# Exploring the Effectiveness of Geomasking Techniques for Protecting the Geoprivacy of Twitter Users

Song Gao[1] and Qunying Huang[1]

Department of Geography, University of Wisconsin, Madison, WI 53706, USA
Email: {song.gao;qhuang46}@wisc.edu

**Abstract.** With the ubiquitous availability of location-based services, large-scale individual-level location data have been widely collected through location-awareness devices. The geoprivacy concerns arise on the issues of user identity de-anonymization and location exposures. In this work, we investigate the effectiveness of geomasking techniques for protecting the geoprivacy of active Twitter users who frequently share geotagged tweets in their home locations. The two-dimensional Gaussian masking technique is found to be more effective than the random perturbation masks.

**Keywords:** Geoprivacy · Geomasking · Social Media.

## 1 Introduction

With the ubiquitous availability of location-based services, large-scale individual-level location data have been collected through the use of mobile phones, GPS devices, geotagged social media, and so on. While such geospatial big data provide new opportunities to study human mobility patterns, complex human-environment interactions, socioeconomic characteristics, urban changes, and business intelligence [11, 14, 6, 12, 8], however, there are also increasing concerns about the social, ethical, legal and behavioral implications of the geoprivacy caused by user identity de-anonymization and location exposures [2, 16, 9].

Generally speaking, geoprivacy refers individual rights to prevent disclosure of a person's locations including but not limit to the location of home, workplace, daily activities, or travel trips [10]. However, the majority of people don't know how the underlying location-related technologies work and what can be inferred from an individual's location record [9] with which people can use various location-based services. As a result, researchers have developed a number of statistical approaches and technical solutions aiming to protect individuals being identified through their location records. A common practice for preserving data confidentiality is aggregation such that detailed individual records are merged into anonymized large-group characteristics. For example, aggregating individual home location into city blocks or census tracts. However, the aggregation reduces the spatial resolution of the analysis that may affect the its effectiveness

[10]. Another family of approaches are called geomasking techniques in which the original location may be hidden or modified for geoprivacy protection but the spatial point patterns are not significantly affected.

There is a rich history of literature on leveraging geographical masking to preserve confidentiality of health records and trajectory data. With child leukemia lymphoma data from North Humberside, England, Armstrong et al. (1999) described and evaluated several types of geographical masks to protect personal privacy as well as allow the conduct of valid spatial analyses [1]. Kwan (2004) examined the effects of random perturbation masks on the results of a spatial analysis using data on lung-cancer deaths [10]. Three different random perturbation masks were implemented with each at three different levels of introduced error. Hampton et al. (2010) extended existing methods of random perturbation by developing an adaptive geomasking technique, known as the donut method [5]. This method guarantees that each geocoded address is not randomly assigned on or too near its original location. Compared with random perturbation method, the performance of k-anonymity using the proposed donut method was at least 42.7% higher in geoprivacy measures and was less than 4.8% in cluster detection measures. Seild et al. (2016) examined the grid masking and random perturbation techniques for anonymizing the GPS trajectory data and test the preservation of both privacy and spatial patterns. They found that as the distance thresholds for grid masking and random perturbation increase, the correlation between density patterns decreases [15].

However, how to use geographical masking methods to prevent the disclosure of important locations of social media users is still not well addressed. To this end, we aim to investigate the effectiveness of geomasking techniques for protecting the geoprivacy of active Twitter users who frequently share geotagged tweets in their home or work location.

## 2    Methodology

In this work, we have applied the two popular geomasking techniques [1]: *Random Perturbation* and *Gaussian Perturbation*, aiming at the preservation of Twitter users' geoprivacy.

***Random Perturbation:*** is a geomasking approach in which each point is displaced in space by a randomly determined distance and direction. A distance threshold is typically added to set the allowed maximum displacement distance in the case of uniform geomasking. As shown in Fig. 2, the original posted locations of the geotagged tweets (Fig. 1) of a Twitter user are randomly displaced within a 1km distance radius.

