

# CESense: Cost-Effective Urban Environment Sensing in Vehicular Sensor Networks

Quan Yuan<sup>1</sup>, Haibo Zhou<sup>1</sup>, Senior Member, IEEE, Zhihan Liu, Jinglin Li<sup>2</sup>, Member, IEEE, Fangchun Yang, Senior Member, IEEE, and Xuemin Shen, Fellow, IEEE

**Abstract**—In vehicular sensor networks, vehicles can act as mobile sensors to monitor the dynamic features of the physical world such as traffic flow, air quality, and temperature. However, the conventional full-coverage sensing approach is neither realizable nor cost-effective since the sensor-equipped vehicles are unevenly distributed and the environmental data are spatiotemporally correlated. To this end, we propose a cost-effective urban environment sensing solution (CESense), that exploits the sensing data correlations to improve the sensing accuracy and efficiency. CESense gathers data only at some specific areas of the whole sensing space and reliably infers the status of unsensed areas. Particularly, CESense uses a probabilistic matrix factorization model to reveal the latent features that impact the environmental status. Then, an appropriate set of sensing areas can be selected by fully taking advantage of these latent features and the sensing resource distribution patterns. In addition, to be adaptive to the dynamic environment, a checkpoint mechanism is designed to supervise the data gathering progress. Extensive experiments, which are based on the real taxicab mobility traces and air quality data collected in Beijing city, demonstrate that CESense can significantly improve the accuracy and efficiency of vehicular sensing.

**Index Terms**—Vehicular sensor networks, urban sensing, sensing quality, cost-effectiveness.

## I. INTRODUCTION

URBAN sensing [1] employs remote sensors to gather various information from urban space, such as traffic flow, crowd density, air quality, and temperature. Although static sensor networks have been widely deployed in cities for urban sensing, their sensing coverage and granularity are limited due to the sparsely installed sensor nodes. To this end, mobile crowdsensing [2], [3] has emerged recently, which empowers ordinary citizens to contribute data sensed from

their mobile devices (e.g., smartphones, wearable devices, and connected vehicles). Undoubtedly, connected vehicles equipped with various sensors are the most appropriate devices for mobile crowdsensing due to their high mobility, sufficient energy, and powerful communication and computation capabilities [4]. With the rising popularity of connected vehicles, vehicular sensing will play an important role in smart city, smart transportation system, and automated driving.

The connected vehicles are organized as a vehicular sensor network (VSN) [5] when participating urban sensing. Compared with existing urban sensing systems, VSNs can exploit the high mobility of vehicles to provide wide-coverage and fine-grained sensing services. Several vehicular sensing platforms [6]–[11] have been designed to collect data generated in urban space. In this paper, a VSN-based urban sensing system is considered, which involves not only a set of sensor-equipped vehicles, but also an urban sensing center and several sensing task publishers. The task publishers are usually organizations who need to monitor the status of specific urban areas, such as traffic management center, meteorological center, and environmental protection agency. The sensing center acts as an intermediary between task publishers and sensor-equipped vehicles, which shields the task publishers from the complexity of sensing task management. It encapsulates the sensing function and provides sensing as a service. Specifically, the sensing center recruits and incentivizes eligible vehicles to perform the sensing tasks submitted by task publishers.

The vehicular sensor data are useful only if their quality is acceptable [12]. The sensing center aims at ensuring the overall quality of sensing service with the least cost/incentive. Generally, the sensing center recruits a minimum set of vehicles to fully cover the sensing space at the required level of granularity. However, *the full-coverage sensing is neither realizable nor cost effective* [13]. On one hand, the vehicles constituting VSNs are usually taxicabs, buses, rental cars, and some private cars. Their real-time trajectories are hardly controllable, which causes the uneven distribution of sensing resources over space and time. To analyze this phenomenon, we use the taxicab trajectory dataset in Beijing and randomly select some taxicabs as sensing resources. Fig. 1a and 1b show the coverage of sensing resources during an off-peak hour and a peak hour, respectively. There exist many vacant areas that are unable to be covered by the sensing resources. Furthermore, Fig. 1c and 1d show the CCDF (Complementary Cumulative Distribution Function) of location vacancy rate and time vacancy rate for different number of sensing resources.

Manuscript received March 22, 2017; revised February 25, 2018; accepted September 20, 2018. Date of publication October 31, 2018; date of current version August 27, 2019. This work was supported in part by the Natural Science Foundation of Beijing under Grant 4181002, in part by the Natural Science Foundation of China under Grant 91638204 and Grant 61876023, and in part by the Natural Sciences and Engineering Research Council of Canada. The Associate Editor for this paper was F. Nashashibi. (Corresponding author: Zhihan Liu.)

Q. Yuan, Z. Liu, J. Li, and F. Yang are with the State Key Laboratory of Networking and Switching Technology, Beijing University of Posts and Telecommunications, Beijing 100876, China (e-mail: yuanquan@bupt.edu.cn; zhihan@bupt.edu.cn; jlli@bupt.edu.cn; fcyang@bupt.edu.cn).

H. Zhou is with the School of Electronic Science and Engineering, Nanjing University, Nanjing 210023, China (e-mail: haibozhou@nju.edu.cn).

X. Shen is with the Department of Electrical and Computer Engineering, University of Waterloo, Waterloo, ON N2L 3G1, Canada (e-mail: sshen@uwaterloo.ca).

Digital Object Identifier 10.1109/TITS.2018.2873112

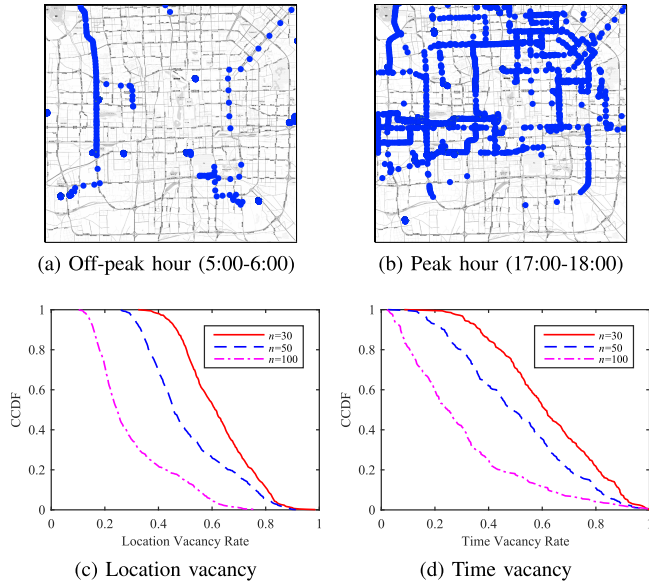


Fig. 1. The distribution of sensing resources within the 4th Ring Road of Beijing (about  $19 \times 18 \text{ km}^2$ ). In (a) and (b), the trajectories of 50 sensor-equipped vehicles are shown. In (c) and (d), the sensing space is partitioned into disjoint sensing areas, where each sensing area is about  $1 \text{ km}^2$  and each sensing cycle is 1 hour. Then, the location and time vacancy for  $n = 30, 50, 100$  sensor-equipped vehicles are shown, respectively.

For instance, when 50 sensor-equipped vehicles are recruited, 70% of the sensing cycles have a location vacancy rate greater than 40%; meanwhile, 60% of the sensing areas have a time vacancy rate larger than 40%. Therefore, it is almost impossible to cover the whole sensing space in every sensing cycle. On the other hand, there are usually spatiotemporal correlations among sensing data. The areas with similar features probably have similar trend of sensing data. Therefore, it is not necessary to gather data at all areas in all sensing cycles.

If the spatiotemporal correlations among sensing data are exploited, the sensing center can improve sensing quality and reduce sensing cost significantly. Sensing center can collect sensing data only at several selected areas during a sensing cycle. As for the unsensed areas, their status can be reliably inferred using the current as well as historical collected data. To this end, we propose a cost-effective urban environment sensing framework in Fig. 2. Specifically, when the task publisher publishes a sensing task, the sensing center decomposes the sensing task into several subtasks which are scattered over the sensing space. Then the sensor-equipped vehicles perform the subtasks and aggregate the sensed data to the sensing center. The sensing center infers the data in the unsensed areas and reports the complete and satisfactory sensing results to the task publisher. In these procedures, the task decomposition and missing data inference are technically challenging. First of all, the missing data inference method is the foundation of the framework. How to uncover and utilize the spatiotemporal correlations among the sensing data to achieve missing data inference is critical. Then, the selected sensing areas have an immediate impact on the overall cost of the sensing task. To be cost effective, the optimal set of subtasks should lie in the areas that are informative enough (i.e., the areas with “○” in Fig. 2). How to quantify the

informativeness of each area before obtaining the sensing data in these areas is difficult. In addition, the subtasks will become infeasible if no sensing resources exist in the corresponding areas. How to determine the subtasks without knowing future distribution of sensing resources is important. Furthermore, it is hard to compromise the subtask optimality with the feasibility to get a suboptimal solution (i.e., the areas with “\*” in Fig. 2). To deal with these challenges, we propose a cost-effective urban environment sensing solution called CESense. CESense can be beneficial to various urban sensing scenarios, provided that the sensing data are with intrinsic spatiotemporal correlations. Our contributions are summarized as follows:

- The latent features of environmental status in both spatial and temporal dimensions are revealed by the probabilistic matrix factorization model. These latent features reflect the spatiotemporal correlations among the sensing data and are used to infer the data in the unsensed areas.
- The latent features are well utilized to quantify the informativeness of each sensing area. The distribution patterns of sensor-equipped vehicles are analyzed and modeled. A greedy algorithm is designed to select potential sensing areas in a batch manner which balances the sensing cost and subtask feasibility.
- The sensing accuracy and efficiency of CESense are evaluated using the real taxicab mobility traces and air quality data in Beijing.

The remainder of this paper is organized as follows. Section II reviews related research works. The system model and problem formulation are presented in Section III. The cost-effective urban environment sensing solution CESense is elaborated in Section IV. The performance of CESense is evaluated in Section V. Finally, the work is concluded in Section VI.

## II. RELATED WORK

In early works [6]–[11], the sensor-equipped vehicles continuously gather data with certain time intervals or distance intervals. Then, the gathered data are aggregated into cloud via Vehicle-to-Infrastructure (V2I) communications; or alternatively, directly shared with other vehicles via Vehicle-to-Vehicle (V2V) or Device-to-Device (D2D) communications [14]–[16]. However, the distribution of sensor-equipped vehicles is uneven over space and time, which usually leads to low quality, low efficiency, and high cost for urban sensing in these vehicular sensing platforms.

