

Quality Estimation for Scarce Scenarios Within Mobile Crowdsensing Systems

Sherif B. Azmy, *Student Member, IEEE*, Nizar Zorba^{IP}, *Senior Member, IEEE*,
and Hossam S. Hassanein, *Fellow, IEEE*

Abstract—Mobile crowdsensing (MCS) is a paradigm that exploits the presence of a crowd of moving human participants to acquire, or generate, data from their environment. As a part of the Internet-of-Things (IoT) paradigm, MCS serves the quest for a more efficient operation of a smart city. Big data techniques employed on this data produce inferences about the participants' environment, the smart city. However, sufficient amounts of data are not always available. Sometimes, the available data are scarce as it is obtained at different times, locations, and from different MCS participants who may not be present. As a consequence, the scale of data acquired may be small and susceptible to errors. In such scenarios, the MCS system requires techniques that acquire reliable inferences from such limited data sets. To that end, we resort to small data (SD) techniques that are relevant for scarce and erroneous scenarios. In this article, we discuss SD and propose schemes to tackle the problems associated with such limited data sets, in the context of the smart city. We propose two novel quality metrics: 1) MAD quality metric (MAD-Q) and 2) MAD bootstrap quality metric (MADBS-Q), to deal with SD, focusing on evaluating the quality of a data set within MCS. We also propose an MCS-specific coverage metric that combines the spatial dimension with MAD-Q and MADBS-Q. We show the performance of all the presented techniques through closed-form mathematical expressions, with which simulation results were found to be consistent.

Index Terms—Data quality, Internet of Things (IoT), IoT architectures, IoT-based services, mobile crowdsensing (MCS), small data (SD).

I. INTRODUCTION

RECENT breakthroughs within the Internet-of-Things (IoT) paradigm have led to an unprecedented integration of sensors across various fields that relate to human life. As a result, infrastructures incorporating sensors are becoming more complex as the quantity and variety of sensors increase, making infrastructures more demanding for schemes that improve the efficiency of a smart city's operation. Such infrastructural demands are a concern for smart city administrators, as more

data could enhance the city's operational efficiency, such as crisis response and transportation. One notable paradigm is mobile crowdsensing (MCS), which exploits the availability of smartphones within a crowd, to benefit from their smart devices' embedded sensors. MCS is capable of transforming a crowd of smartphone users into an extended instrument, MCS participants, for the benefit of the smart city [1].

MCS has a wide range of applications that permit observing the social and physical dynamics of a smart city, which makes MCS a useful tool for smart city management [4]. In MCS, administrators assign tasks to participants in an autonomous manner [4]. This autonomy of sensing can be classified into *participatory* sensing or *opportunistic* sensing [4]. In participatory sensing, the MCS system requests the participant to voluntarily engage in the execution of a task, for example, using the phone's microphone to acquire noise levels [5] or the phone's camera to take a photograph [6]. Opportunistic sensing, on the other hand, waits until the participant satisfies the conditions necessary for the task (such as time, place, sensor capabilities, etc.), executing the task with no active intervention from the participant [7]. So far, MCS management frameworks, such as 4W1H [8] and platforms, such as ParticipAct [6] were developed to address the real-time implementation of MCS. Others have developed MCS models with the objective of reduced mathematical complexity, such as modeling the MCS sensing problem over both space and time, i.e., *spatiotemporal cells* [9].

The previous research efforts [6], [10] focus on large data sets, proposing techniques that assume abundant data is always available, often in the big data (BD) scale. While BD techniques are generally superior when sufficient data are available, small data (SD) techniques are provided insight when the data are scarce, insufficient, and erroneous, in spite of the challenge posed by these three characteristics. Nevertheless, SD can build upon the output of BD which is an inference resulting from the analysis of a large amount of data, for example, a categorical classification. SD can also complement BD by ensuring the resilience of the data sources at the beginning of the data collection pipeline. In other words, SD techniques aim to maintain the quality of the data collected at a local scale, whereas BD techniques analyze at a global scale. The combination of both SD and BD allows the exploitation of the global versus local contrast [11].