***Gaussian Perturbation:*** uses a two-dimensional isotropic (i.e. circularly symmetric) Gaussian kernel to control the random displacement process of a point set such that the distribution of those displaced points follow a 2D Gaussian ("bell-shaped") form:

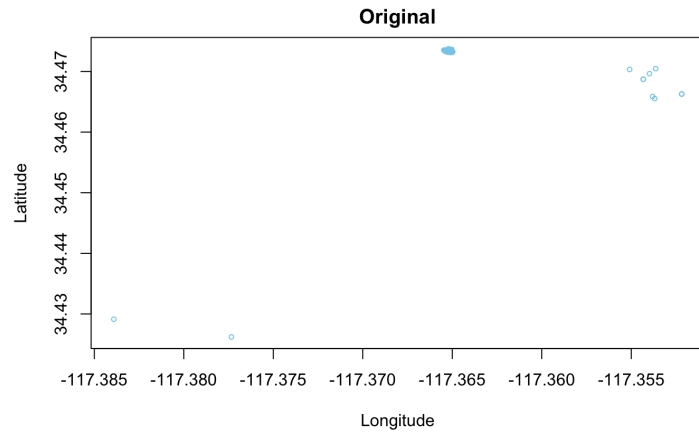$$G(x,y) = (1/2\pi\delta^2)e^{-(x^2+y^2)/2\delta^2} \qquad (1)$$

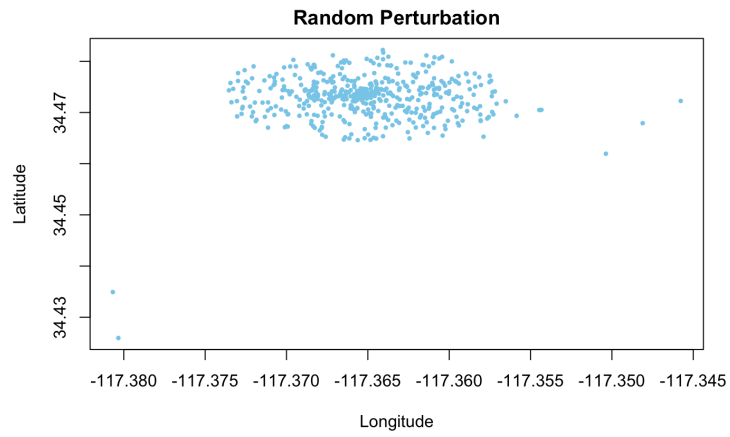**Fig. 1.** The original spatial distribution of geotagged tweets from a user.
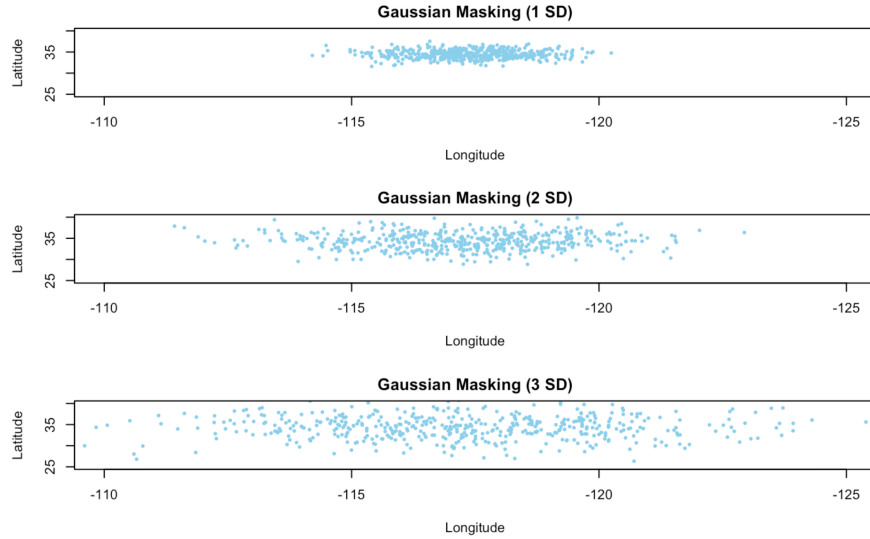


**Fig. 2.** The random perturbation of a user's geotagged tweets.

Where (x,y) is the 2D coordinates of each location after displacement, and $\delta$ specifies the standard deviation of the positional error. As shown in Fig. 3, with the increment of the standard deviation, the displacement of those points spreads more widely.