### A. Participants Selection

Several studies make use of participant selection method to schedule the sensing resources. Hamid *et al.* [17] propose a trajectory-based recruitment scheme which selects a minimum set of vehicles to achieve a required level of coverage for the sensing space. He *et al.* [18] take into account the constraints on sensing quality and the time budgets of participants when allocating sensing tasks. However, these approaches either assume that the future trajectories of participants are known, or only consider the current locations of participants. To overcome these shortcomings, He *et al.* [19] aim at optimizing both the spatial and temporal coverage using predictable

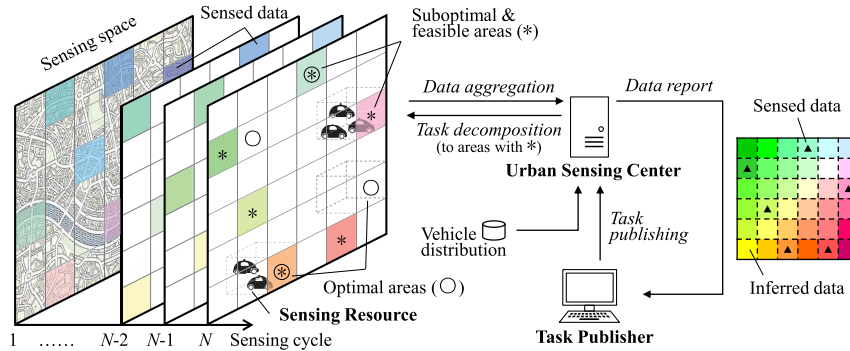


Fig. 2. Framework of cost-effective urban environment sensing using VSN.

mobility of participants. In addition, Ji *et al.* [20] maximize the spatial and temporal coverage without breaking the original commuting plans of participants. However, the participants selection method is applicable only when the sensing resources are abundant. It cannot ensure sensing quality for the areas where sensing resources are scarce.

### B. Missing Data Inference

As sensor-equipped vehicles cannot cover all roads for all time, a lot of spatiotemporal vacancies exist in the collected sensing data. Such missing data problem significantly impedes fine-grained urban sensing applications. The spatial and temporal correlations among sensing data are the foundations of existing missing data inference methods. Several statistic-based methods can be used, such as mean imputation and multiple imputation [21]. Additionally, there are many data-driven methods.  $K$ -Nearest-Neighbor (KNN) [22] is a classic method which estimates a missing value using the weighted average of its  $K$  nearest neighbors. Matrix completion [23] is a new technique to infer the empty entries in a low-rank matrix. Several recent studies [24]–[27] have adopted matrix completion to recover missing or corrupted sensing data. STCDG [24] combines the low-rank and short-term stability features in matrix completion. Kong *et al.* [25] further combine the feature of spatial correlation in matrix completion. Du *et al.* [27] exploit the intrinsic relationship between the entropies of samples and the matrix completion error. Then they propose a novel sampling rule based on this relationship to improve the sensing accuracy. Machine learning is a competitive approach to fill missing data. U-Air [28] uses a semi-supervised learning method to infer the air quality of unsensed areas, based on cross-domain data correlations. ST-MVL [29] is a multi-view-based learning method to fill missing values from both spatio-temporal and global-local perspectives.

### C. Efficient Sensing Approaches

The spatiotemporal correlations among sensing data are also exploited to improve sensing efficiency. Compressive sensing theory [30] makes it possible to recover certain signals from far fewer samples. Some studies [31]–[34] use compressive sensing to collect urban data, which largely reduce communication cost while guaranteeing required sensing accuracy. In [31], each vehicle computes the sparse representation of its original sensing data and transmits the sparse data to an

aggregator for recovery. In [32], the vehicles are clustered and the cluster heads are responsible for in-network data compression. Wang *et al.* [33] design a compressive sensing based monitoring method in a V2V opportunistic scenario. Xu *et al.* [34] apply compressive sensing to efficiently gather the data with multi-dimensional nature. Matrix completion, as an extension of compressive sensing, has been used to collect sensing data at required accuracy with reduced sensing and communication cost [12], [35]–[37]. Our previous work AC-Sense [12], utilizes the temporal uncertainty and the spatial similarity of sensing data to adaptively assign sensing tasks. In [35], Xie *et al.* propose an online data gathering scheme in wireless sensor networks, which adaptively samples different locations according to the environmental conditions. CCS-TA [36] combines matrix completion, Bayesian inference, and active learning techniques to select a minimum set of areas for sensing. Meng *et al.* [37] develop an integrated framework which employs matrix factorization and truth discovery to tackle the redundancy and sparsity problem. However, these methods do not take the dynamic distribution of sensing resources into consideration.

## III. PROBLEM FORMULATION

In this section, some definitions used throughout the paper are presented, and then a cost-effective urban sensing problem is formulated.

**Definition 1 (Sensing Task Requirements):** Sensing task requirements, which are proposed by a task publisher, indicate targeted sensing space, spatiotemporal granularity, and the accuracy of sensing results. To carry out fine-grained sensing, the targeted sensing space is partitioned into disjoint and uniform grids/areas (as shown in Fig. 2). The spatial granularity determines the scale of an area, while the temporal granularity specifies the length of a sensing cycle.

**Definition 2 (Environmental Data):** Environmental data refer to the status of sensing area, such as traffic flow, crowd density, air quality, and temperature.

A sensing task aims at acquiring the environmental data which meet its sensing task requirements. The cost-effective sensing is based on the correlations among environmental data. To explore and then utilize the correlations, we first define a data structure to organize the environmental data.

**Definition 3 (Environmental Matrix):** Environmental matrix is a matrix that holds environmental data within a specific spatiotemporal scope.



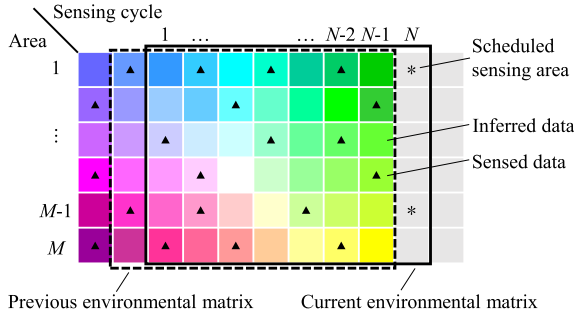


Fig. 3. An example of environmental matrix.

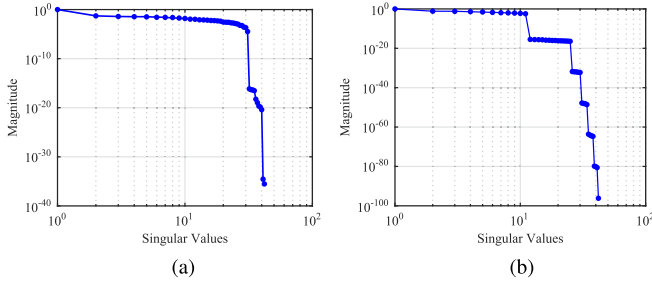


Fig. 4. Magnitude of singular values. (a) PM10 dataset. (b) NO2 dataset.

In this paper,  $\mathbf{X} \in \mathbb{R}^{M \times N}$  is an environmental matrix holding the environmental data for  $M$  areas in recent  $N$  sensing cycles. Specifically, the  $i$ th row of  $\mathbf{X}$  is the environmental data sequence in area  $i$ , while the  $j$ th column of  $\mathbf{X}$  contains the environmental data of all areas in sensing cycle  $j$ .  $\mathbf{X}$  is a dynamic matrix similar to a sliding window and its  $N$ th column always represents current sensing cycle. As shown in Fig. 3, the dashed and solid rectangles are two successive environmental matrices. The index  $(i, j)$  is called a spatiotemporal cell (or cell for short). Denoting the complete set of spatiotemporal cells as  $\mathcal{T} = \{(i, j) \mid 1 \leq i \leq M, 1 \leq j \leq N\}$ . A small set  $\mathcal{S} \subseteq \mathcal{T}$  consists of the spatiotemporal cells with sensed environmental data  $\mathbf{X}_{\mathcal{S}}$ . As for the spatiotemporal cells in  $\mathcal{T} \setminus \mathcal{S}$ , their environmental data  $\mathbf{X}_{\mathcal{T} \setminus \mathcal{S}}$  are inferred by a missing data completion method.

A special matrix  $\tilde{\mathbf{X}} \in \mathbb{R}^{M \times N}$  is used to hold the ground truth of  $\mathbf{X}$ . The vehicles are assumed to be equipped with high-quality sensors, such that the environmental data gathered from vehicles exactly reflect the ground-truth data (i.e.,  $\mathbf{X}_r = \tilde{\mathbf{X}}_r$  for any spatiotemporal cell  $r \in \mathcal{S}$ ). Additionally, it is assumed that one measurement for a spatiotemporal cell is sufficient to obtain the true environmental data.

Once the environmental matrix has been constructed, the spatiotemporal correlations among environmental data can be explored. Singular value decomposition (SVD) is effective to reveal the data correlations in a matrix. As a case study, the SVD is computed for a PM10 dataset and a NO2 dataset, respectively. The magnitude (ratio to the maximum) of singular values is shown in Fig. 4. It can be seen that most of the energy (i.e., the Frobenius norm of the matrix) is constrained by the first few principal components, which will lead the environmental matrix to a low rank. The low-rank feature is the result of spatiotemporal correlations among the environmental data and provides a basis for missing data inference.

Sensing accuracy and sensing cost are the main concerns of this paper. Sensing accuracy is specified in sensing task requirements and quantified by sensing error.

*Definition 4 (Sensing Error):* For sensing cycle  $j$ , sensing error is defined as the root-mean-square error (RMSE) between the environmental data in the  $j$ th column of  $\mathbf{X}$  and their corresponding ground truth in  $\tilde{\mathbf{X}}$ ,

$$\varepsilon_j = \sqrt{\frac{1}{M} \sum_{i=1}^M (\mathbf{X}_{ij} - \tilde{\mathbf{X}}_{ij})^2}. \quad (1)$$

The sensing task requirements usually specify an upper bound  $\epsilon$  for  $\varepsilon_j$  in each sensing cycle, i.e.,  $\varepsilon_j \leq \epsilon$ .

*Definition 5 (Sensing Cost):* Sensing cost is the incentive provided by the sensing center to reward the sensing data contribution from participants. Generally, the amount of incentive is proportional to the effort made by the participants. However, the incentive mechanism can be very complex, especially when game theory based methods are used. For illustration purpose, a simple incentive mechanism is assumed in this paper: constant and pay-per-data incentive. Therefore, the sensing cost is only related to the number of measurements that the sensing center needs.

