

Leveraging crowd to improve data credibility for mobile crowdsensing

Tongqing Zhou, Zhiping Cai, Ming Xu, Yueyue Chen
College of Computer

National University of Defense Technology
Changsha, China

Email:{zhoutongqing, zpcai, xuming}@nudt.edu.cn, chen Yueyue_91@163.com

Abstract—Mobile crowdsensing (MCS) is a new paradigm which takes advantage of pervasive mobile devices to collaboratively collect data and analyze physical phenomenon. As mobile devices are owned and controlled by individuals with various capabilities and intentions, a main challenge MCS applications face is to ensure the credibility of the crowd contributed data. Existed works attempt to increase confidence level of the sensory measurements by validating the location. However, the required infrastructure or neighbor support may not always be available, and the unreliable form containing false sensory data with a valid location is implicitly ignored. In this paper, we propose a novel Crowd-based Credibility Improving Scheme (CCIS) to improve the credibility of data in possible false forms leveraging crowd data property and crowd participants' reputation. Based on the data clusters generated using a lightweight fixed-width clustering algorithm, CCIS is able to adequately identify and filter out the clusters constituted mainly by false data using reputation information as the classifier. We conduct simulations on a publicly available trace with crowd contributed temperature measurements, the results show that CCIS yields an improvement of overall data credibility of around 1.2 with clustering accuracy over 96%.

Index Terms—Mobile crowdsensing; data credibility; reputation score; clustering algorithm.

I. INTRODUCTION

The demand for more pervasive sensing of physical world and the proliferation of human-carried smart devices with rich set of embedded sensors have given rise to a new sensing-oriented application paradigm, known as mobile crowdsensing (MCS) [1]. In a MCS application, ubiquitous sensors of public crowd instead of specialized mote-class sensors perform participatory or opportunistic sensing, and collect interested physical data for further aggregation and analysis in cloud-based platform. A broad spectrum of applications have been developed based on MCS, including environment monitoring (e.g. noise pollution), smart traffic (e.g. congestion monitoring), city management (e.g. network measurement) [2], etc.

One major challenge for the adoption of a MCS application is the assurance for credibility of the contributed sensory data [3]. As an important dimension of data quality, data credibility stands for the extent to which the good faith of data sources can be relied upon to ensure the data really represents what it is supposed to be. The threat level of data falsification is high in MCS as sensors are owned by individuals [4], and

typically the issue exists due to three factors: (1) the openness characteristics of network, (2) human greedy that drives them to maximize their profit while minimizing their efforts [5], and (3) potential malicious intention of participants. Specifically, participants have the motivation and capability to contribute false data to earn money without actually executing the task or merely to mislead the conclusion of MCS applications by injecting artificial data. For instance, an Internet Service Provider may generate false measurement to degrade its competitor's performance evaluation while increasing their own profit. Such misbehavior will deviate aggregation result and hinder the global awareness of phenomenon of interest.

Many works have been done to deal with the flawed ingredient in collected data and improve the credibility. In traditional Wireless Sensor Networks (WSNs), inner-cluster endorsement [6] and statistic analysis [7] are introduced to detect false data injection attack launched by compromised sensors and aggregators, respectively. However, the endorsement scheme requires a prior knowledge on the number of malicious sensors, while the statistic scheme actually focuses on solving a MITM (Man-In-The-Middle) attack. Another category of solution attempts to increase the credibility using location attestation obtained from infrastructure-based [8] or neighbor-assisted [9] verification. Unfortunately, the requirement of infrastructure and neighbor support may not be feasible, while the real-time location verification process is time-consuming, and most importantly, these solutions implicitly ignore a kind of contribution containing false sensory data and valid location. In [10], a reputation framework is proposed to estimate the trustworthiness of contributions for social participatory sensing systems, the heuristic model and empirical thresholds it use constrain its scalability.

This work proposes CCIS, a crowd-based credibility improving scheme to assure the data credibility in MCS applications without third-party involvement or extra-knowledge requirement. Clustering algorithm is performed on the location-intensive crowd-contributed data to formally group false and normal ingredient into different groups. And crowd participants' reputation are introduced to identify and filter out the false ones, improving overall credibility for the crowd data. The main contributions of our work are three-fold:

- 1) Data credibility issue in MCS applications is described through analyzing the state space of data validity, and for

the first time, the type of false data with invalid sensory measurement and valid location is considered;

- 2) We propose a data credibility improving scheme, CCIS, leveraging lightweight clustering algorithm and participant reputation information as two building blocks to filter out the false ingredient from the collected data;
- 3) We validate the scheme using synthetic data, and show the credibility improvement against location attestation-based scheme and good performance on clustering accuracy.

The rest of this paper is organized as follows. Related works are summarized in section II. In Section III, the problem will be formulated and adversary models will be described. We outline the key components of scheme CCIS, and then introduce how CCIS facilitates better data credibility in section IV. In Section V, simulation results that indicate the effectiveness of the scheme will be provided. Finally, conclusions will be presented in Section V.