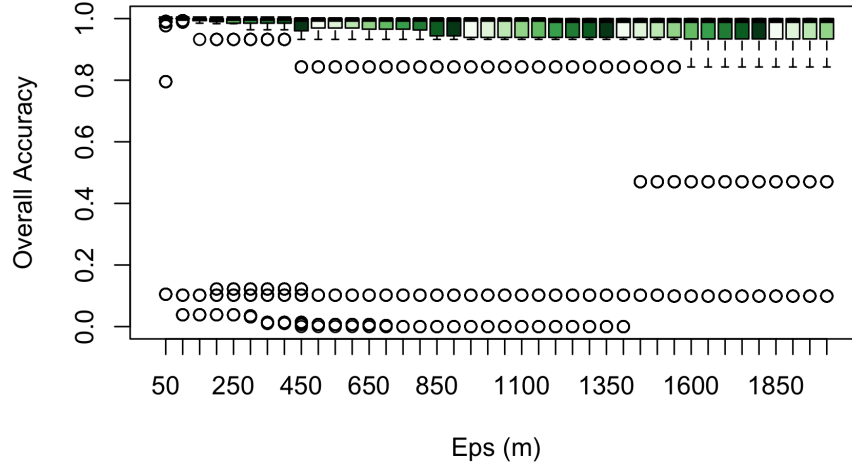**Fig. 3.** The Gaussian perturbation of a user's geotagged tweets.

After the perturbation processing of the original locations of geotagged tweets, we need to further analyze whether users' home location (one of the most sensitive places for an individual's geoprivacy concern) can still be identified through state-of-the-art home location detection algorithms. Specifically, we explored different parameter calibrations for the density-based spatial clustering with noise (DBSCAN) [3, 7] that has been widely used in spatial clustering and the identification of significant human activity places. The DBSCAN algorithm requires two parameters: the searching radius of a cluster (Eps) and the minimum number of points (MinPts) within a cluster. The different combinations of Eps and MinPts values may get different spatial clustering results [6, 13]. In the case of detecting Twitter users' home location, the parameter calibration may generate different candidate clusters or distance shifts from the actual location. Therefore, we have explored different scenarios with varying parameter values for the perturbation and the spatial clustering steps. Some preliminary analysis results are discussed in the following section.

## 3   Results

In this work, we selected 20 most active Twitter users who have frequently posted geotagged tweets in a U.S. metropolitan area from mobile phone devices such that their location information is most accurate for human mobility study [4]. We also manually identified their home locations by overlaying their nighttime (8pm-7am) geotagged tweets onto the high-resolution (about 2m) Digital Globe aerial images as the ground-truth. One limitation of such an approach lies on the uncertainty of users' actual home location without their interview confirmation.