Based on these definitions, the cost-effective urban environment sensing problem is formulated as follows. Given a sensing task and its requirements, the sensing center decomposes it into a set of subtasks scattered over the spatiotemporal cells  $\mathcal{S}$ . With the help of a missing data inference method, the sensing cost of these subtasks should be minimized while the sensing accuracy is ensured,

$$\min |\mathcal{S}|, \quad \text{s.t. } \varepsilon_j \leq \epsilon, \quad \forall j \in \{1, 2, \dots, N\}. \quad (2)$$

However, a subtask is valid only when sensing resources exist in the corresponding spatiotemporal cell. As the validity of a subtask is not known until its sensing cycle has finished, it is necessary to balance the sensing cost and the probable validity of subtasks,

$$\begin{aligned} \min \quad & \sum_{i=1}^M \sum_{j=1}^N \frac{B_{ij}}{R_{ij}}, \\ \text{s.t. } \quad & \varepsilon_j \leq \epsilon, \quad \forall j \in \{1, 2, \dots, N\}, \end{aligned} \quad (3)$$

where  $B_{ij}$  is the indicator function which is equal to 1 if  $(i, j) \in \mathcal{S}$  and is 0 otherwise,  $R_{ij}$  is the probability that sensing resources exist in the spatiotemporal cell  $(i, j)$ . In other words, the subtasks tend to be distributed in the areas that are informative and are likely to have sensing resources.

#### IV. COST-EFFECTIVE URBAN ENVIRONMENT SENSING

In this section, our cost-effective urban environment sensing solution is elaborated. First of all, the spatiotemporal correlations among environmental data are exploited to infer the missing data. Then, taking into account the distribution of sensing resources, we design a method to assign subtasks to appropriate spatiotemporal cells in a batch manner.

### A. Latent Features Revelation and Missing Data Inference

The time-varying environmental status of an area is deeply influenced by a set of features. If these features are known, the spatiotemporal correlations among environmental data can be found out. The features are specific to the domain of environmental data. For instance, the air quality of an area is likely to be influenced by the land use and the function of the area (e.g., residential or commercial areas, parks) as well as the traffic patterns in the area [28]. In addition, the traffic conditions in an area are probably influenced by the function of the area, the road network structure and the road surface conditions in the area. Generally, these features are not obvious and it is difficult to know how these features affect the environmental status. However, once the latent features and their influence are revealed, they can be used to infer the missing data in environmental matrix.

Matrix factorization is a commonly used technique to reveal the latent features and to infer the missing data in the environmental matrix. The environmental matrix  $\mathbf{X}$  can be approximately factorized into the product of an area feature matrix  $\mathbf{U}$  and a sensing cycle feature matrix  $\mathbf{V}$ ,

$$\mathbf{X} \approx \mathbf{U}^T \mathbf{V}, \quad (4)$$

where  $\mathbf{U} \in \mathbb{R}^{D \times M}$ ,  $\mathbf{V} \in \mathbb{R}^{D \times N}$ , and  $D \ll \min\{M, N\}$ . Both areas and sensing cycles are mapped to a joint latent feature space of dimensionality  $D$ . Here  $D$  is just the upper bound instead of the exact number of features. A column of  $\mathbf{U}$  measures the extent to which the area possesses those features. A column of  $\mathbf{V}$  measures the extent of influence the corresponding features have on the environmental status. The matrix factorization can automatically reveal the latent features and their respective influence on the environmental status. Particularly, it's not necessary to explain the meanings of the features. In this paper, a probabilistic matrix factorization (PMF) model [38] is applied on the environmental matrix  $\mathbf{X}$ ,

$$p(\mathbf{X}|\mathbf{U}, \mathbf{V}, \sigma^2) = \prod_{i=1}^M \prod_{j=1}^N \left[ \mathcal{N}(\mathbf{X}_{ij} | \mathbf{U}_i^T \mathbf{V}_j, \sigma^2) \right]^{B_{ij}}, \quad (5)$$

$$p(\mathbf{U}|\alpha^2) = \prod_{i=1}^M \mathcal{N}(\mathbf{U}_i | \mathbf{0}, \alpha^2 \mathbf{I}), \quad (6)$$

$$p(\mathbf{V}|\beta^2) = \prod_{j=1}^N \mathcal{N}(\mathbf{V}_j | \mathbf{0}, \beta^2 \mathbf{I}), \quad (7)$$

where  $\mathcal{N}(x|\mu, \sigma^2)$  denotes the Gaussian distribution with mean  $\mu$  and variance  $\sigma^2$ ,  $\mathbf{0}$  is a  $D \times 1$  zero vector and  $\mathbf{I}$  is a  $D \times D$  order identity matrix. The PMF model can be efficiently trained by using steepest descent to find point estimates of model parameters and hyperparameters. Taking advantage of the latent features, we can approximate the unsensed entry  $(i, j) \in \mathcal{T} \setminus \mathcal{S}$  by  $\mathbf{X}_{ij} \approx \mathbf{U}_i^T \mathbf{V}_j$ .

The latent features reflect the spatiotemporal correlations among environmental data. In this paper, the spatiotemporal correlations in a single data source are considered. Specifically, the areas with similar features tend to have similar changing patterns of environmental status; and the sensing cycles with

similar features usually share similar status distribution across the sensing areas. Furthermore, the correlation between two spatiotemporal cells can be described using the latent features. According to Equation (5), it is reasonable to consider that the environmental data  $\mathbf{X}_{ij}$  is drawn from a Gaussian process (GP), where the area feature vector  $\mathbf{U}_i$  and the sensing cycle feature vector  $\mathbf{V}_j$  are associated covariates of  $\mathbf{X}_{ij}$ . Therefore,  $\mathbf{X}_{ij}$  and  $\mathbf{X}_{i'j'}$  are correlated only when  $\mathbf{U}_i$  and  $\mathbf{V}_j$  are similar to  $\mathbf{U}_{i'}$  and  $\mathbf{V}_{j'}$ , respectively. The covariance function of this GP can be defined by a Mercer kernel  $\mathcal{K}(\cdot, \cdot)$ , and the standard radial basis function (RBF) kernel is used,

$$\begin{aligned} \mathcal{K}(\mathbf{X}_{ij}, \mathbf{X}_{i'j'}) \\ = \exp\left(-\frac{\|\mathbf{U}_i - \mathbf{U}_{i'}\|_2^2}{2\eta_1^2}\right) \cdot \exp\left(-\frac{\|\mathbf{V}_j - \mathbf{V}_{j'}\|_2^2}{2\eta_2^2}\right), \end{aligned} \quad (8)$$

where  $\eta_1$  and  $\eta_2$  are the bandwidth parameters of RBF, which are set to be the median of pairwise distance between data points (i.e., median trick).  $\mathcal{K}(\cdot, \cdot)$  is regarded as a measure of the correlation between the environmental data in two spatiotemporal cells.

### B. Informative Sensing Areas Selection

With the spatiotemporal correlations among environmental data, we can decompose the sensing task only to a small subset of areas to reduce sensing cost while keeping required sensing accuracy. As for the unsensed areas, their environmental data can be approximated using the matrix factorization based on the latent features. A straightforward strategy is to gather sensing data from a random set of areas. Obviously, this strategy supposes that the sensing data from different areas contribute equally to the overall sensing accuracy. However, the sensing areas may have a large impact on the sensing accuracy. Intuitively, some particular sensing areas are more *informative/representative* and the informativeness degree of an area is time-varying. Constrained by the required sensing accuracy, we can minimize sensing cost by selecting the most informative set of sensing areas in each sensing cycle.

The matrix factorization can be considered a process of reducing the amount of uncertainty in the unsensed data. The entropy in information theory is widely adopted to quantify the uncertainty. Here the marginal entropy  $H(\mathbf{X}_{\mathcal{T} \setminus \mathcal{S}})$  represents the amount of uncertainty in the unsensed data, while the conditional entropy  $H(\mathbf{X}_{\mathcal{T} \setminus \mathcal{S}} | \mathbf{X}_{\mathcal{S}})$  measures the amount of uncertainty remaining in the unsensed data after knowing the sensed data. Then the mutual information  $I(\mathbf{X}_{\mathcal{S}}; \mathbf{X}_{\mathcal{T} \setminus \mathcal{S}}) = H(\mathbf{X}_{\mathcal{T} \setminus \mathcal{S}}) - H(\mathbf{X}_{\mathcal{T} \setminus \mathcal{S}} | \mathbf{X}_{\mathcal{S}})$  is the amount of uncertainty in the unsensed data which is removed by knowing the sensed data. Therefore, the most informative sensing areas can most significantly reduce the uncertainty about the unsensed areas, i.e., maximizing  $I(\mathbf{X}_{\mathcal{S}}; \mathbf{X}_{\mathcal{T} \setminus \mathcal{S}})$ . In other words, the correlation between the environmental data in  $\mathbf{X}_{\mathcal{S}}$  and  $\mathbf{X}_{\mathcal{T} \setminus \mathcal{S}}$  is maximized, while the intra-correlation among the environmental data in  $\mathbf{X}_{\mathcal{S}}$  is minimized.

As data gathering is a cycle-by-cycle process, only the uncertainty in unsensed data of current sensing cycle is concerned. Denoting the set of spatiotemporal cells in sensing cycle  $j$  as  $\mathcal{T}_j$  and the set of sensed cells in sensing

cycle  $j$  as  $\mathcal{S}_j$ ,  $\mathcal{S}_j \subseteq \mathcal{T}_j$ . When a new sensing cycle begins, the  $k$  most informative cells  $\mathcal{S}_N^*$  are selected from  $\mathcal{T}_N$ ,

$$\mathcal{S}_N^* = \arg \max_{\substack{\mathcal{V} \subseteq \mathcal{T}_N \\ |\mathcal{V}|=k}} I(\mathbf{X}_{\mathcal{S} \cup \mathcal{V}}; \mathbf{X}_{\mathcal{T}_N \setminus \mathcal{V}}), \quad (9)$$

where  $I(\mathbf{X}_{\mathcal{S} \cup \mathcal{V}}; \mathbf{X}_{\mathcal{T}_N \setminus \mathcal{V}})$  is the mutual information between the environmental data in cells  $\mathcal{S} \cup \mathcal{V}$  and cells  $\mathcal{T}_N \setminus \mathcal{V}$ , and  $\mathcal{V}$  is the possible sensing areas to be selected. Solving this optimization problem is NP-complete, but it can be proved that the mutual information is submodular [39], [40]. That is, for all  $\mathcal{S}' \subseteq \mathcal{S} \subseteq \mathcal{T}$ ,  $\mathcal{S}'_N \subseteq \mathcal{S}_N \subseteq \mathcal{T}_N$ , and  $r \in \mathcal{T}_N \setminus \mathcal{S}_N$ , it holds that  $I(\mathbf{X}_{\mathcal{S}' \cup \{r\}}; \mathbf{X}_{\tilde{\mathcal{S}}'_N}) - I(\mathbf{X}_{\mathcal{S}'}; \mathbf{X}_{\mathcal{T}_N \setminus \mathcal{S}'_N}) \geq I(\mathbf{X}_{\mathcal{S} \cup \{r\}}; \mathbf{X}_{\tilde{\mathcal{S}}_N}) - I(\mathbf{X}_{\mathcal{S}}; \mathbf{X}_{\mathcal{T}_N \setminus \mathcal{S}_N})$ , where  $\tilde{\mathcal{S}}_N = \mathcal{T}_N \setminus (\mathcal{S}_N \cup \{r\})$  and  $\tilde{\mathcal{S}}'_N = \mathcal{T}_N \setminus (\mathcal{S}'_N \cup \{r\})$ . It implies that adding  $r$  to a small cell set  $\mathcal{S}'$  helps more than adding  $r$  to a large cell set  $\mathcal{S}$ .