II. RELATED WORK

As a new sensing paradigm, MCS is a particular subset of the traditional WSNs with the support of widely distributed mobile crowd participants. The reliability issue of collected data also exists in WSNs applications, and various approaches have been proposed in literature. In this section, we summarize the related research works on data credibility assurance from both the perspective of WSNs and MCS.

A. False data detection in WSNs

In the field of WSNs, the potential sources of false data (outlier) include noise, events, and malicious attack [11]. In [6], an interleaved hop-by-hop authentication scheme is proposed to detect injected false data packets by checking endorsements of the co-located nodes. A sensory report is determined to be trusted only when the number of endorsements exceeds the number of possible malicious nodes. And an approach of aggregate-commit-prove is proposed to secure information aggregation by constructing efficient random sampling and interactive proofs in [7]. Clustering algorithm is used to detect anomaly in WSNs in [12]. While the inter-cluster distance is chosen to be the classifier in [12], we propose to use overall reputation information of a sensory data cluster to determine whether it is normal or not.

Another form of WSN that is vulnerable to false data threaten is Smart Grid, which relies on widespread sensors to collect power system information for state estimation. The adversary may inject false measurement reports to mislead the state estimation process through compromising meters and sensors, and result in disruption of the energy distribution [13]. In [14], an efficient injection detection scheme is proposed through exploiting spatial-temporal correlations pattern between grid components. Spatial correlation is also used in our work, while for a different purpose of clustering the sensory data.

Solutions in WSNs are instructive to understand the situation of data credibility in MCS, however, characteristics

of MCS (e.g. human-involvement) must be considered to effectively guide the data analysis process.

B. Approaches in MCS

We emphasize that credibility issue of data in MCS becomes more crucial as sensors are carried and owned by individuals with various capabilities and unknown intentions as it makes every participant a potential threat source.

In work [15] and [16], Trusted Platform Module (TPM) is adopted to assure data trustworthiness. However, such embedded trust module is not yet available for most mobile devices, and malicious participants can still cause distortion of the measurements by deliberately initiating sensing action.

On the other hand, as location being a common tag for sensory measurements, validating location can achieve a certain degree of reliability on the sensed data [17]. Based on this concept, a series of approaches have been proposed, and can be classified as infrastructure-based scheme [18][8] and peer-assisted scheme [9]. In [8], a lightweight protocol named Echo is presented, in which the location of user is successfully proved when it is able to return the challenge packet from the verifier (typically wireless infrastructure) in constrained time. A peer-assisted scheme is proposed in [9] to validate user's location based on the verification from co-located users connected by bluetooth. Such solutions require either infrastructure support or neighbours' involvement, and add a high overhead on sensor devices during proving themselves. Moreover, the case that dishonest participants submit faked data from valid location is not properly handled in those schemes. The proposed scheme CCIS attempts to detect and identify false data from the perspective of sensory data instead of location, which is general and scalable.

III. PROBLEM STATEMENT

Overview of the workflow of MCS applications we consider in this paper is illustrated in Fig. 1. Typically, it consists of a cloud-based platform and a set of participants $U = \{u_1, \dots, u_n\}$ who perform sensing task T at location L . Specifically, we focus on location-based environmental-centric sensing tasks that collect numerical data for further analysis (e.g. temperature, noise level). During the execution of task T , individual u_i collects a series of measurements denoted as s_i , and submits it together with the corresponding location l_i to the platform before pre-defined time deadline, thus the submitted data of u_i is a tuple in key-value syntax, denoted as $d_i = \langle l_i, s_i \rangle$.

By the end of T , the centralized platform will obtain a data set $D = \{d_1, \dots, d_n\}$ (grouped by participant id), based on which some aggregation function f will be performed to provide statistical result for publish. However, as some false or erroneous ingredient exists in D as shown in Fig. 1, the aggregation result may deviate from the expected true value. Here we consider data d_i as trustworthy (normal) if and only if its location component l_i and sensory measurement component s_i are both valid, in other words, d_i is false when either l_i or s_i is invalid. And the validity of l_i and s_i is defined as:

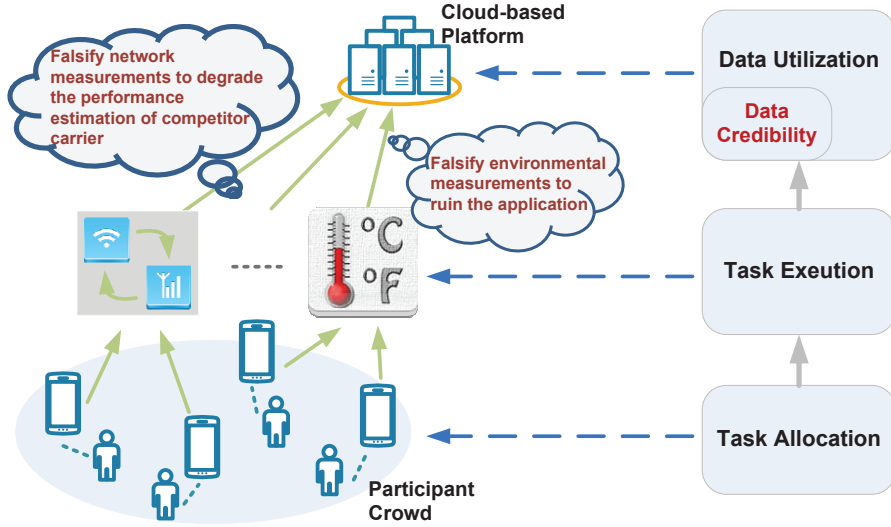


Fig. 1. Overview of the workflow of typical MCS applications, and possible threats to data credibility. Typically, an application is parsed as sensing task which is further carried out in three stages: 1) Task allocation to the participant crowd, 2) Task execution with devices of the crowd, and 3) Crowd data utilization by the application server, where data credibility is considered as an crucial part of the utilization stage, and data falsification threats arise as the crowd participants representing various capabilities and purposes. Specifically, network measurement and temperature monitoring are depicted as two illustrative example applications with potential credibility threats listed

- 1) *Validity of Location*. Announced location l_i of u_i is valid if it is within an acceptable distance from location l_c where u_i currently at.
- 2) *Validity of Sensory Measurement*. Sensory measurement s_i is valid if it reflects the ground truth of physical phenomenon of the corresponding location l_i .

According to the validity of component l_i and s_i , the state space of validity of data d_i can be represented with four categories as shown in Table I, where symbol T means the value is true and symbol F means false. Generally, a location attestation-based scheme attempts to improve data credibility through picking out data with invalid location in category B¹ and category C, ignoring possibly false data in category A, which is more common in a MCS application especially when malicious intention is considered.

TABLE I

SPACE STATE FOR THE VALIDITY OF CONTRIBUTION DATA IN REGARD TO VALIDITY OF LOCATION l_i AND VALIDITY OF SENSORY MEASUREMENT s_i

$l_i \backslash s_i$	T	F
T	Normal data	Category A
F	Category B	Category C

Finally, we assume that location L refers to an area of interest within certain distance of L instead of a specific spot as physical measurements are usually spatial correlated, and

¹Indeed, providing valid sensory data at invalid location (category B) is not feasible as an user is unable to contribute valid data of one location when he is not really there. Otherwise, the contribution is treated as valid considering a possible situation of submission delay.

all sensory data are aligned on measurement features. And two types of adversary model are considered:

- 1) *Random Falsification*. Participants submit measurement data with random value to minimize their efforts, or tamper the measurement to facilitate a misleading effect. For the latter intention, dishonest participants would try to deviate the aggregation result as much as possible.
- 2) *Falsification with Conspiracy Cooperation*. A special case of the random falsification model, in which a group of participants collude with each other to intentionally induce the final aggregation result to a wrong value by submitting false data with similar value. Moreover, in order to avoid being identified by statistical analysis-based abnormal detection method, the dishonest group is able to fabricate and submit data obeying normal distribution.

Collusion among participants would result in a more significant deviation, and the injected artificial data cannot be easily picked out. Taking average function f_{avg} as an example of the aggregation function, if we have a crowd contributed data set $D_{eg} = \{d_n^1, \dots, d_n^N, d_f^1, \dots, d_f^N\}$, where $d_n^i = \langle L, M \rangle$ denotes a normal measurement, and $d_f^i = \langle L, 2M \rangle$ denotes a false measurement, then we will have $f_{avg}(D_{eg}) = 1.5M$ which is 1.5 times larger than the actual value M . Additionally, we do not make any assumption or set any limitation on the number of dishonest participants in U .

IV. DESIGN OF CCIS

This work tries to improve the overall credibility of crowd data in MCS through identifying and discarding the corrupted part with invalid sensory measurements, which refers to the data belonging to category A and C in Table I. Theoretically,

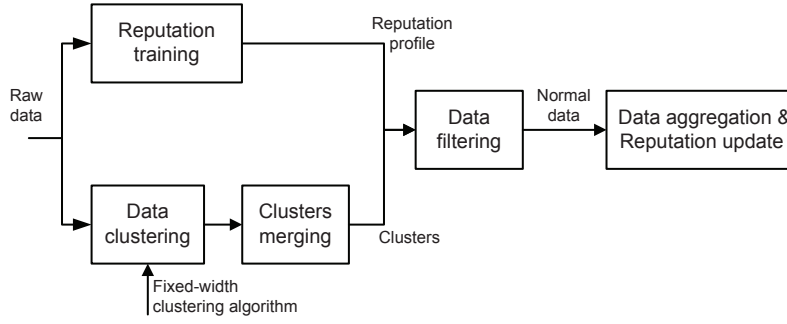


Fig. 2. Framework of CCIS

only normal part in the collected data remains after the processing as data in category B do not pose a threat.

Framework of the proposed scheme CCIS is shown in Fig. 2, which consists of five functional components: a) reputation training, b) data clustering, c) clusters merging, d) data filtering, and e) data aggregation. Based on those components, CCIS is carried out in two phases sequentially, the initialization phase where reputation profile for participants are trained and obtained, and the filtering phase where raw sensory data is classified as false or normal with clustering algorithm and available reputation information as two building blocks. Finally, clusters identified as false data container are filtered out and aggregation function is performed based on the filtering result.