In MCS, if the participants' sensor quality in a spatiotemporal cell is low, then only a subset of the measured values can be considered, in such a case, SD is to be used. SD

Manuscript received January 10, 2020; revised April 11, 2020; accepted April 27, 2020. Date of publication May 14, 2020; date of current version November 12, 2020. This work was supported by the Natural Sciences and Engineering Research Council of Canada under Grant RGPIN-2019-05667. (Corresponding author: Nizar Zorba.)

Sherif B. Azmy is with the Department of Electrical and Computer Engineering, Queen's University, Kingston, ON K7L 3N6, Canada (e-mail: sherif.azmy@queensu.ca).

Nizar Zorba is with the College of Engineering, Qatar University, Doha, Qatar (e-mail: nizarz@qu.edu.qa).

Hossam S. Hassanein is with the School of Computing, Queen's University, Kingston, ON K7L 3N6, Canada (e-mail: hossam@cs.queensu.ca).

Digital Object Identifier 10.1109/JIOT.2020.2994556

also serves MCS from an economic point of view, as the administrator may be asked to reduce costs by reducing the number of sensors, in that case, a deliberate choice of a small number of sensors in combination with SD techniques would reduce the MCS sensing cost, by trading-off quality. In order to optimize such tradeoff, the MCS administrator needs a metric that relates sensing quality to the number of participants for each spatiotemporal cell in order to minimize the cost, further motivating the need for SD.

MCS system costs can have a few general forms, such as the provision of a service [4], incentive payments [7], or data consumed [12], to name a few. An example of an MCS system is one that asks users to perform a certain task, e.g., measure the noise pollution [5]. An incentive payment type of MCS cost within such a system could be a discount coupon to a specific shop or a restaurant, redeemable points in a loyalty program, credit for an online store (such as Google Play credits). Another type of MCS costs is the provision of services, such as tailored recommendations, entertainment, or even a rank in a gamified MCS scoring system. From an MCS administrator's point of view, costs could be data traffic, or the collaboration with entities that provide services to the participants. Ultimately, the costs scale up with the number of participants from the point of view of the system.

The definition of quality [2], [3], [13] depends on several parameters, and on the stakeholders' needs, where for a small number of MCS participants, the accuracy of reports is affected by the presence of heterogeneity among the sensing devices in an MCS system [4]. Such heterogeneity is a source of discrepancy in the readings reported for a specific area of interest, as the accuracy and the precision of the embedded sensors vary from one smart device to another. A sensor estimates the true value of the sensed quantity, however, its reading belongs to a sample distribution whose mean is an estimate of the true value (e.g., the normal distribution that is often observed in physical quantities) [14]. The resulting heterogeneity renders the system vulnerable to inference errors. With assumptions of a large enough data scale, these errors have a minuscule impact. However, in a *spatiotemporal space* or a cell, which may contain a small number of participants, the impact of a single error could be detrimental. Therefore, specific techniques for SD scenarios are needed for the MCS implementation in IoT systems with heterogeneous devices.

In a practical MCS scenario, the MCS system is blind to the true value and thus relies on MCS participants as a proxy to the true value. Furthermore, the system is also blind to how accurate the participants' sensors are. In combination with the scale of the data in SD for a spatiotemporal cell, the truth estimate's susceptibility to errors is greatly increased as a single incorrect measurement, an *abnormality* or an *outlier* could impair the inference, and thus the MCS system [2]. In such SD scenarios, the MCS administrator needs a metric that permits proper classification of readings under the stringent conditions of SD, and characterizes the reliability, and thus the quality, of an MCS cell's readings as a whole. It is imperative that any MCS system is capable of automatic detection and isolation of a potential error before any analysis. In this regard,

we propose novel quality metrics for MCS systems under the SD scenarios, as well as their extension to an MCS-specific coverage metric for the whole MCS area of interest. Such an MCS-specific coverage metric allows better characterization of MCS participants' presence and quality over the area of interest.

Qu *et al.* [15] tackled the same problem as this article however from the perspective of the task, rather than the perspective of the workers and their reported data, and with the assumption that quality can only be ensured by increasing the number of participants. The approach proposed in [15] aims to solve the problem of selecting prices in a posted price model using chance-constrained optimization, to minimize the total cost while maintaining robustness. However, their approach requires an abundance of data as it necessitates the presence of multiple participants performing a task, which limits its applicability to scarce scenarios.