With the submodular property, an approximation algorithm can be used to greedily select the candidate sensing areas. A candidate  $r \in \mathcal{T}_N \setminus \mathcal{S}_N$  is found by maximizing the gain of the mutual information,

$$\begin{aligned} & I(\mathbf{X}_{\mathcal{S} \cup \{r\}}; \mathbf{X}_{\tilde{\mathcal{S}}_N}) - I(\mathbf{X}_{\mathcal{S}}; \mathbf{X}_{\mathcal{T}_N \setminus \mathcal{S}_N}) \\ &= [H(\mathbf{X}_{\tilde{\mathcal{S}}_N}) - H(\mathbf{X}_{\tilde{\mathcal{S}}_N} | \mathbf{X}_{\mathcal{S} \cup \{r\}})] \\ &\quad - [H(\mathbf{X}_{\tilde{\mathcal{S}}_N \cup \{r\}}) - H(\mathbf{X}_{\tilde{\mathcal{S}}_N \cup \{r\}} | \mathbf{X}_{\mathcal{S}})] \\ &= [H(\mathbf{X}_{\tilde{\mathcal{S}}_N}) - H(\mathbf{X}_{\mathcal{S} \cup \tilde{\mathcal{S}}_N \cup \{r\}}) + H(\mathbf{X}_{\mathcal{S} \cup \{r\}})] \\ &\quad - [H(\mathbf{X}_{\tilde{\mathcal{S}}_N \cup \{r\}}) - H(\mathbf{X}_{\mathcal{S} \cup \tilde{\mathcal{S}}_N \cup \{r\}}) + H(\mathbf{X}_{\mathcal{S}})] \\ &= H(\mathbf{X}_r | \mathbf{X}_{\mathcal{S}}) - H(\mathbf{X}_r | \mathbf{X}_{\tilde{\mathcal{S}}_N}). \end{aligned} \quad (10)$$

As the environmental data  $\mathbf{X}_{ij}$  is drawn from a GP,  $p(\mathbf{X}_r | \mathbf{X}_{\mathcal{S}})$  is a Gaussian distribution whose conditional variance  $\sigma_{r|\mathcal{S}}^2$  is given by

$$\sigma_{r|\mathcal{S}}^2 = \mathcal{K}(\mathbf{X}_r, \mathbf{X}_r) - \mathbf{C}_{r\mathcal{S}} \mathbf{C}_{\mathcal{S}\mathcal{S}}^{-1} \mathbf{C}_{\mathcal{S}r}, \quad (11)$$

where  $\mathbf{C}_{\mathcal{S}\mathcal{S}}$  is the covariance matrix whose entry for  $s, s' \in \mathcal{S}$  is  $\mathcal{K}(\mathbf{X}_s, \mathbf{X}_{s'})$ ,  $\mathbf{C}_{r\mathcal{S}}$  is the covariance vector whose entry for  $s \in \mathcal{S}$  is  $\mathcal{K}(\mathbf{X}_r, \mathbf{X}_s)$ , and  $\mathbf{C}_{\mathcal{S}r} = \mathbf{C}_{r\mathcal{S}}^T$ . Then the conditional entropy of  $\mathbf{X}_r$  given  $\mathbf{X}_{\mathcal{S}}$  is calculated as

$$\begin{aligned} H(\mathbf{X}_r | \mathbf{X}_{\mathcal{S}}) &= \frac{1}{2} \log(2\pi e \sigma_{r|\mathcal{S}}^2) \\ &= \frac{1}{2} \log \sigma_{r|\mathcal{S}}^2 + \frac{1}{2} (\log(2\pi) + 1), \end{aligned} \quad (12)$$

and  $H(\mathbf{X}_r | \mathbf{X}_{\tilde{\mathcal{S}}_N})$  can be calculated similarly. Thanks to the latent features and the GP properties for the data in  $\mathbf{X}$ , we can calculate the mutual information without knowing the exact data of the unsensed cells. To maximize the gain of the mutual information in Equation (10), the next most informative sensing area  $r$  can be selected by computing

$$\max_{r \in \mathcal{T}_N \setminus \mathcal{S}_N} \delta_r = \frac{\sigma_{r|\mathcal{S}}^2}{\sigma_{r|\tilde{\mathcal{S}}_N}^2} = \frac{\mathcal{K}(\mathbf{X}_r, \mathbf{X}_r) - \mathbf{C}_{r\mathcal{S}} \mathbf{C}_{\mathcal{S}\mathcal{S}}^{-1} \mathbf{C}_{\mathcal{S}r}}{\mathcal{K}(\mathbf{X}_r, \mathbf{X}_r) - \mathbf{C}_{r\tilde{\mathcal{S}}_N} \mathbf{C}_{\tilde{\mathcal{S}}_N \tilde{\mathcal{S}}_N}^{-1} \mathbf{C}_{\tilde{\mathcal{S}}_N r}}. \quad (13)$$

It can be seen that the latent features are fully utilized to calculate the conditional entropy and to select the most informative sensing areas. The candidate cell selected using Equation (13) is strongly correlated with the unsensed cells and weakly correlated with the already sensed cells.

By iteratively computing Equation (13), the next  $k$  most informative areas in a sensing cycle can be found, as shown

in Algorithm 1. First, the algorithm is initialized and PMF is used to reveal the latent features in  $\mathbf{X}$  (line 1, 2). Next, the whole covariance matrix  $\mathbf{C}_{\mathcal{A}\mathcal{A}}$  is computed (line 3). Then, the algorithm iteratively selects the unsensed area  $r$  with the largest  $\delta_r$  into the set of candidate sensing areas for current sensing cycle (line 4-10). Specifically, the covariance matrices for computing  $\delta_r$  in line 6 are directly extracted from  $\mathbf{C}_{\mathcal{A}\mathcal{A}}$ . Finally, the algorithm terminates when  $k$  areas have been added to the set. It is worth noting that the algorithm is able to compare the informativeness of sensing areas without really knowing the sensing data. Therefore, the most informative sensing areas are selected in a batch manner, which is practical and applicable in the highly dynamic vehicular sensing scenario.

---

#### Algorithm 1 Selecting the Most Informative Sensing Areas

---

##### Input:

$\mathbf{X} \in \mathbb{R}^{M \times N}$ : environmental matrix.

$k$ : the required number of sensing areas.

##### Output:

$\mathcal{S}_N^* \subseteq \mathcal{T}_N \setminus \mathcal{S}_N$ : the next  $k$  most informative sensing areas in sensing cycle  $N$ .

- 1:  $\mathcal{S}_N^* \leftarrow \emptyset$
  - 2: Use PMF to factorize  $\mathbf{X}$  into  $\mathbf{U}$  and  $\mathbf{V}$ .
  - 3: Compute  $\mathbf{C}_{\mathcal{A}\mathcal{A}}$ , where  $\mathcal{A} = \mathcal{S} \cup (\mathcal{T}_N \setminus \mathcal{S}_N)$ .
  - 4: **while**  $|\mathcal{S}_N^*| < k$  **do**
  - 5:   **for all**  $r \in \mathcal{T}_N \setminus (\mathcal{S}_N \cup \mathcal{S}_N^*)$  **do**
  - 6:      $\delta_r = \frac{\mathcal{K}(\mathbf{X}_r, \mathbf{X}_r) - \mathbf{C}_{r\mathcal{B}} \mathbf{C}_{\mathcal{B}\mathcal{B}}^{-1} \mathbf{C}_{\mathcal{B}r}}{\mathcal{K}(\mathbf{X}_r, \mathbf{X}_r) - \mathbf{C}_{r\tilde{\mathcal{B}}_N} \mathbf{C}_{\tilde{\mathcal{B}}_N \tilde{\mathcal{B}}_N}^{-1} \mathbf{C}_{\tilde{\mathcal{B}}_N r}}$ ,  
       where  $\mathcal{B} = \mathcal{S} \cup \mathcal{S}_N^*$  and  $\tilde{\mathcal{B}}_N = \mathcal{T}_N \setminus (\mathcal{S}_N \cup \mathcal{S}_N^* \cup \{r\})$ .
  - 7:   **end for**
  - 8:    $r^* \leftarrow \arg \max_{r \in \mathcal{T}_N \setminus (\mathcal{S}_N \cup \mathcal{S}_N^*)} \delta_r$
  - 9:    $\mathcal{S}_N^* \leftarrow \mathcal{S}_N^* \cup \{r^*\}$
  - 10: **end while**
- 

The complexity of Algorithm 1 is analyzed. Let  $|\mathcal{S}| = m$ , it takes  $O(m)$  operations to solve the PMF model using steepest descent. The time complexity for computing the whole covariance matrix  $\mathbf{C}_{\mathcal{A}\mathcal{A}}$  is  $O(m^2)$ . Computing each  $\delta_r$  requires  $O(m^2)$  operations, but it only requires  $O(m)$  operations to update each  $\delta_r$ . Since only one cell is added to  $\mathcal{B}$  and removed from  $\tilde{\mathcal{B}}_N$ ,  $\delta_r$  can be updated by only re-calculating the changed component. For the  $k$  iterations, the first iteration requires  $O(Mm^2)$  operations, and the remaining  $k - 1$  iterations requires  $O(kMm)$  operations. Therefore, the overall complexity of Algorithm 1 is  $O(Mm^2)$ , which is acceptable for the sensing center with powerful computing resources and can be accelerated by parallel computing for matrix multiplication.

#### C. Adapt Subtasks to Sensing Resource Distribution

As described above, the trajectories of sensor-equipped vehicles are hardly controllable. If the sensing center just waits for vehicles to travel through and sense at the areas selected using Algorithm 1, some of the areas may never gather sensing data as planned. Without an adequate number of sensing data, the required sensing accuracy cannot be satisfied.

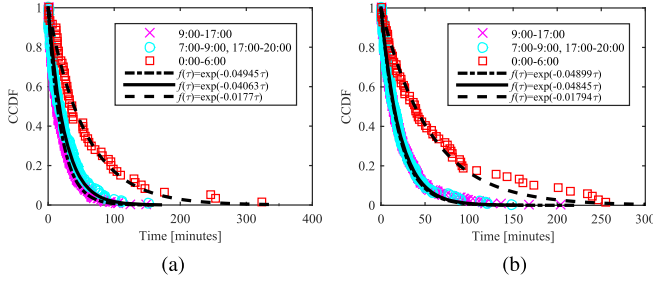


Fig. 5. The CCDF of inter-arrival time of 50 sensor-equipped vehicles in Beijing during peak hours (7:00-9:00 and 17:00-20:00), mid-peak hours (9:00-17:00), and off-peak hours (0:00-6:00), respectively. (a) Commercial area (Guomao). (b) Residential area (Anzhen).