A. Initialization phase

In order to obtain basic knowledge about the participants' reputation, we propose to introduce an initialization phase for reputation training. During this phase, the distance between each piece of sensory data s_i and real world measurement s_r is first calculated as,

$$\text{dist}(s_i, s_r) = \sqrt{\sum_{j=1}^{N_f} (s_i(j) - s_r(j))^2} \quad (1)$$

where $s(j)$ denotes the value of the j -th feature for measurement s , and N_f denotes the size of feature space. The validity of s_i is determined by comparing $\text{dist}(s_i, s_r)$ with a predefined threshold d_{th} . Then individual reputation scores are calculated based on the validity of contribution in an AIMD (additive increase and multiplicative decrease) way in order to clearly set reputation of participant with different behavior apart. The process is illustrated in Algorithm 1, and it ends up with reputation evaluation result $REP = \{R_1, R_2, \dots, R_n\}$. For simplicity, reputation scores in REP are further normalized using:

$$R_i = \frac{R_i - \min(REP)}{\max(REP) - \min(REP)} \quad (2)$$

Note that this training process is only required in the very start of a MCS application, corresponding knowledge can be used as a metric in both the continued tasks and other MCS applications. User profiles from social network are feasible

alternative to obtain reputation information, but it is not in the scope of our work.

Algorithm 1 Participant Reputation Score Calculation

R_i : reputation score of the i -th participant with same initially value r_I
 r_a : additive increase factor
 r_m : multiplicative decrease factor
 n : the number of participants
 m_i : the number of contributions from participant i

```

for  $i = 1 \rightarrow n$  do
  for  $j = 0 \rightarrow m_i$  do
    if  $\text{dist}(s_i^j, s_r^j) > d_{th}$  then
       $R_i += r_a$ 
    else
       $R_i /= r_m$ 
    end if
  end for
end for
  
```

B. Filtering phase

It is reasonable to say that contribution from participant with better reputation tends to be more reliable. However, using reputation score alone as the credibility metric is ineffective, because setting a proper threshold to distinguish corrupted part from the crowd data is hard and inflexible. Moreover, such straightforward strategy basically denies the possibility for participants with low reputation to submit normal data, and neglects accidentally low quality contribution from participants with high reputation. On the other hand, clustering algorithm is effective for detecting outlier of a set of data, while not adequate in handling samples containing fabricated false data with ambiguous amount. Additionally, the choice of a proper cluster width parameter for a clustering algorithm is non-trivial. To overcome these drawbacks while retaining the advantages, a hybrid approach is proposed to identify the clusters with well-defined width as normal or false utilizing participants reputation as the classifier.

Generally, physical measurements act as signatures that characterize a place of interest, which implies that measurements for the same location are correlated with each other. Meanwhile, the collected data are mainly exploited at a

community scale which provides sufficient participant density support for clustering the data around a specific sensing location. Hence, a fixed width clustering algorithm is first performed on D to group similar data instances into clusters with similar property. The first data is assigned to be the centroid of the first cluster. Then for every subsequent data d_i in D , Euclidean distance between centroid of each cluster and data d_i is calculated. If its distance to one cluster is less than the cluster width² ω , then it is added to that cluster and the centroid of that cluster is updated accordingly. Otherwise a new cluster is formed with that data as the initial centroid. Here, we novelly define ω as half of the minimum expected deviation from the true aggregation result for a potential falsify behavior among the crowd data, i.e.,

$$\omega = 1/2 \cdot \min_i \left(\left| f(D(i)) - f(\tilde{D}) \right| \right) = 1/2 \cdot \sigma_{dev} \quad (3)$$

where $D(i)$ represents one of the possible collected data sets containing corrupted ingredient, whose value and amount are both unknown, \tilde{D} denotes the set of normal data, and σ_{dev} denotes the minimum expected misleading degree. The value of parameter σ_{dev} is adjusted adaptively according to the application context, e.g., a dishonest participant may prefer to consider a deviation of at least 4°C as effective for a task that measures city temperature, while 10dB may be a meaningful value in an application of received signal strength measurement. The clustering operation produces a set of fixed width clusters $C = \{C_1, \dots, C_n\}$ in the feature space. The advantage of this simple approach is that only one pass is required, and the complexity of the algorithm is $O(n_c n_d)$, where n_c and n_d is the number of clusters and data points, respectively.

The clusters are then labelled as normal or false with the aid of their overall reputation. In this way, both humanity dimension (participant reputation) and data dimension (cluster property) are taken into consideration. The reputation of a cluster is defined as the average reputation of the contributors of every piece of data in that cluster:

$$R(C_i) = 1/N_{C_i} \cdot \sum_j R_{u(s_j)} \quad (4)$$

subject to $R_{u(s_j)} \in REP, s_j \in C_i, j \in [1, N_{C_i}]$

where N_c denotes the number of sensory measurements in C_i , and $u(s_j)$ is the participant id of measurement s_j . Metric $R(C_i)$ reveals the credibility of C_i with respect to the overall reputation condition within C_i . Hence, the cluster that represents the highest reputation is first selected and determined to be the container of normal data, i.e.,

$$C^* = \arg \max_{C_i} \{R(C_i)\} \quad (5)$$

and the data points inside are remained.