Another related work attempted to ensure the quality by means of recruiting a cross-validating crowd, to validate the crowd sensed data's validity and rate the participants [16]. The work in [16] addresses the data credibility problem by means of recruiting a cross-validating crowd, based on their social profiles and technical expertise, for a crowdsourcing campaign on top of a crowdsensing campaign. The framework and mechanism proposed in [16], while mildly susceptible to the well-known mainstream bias, are capable of increasing data quality by means of using the cross-validating crowd, to update the crowd sensed data distribution and evaluating the participants' contributions, and thus affecting the received incentives. The contribution of [16] is interesting and practically applicable, however, it does not cater to situations where the SD problem in question, in addition to it relying on experts' subjective opinions in the evaluation of the crowd sensed data.

The contributions of this article address the event when the amount of data is too small for proper inference, or some readings come from poor or erroneous sensors. The first metric, the MAD quality metric (MAD-Q), allows the discernment of quality for a small number of participants as low as 11. The second metric, the MAD bootstrap quality metric (MADBS-Q), is more complex, but it can discern the quality for an even smaller number of participants as small as eight. Finally, we address the coverage problem, in which the quality of each cell is extended to define an overall regional quality, allowing the MCS administrator to target the recruitment process toward specific cells to achieve coverage uniformity over the whole region.

The proposed quality metrics are developed mathematically in order to show the impact of the involved parameters. Such quality metrics give the administrator a degree of control over the outlier sensitivity and outlier tolerance of the MCS system. This ultimately controls the MCS system's outlier sensitivity. The control on both parameters, outlier sensitivity and tolerance, is important for commercial systems, as their tuning impacts the relation between the desired quality, the allowed cost, and the number of participants needed within an MCS system. Moreover, this method is applicable

for MCS applications whose readings come from a symmetric distribution.

This article is structured as follows. Section II introduces the mathematical essentials required throughout this article, in particular, the median absolute deviation-based mean (MAD-mean), the statistical bootstrap, the bootlier method, and the MCS spatiotemporal cell model. Section III addresses the cell-specific quality metrics, MAD-Q, and MADBS-Q. Section IV develops an MCS-specific coverage metric while Section V describes the usage of the three proposed techniques: 1) MAD-Q; 2) MADBS-Q; and 3) the coverage metric in a simulated scenario. Section VI concludes with an overview of SD and developed techniques.

II. MATHEMATICAL ESSENTIALS AND THE SPATIOTEMPORAL MCS MODEL

A. MAD-Trimmed Mean

For a sample, $X = \{x_1, x_2, \dots, x_N\}$, the dispersion of the population can be estimated via the standard error. However, it is recommended for a large sample set and cannot be used for smaller samples, as the sample standard error is a nonrobust measure of dispersion which is vulnerable to outliers. The vulnerability of a statistic θ to outliers is expressed in terms of the *breakdown point* b_θ . The breakdown point is the proportion of outliers within a sample at which the statistic becomes blind to outliers and not robust anymore. The standard error $\bar{\sigma}_{\bar{x}}$ is based on the mean \bar{x} which has a breakdown point $b_{\bar{x}} = 0$, at which both measures, $\bar{\sigma}_{\bar{x}}$ and \bar{x} , are vulnerable to outliers unless the sample is large enough [17]. Therefore, when the sample is small, we resort to the usage of robust statistics, such as the sample median and the median absolute deviation (MAD), which both have a breakdown point of $b_{\bar{x}} = 0.5$; a point beyond which the outliers are populating the sample itself. In practical systems, it is nonrealistic to have more than 50% outliers. This implies that the MAD is a reliable practical measure of dispersion. The MAD is defined as the median of the absolute deviations from the sample median [2], expressed as

$$\text{MAD} = \text{median} \left\{ \left| x_i - \overbrace{\text{median}(X)}^{\bar{x}} \right| \right\} \quad (1)$$

where x_i is the i th sample in X , and the median is denoted by \bar{x} .

