , geographic differential privacy still holds. Then, based on *Lemma 2*, we can know that the  $\hat{P}$  satisfying Eq. 29 and 30 can lead to the upper bound of the objective value.  $\square$

### Detailed Proof of Lemma 5

According to the geographic differential privacy constraints, we have

$$e^{\epsilon d(l, l_t)} P(\hat{l}^*|l) \geq P(\hat{l}^*|l_t), \quad l \notin \mathbb{L} \text{ and } l_t \in \mathbb{L} \quad (34)$$

$$\Rightarrow \pi(l_t) e^{\epsilon d(l, l_t)} P(\hat{l}^*|l) \geq \pi(l_t) P(\hat{l}^*|l_t), \quad l \notin \mathbb{L} \text{ and } l_t \in \mathbb{L} \quad (35)$$

$$\Rightarrow \sum_{l_t \in \mathbb{L}} \pi(l_t) e^{\epsilon d(l, l_t)} P(\hat{l}^*|l) \geq \sum_{l_t \in \mathbb{L}} \pi(l_t) P(\hat{l}^*|l_t), \quad l \notin \mathbb{L} \quad (36)$$

$$\Rightarrow P(\hat{l}^*|l) \geq \frac{\sum_{l_t \in \mathbb{L}} \pi(l_t) P(\hat{l}^*|l_t)}{\sum_{l_t \in \mathbb{L}} \pi(l_t) e^{\epsilon d(l, l_t)}}, \quad l \notin \mathbb{L} \quad (37)$$

Then,

$$\frac{\sum_{l_t \in \mathbb{L}} \pi(l_t) P(l^*|l_t)}{\sum_{l \in \mathcal{L}} \pi(l) P(l^*|l)} \quad (38)$$

$$= \frac{\sum_{l_t \in \mathbb{L}} \pi(l_t) P(l^*|l_t)}{\sum_{l_t \in \mathbb{L}} \pi(l_t) P(l^*|l_t) + \sum_{l \notin \mathbb{L}} \pi(l) P(l^*|l)} \quad (39)$$

For the ease of presentation, we denote  $C = \sum_{l_t \in \mathbb{L}} \pi(l_t) P(l^*|l_t)$ ,

$$\frac{\sum_{l_t \in \mathbb{L}} \pi(l_t) P(l^*|l_t)}{\sum_{l \in \mathcal{L}} \pi(l) P(l^*|l)} \quad (40)$$

$$= \frac{C}{C + \sum_{l \notin \mathbb{L}} \pi(l) P(l^*|l)} \quad (41)$$

$$\leq \frac{C}{C + \sum_{l \notin \mathbb{L}} \pi(l) \frac{C}{\sum_{l_t \in \mathbb{L}} \pi(l_t) e^{\epsilon d(l, l_t)}}} \quad (42)$$

$$= (1 + \sum_{l \notin \mathbb{L}} \sum_{l_t \in \mathbb{L}} \frac{\pi(l)}{\pi(l_t)} e^{-\epsilon d(l, l_t)})^{-1} \quad (43)$$

$\square$

### Mobility Profiling

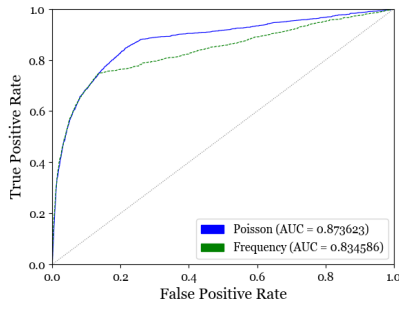
We consider two popular mobility profiling methods used in literature, and choose the better one in our experiments.

(1) **Frequency** (Guo et al. 2017). This method counts daily (or weekly) frequency that a user visits a location in her/his historical mobility records. For example, suppose we have a user's 7-day mobility history and s/he visits a location  $l_i$  in 5 days, then the daily visiting probability is 5/7.

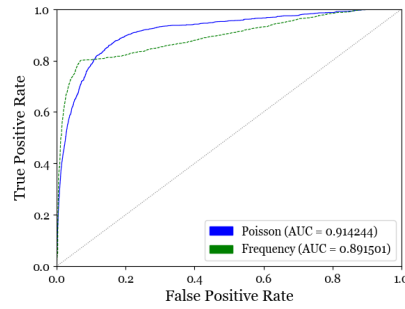
(2) **Poisson** (Xiong et al. 2016). Given a user  $u$ 's average daily (or weekly) visiting times to location  $l_i$  in the past, denoted as  $\lambda_{u,i}$ , then the Poisson process estimates that  $u$  visits  $l_i$  at least once in one day (week) is:

$$p_{u,i} = 1 - e^{-\lambda_{u,i}} \quad (44)$$

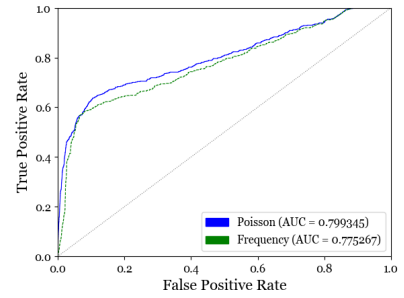
Figure 10 plots the receiver operating characteristics (ROC) curves (Fawcett 2006) and shows the area under the ROC curve (AUC) values for both profiling methods on the FS, CMCC, and D4D datasets, respectively. The larger AUC value implies better performance in predicting a user's future mobility patterns. From the results, we see that



(a) New York (FS)



(b) Hangzhou (CMCC)



(c) Abidjan (D4D)

Figure 10: ROC curves of mobility profiling.

Poisson-based mobility profiling method beats Frequency-based method, and thus we use the Poisson-based method in all the experiments.

In addition, we observe that the mobility prediction on the D4D dataset is more difficult than on the other two datasets, as it gets a lower AUC value. As expected, our experiments in the paper (Figure 9) show that the selected users on the D4D dataset achieve a lower coverage probability than the other two datasets with the same user selection mechanism.