

Privacy in Participatory Sensing

Lab Based Report



Department of Electronics and Communication
Engineering

Indian Institute of Technology Roorkee

Submitted by

Palak (17116050)

Rahul (17116055)

Submitted to

Dr. Dheeraj Kumar

Privacy in Participatory Sensing

Palak Goenka
17116050
pgoenka@ec.iitr.ac.in

Rahul Patil
17116055
rpatil@ec.iitr.ac.in

June 3, 2020

Abstract

Nowadays mobile phones have evolved from merely being phones to full-fledged computing, sensing, and communication devices. The presence of multimodal sensors on mobile phones and their reachability has profoundly opened a great way for participatory sensing. Participatory sensing is the concept of communities contributing sensory information to form a body of knowledge. A well-known example of participatory sensing is the *Aarogya Setu Mobile App*. Now, it is pretty clear that participatory sensing is important for us but it comes with great cost in terms of privacy. In this report, we attempt to provide a detailed discussion of the various techniques to protect privacy in mobile participatory sensing. Furthermore, we have shown the results of our baseline model on environment centric participatory sensing.

1 Introduction

The traditional technique of environment data collection involved the deployment of a large number of static wireless sensor devices. This technique suffered from both less coverage as well as installation cost overhead. With the advent of participatory sensing, a technique to empower ordinary citizens to collect and share sensed-data from their surrounding environments using their mobile phones, these two major issues with the previous system have been resolved. In participatory sensing, the use of existing sensors (mobile phones) provides a zero deployment cost as well as unique spatiotemporal coverage.

A plethora of novel and exciting participatory sensing applications have emerged in recent years, which are used to collect information about things like traffic, quality of en-route WiFi access points, urban air quality, noise pollution, cyclist experiences, diets, and etc. A very common example of participatory sensing is Google maps. Their traffic recommendations are based on two major factors:

1. Historical data about the average time taken to travel a particular section of the road at specific times on specific days.
2. Real-time traffic data sent by sensors and smartphones.

In Fig 1, a man generated fake real-time data using 99 smartphones with sim cards and fooled Google maps by making a fake traffic jam following him.



Figure 1: Google Maps: An example of participatory sensing

1.1 Privacy aspect of participatory sensing

Current participatory sensing applications are primarily focused on the collection of data on a large scale which makes it reluctant towards privacy attacks. Although the subjects of interests of environment-centric applications are not the participants themselves, all considered applications monitor the spatiotemporal context of the participants and therefore represent a danger to their privacy. Furthermore, additional captured sensing modalities may provide further insights into the participants. As a result, environment-centric applications can endanger the privacy of the participants, even if the threats are less perceptible at first sight than in the case of people-centric applications. Hence, concluding sensing **privacy threats represent an inherent problem of any participatory sensing application.**

The participatory system is all based on the user's contribution and if the users become reluctant to share the data this whole technique will fail. So a lot of efforts have been made to ensure user's privacy in participatory sensing. In the subsequent sections, we have tried to present a summary of various privacy-preserving techniques out there and our plan for the same.

2 Privacy-preserving Techniques for participatory sensing

Some researchers have opted for graph theory (Gao, Ma, Shi, Zhan, and Sun (2013)), cell ID (a unique identifier for the cell tower that the mobile device is currently connected to) instead of GPS for location-based information (Fanourakis (2018)), group anonymization techniques (Gisdakis, Giannetsos, and Papadimitratos (2014)) to protect user's location and trajectory breaches. K-anonymity based solutions (Vu, Zheng, and Gao (2012)) are found

to be the most general opted technique. These include suppression based (Replace some of the attributes with asterisks and b) as well as generalization technique (Replace values with more generalized ones.). Letting people have the option to decide extent of granularity and hence the privacy protection level they desire (Agir, Papaioannou, Narendula, Aberer, and Hubaux (2014)), is also a customary way for privacy protection. A detailed description of some of the approaches has been presented.