The MAD observes the deviations from the median, unlike the standard error which considers the square of the deviations from the mean. However, they are related as the MAD is a robust consistent estimator of the standard deviation [17], whereas the standard error is a nonrobust consistent estimator of the standard deviation. The MAD can be used to estimate the standard deviation as

$$\hat{\sigma}_{\text{MAD}} = \frac{1}{\Phi^{-1}(3/4)} \text{MAD} = \frac{1}{1.4826} \text{MAD} \Big|_{f(x)=N(\mu, \sigma)} \quad (2)$$

where $\hat{\sigma}_{\text{MAD}}$ is the MAD-based standard deviation estimator, $f(x)$ is the probability distribution followed by the population, $N(\mu, \sigma)$ is a normal distribution centered at μ with a spread of

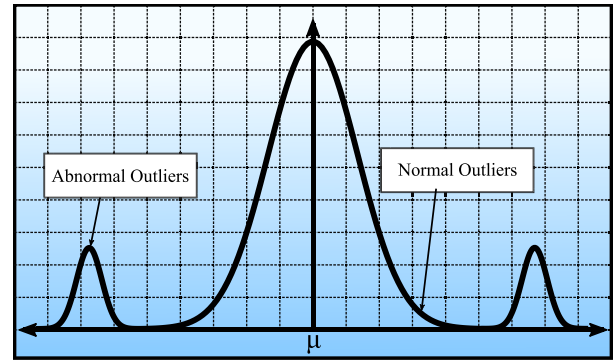


Fig. 1. Population with abnormal outliers (lower and upper).

σ , and $\Phi^{-1}(3/4)$ is the quantile function (which is the inverse cumulative distribution function) at 75%.

The MAD covers the median distance from the sample median \bar{x} , which is located at the center having it spans the portion between 25% and 75%, and thus the reason why the quantile function is evaluated at that point. For any symmetric distribution, the MAD covers 50% spanning from the left to the right of \bar{x} . For any normal distribution $N(\mu, \sigma)$, $\Phi^{-1}(3/4) = 1.4826$. As a result, the relation of the MAD and the standard deviation was derived, nevertheless, this relation is valid for any well-defined distribution $f(x)$ [17].

The MAD inherits the robustness of the median, which makes it useful for the detection and removal of outliers. However, the definition of outliers is vague and depends on the application in hand. For MCS systems, physical sensor measurements tend to follow an even symmetric distribution (i.e., defined by an even function), such as the normal distribution. The normal distribution includes *normal*¹ outliers that rise from the extremities of the distribution. However, there are *abnormal outliers* which lie far outside the three-sigma range, defined by the three-sigma rule ($\pm 3\sigma$) which is valid for all symmetric unimodal distributions [18]. In case the distribution is an even multimodal distribution, other criteria would be required for the definition of expected outliers, such as those based on the interquartile range. These abnormal outliers come from an *unexpected phenomenon* that is not modeled by the normal distribution such as the tail modalities in Fig. 1. The distribution in Fig. 1 is a multimodal distribution which is a superposition of abnormal distributions representing an unexpected phenomenon and an expected distribution of a sensed physical measurement.

In an MCS system, these abnormal outliers need to be isolated as they pose a hindrance that impairs the MCS system. The MAD allows the removal of outlier samples present in X , producing a MAD-trimmed sample, denoted X_{MAD} , where only values from within the range of the normal distribution are considered [19], expressed as

$$X_{\text{MAD}} = \left\{ X : x_j \in \left[\bar{x} \pm \frac{\lambda \text{MAD}}{\delta \hat{\sigma}_{\text{MAD}}} \right] \right\} \quad (3)$$

¹Normal here refers to expected outliers coming from a distribution itself, not necessarily a normal distribution.

