

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/269297802>

# IoT-Privacy: To be private or not to be private

Conference Paper in Proceedings - IEEE INFOCOM · April 2014

DOI: 10.1109/INFCOMW.2014.6849186

CITATIONS

41

READS

1,324

3 authors:



[Arijit Ukil](#)

Tata Consultancy Services Limited

61 PUBLICATIONS 481 CITATIONS

SEE PROFILE



[Soma Bandyopadhyay](#)

Tata Consultancy Services Limited

51 PUBLICATIONS 518 CITATIONS

SEE PROFILE



[Arpan Pal](#)

Tata Consultancy Services Limited

189 PUBLICATIONS 841 CITATIONS

SEE PROFILE

Some of the authors of this publication are also working on these related projects:



Home Infotainment Platform [View project](#)



Human Sensing for wellness and behavior modeling [View project](#)

# IoT-Privacy :To Be Private or Not To Be Private

Arijit Ukil  
Innovation Lab  
Tata Consultancy Services  
Kolkata, India  
[arijit.ukil@tcs.com](mailto:arijit.ukil@tcs.com)

Soma Bandyopadhyay  
Innovation Lab  
Tata Consultancy Services  
Kolkata, India  
[soma.bandyopadhyay@tcs.com](mailto:soma.bandyopadhyay@tcs.com)

Arpan Pal  
Innovation Lab  
Tata Consultancy Services  
Kolkata, India  
[arpan.pal@tcs.com](mailto:arpan.pal@tcs.com)

**Abstract**—Privacy breaching attacks pose considerable challenges in the development and deployment of Internet of Things (IoT) applications. Though privacy preserving data mining (PPDM) minimizes sensitive data disclosure probability, sensitive content analysis, privacy measurement and user’s privacy awareness issues are yet to be addressed. In this paper, we propose a privacy management scheme that enables the user to estimate the risk of sharing private data like smart meter data. Our focus is to develop robust sensitivity detection, analysis and privacy content quantification scheme from statistical disclosure control aspect and information theoretic model. We depict performance results using real sensor data.

**Keywords**—privacy; smart meter; statistical disclosure; sensitivity; Wasserstein distance;

## I. INTRODUCTION

IoT applications like smart home, smart energy management render innumerable benefits to human society. Sensors like smart meters collect sensitive personal information like detailed household energy consumption profile. However when such data is released to third parties, possibility of unintentional or malicious privacy breach like activity detection of the users is very high. Though released data can be privacy protected by standard privacy preservation techniques like noise addition, suppression; user or private data owner should also be aware of the privacy content of his/her sensor data. Such analysis also acts as a precursor in enforcing optimal privacy preservation. In this paper we present a privacy management scheme that detects, analyzes sensitive content of sensor data and measures the amount of privacy. We consider time-series sensor data, though proposed scheme is generic in nature. Without loss of generality, we consider smart meter data for analysis and assume appliance load change recovery and peak load monitoring as privacy threats and main privacy violators. Our scheme has two distinct components. First we detect and analyze the sensitivity in typical sensor dataset using robust statistical method. Then we measure the privacy content of the sensor dataset using information theoretic model that is consistent with analytical reasoning. We perform experimental tests using real sensor data [1] and compare our results with relevant techniques [2-3].

## II. SENSITIVITY DETECTION AND ANALYSIS

We define sensitivity as statistical anomalies in the sensor data that describe the presence of unusual or unanticipated events. Most of the relevant proposed schemes are either completely dependent on the application use case (intrusive) [4] or totally independent [2 – 3] (rudimentary). We propose

an unobtrusive sensitivity detection and analysis algorithm that discovers the anomaly points in sensor dataset ( $\mathcal{S}_t$ ) while optimizing the masking and swamping effects. The pattern of  $\mathcal{S}$  over a large time period may show periodic and predictable pattern, few anomalies would provide hint of existence of private events. Firstly, the distribution pattern of  $\mathcal{S}_t, t = n, n = 1, 2, \dots$  is derived using fourth order statistical moment (kurtosis,  $\kappa$ ). When  $\kappa > 3$  (leptokurtic), Rosner filtering is executed to minimize swamping effect [5 – 6]. Unlike in traditional outlier detection tests or clustering algorithms, we need not specify the number of sensitive points a priori. Given the upper bound,  $\Phi$ ,  $\Phi$  number of backward selection tests are performed to identify initial anomalies,  $v_i, i = 1, 2, \dots, \Phi$  with respect to median deviation. Accordingly,  $\Phi$  critical values are computed as:

$$\epsilon_i = \frac{(\Phi-i)s_{p,n-i-1}}{\left(\left(\Phi-i-1\right)s_{p,n-i-1}^2\right)^{1/2}}, \quad s_{p,v} \text{ is the } 100p$$