2.1 AnonySense: a system for anonymous opportunistic sensing (Shin et al. (2011))

2.1.1 Description

- A privacy-aware system that allows applications to submit Sensing tasks, distributed across participating mobile devices.
- Receive verified and anonymized data report from the field.

2.1.2 Key Features

- **Framework:** A general-purpose framework for anonymous opportunistic tasking and reporting.
- Privacy of users as well as the integrity of reports.
- **Authentication:** Allows user or application authentication.
- **Efficient:** Consume little CPU time, network bandwidth, and battery energy.
- **Flexibility Feasibility Two applications:** RougeFinder ObjectFinder
- **AnonyTL:** A Simple and expressive language called AnonyTL for applications to specify their tasks. Instead of using an existing language such as SQL or XQuery, they designed AnonyTL with Lisp-like syntax to allow concise task descriptions, a small interpreter, easy portability to embedded platforms, and a clean fit with the sensing/reporting semantics.

2.1.3 Threats:

- **Threats to carrier privacy (mobile users)**
 - **Narrow Tasking:** A malicious tasking entity may submit tasks with restrictive acceptance conditions including, e.g., a rare sensor type or specific mobile phone brand, known as Narrow Tasking. This attack may allow the attacker to de-anonymize the mobile phones accepting these highly device-specific tasks, as only one or a few mobile phones share these restrictive conditions. Eg. Those who carry a heart-rate sensor paired with an iPhone.
 - **Tasking de-anonymization**

- Reporting de-anonymization
- **Selective Tasking:** A malicious tasking entity may attempt to differentiate and identify anonymous participants by launching selective tasking attacks, where the tasking entity distributes a task to only a restricted pool of mobile phones.
- **Threats to Report Integrity:**
 - Threats Report Tampering
 - Report Relay

2.1.4 Architecture Protocols:

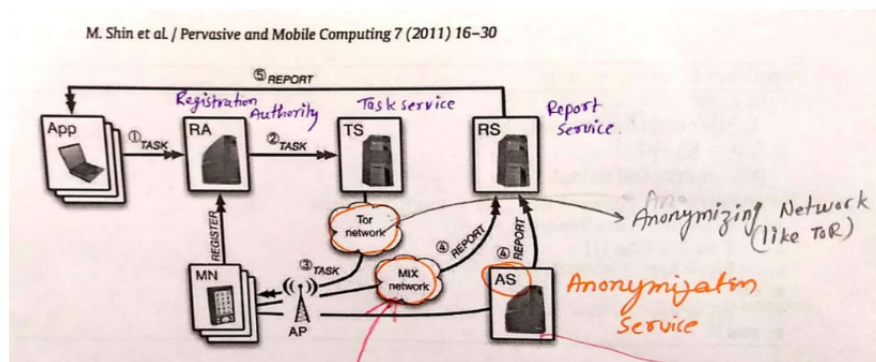


Figure 2: System Architecture

- **Registration Authority:**
 - The root of trust in a sense every other component trusts the certificates issued by RA.
 - Installation of AnonySense Software in MN, hence to calibrate all sensors in MN.
 - Providing public key of TS, RS, and AS.
 - Releasing task securely.
- **Anonymous Network (TOR)**
 - Serves to protect the network identity and location of the MN when it connects to the TS to download new tasks. Tor allows clients to anonymously connect to servers through multiple Tor relays using onion encryption. Traffic is delivered as quickly as possible and the location (IP address) of the client is hidden from the server. Each MN connects to the TS with a randomized interval, making it harder for the adversary to link between tasking connections.
- **Mixed Network**
 - An anonymizing channel between MNs and the RS to submit a report.

- Uses Mixmaster, the most popular MIX in use today.
- Although any remailer-type MIX network supporting SMTP email protocols would suffice.
- **Anonymization Service (AS)**
 - Blurring Techniques.