Algorithm 1 Computing the MAD-Mean**Input:** A sample: $X = \{x_1, x_2, \dots, x_n\}$ **Output:** MAD-mean: \bar{x}_{MAD} Initialize : λ

```

1: MAD(X) = median(median(X) - X)
2: for all  $x_i$  do
3:   if  $x_i \notin [\text{median}(X) \pm \lambda \text{MAD}]$  then
4:      $X_{\text{MAD},o} = \text{append}(x_i, X_{\text{MAD},o})$ 
5:   else
6:      $X_{\text{MAD}} = \text{append}(x_i, X_{\text{MAD}})$ 
7:   end if
8: end for
9: return  $\bar{x}_{\text{MAD}} = \text{mean}(X_{\text{MAD}})$ 

```

where λ is how many multiples of MAD far from the median are the nonoutlier samples. The product λMAD can be related to the consistent estimation of the standard deviation $\delta \hat{\sigma}_{\text{MAD}}$, to define it as multiples of σ deviations from the mean, where $\delta = \lambda \Phi^{-1}(3/4)$. λ is an important parameter as it controls the range for which values, even if outliers, are indeed in X_{MAD} .

Nevertheless, normal outliers are *expected*, unlike *abnormal outliers* which are beyond the three-sigma range. Thus, for the administrator to ensure proper estimation of the true value μ that considers normal outliers, λ , such as $\lambda = 4$ corresponding to the three-sigma range, for example, could be selected. As a result, only values that belong to the interval $\bar{x} \pm \lambda \text{MAD}$ will be considered in X_{MAD} . Therefore, the MAD-mean can be defined as

$$\begin{aligned} \bar{x}_{\text{MAD}} &= \frac{1}{N - N_{\text{MAD},o}} \sum_{i=1}^{N - N_{\text{MAD},o}} x_{\text{MAD},i} \\ &= \frac{1}{N_{\text{MAD}}} \sum_{i=1}^{N_{\text{MAD}}} x_{\text{MAD},i} \end{aligned} \quad (4)$$

where N is the sample size, $N_{\text{MAD},o}$ is the number of outliers outside the $\bar{x} \pm \lambda \text{MAD}$ range, $x_{\text{MAD},i}$ is the i th element in the MAD-trimmed sample X_{MAD} , and $N_{\text{MAD}} = |X_{\text{MAD}}|$. Algorithm 1 shows the steps in which the MAD-mean is computed.

Notice that the value of λ provides a degree of freedom for system administrators to decide the range of the measurements and the consideration of outliers. λ is a parameter that controls the degree of outlier tolerance. For absolute outlier intolerance, $\lambda = 2$ is sufficient as the resulting $\delta < 3$ is within the $3\text{-}\sigma$ range. On the other hand, for outlier tolerance, a choice of $\lambda \geq 4$ is recommended as $\delta > 3$ is beyond the $3\text{-}\sigma$ range. Since λ has no upper bound, it needs to be carefully selected by the administrator to avoid including abnormal outliers. However, the breakdown point of the MAD-mean $b_{\bar{x}_{\text{MAD}}}$, is impacted by the choice of λ ; in particular, $b_{\bar{x}_{\text{MAD}}} < 0.5$ for $\lambda > 4$, as the MAD-mean is based on the median. The MAD-based outlier detection is a technique that is especially useful for small-sample sizes, which makes it useful for the small-sample scenarios present in MCS systems' spatiotemporal cells [20]. The MAD estimator is a very good option that is further characterized by low complexity. Its parameter λ

provides an interesting tradeoff between accuracy and robustness, which is very important for commercial systems. Later, simulations show deeper considerations of this tradeoff.

B. Nonparametric Bootstrap

The nonparametric bootstrap [21] is a population-agnostic method, which allows the construction of sample distributions without prior assumptions about the population's distribution. Its numerical nature allows it to achieve its target by means of sampling with replacement a large number of times B . The nonparametric bootstrap resamples the original sample $X = \{x_1, x_2, \dots, x_N\}$, generating B resamples, $X_b = \{x_{b1}, x_{b2}, \dots, x_{bN}\}$, where x_{bi} can appear more than once in X_b , and x_{bi} is a random variable that samples uniformly from X . From each bootstrapped sample set X_b , a bootstrap statistic θ_b can be acquired. To construct the sample distribution of a statistic θ , each X_b is employed to compute B θ_b bootstrap statistics, collected in a vector of bootstrap statistics $\theta^* = \{\theta_1, \theta_2, \dots, \theta_B\}$, whose histogram represents the sample distribution.