Note that data points in C^* only cover a fraction of the normal ingredient of \tilde{D} . In order to improve clustering

²We use cluster width to describe the maximum acceptable distance to the centroid for one data point to be added into that cluster, and use radius to describe the actual maximum distance to the centroid for all the data points in that cluster.

accuracy, we further introduce a merging stage to combine clusters similar to C^* with C^* to form a new cluster C_m^* . The merging process will enlarge the radius of C^* as more data are integrated into it, thus a new cluster width ω' need to be chosen first to set an upper bound for the distance between a data point and the centroid of C^* . Here, we set parameter ω' to $\omega + \sigma_{C^*}$, where σ_{C^*} denotes the standard deviation of measurements in C^* . This operation equals to transferring the centroid of C^* to a circle with radius σ_{C^*} . σ_{C^*} is here chosen to be the increment because it reflects the deviation level of C^* , thus a centroid stays within the distance of σ_{C^*} to C^* can be considered as acceptable.

Proposition 1: If cluster C_i satisfies $dist(C^*, C_i) < \sigma_{C^*}$, then $\forall d_j \in C_i, d_j$ is qualified to join C_m^* .

Proof: According to the definition of ω , the radius of C_i satisfies $r_i < \omega$, then

$$\begin{aligned} \max_j (dist(C^*, d_j)) &= dist(C^*, C_i) + r_i \\ &< \sigma_{C^*} + \omega = \omega' \end{aligned} \quad (6)$$

such that all points in C_i are within the distance constraint to be added into C_m^* . ■

The merging stage is carried out as follows:

- 1) For each cluster C_i in C , the distance $dist(C^*, C_i)$ between C_i and C^* is calculated as the Euclidean distance between their centroid.
- 2) According to Proposition 1, the distances are then compared with σ_{C^*} , and cluster C_i is classified as normal and integrated with C^* if $dist(C^*, C_i) < \sigma_{C^*}$.

The process ends up with a output C_m^* , and data points not belonging to C_m^* form another cluster C_f labelled as false.

Finally, data in C_f is regarded as corrupted ingredient and filtered out from D to improve the overall credibility. And participants' reputation scores are dynamically updated using Algorithm 1 (increase the reputation of the contributors of data in C_m^* , and decrease the reputation of those in C_f) to involve the evaluation of participants' contribution during this task.

V. EVALUATION

A. Settings and Metrics

We consider a typical application in this section, say, leveraging MCS for environmental monitoring. In such applications, portable sensors are equipped with mobile participants for collecting physical information. Specifically, we evaluate the proposed scheme on an open source temperature measurement traces obtained from the CRAWDAD data set [19], which contains 5030 measurement items from 289 active taxicabs collected around the GPS location (41.9, 12.5) in Rome. Meanwhile, some items are modified to simulate the dishonest behaviors that falsifies sensory data. Here, we consider the adversary model of falsification with conspiracy cooperation as it is harder to detect. Typically, participants with lower reputation are more likely to submit faked data than their counterparts, so we randomly replace their temperature measurements with random values generated from a normal distribution with mean parameter μ equalling to the value of misleading target S_{err}

TABLE II
PARAMETERS SETUP FOR CCIS

Parameter	Additive Factor	Multiplicative Factor	Reputation Threshold	Cluster Width
Value	1	2.5	0.2	2

and standard deviation parameter $\sigma = 1$. We generate falsification data obeying normal distribution to imitate smart collusion among a dishonest group. Moreover, the synthetic data set is divided into two sets according to the submission time of the contribution to conduct two experiments independently.

The corresponding setup parameters are presented in Table II. We consider the minimum possible deviation caused by data falsification to be 4° , so the cluster width is set to be 2° . And the mean temperature measurement for the two time period is 8.85° and 14.05° . In order to effectively mislead the aggregation result, the falsification target S_{err} for time period 1 and time period 2 is set to 14° and 8° , respectively. Two sets of faked measurement are then generated and used to replace the measurements of selected data items in the original set.

For the evaluation part, the performance of clustering algorithm is evaluated using overall accuracy, which is defined as

$$A_{overall} = \frac{\sum_{i=1}^{N_{cluster}} TP_{C_i}}{N_{data}} \quad (7)$$

where TP_{C_i} refers to True Positive and equals to the number of correctly classified data in cluster C_i (e.g. Normal data that classified into a false data cluster which will be discarded do not belong to TP_{C_i}), and $N_{cluster}$ and N_{data} denotes the number of clusters and data points, respectively. Meanwhile, the effectiveness of CCIS is evaluated using credibility metric \mathfrak{R}_D , given by

$$\mathfrak{R}_D = 1 - \left(\frac{|f(D) - f(\tilde{D})|}{\min(f(D), f(\tilde{D}))} \right) \quad (8)$$

where $\tilde{D} = D - D_f$, and D_f is the set of false data. Obviously, the less false data in D , the more similar $f(D)$ and $f(\tilde{D})$ will be, and D will have a higher credibility. Finally, without loss of generality, we adopt average function as the aggregation function f during data analysis.