*The impact of Eps and MinPts:* Before applying the geomasks, we first tested how the choice of *MinPts* and *Eps* in DBSCAN would affect the effectiveness of identifying the home clusters. We chose the *MinPts* ranging from 4 points to the square root of the total number of nighttime tweets, and the search radius *Eps* in a range of 50m to 2000m with a step of 50m. As shown in Fig. 4, we grouped the home cluster detection results based on the *Eps*, each sub-boxplot represents the overall accuracy with varying *MinPts* in DBSCAN. The overall accuracy is the ratio between number of correctly identified clusters (including both true positive and true negative cases) and total of number of clusters. Not surprisingly, the mean of overall accuracy is over 0.98 and keeps very high values (over 0.9) regardless of the parameter choices. It also indicates the potential risk of location exposures of those active users as their home locations are easily identified without parameter calibration.
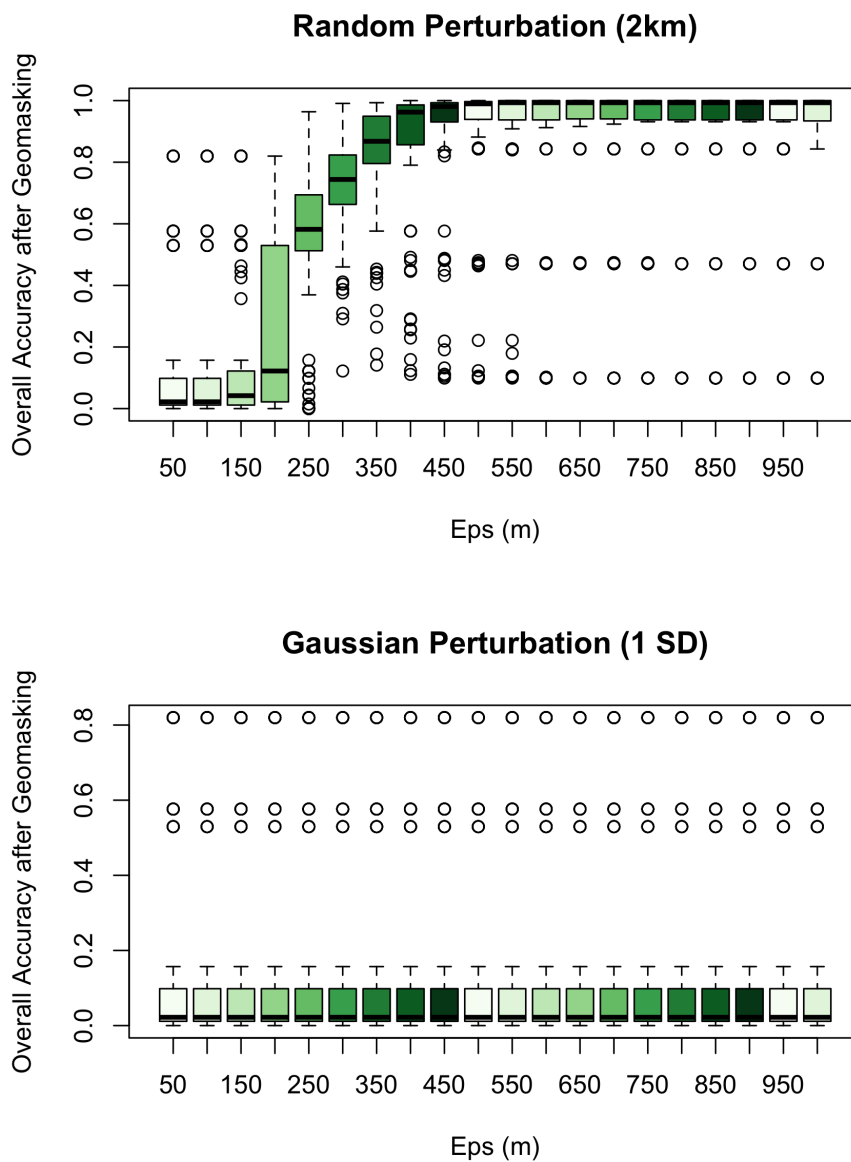
*Comparing the effectiveness of different geomasking techniques:* First, we explored the impact of the random perturbation geomasking with different thresholds. Existing studies have found that the choice of Eps=200m to 300m could generate good spatial clustering results for urban areas of interest or human activity zones [6, 8]. Therefore, we are interested in whether the geomasking process could protect users' home location privacy in such a range. However, our experiments show that small-distance (such as within 1km) random perturbations don't help the protection of users' geoprivacy because their home location clusters can still be correctly identified with over 90% overall accuracy. Note that our study area in this work is different from [6, 8]. Even when the displacement threshold reaches 2km, the mean of overall accuracy is still over 70%. But the 2km random perturbation mask is effective for protecting users' home location within the searching radius of 200m; and only less than 20% of overall accuracy can be achieved for identifying the users' home location (in Fig. 5). And the mean of distance shifts between the true home location and the medoids of home clustering results is about 480m. As for the Gaussian perturbation, we found that it is very effective for protecting users' home location. The mean of overall accuracy for identifying users' home clusters with 2D Gaussian kernels is less than 12% in average regardless of the DBSCAN parameter choices. In addition to accuracy, we also explored other clustering assessment measures including sensitivity, balanced accuracy, F1-score, and distance shifts of centroids, which will be further described in a full paper.