Despite its numerical complexity for a large B , the bootstrap is useful for crowdsensing applications where the scenario is sparse [1]. It is of particular usefulness for small-sample cases due to the fact that its fair resampling has a low probability P of selecting a homogeneous sample (i.e., all samples in X_b being exactly the same); which is obtained as

$$P(X_b = \{x_i, x_i, \dots, x_i\}) = [1 - (1 - 1/N)^N]^N \quad (5)$$

which is the probability of a binomial case. Inspecting (1), out of B resamples, only 3.45% will be extremely biased for a sample of size $N = 8$.

The usefulness of the bootstrap in obtaining the quality of small-sample scenarios comes from the presence of outliers. Since resampling is uniform, the outliers presence is promoted to have a probability of $1/N$, like any other sample x_i in X . Further analysis of the binomial probability shows that the probability of an element's inclusion is

$$P(x_i \in X_b) = 1 - (1 - 1/N)^N \quad (6)$$

which converges for a large N to 67%, i.e., each sample is present in 67% of the B X_b sample sets.

Furthermore, since B is a sufficiently large number, the properties of the central limit theorem (CLT) are also applicable to B θ_b statistics obtained. This property is of particular usefulness in the discussion of central measures of tendency (mean, median, mode, etc.), as it causes the bootstrap for the mean to follow the normal distribution. However, due to the robustness of the MAD-mean, the resulting distribution is inherently multimodal, even if it is seemingly normal, as the employment of the median in the MAD-mean's trimming process generates multimodalities. It limits the set of *medians* to be selected to a discrete number of medians (for an odd N , there are N possible medians, and for an even N , there are $N^2 - \sum_{i=0}^{N-1} i$ possible medians).

C. Bootlier

The bootlier, a graphical tool developed in [22], exploits the outlier promotion flaw in the bootstrap to detect outliers by computing the difference statistic between the mean \bar{x} and the trimmed mean \bar{x}_k . The trimmed mean is a robust estimate of centrality with a breakdown point $b_{\bar{x}_k} = k\%$. In the bootlier, the difference between the mean and the trimmed mean, $\bar{x} - \bar{x}_k$ is bootstrapped to construct a sample distribution. The histogram of the resulting sample distribution is called a bootlier plot, which shows the impact of outliers present in a sample. Singh and Xie [22] investigated the bootlier's multimodality and smoothness, and developed "the bootlier index" as a measure to quantify its smoothness. However, their work heavily depends on human intuition for assessing the quality of a sample, and for the detection of outliers, in addition to the lack of a straight metric for the evaluation of a sample set's quality.

We analyzed the bootlier plot, by considering an infinite sample size ($N \rightarrow \infty$), and found that its ideal reference is in fact the distribution $N(0, 0)$, which is best described as the Dirac delta impulse distribution [23]. We also found that the bootlier, from a signal processing perspective, is a *superposition of leaking impulses*, which are very distinct when an outlier is present as it introduces another impulse at $x_o - \bar{x}_k$, where x_o is an outlier sample and \bar{x}_k is the trimmed mean of the sample. The source of the leakage is the variation within the original sample, which causes the resulting sample distribution, i.e., the bootlier, to be smooth. The bootlier statistic hides the resulting multimodalities which combine to form a seemingly normal distribution around 0, that should ideally—for an absolutely perfect sample—be an impulse. Section III introduces the mean-MAD mean-trimmed mean (MMTM) statistic and discusses the modalities in its distribution based on the trimmed mean and the MAD-mean.