2.2 HealthSense: Classification of Health-related Sensor Data through User-Assisted Machine Learning (Stuntebeck, Davis, Abowd, and Blount (2008))

2.2.1 Description

- Despite the development of highly precise sensors to detect certain medical conditions, there are still few health-care conditions being directly undetectable like pain and depression. Believing this detectable bio-physical data correlated to events related to health conditions, a machine learning model can be used to predict such directly undetectable health events.
- **Example:**
 - Remote health monitoring scenarios include applying pulse oximeters to monitor blood oxygen levels.
 - Using glucometers for diabetic patients
 - Detecting cardiac events with portable EKG monitors
 - There are networked weight sensors for detecting weight loss and gain
 - Portable EEG sensors for monitoring epilepsy and EMG sensors for detecting muscle dynamics.
 - Additional sensors are on the horizon for application to a wide variety of conditions and will change the relationship between the health-care provider and patient.

2.2.2 Key Features

- A framework for real-time tagging of health-related sensor data.
- HealthSense transmits sensor data from the patient to a server for analysis via machine learning techniques.
- The system uses patient input to assist with the classification of unusual events (e.g., pain or itching).
- Occasionally asking the patient whether they are experiencing the condition being monitored (predicted) or not to improve the ML model.

2.3 Privacy preservation for participatory sensing data (LOCATE) (Boutsis and Kalogeraki (2013))

2.3.1 Description

- **LOCATE** allows users to locally sense and store data, as well as issue queries on data stored across various systems. In this case, the data is not uploaded to the centralized server.
- **Example Query:** "Give me the trajectories from location A and B".

2.3.2 System detail

- Each user maintains a local database.
- The local database can have two types of data, the first one is its own data and the second one is data of other users, which they have collected via querying.
- The queries are received through WiFi and cellular data (3G/4G).
- Proposed a data exchange approach via distributing the user data trajectories among multiple user databases, based on the local entropy.

2.3.3 Attack detail

- Queries are not known, so while executing queries filters are needed. Regex based filters have been implemented to remove malicious queries. Example of a malicious query, 'When a system (mobile phone in this case) queries another to share the personal information'.
- Now, *intents* are responsible for executing queries or say msg transfer in mobile devices. So, attackers might be able to access any private components in the phone.

2.3.4 Implementation detail

- This paper is very task-specific. As mentioned in the paper, each user stores 6 fields to their local storage. "*latitude, longitude, altitude, timestamp, point type, traj_id*"
 - **Point type:** Characterizes the point as starting point, destination point, moving point. Point types are defined in a specific way. For example, if the location coordinate doesn't change then the point type is marked to be *stop*.

2.4 PAMPAS: Privacy-Aware Mobile Participatory Sensing Using Secure Probes (That, Popa, Zeitouni, and Borcea (2016))

2.4.1 Description

- A privacy-aware mobile distributed system for efficient data aggregation in mobile participatory sensing and a solution for querying large numbers of users which protects user location privacy and works in real-time.

2.4.2 System Detail

- Mobile devices have secure probes (SPs) that perform distributed query processing while preventing users from accessing other users' data.
- A supporting server infrastructure (SSI) coordinates the inter-SP communication in an encrypted format and the computation tasks are executed on SPs.
- PAMPAS ensures that SSI cannot link the location reported by SPs to the user identities even if SSI has additional background information. Communication between SPs and SSI is anonymous, e.g., by using a proxy forwarder or anonymization network (e.g., Tor).
- Privacy-aware location-based aggregation and adaptive spatial partitioning of SPs that work efficiently on resource-constrained SPs.