Therefore, the validity of subtasks should be integrated into the process of selecting the candidate sensing areas.

If the distribution patterns of sensing resources can be discovered, the validity of subtasks can be well estimated. For a specific area, the arrival process of vehicles is considered to follow a Poisson process [41]–[43]. Under the assumption that sensor-equipped vehicles are uniformly distributed among vehicles on the road, the arrival of sensor-equipped vehicles can be regarded as a thinned Poisson process, and we have validated it by analyzing the real trajectory data in Beijing. In Fig. 5, the inter-arrival time distribution of sensor-equipped vehicles is shown and compared with the exponential distribution. For different areas (i.e., commercial and residential area) and different periods of a day (i.e., peak hours, mid-peak hours, and off-peak hours), the Poisson arrival property for the sensor-equipped vehicles is reasonable. Therefore, this property can be utilized to select sensing areas which compromise between cost-effectiveness and subtask feasibility.

Denoting the current arrival rates of sensor-equipped vehicles in all  $M$  areas as  $\mathcal{L} = \{\lambda_1, \lambda_2, \dots, \lambda_M\}$ . The probability of sensor-equipped vehicles being available in area  $i$  for a time period  $\tau$  can be computed as

$$R_i(\tau) = 1 - e^{-\lambda_i \tau}. \quad (14)$$

In addition, as shown in Fig. 5, the arrival rate in an area varies with time. For example, during peak and mid-peak hours, there will be more sensor-equipped vehicles coming into the sensing areas. Therefore, the arrival rate for each area should be updated periodically. An exponentially weighted moving average is used to update  $\lambda_i$ ,

$$\lambda_i \leftarrow \omega \cdot v_i + (1 - \omega) \cdot \lambda_i, \quad (15)$$

where  $v_i$  is the arrival rate in area  $i$  during the most recent observing window, and  $\omega$  is the weight of  $v_i$ . The length of the observing window is set to 1 hour, and  $\omega$  is set to 0.6.

Since it is unknown in advance whether the candidate areas can be sensed or not, the optimization problem in Equation (9) is not applicable. Instead, by considering the distribution of sensing resources, the expectation of the mutual information is used,

$$\begin{aligned} \mathcal{S}_N^* &= \arg \max_{\substack{\mathcal{V} \subseteq \mathcal{T}_N \\ |\mathcal{V}|=k}} \mathbb{E}[I(\mathbf{X}_{\mathcal{S} \cup \mathcal{V}}; \mathbf{X}_{\mathcal{T}_N \setminus \mathcal{V}})] \\ &= \arg \max_{\substack{\mathcal{V} \subseteq \mathcal{T}_N \\ |\mathcal{V}|=k}} \sum_{\mathcal{W} \subseteq \mathcal{V}} \Pr(\mathcal{W}) \cdot I(\mathbf{X}_{\mathcal{S} \cup \mathcal{W}}; \mathbf{X}_{\mathcal{T}_N \setminus \mathcal{W}}), \quad (16) \end{aligned}$$

where  $\mathbb{E}[\cdot]$  is the expectation in terms of the probability that sensing resources are available at candidate areas,  $\langle \mathcal{V} \rangle = \sum_{(i,N) \in \mathcal{V}} R_i(t)$  is the expected value for the number of valid subtasks in  $\mathcal{V}$ ,  $\lfloor \cdot \rfloor$  is the floor function,  $k$  is the required number of sensing areas,  $\Pr(\mathcal{W}) = \prod_{(i,N) \in \mathcal{W}} R_i(t) \cdot \prod_{(i',N) \in \mathcal{V} \setminus \mathcal{W}} [1 - R_{i'}(t)]$  is the probability of successfully gathering data from  $\mathcal{W}$  and not from  $\mathcal{V} \setminus \mathcal{W}$ , and  $t$  is the remaining time in current sensing cycle. Then, the area which has the largest expected gain of the mutual information can be selected greedily,

$$\begin{aligned} &\mathbb{E}[I(\mathbf{X}_{\mathcal{S} \cup \{r\}}; \mathbf{X}_{\mathcal{S}_N}) - I(\mathbf{X}_{\mathcal{S}}; \mathbf{X}_{\mathcal{T}_N \setminus \mathcal{S}_N})] \\ &= R_i(t) \cdot [H(\mathbf{X}_r | \mathbf{X}_{\mathcal{S}}) - H(\mathbf{X}_r | \mathbf{X}_{\mathcal{S}_N})] \\ &= R_i(t) \cdot \frac{1}{2} \log \frac{\sigma_{r|\mathcal{S}}^2}{\sigma_{r|\mathcal{S}_N}^2} \\ &= \frac{1}{2} \log \delta_r^{R_i(t)}, \quad (17) \end{aligned}$$

where  $r = (i, N) \in \mathcal{T}_N \setminus \mathcal{S}_N$ . As it is too costly to consider the uncertainty from the previous selected candidate areas, the areas in  $\mathcal{S}$  are deemed to be determinate, and the checkpoint mechanism in the next subsection will make this simplification acceptable by adjusting the candidate areas at suitable time. Therefore, the optimization problem in Equation (13) is evolved to

$$\max_{r=(i,N) \in \mathcal{T}_N \setminus \mathcal{S}_N} \delta_r^+ = \delta_r^{R_i(t)}. \quad (18)$$

Then, Algorithm 2 is proposed for selecting candidate sensing areas which considers the distribution of sensing resources.

---

#### Algorithm 2 Selecting the Candidate Sensing Areas

---

##### Input:

- $\mathbf{X} \in \mathbb{R}^{M \times N}$ : environmental matrix.
- $k$ : the required number of sensing areas.
- $\mathcal{L}$ : sensing resources' arrival rates.
- $t$ : the remaining time in the sensing cycle.

##### Output:

- $\mathcal{S}_N^* \subseteq \mathcal{T}_N \setminus \mathcal{S}_N$ : the candidate sensing areas in sensing cycle  $N$ , where the expected value for the number of valid subtasks is  $k$ .
  - 1:  $\mathcal{S}_N^* \leftarrow \emptyset$
  - 2: Use PMF to factorize  $\mathbf{X}$  into  $\mathbf{U}$  and  $\mathbf{V}$ .
  - 3: Compute  $\mathbf{C}_{\mathcal{A}\mathcal{A}}$ , where  $\mathcal{A} = \mathcal{S} \cup (\mathcal{T}_N \setminus \mathcal{S}_N)$ .
  - 4: **while**  $\langle \mathcal{S}_N^* \rangle < k$  **do**
  - 5:   **for all**  $r = (i, N) \in \mathcal{T}_N \setminus (\mathcal{S}_N \cup \mathcal{S}_N^*)$  **do**
  - 6:     
$$\delta_r^+ = \left( \frac{\mathcal{K}(\mathbf{X}_r, \mathbf{X}_r) - \mathbf{C}_{r\mathcal{B}} \mathbf{C}_{\mathcal{B}\mathcal{B}}^{-1} \mathbf{C}_{\mathcal{B}r}}{\mathcal{K}(\mathbf{X}_r, \mathbf{X}_r) - \mathbf{C}_{r\mathcal{B}_N} \mathbf{C}_{\mathcal{B}_N \mathcal{B}_N}^{-1} \mathbf{C}_{\mathcal{B}_N r}} \right)^{R_i(t)},$$
    - where  $\mathcal{B} = \mathcal{S} \cup \mathcal{S}_N^*$  and  $\mathcal{B}_N = \mathcal{T}_N \setminus (\mathcal{S}_N \cup \mathcal{S}_N^* \cup \{r\})$ .
  - 7:   **end for**
  - 8:    $r^* \leftarrow \arg \max_{r \in \mathcal{T}_N \setminus (\mathcal{S}_N \cup \mathcal{S}_N^*)} \delta_r^+$
  - 9:    $\mathcal{S}_N^* \leftarrow \mathcal{S}_N^* \cup \{r^*\}$
  - 10: **end while**
- 

Here the mutual information gain  $\delta_r$  has been substituted with the expected mutual information gain  $\delta_r^+$  (line 6). In addition, the algorithm terminates when the expected value for the number of valid subtasks (i.e.,  $\langle \mathcal{S}_N^* \rangle$ ) becomes greater



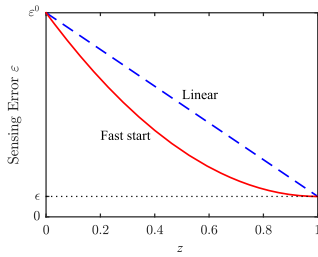


Fig. 6. Sensing error objective curve.

than or equal to  $k$ . Obviously, the overall time complexity is also  $O(Mm^2)$ . The candidate areas selected using Algorithm 2 are not the most informative ones, but the areas that are balanced between the informativeness and feasibility.

#### D. Determine the Number of Subtasks

Given a missing data inference method and a sensing areas selection algorithm, the number of gathered sensing data directly affects the sensing accuracy. To facilitate the sensing areas selection, the required number of subtasks should be determined at the very beginning of each sensing cycle. However, it is impossible to exactly determine the number when the environmental data in the new sensing cycle are almost unknown. Therefore, a rough number is used at beginning and the number is adjusted during the data gathering procedure. As the environmental data will not change severely in successive sensing cycles, the number of subtasks in previous sensing cycle is a good estimation of that in current sensing cycle.

In addition, after gathering some data, the feature of current sensing cycle becomes more clear. Thus the calculated most informative sensing areas may need some revision. The sensing center should supervise the data gathering process and adjust the sensing strategy at the suitable time. The current sensing error is the major indicator to determine whether it is necessary to adjust the sensing strategy. Therefore, an objective curve is defined and used to supervise the data gathering progress in a real-time way. Let  $0 \leq z \leq 1$  be the fraction of time elapsed in a sensing cycle, then the objective curve  $g(z)$  indicates the upper bound of sensing error in each timestamp. When the current sensing error becomes larger than its upper bound (i.e.,  $\epsilon > g(z)$ ), the sensing center should increase the number of subtasks and adjust candidate sensing areas. According to the expected rate of data gathering, the objective curve can be linear or fast start, as shown in Fig. 6. The linear objective curve is defined as  $g(z) = (\epsilon^0 - \epsilon)(1 - z) + \epsilon$ , which requires that the sensing error decreases to  $\epsilon$  steadily. The fast start objective curve is quadratic and defined as  $g(z) = (\epsilon^0 - \epsilon)(1 - z)^2 + \epsilon$ , where the sensing error decreases faster at beginning.  $\epsilon^0$  is the initial sensing error in current sensing cycle. In addition, several checkpoints are set and the sensing progress is checked only at these checkpoints to improve the solution efficiency.