percentage point from the student’s  $t$  distribution with  $\nu$  degrees of freedom and  $p = 1 - \frac{\theta}{2(\Phi-i-1)}$ . Sensitive points in  $\mathcal{S}_t$  are the largest  $i$  such that  $v_i > \epsilon_i$ . Appropriate values of  $\Phi$  and  $\theta$  are chosen, which are derived from application metadata to preserve the generic nature, while capturing the properties of particular sensor application. For example, we choose  $\Phi = \mathcal{B} \times |\mathcal{S}_t|$ , which is a very conservative estimate. When  $\kappa \leq 3$ , Hampel filter is employed to detect sensitive points so that masking effect is minimized. We test sensitivity detection and analysis outcome using real smart meter data, REDD dataset [1] and the result is shown in figure 1.

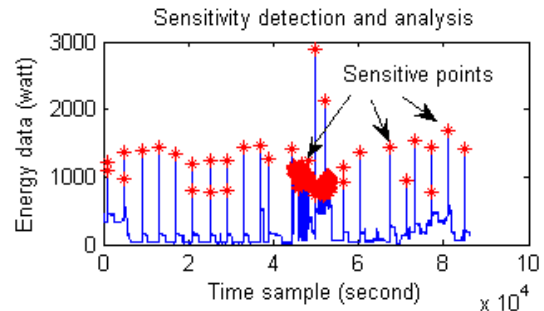


Fig. 1. Sensitivity detection and analysis

We can prove that our proposed scheme has better detection capability with optimal masking and swamping effects when compared with [2– 3] using different measures like Mahalanobis distance, KullbackLeibler (KL) divergence.

For example, KL divergence of our scheme is more than that of [2–3] and Sanov’s theorem from large deviation theory shows more KL divergence means more detection capability.

### III. PRIVACY MEASUREMENT AND QUANTIFICATION

Privacy measurement and quantification scheme is computed based on statistical and information theoretic model. Our objective is to derive privacy measure from fundamental principle for disambiguation of privacy measure among different privacy preservation guarantees like  $k$ -anonymity,  $l$ -diversity. Below we describe our proposed scheme. Functional block diagram is shown in figure 2.

#### A. Privacy measurement

Consider  $v$  be the sensitive part of  $\mathcal{S}$  and  $\lambda$  be the non-sensitive part;  $\mathcal{S} = v \cup \lambda$ . We define privacy measure ( $\rho_M$ ) as the amount of difficulty to infer  $v$ , i.e. how much the probability of finding  $v$ , given  $\mathcal{S}$ , i.e. the information leakage transfer function  $Y_{\mathcal{S},v}$ :  $\mathcal{S} \rightarrow v$  and  $Y_{\mathcal{S},v} = \rho_M = \frac{\sum_{i=1}^{|\mathcal{V}|} Pr(v_i) \log_2 \frac{1}{Pr(v_i)}}{\sum_{i=1}^{|\mathcal{S}|} Pr(S_i) \log_2 \frac{1}{Pr(S_i)}}$ . In [7], such metric is derived using mutual information  $I(\mathcal{S}, v)$ . As  $v \in \mathcal{S}$ , leakage function  $\mathcal{L}_{\mathcal{S},v} = I(\mathcal{S}, v)$  would indicate the maxima.

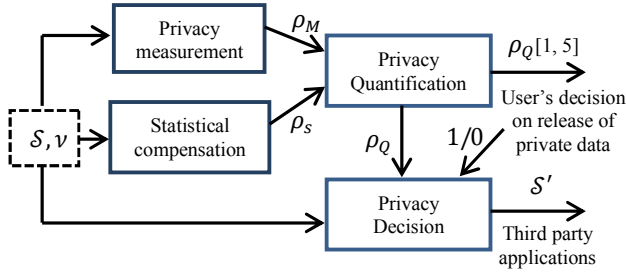


Fig. 2. Privacy measurement, quantification, user decision

#### B. Statistical compensation

In order to enhance privacy measurement and quantification accuracy, statistical compensation due to statistical relation between  $\mathcal{S}$  and  $v$  is computed. Here, two-sample Kolmogorov-Smirnov (KS) test of  $\mathcal{S}$  and  $v$  is performed. KS test is a nonparametric hypothesis test that evaluates the difference between the Cumulative Distribution Functions (CDFs). It computes under the null hypothesis that  $\mathcal{S}$  and  $v$  are drawn from the same distribution. When KS-test accepts null hypotheses, statistical compensation  $\rho_S = 1$ . When KS-test rejects null hypotheses, we propose L1-Wasserstein metric between  $\mathcal{S}, v$  ( $w_{\mathcal{S},v}$ ) to estimate statistical misfit or compensation  $\rho_S = w_{\mathcal{S},v}$ . Wasserstein distance quantifies the numerical cost with respect to distribution dissimilarity between pair of distributions, defined for  $\mathcal{S}, v \in \Omega$ :

$$w_{\mathcal{S},v} := \inf_{\mu \in \Omega(\mathcal{S},v)} \int_{\Omega} |x - y| d\mu(x, y), \quad x \in \mathcal{S}, \quad y \in v$$