**Fig. 4.** The boxplot of overall accuracy changes of home cluster detection with different DBSCAN parameters.

## 4   Conclusion and Discussion

In this work, we have explored the effectiveness of two geomasking techniques for protecting the geoprivacy of active Twitter users who frequently share geotagged tweets in their home locations. Based on our preliminary experiments, the two-dimensional Gaussian masking is found to be more effective on hiding or shifting user's home location than the random perturbation masks. One discussion issue would be what is the impact of these geoprivacy enhancements on the user experience comparing with simply removing the benefit to the user of posting geotagged tweets. However, this conclusion is drawn through a limited number of active users with ground-truth location; and the current home-detection method only relies on DBSCAN for nighttime geotagged tweets. Other approaches also exist such as the detection of home-job locations using recurring trips. We would like to extend our workflow to a larger group of users and in other cities to test whether our conclusion is extensible or not in future work.

**Random Perturbation (2km)**



**Gaussian Perturbation (1 SD)**



**Fig. 5.** The boxplot of overall accuracy changes of home cluster detection with different DBSCAN parameters with geomasking.

# References

1. Armstrong, M.P., Rushton, G., Zimmerman, D.L., et al.: Geographically masking health data to preserve confidentiality. Statistics in medicine **18**(5), 497–525 (1999)
2. Beresford, A.R., Stajano, F.: Location privacy in pervasive computing. IEEE Pervasive computing **2**(1), 46–55 (2003)
3. Ester, M., Kriegel, H.P., Sander, J., Xu, X., et al.: A density-based algorithm for discovering clusters in large spatial databases with noise. In: Kdd. vol. 96, pp. 226–231 (1996)
4. Gao, S., Yang, J.A., Yan, B., Hu, Y., Janowicz, K., McKenzie, G.: Detecting origin-destination mobility flows from geotagged tweets in greater los angeles area. In: Eighth International Conference on Geographic Information Science (GIScience'14). Citeseer (2014)
5. Hampton, K.H., Fitch, M.K., Allshouse, W.B., Doherty, I.A., Gesink, D.C., Leone, P.A., Serre, M.L., Miller, W.C.: Mapping health data: improved privacy protection with donut method geomasking. American journal of epidemiology **172**(9), 1062–1069 (2010)
6. Hu, Y., Gao, S., Janowicz, K., Yu, B., Li, W., Prasad, S.: Extracting and understanding urban areas of interest using geotagged photos. Computers, Environment and Urban Systems **54**, 240–254 (2015)
7. Huang, Q., Wong, D.W.: Modeling and visualizing regular human mobility patterns with uncertainty: An example using twitter data. Annals of the Association of American Geographers **105**(6), 1179–1197 (2015)
8. Huang, Q., Wong, D.W.: Activity patterns, socioeconomic status and urban spatial structure: what can social media data tell us? International Journal of Geographical Information Science **30**(9), 1873–1898 (2016)
9. Keßler, C., McKenzie, G.: A geoprivacy manifesto. Transactions in GIS **22**(1), 3–19 (2018)
10. Kwan, M.P., Casas, I., Schmitz, B.: Protection of geoprivacy and accuracy of spatial information: how effective are geographical masks? Cartographica: The International Journal for Geographic Information and Geovisualization **39**(2), 15–28 (2004)
11. Liu, L., Andris, C., Ratti, C.: Uncovering cabdrivers' behavior patterns from their digital traces. Computers, Environment and Urban Systems **34**(6), 541–548 (2010)
12. Liu, Y., Liu, X., Gao, S., Gong, L., Kang, C., Zhi, Y., Chi, G., Shi, L.: Social sensing: A new approach to understanding our socioeconomic environments. Annals of the Association of American Geographers **105**(3), 512–530 (2015)
13. Mai, G., Janowicz, K., Hu, Y., Gao, S.: Adcn: An anisotropic density-based clustering algorithm for discovering spatial point patterns with noise. Transactions in GIS **22**(1), 348–369 (2018)
14. McKenzie, G., Janowicz, K., Gao, S., Yang, J.A., Hu, Y.: Poi pulse: A multi-granular, semantic signature–based information observatory for the interactive visualization of big geosocial data. Cartographica: The International Journal for Geographic Information and Geovisualization **50**(2), 71–85 (2015)
15. Seidl, D.E., Jankowski, P., Tsou, M.H.: Privacy and spatial pattern preservation in masked gps trajectory data. International Journal of Geographical Information Science **30**(4), 785–800 (2016)
16. Song, C., Qu, Z., Blumm, N., Barabási, A.L.: Limits of predictability in human mobility. Science **327**(5968), 1018–1021 (2010)