2.4.3 Hybrid Architecture : (Decentralized + Centralized SSI)

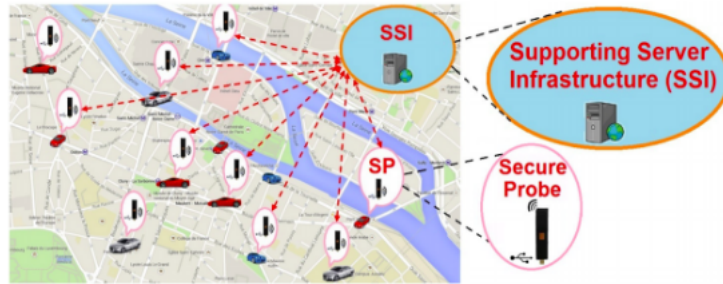


Figure 3: Hybrid architecture

2.4.4 Threat

- The attackers in PAMPAS could be either users or the owners of SSI. SSI is honest-but-curious, which means it will follow the protocol but may try to infer things from the data. Although, IP anonymity has been implemented using TOR networks. But IP anonymity is not enough to protect the user's privacy in MPSS because identity information could be determined from the location and sensing data also.

2.5 Enhancing Privacy in Participatory Sensing Applications with Multidimensional Data (Groat, Edwards, Horey, He, and Forrest (2012))

2.5.1 Description

- Negative survey-based approach on multivariate categorical data to protect privacy in participatory sensing applications. Along with a novel efficient reconstruction algorithm.

2.5.2 Key Features

- Perturbation using negative surveys, modified data from each user is sent to the central server.
- The server then reconstructs the probability density functions of the original distributions of sensed values, without knowing the participants' actual data.

2.6 Mutual Privacy-Preserving k-Means Clustering in Social Participatory Sensing (Xing, Hu, Yu, Cheng, and Zhang (2017))

2.6.1 Description

- A technique that can mutually protect the privacy of both the participants and the community, i.e., a technique that allows the data analyst (the social application server) to extract information about the community without accessing any user's private data, while no participatory participant can obtain any information about other participants and the community.

2.6.2 System Detail

- Consists of two privacy-preserving algorithms called at each iteration of the k-means clustering.
 - **Stage1:** Assign Participants to Their Nearest Centers
 - **Stage2:** Update the Cluster Centers. A public-key-based additive homomorphic encryption scheme is adopted for this.
- In the execution of these 2 algorithm data analysts are taken into consideration.

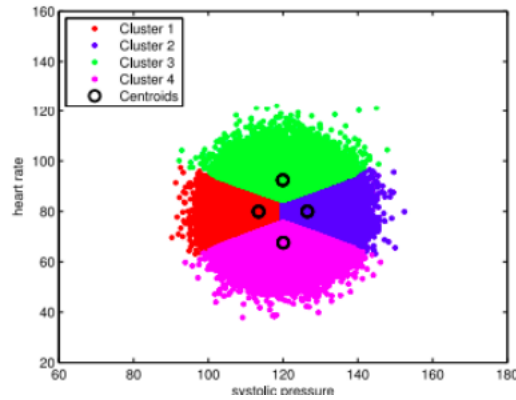


Figure 4: Clustering results

2.6.3 Privacy at participant and data analyst level

- Participants don't get the following information
 - Cluster centers
 - Other participants in the same cluster
 - Private data of other participants.
- Data analyst:
 - Knows which participant belongs to which cluster.
 - Don't have access to participant private data as well as the distance of the participant from the cluster center it is associated with.

2.7 Virtual Individual Servers as Privacy-Preserving Proxies for Mobile Devices (Cáceres, Cox, Lim, Shakimov, and Varshavsky (2009))

2.7.1 Description

- An approach for Privacy-aware data storage: In this kind of technique, people have control over what/with whom they want to share. Mostly in this type of approach people share data, but have control over what they want to share, either they sign a contract (using blockchain) or are asked to grant permission for specific data.
- Proposed the concept of using personal vaults to store the data from mobile phones. Such that the **owner can have control over his/her data.**

2.8 NoiseTube: Measuring and mapping noise pollution with mobile phones (Maisonneuve, Stevens, Niessen, and Steels (2009))

2.8.1 Description

- A technique for privacy protection in participatory sensing via sharing data processing between mobile phones and the application server. For example, if the central

server needs only loudness as a feature then instead of sending the voice signal a signal processing algorithm could be run on the phone itself - this will save both privacy as well as the bandwidth required to send the data.