However, in practice the ground-truth data of the area without sensing data is not known, that is, Equation (1) cannot be obtained to evaluate the sensing error. Therefore, the error is evaluated approximately using original sensing data and their

corresponding inference results,

$$\hat{\epsilon}_j = \sqrt{\frac{1}{|\mathcal{S}_j|} \sum_{(i,j) \in \mathcal{S}_j} (\mathbf{X}_{ij} - \mathbf{U}_i^T \mathbf{V}_j)^2}. \quad (19)$$

Some empirical study results will be used in Section V to show that  $\hat{\epsilon}$  is a good approximation of  $\epsilon$ . Additionally, when  $l$  successive data gatherings satisfy  $\hat{\epsilon} \leq \epsilon$  (i.e., stopping condition), the sensing results are considered to meet the requirement of sensing accuracy.

The cost-effective urban environment sensing solution is summarized in Algorithm 3. First, the algorithm sets the required number of subtasks according to the previous sensing cycle (line 1). Next,  $k^0$  initial sensing data are gathered to bootstrap the data gathering process (line 2). These initial sensing data can solve the cold start problem for matrix factorization. Since the feature of current sensing cycle is not known, a heuristic strategy called most-hungry-area-first is used to gather data from the areas with sensing resources. Specifically, the most hungry area is the area that has not acquired sensing data for the longest time. As a side benefit of the most-hungry-area-first strategy, the probability of missing an entire row of  $\mathbf{X}$  is reduced. Then, the algorithm selects the initial candidate sensing areas and sets the objective curve (line 3, 4). After that, the sensing center gathers sensing data from the candidate sensing areas  $\mathcal{S}_N^*$  until meeting the stopping condition (line 5-13). Additionally, in periodical checkpoints, when data gathering progress becomes slower than expectation (i.e.,  $\hat{\epsilon}_N > g(z)$ ), extra subtasks should be added to  $\mathcal{S}_N^*$  to boost the data gathering. The number of extra subtasks can be calculated as  $\Delta k = v \cdot \frac{\hat{\epsilon}_N - g(z)}{\Delta \hat{\epsilon}_N}$ , where  $v$  is the number of data gathered between the latest two checkpoints, and  $\Delta \hat{\epsilon}_N$  is the reduction of the sensing error between the latest two checkpoints. Finally, the unsensed data are inferred and the sensing resources' arrival rates are updated (line 14, 15).

## V. PERFORMANCE EVALUATION

Extensive experiments have been performed for evaluating the performance of our proposed cost-effective urban environment sensing solution. In the following, the experimental setup is presented, and then the compared methods are introduced. Finally, performance results are presented and discussed.

#### A. Experimental Setup

The experiments are driven by real taxicab mobility traces and air quality data in Beijing. The taxicab traces contain the GPS trajectory recorded by over 12,000 taxicabs in November of 2012. The air quality data are obtained from the U-Air project [28], [44] undertaken by Microsoft Research in March of 2015. These two datasets are combined by a time-shifted and space-aligned mapping. Moreover, some taxicabs are considered to have the ability to sense air quality. Specifically, the concentrations of PM10 and NO2 are separately measured in two sensing tasks and used in two sets of experiments. As shown in Fig. 7, the sensing space is a part within the 3rd Ring Road of Beijing, and stretches



**Algorithm 3** Cost-Effective Urban Environment Sensing**Input:**

$\mathbf{X} \in \mathbb{R}^{M \times N}$ : environmental matrix.  
 $\mathcal{L}$ : sensing resources' arrival rates.

**Output:**

Complete and satisfactory sensing results for the  $N$ th column of  $\mathbf{X}$ .

- 1: Set the expected number of subtasks  $k$  as that in previous sensing cycle.
- 2: According to the current distribution of sensing resources, gather  $k^0$  initial sensing data using the most-hungry-area-first strategy.
- 3: Use Algorithm 2 to select candidate sensing areas  $\mathcal{S}_N^*$ .
- 4: Use Equation (19) to estimate the initial sensing error  $\varepsilon_N^0$ , and set the objective curve  $g(z)$ .
- 5: **while** Not meet the stopping condition **do**
- 6:   Wait sensing data at cells  $\mathcal{S}_N^*$ .
- 7:   **if** Checkpoint **then**
- 8:     Calculate current  $\hat{\varepsilon}_N$ .
- 9:     **if**  $\hat{\varepsilon}_N > g(z)$  **then**
- 10:      Add  $\Delta k$  areas, adjust  $\mathcal{S}_N^*$  with Algorithm 2.
- 11:     **end if**
- 12:   **end if**
- 13: **end while**
- 14: For the unsensed area  $(i, j) \in \mathcal{T}_N \setminus \mathcal{S}_N$ ,  $\mathbf{X}_{ij} \approx \mathbf{U}_i^T \mathbf{V}_j$ .
- 15: Update  $\mathcal{L}$  for next sensing cycle using Equation (15).

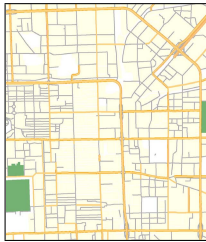


Fig. 7. Sensing space within the 3rd Ring Road of Beijing, which is partitioned into  $7 \times 6$  disjoint grids.

about 7 km from north to south and 6 km from west to east. The area is divided into 42 ( $= 7 \times 6$ ) disjointed grids with each grid about  $1 \text{ km}^2$ . The air quality dataset holds one PM10 and one NO2 record per hour for each grid. A complete subset of the original air quality data for 14 successive days is extract. In addition, 50 or 100 taxicabs are randomly selected as vehicles that are equipped with sensors.

The environmental matrix  $\mathbf{X}$  is of the size  $M = 42$  and  $N = 48$  ( $= 24 \times 2$ ). The data from the first 2 days are used to initialize the experiment; and the data from the remaining 12 days (i.e.,  $12 \times 24 = 288$  sensing cycles) are used to show the performance. The upper bound of sensing error  $\varepsilon$  is set to 10. In Algorithm 3, the stopping condition  $l$  is set to 3, the number of initial sensing data  $k^0$  is set to 5, and the checkpoint is performed every 10 minutes. The parameters in PMF model are set based on a group of parameter test.

We perform a competitive study, comparing our CESense solution with other alternative solutions that will be introduced

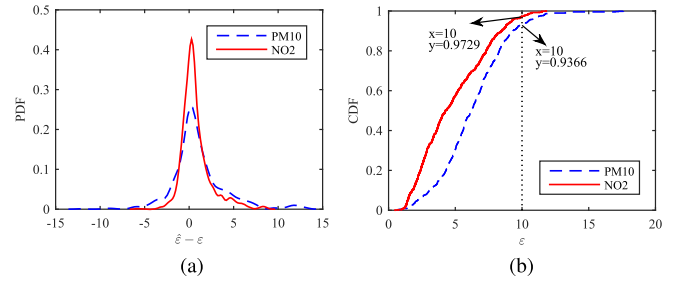


Fig. 8. The correlation between the estimated and actual error. (a) The PDF of  $\hat{\varepsilon} - \varepsilon$ . (b) The CDF of  $\varepsilon$  when  $\hat{\varepsilon}$  meets the requirement (i.e.,  $\hat{\varepsilon} \leq \varepsilon$ ,  $\varepsilon = 10$ ).

in following subsection. As there is no tightly-coupled relationship between the distribution of sensing resources and the environmental status to be sensed, the experiments based on the time-shifted and space-aligned datasets are convincing. Sensing accuracy and efficiency are performance metrics for evaluation and comparison. A better urban environment sensing solution is the one which meets the required sensing accuracy with a fewer number of subtasks.

### B. Compared Methods

Our cost-effective sensing solution is compared with three other methods.

1) *ST-Interp*: The method assigns the same number of subtasks in each sensing cycle. It simply gathers data once an unsensed area is available to be sensed until reaching the predefined number of subtasks. As for the missing data inference, a spatiotemporal KNN algorithm with the inverse distance weights is used. The algorithm assigns a weight to each sensed data of the  $K_{ST}$  nearest cells according to their inverse distance to the target cell. Then the missing data in cell  $r$  is estimated as  $\mathbf{X}_r \approx \frac{\sum_{r'} d_{rr'}^{-1} \mathbf{X}_{r'}}{\sum_{r'} d_{rr'}^{-1}}$ , where  $r'$  is from the  $K_{ST}$  neighbours of  $r$ , and  $d_{rr'}$  is the Euclid distance between  $r$  and  $r'$ .  $K_{ST}$  is set to 20 in following experiments.

2) *SVT*: The method uses the same way with *ST-Interp* to gather data. However, it infers the missing data with a matrix completion algorithm called singular value thresholding (SVT) [45], which minimizes the nuclear norm of the environmental matrix. The parameters in SVT are determined through a group of parameter test.

3) *SiSense*: It is a simplified version of CESense. SiSense is almost the same with CESense except that it does not use Algorithm 2 to select the candidate sensing areas. Instead, SiSense considers all areas with the same informativeness and selects those with largest  $R_i(t)$  as candidate sensing areas.

### C. Experimental Results

First of all, to demonstrate that the estimated sensing error defined in Equation (19) can be a substitute for the actual one defined in Equation (1), their correlation is shown in Fig. 8. The PDF (Probability Density Function) of the differences between  $\hat{\varepsilon}$  and  $\varepsilon$  is shown in Fig. 8a. It can be found that, for both PM10 and NO2 dataset, the differences have a narrow distribution around 0, which implies great consistency between  $\hat{\varepsilon}$  and  $\varepsilon$ . Moreover, the CDF (Cumulative Distribution Function) of  $\varepsilon$  when  $\hat{\varepsilon} \leq \varepsilon$  is shown in Fig. 8b. It can be seen

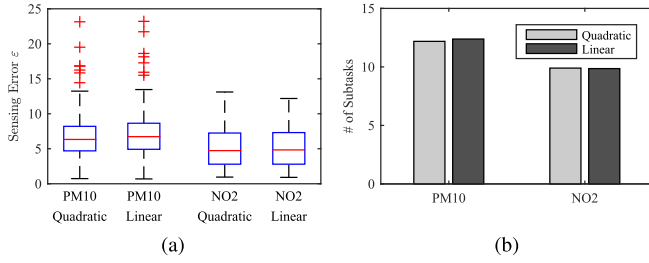


Fig. 9. The performance of CE Sense using different objective curves, where 50 sensor-equipped vehicles are involved. (a) Box plots of sensing errors. (b) Average sensing cost.