As  $w_{\mathcal{S},v}$  is not straightforward for implementation, we choose closed solution considering CDFs of  $\mathcal{S}, v$  as [8]:

$$w_{\mathcal{S},v} = \int_0^1 |\mathcal{P}^{-1}(q) - \mathcal{Q}^{-1}(q)| dq, \quad \text{where } \mathcal{P}, \mathcal{Q} \quad \text{be} \quad \text{the} \quad \text{distribution} \quad \text{functions} \quad \text{of} \quad \mathcal{S}, v.$$

#### C. Privacy quantification

Logically, privacy quantification  $\rho_Q = \rho_M \wedge \rho_S$ . Algebraically,  $\rho_Q = \rho_M \times \rho_S$ . With  $\rho_Q \in [0, 1]$ , we scale  $\rho_Q$  as  $\rho_Q \mapsto \lceil \rho_Q \times 5 \rceil$ :  $\rho_Q = [1, 5]$ , with high magnitude of  $\rho_Q$  signifies more privacy risk probability in  $\mathcal{S}$ .

We depict in figure 3 the outcome of privacy quantification of our method comparing with [2–3] using REDD dataset. We compute privacy measure of four equal parts of the day. However, efficacy of the proposed scheme can be established by measuring the privacy risk probability when an attack with standard disaggregation or NILM (Non-Intrusive Load Monitoring) is launched, which is our future research scope.

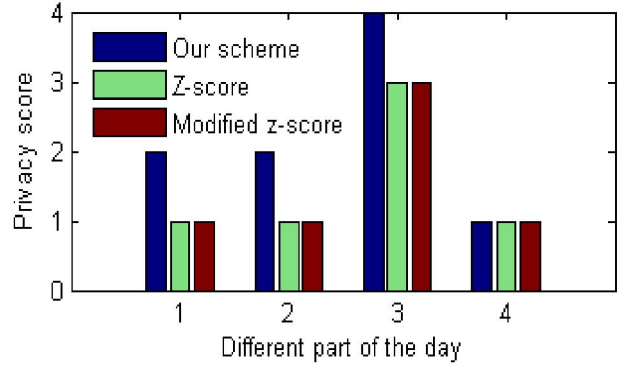


Fig. 3. Privacy measure outcome

#### D. Privacy decision

Privacy quantification enables the users to decide whether to share its private sensor data to avail third party applications. When affirmative,  $\mathcal{S}'$  is computed and shared to third parties.

Here,  $\mathcal{S} \xrightarrow{\rho_Q} \mathcal{S}'$  using standard PPDM, where  $\mathcal{S}'$  is the privacy preserved sensor data. Informally, the strength of perturbation or generalization of PPDM scheme is proportional to  $\rho_Q$ . For example, if  $\mathcal{S}' = \mathcal{S} + \mathcal{N}$ ,  $\mathcal{N} = f(\rho_Q)$  or  $l$ -diversity of  $\mathcal{S}$  would be  $l = g(\rho_Q)$ .

### REFERENCE

- [1] Z. Kolter, and M. J. Johnson, "REDD: A public data set for energy disaggregation research," SustKDD, 2011.
- [2] R. Rao, S. Akella, G. Guley, "Power Line Carrier (PLC) Signal Analysis of Smart Meters for Outlier Detection," IEEE SmartGridComm, pp. 291 - 296, 2011.
- [3] R. M. Nascimento, et al., "Outliers' Detection and Filling Algorithms for Smart Metering Centers," IEEE PES, pp.1 - 6, 2012.
- [4] W. Yang, et al., "Minimizing Private Data Disclosures in the Smart Grid," ACM CCS, pp. 412- 427, 2012.
- [5] B. Rosner, "Percentage points for a generalized ESD many-outlier procedure," Technometrics, vol. 25, issue. 2, pp. 165 - 172, 1983.
- [6] R. Serfling, and S. Wang, "General Foundations for Studying Masking and Swamping Robustness of Outlier Identifiers," Elsevier Statistical Methodology, August 2013.
- [7] L. Sankar, S.R. Rajagopalan, S. Mohajer, and H.V. Poor, "Smart Meter Privacy: A Theoretical Framework," IEEE Transactions on Smart Grid, vol. 4, issue. 2, pp. 837 - 846, 2013.
- [8] A. Halder, and R. Bhattacharya, "Further results on probabilistic model validation in Wasserstein metric," IEEE Annals of the IEEE Conference on Decision and Control, pp. 5542 - 5547, 2012.