D. MCS Spatiotemporal Model

In order to achieve efficient characterization of an area of interest under an MCS system, it needs to be appropriately divided into *geofences* or *cells* to which MCS participants are assigned. In addition to such a spatial division, a temporal division is also required for proper real-time sensing. These divisions in space and time can be represented as a *spatiotemporal diagram* as shown in Fig. 2, based on that conceived in [9]. An MCS administrator could divide the spatiotemporal space in a manner that satisfies the objective of the system while being consistent with the spatial and temporal versions of the sampling theorem [24]. MCS participants are assigned tasks according to their availability, for example, if the MCS administrator requires temperature values to be sensed by MCS participants, the m th cell will have a set of readings, X_m

$$X_m = \{x_{m,1}, x_{m,2}, \dots, x_{m,N_m}\} \quad (7)$$

where the reading obtained by the i th participant is modeled as a random variable $x_{m,i}$, and N_m is the number of participants who executed the task in the m th spatiotemporal cell. Each cell corresponds to a 3-tuple (a, b, c) , where a maps to the x -location, b maps to the y -location, and c maps to the c th sensing cycle (temporal location).

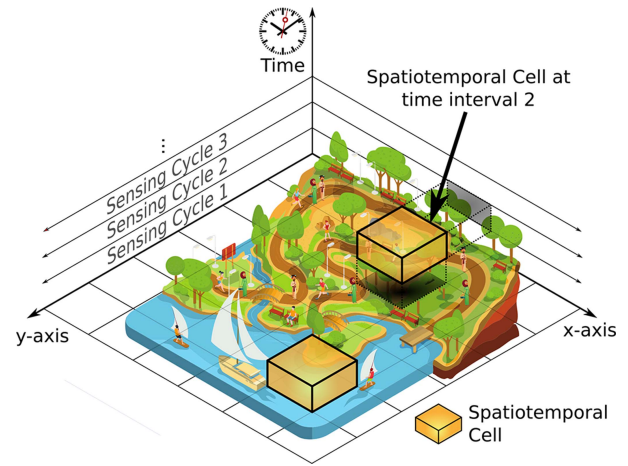


Fig. 2. Spatiotemporal diagram.

The true value of the sensed quantity in the m th cell μ_m is estimated from the sample obtained in (7) by computing the mean, \bar{x}_m

$$\hat{\mu}_m = \bar{x}_m = \text{mean}(X) = \frac{1}{N_m} \sum_{i=1}^{N_m} x_{m,i}. \quad (8)$$

The standard deviation for the m th cell can be estimated by means of the sample error

$$\hat{\sigma}_m = \sqrt{\frac{\sum_{i=1}^{N_m} (x_{m,i} - \hat{\mu}_m)^2}{N_m - 1}}. \quad (9)$$

However, this is only valid when the sample is large enough that $N_m \geq 30$ samples are suitable to represent the population. At such a size, the CLT is viable as the combination of all $x_{m,i}$ random variables would ultimately conjure a normal distribution [25]. While this assumption is very useful, it is not the case when $N \ll 30$, when there are barely enough measurements within a spatiotemporal cell. The problem, then, is of the SD scale. An MCS system should be capable of inferring as much as possible under the stringent conditions of SD. Equations (8) and (9) become less accurate as the presence of $N_{m,o}$ outliers $x_{m,o}$ would throw off the estimation.

Outlier values $x_{m,o}$ can be viewed as an estimation that is offset from μ_m

$$x_{m,o} = \mu_m + \Theta_{m,o} \sigma_m \quad (10)$$

where $\Theta_{m,o} \in \mathbb{R}$ is the outlier deviation factor (ODF), which we define as the multiple of standard deviations σ_m , the o th outlier value $x_{m,o}$ is far from the mean μ_m . This allows the expression of (10) as

$$\hat{\mu}_m = \mu_m + \bar{\Theta}_m \sigma_m \quad (11)$$

where $\bar{\Theta}_m = \sum_{o=1}^{N_{m,o}} \Theta_{m,o} / N_{m,o}$ is the average ODF (AODF).

The mean \bar{x}_m can then be rewritten as

$$\begin{aligned} \bar{x}_m &= \frac{1}{N_m} \sum_{i=1}^{N_m} x_{m,i} \\ &= \frac{1}{N_m} \left[\sum_{i=1}^{N_m - N_{m,o}} x_{m,i} + \sum_{j=N_{m,o}}^{N_m} x_{m,j} \right] \end{aligned}$$