2.8.2 System Detail

- NoiseTube is installed on GPS equipped phones, which collects data like time, GPS coordinates, noise to the NoiseTube server. In order to measure loudness in real-time, they have used the signal processing algorithm of mobile phones.

3 Proposal: A Federated learning approach for training a machine learning model on participatory sensing data (HealthCare)

3.1 Problem Statement

Even though new sensors are being developed, there are certain medical conditions for which automated detection is not feasible. These categories of health-care conditions are referred to as being directly undetectable. Heart rate falls under the category of directly detectable conditions while anxiety/depression falls under directly undetectable conditions. It has been believed that detectable conditions along with patients' feedback could help in determining undetectable conditions Stuntebeck et al. (2008). However, collecting these personal data may introduce privacy issues. **How can we overcome these privacy issues?**

3.2 Background

Standard machine learning approaches require centralizing the training data on one machine or in a data center under the assumption that servers are non-malicious. While previous studies have shown that malicious servers are the major cause of privacy breaches in participatory sensing. Unlike standard machine learning algorithms, **Federated learning** proposes a new way in which the raw data of users need not be sent to the central server/cloud.

3.2.1 Example of an existing product what works on Federated learning:

Google Keyboard

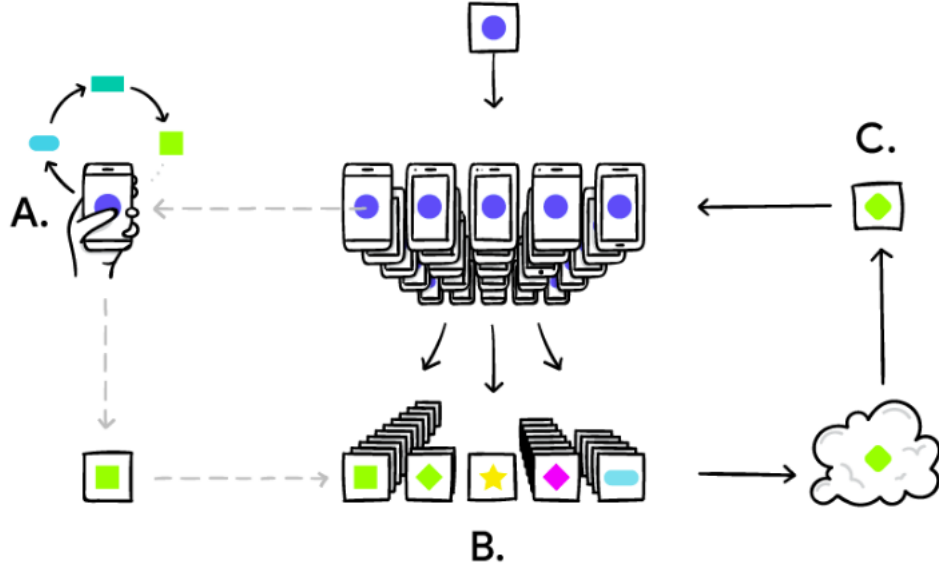


Figure 5: Federated learning

Working (see fig 5): Each device locally downloads the currently trained model (say for language modeling) from the server, this is the initial model for any device. After this, devices train the initial model locally and after a fixed amount of time (for eg. say 1 month), summarizes the changes as a small focused update. Only this update to the model is sent to the cloud, using encrypted communication. When a small patch/update is received at the server-side, small updates are aggregated together for the formation of an improved shared model.

Federated Learning allows for smarter models, lower latency, and less power consumption, all while ensuring privacy. Moreover, this approach has another immediate benefit, in addition to providing an update to the shared model, the improved model on the phones can also be used immediately, powering experiences personalized by the way one uses phones.

