

Preserving Differential Privacy in Publication of Trajectory Data

Kulkarni Sourabh Shrinivasrao, 21114053

Abstract— For development and improvement of many modern day applications, collecting and publishing the trajectory data is very important. But publishing it right away possesses different challenges such as compromise in user’s privacy. Also, traditional trajectory merging algorithms are too slow. We try to solve both of these problems using our database publication algorithm and trajectory merging algorithm. Theoretical analysis shows that our privacy loss is much less.

Index Terms—

1. Differential Privacy
2. Trajectory Data
3. Database Publication
4. Laplace Noise

INTRODUCTION

With the revolution in the devices able to accommodate the GPS services , generating and collecting the trajectory data is becoming increasingly easier for users of those devices. Of course, This data is very much useful for inventors, researchers and solving real life problems. But this also comes with a cost. Publishing the trajectory data as it is exposes sensitive details of users such as their routines, habits, personal information etc.

Differential privacy is an elegant and modern approach to protect each individuals privacy while using their data in the database. Third party with access to the database cannot reach a conclusion whether the data belongs to an user or not. Using Laplace mechanism is the most commonly used method to achieve differential privacy adding noise in the **true count** (trajectory count on a road). While adding this noise, one must look after the boundedness of the noise else making release meaningless. Many work is available on this topic such as first applied differential privacy to sequential data, n-gram model etc.

TERMS AND DEFINITIONS

1.Differential Privacy: If the output of a query does not depend on the presence of a single record, then we can say this database possesses differential privacy.

2.Sensitivity: for a function $f: D \rightarrow \mathbb{R}^d$, the sensitivity $df = \max \|f(D1)-f(D2)\|$ where $D1, D2$ differ in atmost 1 record.

3.Laplace mechanism: $\text{pdf}(x)=(1/2b) e^{-|x-m|/b}$ where m is mean and $b=df/E$ where E is privacy budget.

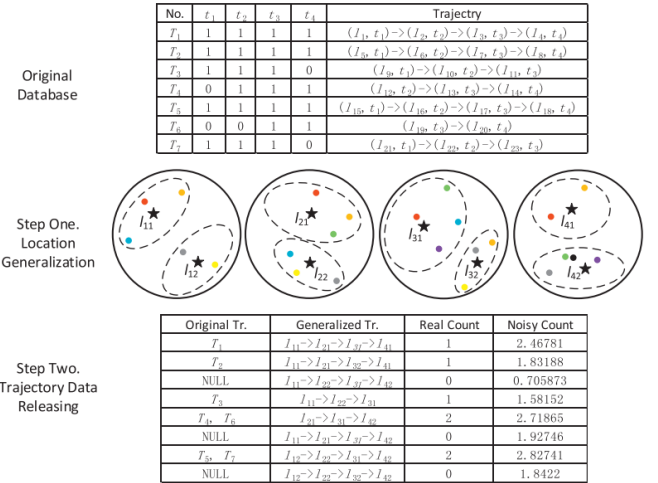
4.Trajectory: a list of location-time pairs: $T = (l_1, t_1) \rightarrow (l_2, t_2) \rightarrow \dots \rightarrow (l_T, t_T)$

5.True count: If a trajectory exists in the database, True count=1. Else 0.

CLAIM:

Unbounded Laplace noise does not yield preservation of utility and privacy. Thus, we will use bounded noises.

PROCEDURE:



Pictorial representation of general idea

The general idea is, we want to add some “fake” (generated) trajectories in the database to the original trajectories. While doing this we will not include all of the original trajectories in order to match the original size of the database. After that, we will add noise to the true counts in order to ensure that these trajectories become indistinguishable. We can achieve this by following these steps:

PROOF THAT DIFFERENTIAL PRIVACY IS PRESERVED

For any sequence r of outcomes $r_i \in \text{Range}(A_i)$, $i = 1, 2$, we write A^r_1 and A^r_2 for algorithm A_1 and A_2 supplied with r_1 and r_2 . The probability of output r from the sequence $A^r_1(D)$ and $A^r_2(D)$ on database D is

$$\Pr[A(D) = r] = \Pr[A^r_1(D) = r_1] \Pr[A^r_2(D) = r_2].$$

Applying definition of differential privacy, we get:

$$\begin{aligned} \Pr[A^r_1(D) = r_1] \Pr[A^r_2(D) = r_2] &\leq (e^{\epsilon_1} \Pr[A^r_1(D') = r_1]) * (e^{[T] \cdot \epsilon_2} \Pr[A^r_2(D') = r_2]) \\ &= e^{\epsilon_1 + [T] \cdot \epsilon_2} \Pr[A^r_1(D') = r_1] \Pr[A^r_2(D') = r_2] \\ &= e^{\epsilon_1 + [T] \cdot \epsilon_2} \Pr[A^r_1(D') = r_1 \wedge A^r_2(D') = r_2] \\ &= e^{\epsilon_1 + [T] \cdot \epsilon_2} \Pr[A(D') = r] \end{aligned}$$

hence, we can achieve differential privacy.

- [2] Source for many images is from the implementation of this paper by me: [GitHub link for code](#)
 [3] [What is differential privacy](#)

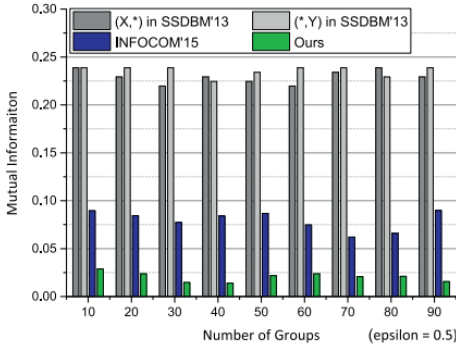
RESULTS

Lets Define Mutual information metric as following:

$$MI(x, y) = \sum_t \sum_{x(t)} \sum_{y(t)} \Pr(x(t), y(t)) \log \frac{\Pr(x(t), y(t))}{\Pr(x(t)) \Pr(y(t))}$$

Which calculates the dependency between two variables X and Y . If X and Y are not dependent, we can know nothing about the second if we know about first and vice versa.

Using this metric, we compare our algorithm of data publication with other as shown in below bar diagram. Clearly the less MI, the better algorithm.



CONCLUSION

In this paper, we saw how it is inappropriate to release trajectory data without any treatment. With given algorithms, one can achieve differential privacy with first applying clustering algorithm and then by adding bounded Laplace noises.

REFERENCES

- [1] "Achieving differential privacy of trajectory data publishing in participatory sensing" by Meng Li, Liehuang Zhu, Zijian Zhang, Rixin Xu. [Link to the paper](#).