B. Results

During the initialization phase, each participant's reputation is evaluated based on the quality of their submitted measurements. We assume the raw data in [19] are all valid, and take the mean value of the temperature measurements in corresponding time period as the real temperature. Then one's reputation is adequately decreased (increased) if the distance between the real value and his sensory measurement exceeds (stays below) predefined threshold 1.5 (this parameter can be adjusted to simulate different amount of dishonest participants). Fig. 3 shows the normalized training results regarding each participant's reputation with initial reputation of each participant set to 1.

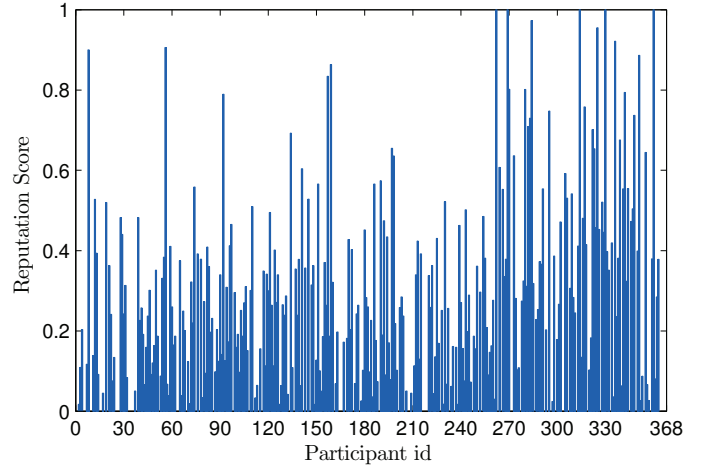


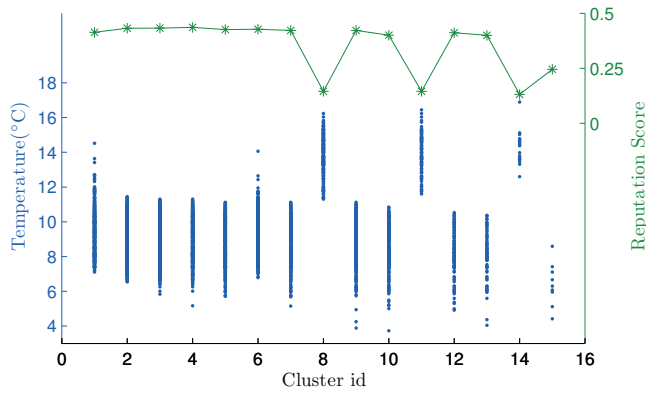
Fig. 3. Reputation of participants after the initialization training phase (participant id ranges from 4 to 368 with 289 valid ids)

During the filtering phase, clustering algorithm in CCIS is first performed on the modified set to gather data points with similar value into the same group. Clustering results for the two time periods are illustrated in Fig 4 with the value of overall reputation for each cluster depicted as well. As we can tell, clusters containing false ingredient turn out to have low reputation in both periods, and cluster 4 (reputation 0.436) and cluster 21 (reputation 0.512) represent the highest overall reputation for each period. Meanwhile, clusters generated for each period share joint area, leaving a space for merging.

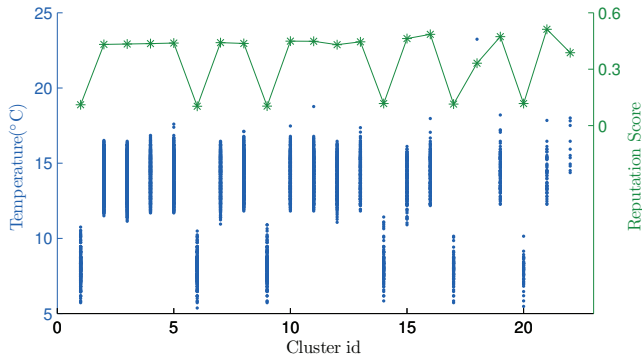
According to the scheme, cluster 4 in the first set and cluster 21 in the second set are first identified and determined to be normal. Without further processing, the aggregation result will be 8.864° and 14.281° with accuracy 58% and 6% for each period. The estimation result is acceptable, while the accuracy is poor. In order to improve clustering accuracy, the merging approach of CCIS is performed among the generated clusters. The cluster width is updated by adding the standard deviation of cluster 4 and cluster 21 (both equal to 1.2) to ω , which renders $\omega' = 3.2$ for both periods. The merging results are illustrated in Fig. 5. For each period, two clusters labelled as normal or false are generated. The normal one consists of 1) data from cluster C^* identified in data clustering stage and described as normal data-1, and 2) data from the neighbor clusters of C^* identified in the cluster merging stage and described as normal data-2, which may introduce a small amount of false data. The remaining data form the false cluster and will be eventually discarded.

In Fig. 6, we take time period 1 as an example, and statistically compare the temperature measurement distribution of the modified data set and the data set after being filtered using CCIS. The falsification operation deviates the mean value by fabricating a faked normal distribution around $14^\circ C$ (Fig. 6(a)). On the other hand, the misleading deviation is successfully removed with CCIS to bring the statistical result back to its real value (Fig. 6(b)).