3.3 Plan

Taking inspiration from how Google keyboard works, we can use FL to train a machine learning model to predict directly undetectable medical conditions. As model training will be on the user's mobile phone itself, their personalized data will not be exposed to any server. With enhanced privacy, more patients will be willing to take part in collaborative model training thus, better inference models can be built.

4 A baseline model

4.1 Dataset

- WSN dataset.

- There are 54 unique positions, whose x,y coordinates are known.
- For each unique position, we have a time series of humidity values at an interval of 10 mins.

4.2 Experiment

Humidity Interpolation using kriging. Kriging is an advanced geostatistical procedure that generates an estimated surface from a scattered set of points with z-value. To make a prediction with the kriging interpolation method, two tasks are necessary:

1. Uncover the dependency rules.
2. Make the predictions.

To realize these two tasks, kriging goes through a two-step process:

1. It creates the variograms and covariance functions to estimate the statistical dependence (called spatial autocorrelation) values that depend on the model of autocorrelation (fitting a model).
2. It predicts the unknown values (making a prediction).

It is because of these two distinct tasks that it has been said that kriging uses the data twice: the first time to estimate the spatial autocorrelation of the data and the second to make the predictions.

4.3 Results

- For each time interval, we have calculated the RMS error averaged over the number of modes.
- Some of the error value and the plot has been mentioned below.