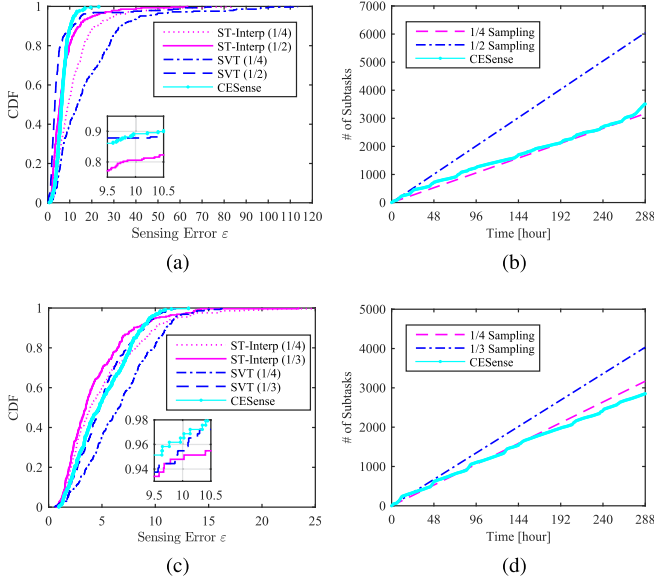


Fig. 10. Performance of different sensing solutions with 50 sensor-equipped vehicles. (a) Sensing error, PM10. (b) Accumulative sensing cost, PM10. (c) Sensing error, NO2. (d) Accumulative sensing cost, NO2.

that more than 93% and 97% of the sensing results meet the requirement for PM10 and NO2 dataset, respectively, when their estimated errors meet the requirement. Therefore, it is reasonable to use  $\hat{\epsilon}$  to supervise the progress of data gathering.

Two kinds of objective curves (i.e., linear and fast start) have been designed to indicate the desired sensing progress. Their influence on the overall performance of CE Sense is shown in Fig. 9. The distribution of sensing errors with these two objective curves are very similar for both PM10 and NO2 dataset (Fig. 9a). Besides, these two objective curves result in almost the same average sensing cost in unit sensing cycle (Fig. 9b). As the type of objective curve has little impact on the performance of CE Sense, only the quadratic curve is used in the following experiments.

In Fig. 10, our CE Sense is compared with ST-Interp and SVT method, where 50 sensor-equipped vehicles are involved. In addition, ST-Interp and SVT are executed with different number of subtasks in each sensing cycle, that is, 1/4, 1/3 and 1/2 of the total number of sensing areas. Fig. 10a and Fig. 10b show the experimental results on PM10 dataset. CE Sense can generate the results that 89% of the sensing cycles meet the required sensing accuracy, the highest of all. Though the sensing cost of CE Sense is just similar to the 1/4 sampling, its sensing accuracy is even better than ST-Interp and SVT

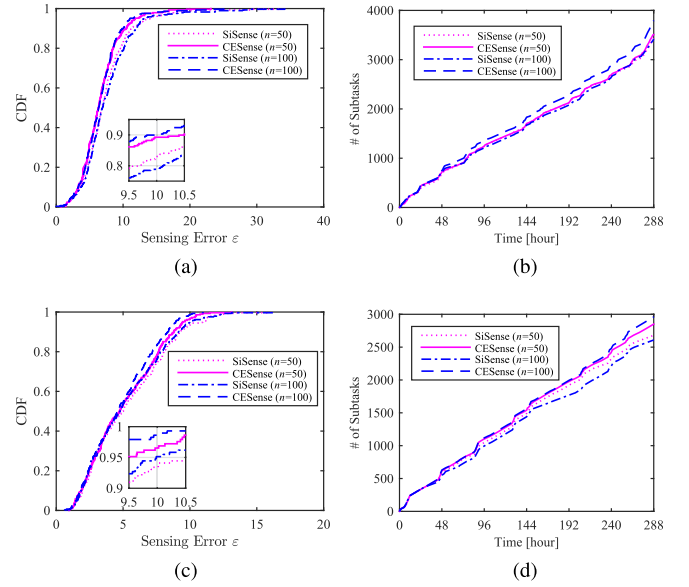


Fig. 11. The performance gain introduced by the informative sensing areas selection algorithm. (a) Sensing error, PM10. (b) Accumulative sensing cost, PM10. (c) Sensing error, NO2. (d) Accumulative sensing cost, NO2.

method with 1/2 sampling. That is, as compared to ST-Interp and SVT, CE Sense can reduce the sensing cost by half. The experimental results on NO2 dataset are presented in Fig. 10c and Fig. 10d. CE Sense also has advantages on NO2 dataset. It provides the results that more than 96% of the sensing cycles meet the requirement of sensing accuracy, the best of all. Though the sensing cost of CE Sense is even lower than 1/4 sampling, its sensing accuracy is better than both ST-Interp and SVT method with 1/3 sampling. In both PM10 and NO2 dataset, it is worth noting that the CDF of sensing errors for CE Sense is partly lower than ST-Interp and SVT method when  $\epsilon < \epsilon$ . This phenomenon demonstrates the high efficiency of CE Sense from another perspective, because CE Sense never excessively pursues sensing accuracy and just seeks to ensure the required accuracy. There are a few sensing cycles that do not meet the sensing requirement. From the previous discussion on Fig. 8b, it can be known that these unsatisfactory results are mainly caused by the gap between the estimated and actual errors. The insufficiency of sensing resources is another reason for the unsatisfactory results.

To clarify the benefits introduced by the informative sensing area selection algorithm, CE Sense is compared with SiSense in Fig. 11. For the PM10 dataset with 50 sensor-equipped vehicles, by assigning subtasks to the informative sensing areas, the proportion of sensing cycles that meet requirement increases from 0.83 to 0.89. When 100 sensor-equipped vehicles are involved, the proportion is increased from 0.79 for SiSense to 0.9 for CE Sense. The sensing costs for CE Sense and SiSense are very similar. Therefore, the informative sensing area selection algorithm plays an important role in CE Sense. Specifically, the sensing accuracy of CE Sense is improved when more sensing resources are available, while the sensing cost does not have a significant change. This happens for two reasons. First, it has the chance to select more informative sensing areas when more sensing resources are available. Moreover, the shortage of sensing resources in

some sensing cycles is relieved. From the experimental results on NO<sub>2</sub> dataset, the similar phenomenon can be seen, and then the same conclusion can be drawn.

## VI. CONCLUSION

In this paper, we have proposed a cost-effective urban environment sensing solution, which is applicable to the data source with spatiotemporal correlations. To reduce sensing cost while considering the dynamic distribution of sensing resources, the proposed solution intelligently decomposes the sensing task into subtasks and assigns the subtasks to the areas of considerable informativeness. We have also designed a checkpoint mechanism to supervise the progress of data gathering. Extensive experimental results, based on real taxicab mobility traces and air quality data in Beijing, have demonstrated that our proposed solution improves sensing quality while keeping low cost.

In the future, we will design a cost-effective urban environment sensing solution with the cooperation among participants. The participants can exchange knowledge using vehicular ad-hoc networks and collaboratively perform sensing tasks assigned by the urban sensing center. This approach should make the data gathering solution more adaptive to the dynamic distribution of sensing resources.