Finally, the evaluation results for our proposed scheme CCIS



(a) Result for time period 1 (from 6 o'clock to 10 o'clock)



(b) Result for time period 2 (from 11 o'clock to 15 o'clock)

Fig. 4. Clustering results (left Y-axis) and overall cluster reputation (right Y-axis) for each cluster

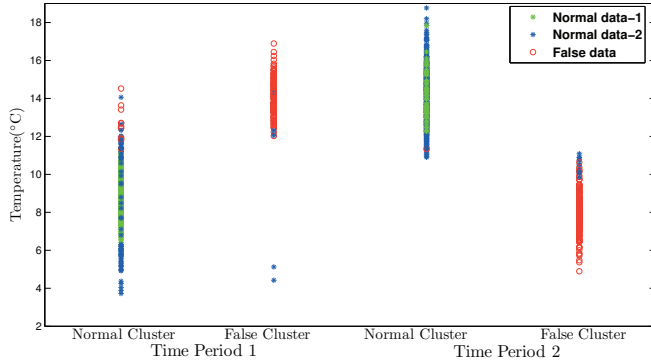
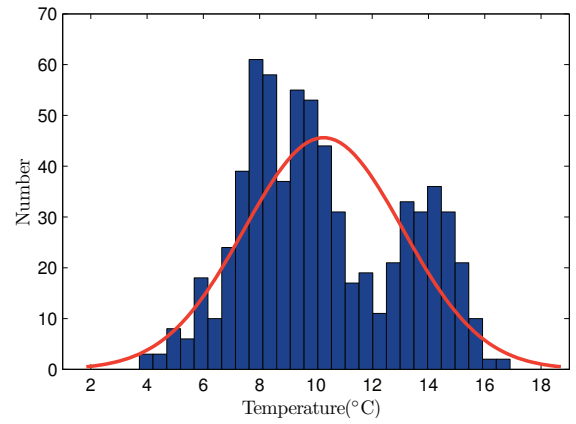


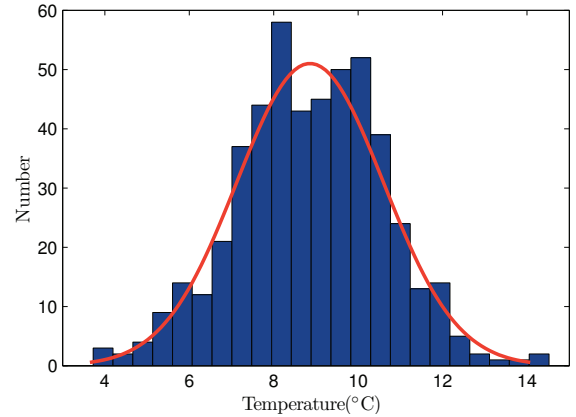
Fig. 5. Clusters merging results of CCIS for two periods of time. For both periods, the crowd data are clustered into two clusters, and labelled as normal or false based on each cluster's reputation

are illustrated in Table III. Temperature estimation is misled to a wrong value (resulting in a higher estimation for period 1 and a lower estimation for period 2) due to the introduction of artificial measurements. Specifically, for time period 1, the scheme improves the overall data credibility from 84% to 99% with clustering accuracy 96%. And the scheme yields an improvement of overall data credibility of 1.19 with clustering accuracy 97% for time period 2.

(Comparison) For the scenario considered in our work, location attestation-based schemes are unable to identify any



(a) Temperature distribution of the modified data set



(b) Temperature distribution of the crowd contributed data after CCIS filtering

Fig. 6. Statistical analysis for data distribution before and after CCIS filtering for time period 1. The Y-axis represents the number of data points regarding different temperature, and the curve in red shows the normal fitting for the data set

false contributions as all the location components in the data set are valid. Hence, we emphasize that the improvement of credibility achieved by CCIS stated in Table III is also relative to the credibility level location attestation-based schemes can achieve.

VI. CONCLUSION

We have presented a crowd-based scheme CCIS to improve data credibility for typical MCS applications. Other than the scenarios considered in location attestation-based scheme, CCIS especially handles the cases that dishonest participants submit false data from valid location. A fixed-width clustering algorithm is introduced to cluster the contributions into groups, and a merging approach is further performed on the groups to improve clustering accuracy. Finally, reputation information is adopted to identify clusters holding false data and filters them out. The simulation results show that CCIS adequately improve the overall credibility under the cases considered. Further works will include reputation profile evaluation with social network information, the consideration for more general form of data.

TABLE III
SUMMARY OF FILTERING RESULT AND EVALUATION RESULT

Metric	Real Result	Misleading Target	Falsification Result	CCIS Result	Accuracy	Credibility (False, CCIS)	Credibility Improvement
Period 1	8.85	14	10.27(↑)	8.86	96%	(0.84, 0.99)	1.17
Period 2	14.05	8	12.05(↓)	14.13	97%	(0.83, 0.99)	1.19

ACKNOWLEDGMENT

This work is supported by National Key Basic Research Program of China (No.2012CB933504), the National Natural Science Foundation of China (NSNF) under Grant Nos. 61379144, 61379145, 61402513, 61363071, 61402275.