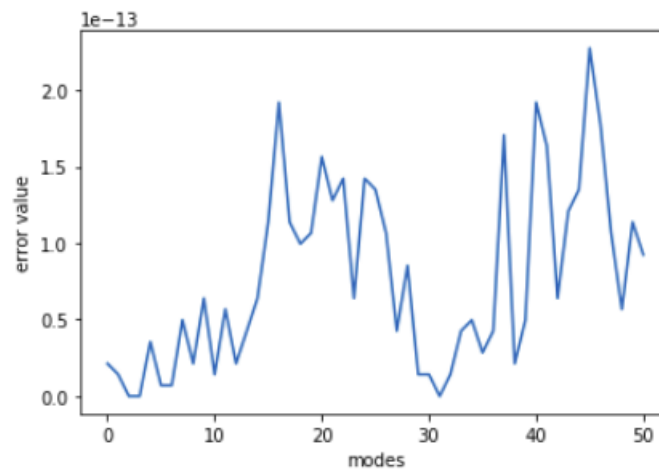


Figure 6: At time = 5th interval, RMS Value: 7.732376830330502e-14

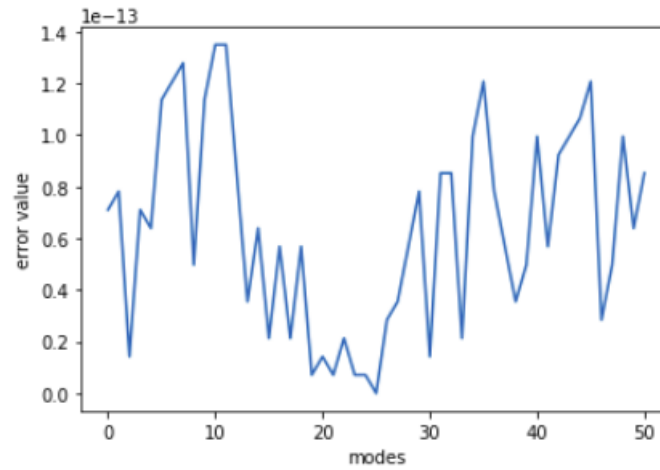


Figure 7: **At time** = 10th interval, **RMS Value:** 6.367020200830702e-14

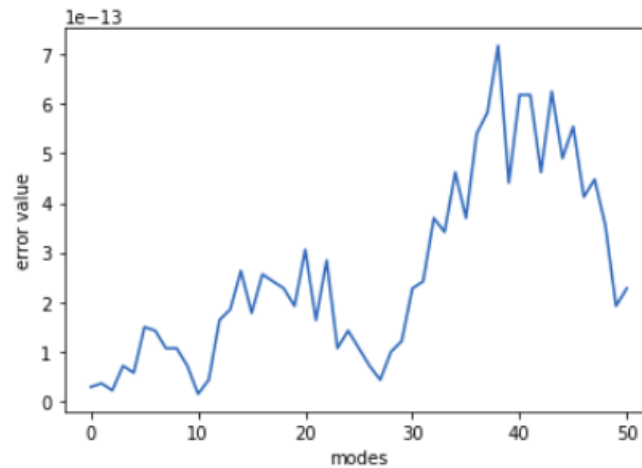


Figure 8: **At time** = 50th interval, **RMS Value:** 2.6053233644537005e-13

References

- Agir, B., Papaioannou, T. G., Narendula, R., Aberer, K., & Hubaux, J.-P. (2014). User-side adaptive protection of location privacy in participatory sensing. *GeoInformatica*, 18(1), 165–191.
- Boutsis, I., & Kalogeraki, V. (2013). Privacy preservation for participatory sensing data. In *2013 IEEE International Conference on Pervasive Computing and Communications (PerCom)* (pp. 103–113).
- Cáceres, R., Cox, L., Lim, H., Shakimov, A., & Varshavsky, A. (2009). Virtual individual servers as privacy-preserving proxies for mobile devices. In *Proceedings of the 1st ACM workshop on networking, systems, and applications for mobile handhelds* (pp. 37–42).
- Fanourakis, M. A. (2018). *On the feasibility and privacy benefits of on-device data mining for opportunistic crowd-sensing and service self-provisioning* (Unpublished doctoral dissertation). University of Geneva.

- Gao, S., Ma, J., Shi, W., Zhan, G., & Sun, C. (2013). Trpf: A trajectory privacy-preserving framework for participatory sensing. *IEEE Transactions on Information Forensics and Security*, 8(6), 874–887.
- Gisdakis, S., Giannetsos, T., & Papadimitratos, P. (2014). Sppear: security & privacy-preserving architecture for participatory-sensing applications. In *Proceedings of the 2014 acm conference on security and privacy in wireless & mobile networks* (pp. 39–50).
- Groat, M. M., Edwards, B., Horey, J., He, W., & Forrest, S. (2012). Enhancing privacy in participatory sensing applications with multidimensional data. In *2012 ieee international conference on pervasive computing and communications* (pp. 144–152).
- Maisonneuve, N., Stevens, M., Niessen, M. E., & Steels, L. (2009). Noisetube: Measuring and mapping noise pollution with mobile phones. In *Information technologies in environmental engineering* (pp. 215–228). Springer.
- Shin, M., Cornelius, C., Peebles, D., Kapadia, A., Kotz, D., & Triandopoulos, N. (2011). Anonymsense: A system for anonymous opportunistic sensing. *Pervasive and Mobile Computing*, 7(1), 16–30.
- Stuntebeck, E. P., Davis, J. S., Abowd, G. D., & Blount, M. (2008). Healthsense: classification of health-related sensor data through user-assisted machine learning. In *Proceedings of the 9th workshop on mobile computing systems and applications* (pp. 1–5).
- That, D. H. T., Popa, I. S., Zeitouni, K., & Borcea, C. (2016). Pampas: privacy-aware mobile participatory sensing using secure probes. In *Proceedings of the 28th international conference on scientific and statistical database management* (pp. 1–12).
- Vu, K., Zheng, R., & Gao, J. (2012). Efficient algorithms for k-anonymous location privacy in participatory sensing. In *2012 proceedings ieee infocom* (pp. 2399–2407).
- Xing, K., Hu, C., Yu, J., Cheng, X., & Zhang, F. (2017). Mutual privacy preserving k -means clustering in social participatory sensing. *IEEE Transactions on Industrial Informatics*, 13(4), 2066–2076.