## REFERENCES

- [1] L. Liu, W. Wei, D. Zhao, and H. Ma, "Urban resolution: New metric for measuring the quality of urban sensing," *IEEE Trans. Mobile Comput.*, vol. 14, no. 12, pp. 2560–2575, Dec. 2015.
- [2] B. Guo *et al.*, "Mobile crowd sensing and computing: The review of an emerging human-powered sensing paradigm," *ACM Comput. Surv.*, vol. 48, no. 1, 2015, Art. no. 47.
- [3] X. Cheng, L. Fang, X. Hong, and L. Yang, "Exploiting mobile big data: Sources, features, and applications," *IEEE Netw.*, vol. 31, no. 1, pp. 72–79, Jan./Feb. 2017.
- [4] S. Abdelhamid, H. Hassanein, and G. Takahara, "Vehicle as a resource (VaaR)," *IEEE Netw.*, vol. 29, no. 1, pp. 12–17, Jan. 2015.
- [5] U. Lee and M. Gerla, "A survey of urban vehicular sensing platforms," *Comput. Netw.*, vol. 54, no. 4, pp. 527–544, Mar. 2010.
- [6] B. Hull *et al.*, "CarTel: A distributed mobile sensor computing system," in *Proc. ACM SenSys*, Boulder, CO, USA, Nov. 2006, pp. 125–138.
- [7] U. Lee, B. Zhou, M. Gerla, E. Magistretti, P. Bellavista, and A. Corradi, "Mobeyes: Smart mobs for urban monitoring with a vehicular sensor network," *IEEE Wireless Commun.*, vol. 13, no. 5, pp. 52–57, Oct. 2006.
- [8] M. D. Dikaiakos, A. Florides, T. Nadeem, and L. Iftode, "Location-aware services over vehicular ad-hoc networks using car-to-car communication," *IEEE J. Sel. Areas Commun.*, vol. 25, no. 8, pp. 1590–1602, Oct. 2007.
- [9] J. Eriksson, L. Girod, B. Hull, R. Newton, S. Madden, and H. Balakrishnan, "The pothole patrol: Using a mobile sensor network for road surface monitoring," in *Proc. ACM MobiSys*, Breckenridge, CO, USA, 2008, pp. 29–39.
- [10] Q. Yuan, Z. Liu, J. Li, J. Zhang, and F. Yang, "A traffic congestion detection and information dissemination scheme for urban expressways using vehicular networks," *Transp. Res. C, Emerg. Technol.*, vol. 47, pp. 114–127, Oct. 2014.
- [11] J. H. Ahn and M. Potkonjak, "VeSense: High-performance and energy-efficient vehicular sensing platform," *Pervasive Mobile Comput.*, vol. 12, pp. 112–122, Jun. 2014.
- [12] Q. Yuan, Z. Liu, J. Li, S. Yang, and F. Yang, "An adaptive and compressive data gathering scheme in vehicular sensor networks," in *Proc. IEEE ICPADS*, Melbourne, VIC, Australia, Dec. 2015, pp. 207–215.
- [13] Y.-C. Wang and G.-W. Chen, "Efficient data gathering and estimation for metropolitan air quality monitoring by using vehicular sensor networks," *IEEE Trans. Veh. Technol.*, vol. 66, no. 8, pp. 7234–7248, Aug. 2017.
- [14] K. Abboud, H. A. Omar, and W. Zhuang, "Interworking of DSRC and cellular network technologies for V2X communications: A survey," *IEEE Trans. Veh. Technol.*, vol. 65, no. 12, pp. 9457–9470, Dec. 2016.
- [15] Z. Haibo *et al.*, "ChainCluster: Engineering a cooperative content distribution framework for highway vehicular communications," *IEEE Trans. Intell. Transp. Syst.*, vol. 15, no. 6, pp. 2644–2657, Dec. 2014.
- [16] X. Cheng, L. Yang, and X. Shen, "D2D for intelligent transportation systems: A feasibility study," *IEEE Trans. Intell. Transp. Syst.*, vol. 16, no. 4, pp. 1784–1793, Jan. 2015.
- [17] S. A. Hamid, G. Takahara, and H. S. Hassanein, "On the recruitment of smart vehicles for urban sensing," in *Proc. IEEE GLOBECOM*, Atlanta, GA, USA, Dec. 2013, pp. 36–41.
- [18] S. He, D.-H. Shin, J. Zhang, and J. Chen, "Toward optimal allocation of location dependent tasks in crowdsensing," in *Proc. IEEE INFOCOM*, Toronto, ON, Canada, Apr./May 2014, pp. 745–753.
- [19] Z. He, J. Cao, and X. Liu, "High quality participant recruitment in vehicle-based crowdsourcing using predictable mobility," in *Proc. IEEE INFOCOM*, Hong Kong, Apr./May 2015, pp. 2542–2550.
- [20] S. Ji, Y. Zheng, and T. Li, "Urban sensing based on human mobility," in *Proc. ACM UbiComp*, Heidelberg, Germany, 2016, pp. 1040–1051.
- [21] R. J. A. Little and D. B. Rubin, *Statistical Analysis With Missing Data*, 2nd ed. Hoboken, NJ, USA: Wiley, 2002.
- [22] T. Cover and P. Hart, "Nearest neighbor pattern classification," *IEEE Trans. Inf. Theory*, vol. IT-13, no. 1, pp. 21–27, Jan. 1967.
- [23] E. J. Candès and B. Recht, "Exact matrix completion via convex optimization," *Found. Comput. Math.*, vol. 9, no. 6, pp. 717–772, 2009.
- [24] J. Cheng, Q. Ye, H. Jiang, D. Wang, and C. Wang, "STCDG: An efficient data gathering algorithm based on matrix completion for wireless sensor networks," *IEEE Trans. Wireless Commun.*, vol. 12, no. 2, pp. 850–861, Feb. 2013.
- [25] L. Kong, M. Xia, X.-Y. Liu, M.-Y. Wu, and X. Liu, "Data loss and reconstruction in sensor networks," in *Proc. IEEE INFOCOM*, Turin, Italy, Apr. 2013, pp. 1654–1662.
- [26] Y. Zhu, Z. Li, H. Zhu, M. Li, and Q. Zhang, "A compressive sensing approach to urban traffic estimation with probe vehicles," *IEEE Trans. Mobile Comput.*, vol. 12, no. 11, pp. 2289–2302, Nov. 2013.
- [27] R. Du, C. Chen, B. Yang, N. Lu, X. Guan, and X. Shen, "Effective urban traffic monitoring by vehicular sensor networks," *IEEE Trans. Veh. Technol.*, vol. 64, no. 1, pp. 273–286, Jan. 2015.
- [28] Y. Zheng, F. Liu, and H.-P. Hsieh, "U-Air: When urban air quality inference meets big data," in *Proc. ACM KDD*, Chicago, IL, USA, 2013, pp. 1436–1444.
- [29] X. Yi, Y. Zheng, J. Zhang, and T. Li, *ST-MVL: Filling Missing Values in Geo-Sensory Time Series Data*. New York, NY, USA: IJCAI, 2016, pp. 2704–2710.
- [30] E. J. Candès and M. B. Wakin, "An introduction to compressive sampling," *IEEE Signal Process. Mag.*, vol. 25, no. 2, pp. 21–30, Mar. 2008.
- [31] X. Yu, H. Zhao, L. Zhang, S. Wu, B. Krishnamachari, and V. O. K. Li, "Cooperative sensing and compression in vehicular sensor networks for urban monitoring," in *Proc. IEEE ICC*, Cape Town, South Africa, May 2010, pp. 1–5.
- [32] C. Liu, C. Chigan, and C. Gao, "Compressive sensing based data collection in VANETs," in *Proc. IEEE WCNC*, Shanghai, China, Apr. 2013, pp. 1756–1761.
- [33] H. Wang, Y. Zhu, and Q. Zhang, "Compressive sensing based monitoring with vehicular networks," in *Proc. IEEE INFOCOM*, Turin, Italy, Apr. 2013, pp. 2823–2831.
- [34] L. Xu, X. Hao, N. D. Lane, X. Liu, and T. Moscibroda, "More with less: Lowering user burden in mobile crowdsourcing through compressive sensing," in *Proc. ACM UbiComp*, Osaka, Japan, Sep. 2015, pp. 659–670.
- [35] K. Xie, L. Wang, X. Wang, J. Wen, and G. Xie, "Learning from the past: Intelligent on-line weather monitoring based on matrix completion," in *Proc. IEEE ICDCS*, Madrid, Spain, Jun./Jul. 2014, pp. 176–185.
- [36] L. Wang *et al.*, "CCS-TA: Quality-guaranteed online task allocation in compressive crowdsensing," in *Proc. ACM UbiComp*, Osaka, Japan, Sep. 2015, pp. 683–694.
- [37] C. Meng, H. Xiao, L. Su, and Y. Cheng, "Tackling the redundancy and sparsity in crowd sensing applications," in *Proc. ACM SenSys*, Stanford, CA, USA, 2016, pp. 150–163.
- [38] R. Salakhutdinov and A. Mnih, "Probabilistic matrix factorization," in *Proc. Neural Inf. Process. Syst.*, Vancouver, BC, Canada, 2007, pp. 1257–1264.



- [39] A. Krause, A. Singh, and C. Guestrin, "Near-optimal sensor placements in Gaussian processes: Theory, efficient algorithms and empirical studies," *J. Mach. Learn. Res.*, vol. 9, no. 2, pp. 235–284, 2008.
- [40] J. Silva and L. Carin, "Active learning for online Bayesian matrix factorization," in *Proc. ACM KDD*, Beijing, China, 2012, pp. 325–333.
- [41] N. Wisitpongphan, F. Bai, P. Mudalige, V. Sadekar, and O. Tonguz, "Routing in sparse vehicular ad hoc wireless networks," *IEEE J. Sel. Areas Commun.*, vol. 25, no. 8, pp. 1538–1556, Oct. 2007.
- [42] L. Huang, H. Jiang, Z. Zhang, and Z. Yan, "Optimal traffic scheduling between roadside units in vehicular delay-tolerant networks," *IEEE Trans. Veh. Technol.*, vol. 64, no. 3, pp. 1079–1094, Mar. 2015.
- [43] J. He, L. Cai, P. Cheng, and J. Pan, "Delay minimization for data dissemination in large-scale VANETs with buses and taxis," *IEEE Trans. Mobile Comput.*, vol. 15, no. 8, pp. 1939–1950, Aug. 2016.
- [44] Microsoft Research. (2017). *Urban Air*. [Online]. Available: <http://urbanair.msra.cn/>
- [45] J.-F. Cai, E. J. Candès, and Z. Shen, "A singular value thresholding algorithm for matrix completion," *SIAM J. Optim.*, vol. 20, no. 4, pp. 1956–1982, 2010.



**Quan Yuan** received the Ph.D. degree in computer science and technology from the Beijing University of Posts and Telecommunications (BUPT), China, in 2018. He is currently a Post-Doctoral Fellow at the State Key Laboratory of Networking and Switching Technology, BUPT. His current research interests include mobile computing, crowdsensing, and vehicular networks.



**Haibo Zhou** received the Ph.D. degree in information and communication engineering from Shanghai Jiao Tong University, China, in 2014. From 2014 to 2017, he was a Post-Doctoral Fellow with the Broadband Communications Research Group, ECE Department, University of Waterloo, Canada. He is currently an Associate Professor with the School of Electronic Science and Engineering, Nanjing University. His current research interests include resource management and protocol design in cognitive radio networks and vehicular networks. He was a recipient of the WCSP 2015 Best Paper Award. He has been serving as a Guest Editor for the *IEEE Communications Magazine*, *IET Communications*, and the *International Journal of Distributed Sensor Networks*.



**Zhihan Liu** received the M.S. degree in computer science and technology from the Beijing University of Posts and Telecommunications (BUPT) in 2004. He is currently a Researcher with the State Key Laboratory of Networking and Switching Technology, BUPT. His research interests include the Internet of Vehicles, Internet of Things, and mobile Internet.



**Jinglin Li** received the Ph.D. degree in computer science and technology from the Beijing University of Posts and Telecommunications (BUPT). He is currently an Associate Professor with the State Key Laboratory of Networking and Switching Technology, BUPT. His research interests are mainly in the areas of network intelligence, mobile Internet, Internet of Things, and Internet of Vehicles.



**Fangchun Yang** received the Ph.D. degree in communications and electronic systems from the Beijing University of Posts and Telecommunications (BUPT). He is currently a Professor with the State Key Laboratory of Networking and Switching Technology, BUPT. His current research interests include network intelligence, service computing, and Internet of Vehicles. He is a fellow of IET.



**Xuemin (Sherman) Shen** (M'97–SM'02–F'09) received the B.Sc. degree from Dalian Maritime University, China, in 1982, and the M.Sc. and Ph.D. degrees from Rutgers University, NJ, USA, in 1987 and 1990, respectively, all in electrical engineering. He is currently a Professor and the University Research Chair with the Department of Electrical and Computer Engineering, University of Waterloo, Canada. He is also the Associate Chair for Graduate Studies. His research focuses on resource management in interconnected wireless/wired networks, wireless network security, social networks, smart grid, and vehicular ad hoc and sensor networks. He is a registered Professional Engineer of Ontario, an Engineering Institute of Canada Fellow, a Canadian Academy of Engineering Fellow, a Royal Society of Canada Fellow, and a Distinguished Lecturer of the IEEE Vehicular Technology Society and the Communications Society. He is an Elected Member of the IEEE ComSoc Board of Governor, and the Chair of the Distinguished Lecturers Selection Committee. He received the Distinguished Performance Award in 2002 and 2007 from the Faculty of Engineering, University of Waterloo, the Premier's Research Excellence Award in 2003 from the Province of Ontario, Canada, the Outstanding Performance Award in 2004, 2007, 2010, and 2014 from the University of Waterloo, and the Excellent Graduate Supervision Award in 2006. He served as the Technical Program Committee Chair/Co-Chair for GLOBECOM 2007, IEEE VTC2010 Fall, IEEE Infocom 2014, and IEEE GLOBECOM 2016, the Symposia Chair for IEEE ICC 2010, the Tutorial Chair for IEEE ICC 2008 and IEEE VTC2011 Spring, the General Co-Chair for QShine 2006, ChinaCom 2007, and ACM MobiHoc 2015, and the Chair for the IEEE Communications Society Technical Committee on Wireless Communications and P2P Communications and Networking. He also serves/served as the Editor-in-Chief for the IEEE NETWORK, *Peer-to-Peer Networking and Application*, *IET Communications*, and the IEEE INTERNET OF THINGS JOURNAL, a Founding Area Editor for the IEEE TRANSACTIONS ON WIRELESS COMMUNICATIONS, an Associate Editor for the IEEE TRANSACTIONS ON VEHICULAR TECHNOLOGY, *Computer Networks*, and *ACM Wireless Networks*, and a Guest Editor for the IEEE JOURNAL ON SELECTED AREAS IN COMMUNICATIONS, the IEEE WIRELESS COMMUNICATIONS, the *IEEE Communications Magazine*, and *ACM Mobile Networks and Applications*.