REFERENCES

- [1] R. K. Ganti, F. Ye, and H. Lei, "Mobile crowdsensing: current state and future challenges," *Communications Magazine, IEEE*, vol. 49, no. 11, pp. 32–39, 2011.
- [2] W. Z. Khan, Y. Xiang, M. Y. Aalsalem, and Q. Arshad, "Mobile phone sensing systems: A survey," *Communications Surveys & Tutorials, IEEE*, vol. 15, no. 1, pp. 402–427, 2013.
- [3] S. Reddy, V. Samanta, J. Burke, D. Estrin, M. Hansen, and M. Srivastava, "Mobisense: mobile network services for coordinated participatory sensing," in *Autonomous Decentralized Systems, 2009. ISADS'09. International Symposium on*. IEEE, 2009, pp. 1–6.
- [4] P. Johnson, A. Kapadia, D. Kotz, N. Triandopoulos, and N. Hanover, "People-centric urban sensing: Security challenges for the new paradigm," *Dept. of Computer Science, Dartmouth College. URL http://www.cs.dartmouth.edu/~dfk/papers/johnson-metrosec-challenges-tr.pdf.-Zugriffsdatum*, vol. 26, p. 2008, 2007.
- [5] J. S. Downs, M. B. Holbrook, S. Sheng, and L. F. Cranor, "Are your participants gaming the system? screening mechanical turk workers," in *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. ACM, 2010, pp. 2399–2402.
- [6] S. Zhu, S. Setia, S. Jajodia, and P. Ning, "An interleaved hop-by-hop authentication scheme for filtering of injected false data in sensor networks," in *Security and privacy, 2004. Proceedings. 2004 IEEE symposium on*. IEEE, 2004, pp. 259–271.
- [7] B. Przydatek, D. Song, and A. Perrig, "Sia: Secure information aggregation in sensor networks," in *Proceedings of the 1st international conference on Embedded networked sensor systems*. ACM, 2003, pp. 255–265.
- [8] N. Sastry, U. Shankar, and D. Wagner, "Secure verification of location claims," in *Proceedings of the 2nd ACM workshop on Wireless security*. ACM, 2003, pp. 1–10.
- [9] M. Talasila, R. Curtmola, and C. Borcea, "Link: Location verification through immediate neighbors knowledge," in *Mobile and Ubiquitous Systems: Computing, Networking, and Services*. Springer, 2012, pp. 210–223.
- [10] H. Amintoosi and S. S. Kanhere, "A reputation framework for social participatory sensing systems," *Mobile Networks and Applications*, vol. 19, no. 1, pp. 88–100, 2014.
- [11] Y. Zhang, N. Meratnia, and P. Havinga, "Outlier detection techniques for wireless sensor networks: A survey," *Communications Surveys & Tutorials, IEEE*, vol. 12, no. 2, pp. 159–170, 2010.
- [12] S. Rajasegarar, C. Leckie, M. Palaniswami, and J. C. Bezdek, "Distributed anomaly detection in wireless sensor networks," in *Communication systems, 2006. ICCS 2006. 10th IEEE Singapore International Conference on*. IEEE, 2006, pp. 1–5.
- [13] Y. Mo, E. Garone, A. Casavola, and B. Sinopoli, "False data injection attacks against state estimation in wireless sensor networks," in *Decision and Control (CDC), 2010 49th IEEE Conference on*. IEEE, 2010, pp. 5967–5972.
- [14] P.-Y. Chen, S. Yang, J. McCann, J. Lin, X. Yang *et al.*, "Detection of false data injection attacks in smart-grid systems," *Communications Magazine, IEEE*, vol. 53, no. 2, pp. 206–213, 2015.
- [15] P. Gilbert, J. Jung, K. Lee, H. Qin, D. Sharkey, A. Sheth, and L. P. Cox, "Youprove: authenticity and fidelity in mobile sensing," in *Proceedings of the 9th ACM Conference on Embedded Networked Sensor Systems*. ACM, 2011, pp. 176–189.
- [16] P. Gilbert, L. P. Cox, J. Jung, and D. Wetherall, "Toward trustworthy mobile sensing," in *Proceedings of the Eleventh Workshop on Mobile Computing Systems & Applications*. ACM, 2010, pp. 31–36.
- [17] M. Talasila, R. Curtmola, and C. Borcea, "Improving location reliability in crowd sensed data with minimal efforts," in *Wireless and Mobile Networking Conference (WMNC), 2013 6th Joint IFIP*. IEEE, 2013, pp. 1–8.
- [18] J. Brassil, R. Netravali, S. Haber, P. Manadhata, and P. Rao, "Authenticating a mobile device's location using voice signatures," in *Wireless and Mobile Computing, Networking and Communications (WiMob), 2012 IEEE 8th International Conference on*. IEEE, 2012, pp. 458–465.
- [19] M. A. Alswailim, H. S. Hassanein, and M. Zulkernine, "CRAWDAD dataset queensu/crowd_temperature (v. 2015-11-20)," Downloaded from http://crawdad.org/queensu/crowd_temperature/20151120, Nov. 2015.