$$\begin{aligned}
&= \frac{N_m - N_{m,o}}{N_m} \hat{\mu}_{\bar{x}_m} + \frac{N_{m,o}}{N_m} (\mu_m + \bar{\Theta}_{\bar{x}_m} \sigma_m) \\
&= \hat{\mu}_{\bar{x}_m} + \underbrace{\frac{N_{m,o}}{N_m} (\mu_m - \hat{\mu}_{\bar{x}_m})}_{\text{sample size error due to outliers}} + \underbrace{\frac{N_{m,o}}{N_m} \bar{\Theta}_{\bar{x}_m} \sigma_m}_{\text{deviation error due to outliers}} \quad (12)
\end{aligned}$$

where $\hat{\mu}_{\bar{x}}$ corresponds to the mean's estimate of the true value μ , and $\bar{\Theta}_{\bar{x}}$ corresponds to the average ODF of the mean.

The second and third terms of (12) correspond to the errors induced by the presence of the outliers. The second is a sample size error due to the presence of outliers that inflates the denominator N_m , while the third is a deviation error due to distance between the outliers and μ . The third term contributes the majority of the error in the nonrobust mean.

Thus, the error of the mean, denoted $\xi_{\bar{x}_m}$, can be defined as the distance between the mean \bar{x}_m and the true value μ_m , where

$$\begin{aligned}
\xi_{\bar{x}_m} &= \bar{x}_m - \mu_m \\
&= \underbrace{(\hat{\mu}_{\bar{x}_m} - \mu_m)}_{\text{estimation error}} + \frac{N_{m,o}}{N_m} \left[\underbrace{(\mu_m - \hat{\mu}_{\bar{x}_m})}_{\text{sample size error due to outliers}} + \underbrace{\bar{\Theta}_{\bar{x}_m} \sigma}_{\text{deviation error due to outliers}} \right]. \quad (13)
\end{aligned}$$

The errors of the trimmed mean \bar{x}_k and the MAD-mean \bar{x}_{MAD} can be derived in a manner similar to that employed in (12) and (13).

Equation (13) can be rewritten as

$$\xi_{\bar{x}_m} - (\hat{\mu}_{\bar{x}_m} - \mu_m) = \frac{N_{m,o}}{N_m} [(\mu_m - \hat{\mu}_{\bar{x}_m}) + \bar{\Theta}_{\bar{x}_m} \sigma] \quad (14)$$

indicating that for a more accurate estimation, the errors on the right-hand side need to be minimized. This can be achieved by two options: 1) reducing the proportion of outliers, $N_{m,o}/N_m \rightarrow 0$, which translates to outlier removal or 2) increasing the sample size by taking $N_m \rightarrow \infty$. The latter case is when the data are large scale and CLT is applicable, while the former case is the problem faced under SD. In SD, the proportion of outliers cannot easily become zero since $N_{m,o} \geq 1$ would render the proportion $N_{m,o}/N_m$ significant. Due to the nature that the MCS data comes from participants whose sensors are different, prone to error, or even maliciousness, its susceptibility to outlier errors is more significant when SD is involved.

III. CELL-SPECIFIC QUALITY METRICS

In the previous section, we saw how the MAD-mean is a robust centrality estimate for the true value μ . We have also covered the basics of the statistical bootstrap. Now, we employ both tools to develop a cell-specific quality metric that accounts for both range and accuracy. For range, a suitable solution is the difference between the mean and the MAD-mean $\bar{x} - \bar{x}_{\text{MAD}}$, as the administrator's choice of λ allows the extension of the MAD's trimming range, either to include or exclude normal outliers. For accuracy, the difference between the mean and the trimmed mean $\bar{x} - \bar{x}_k$, provides a degree of outlier intolerance, as the trimmed mean suppresses all

outliers, including normal outliers, because it trims at least a single outlier as long as $k \neq 0$. However, in case the case that $kN < N_o$, some outliers might be present in the trimmed mean term. However, for the scale of SD, this is unlikely to happen as the proportion of k is larger for SD ($k = 10\%$ trims at most 1–3 samples for the range 8:30, this implies that outliers will be filtered for $k > N_o/N$). Therefore, in this section, we introduce the hybrid MMTM statistic and use it as the basis upon which we formulate two quality metrics for SD in MCS. One is computationally simple for samples as small as 11 by measuring the closeness of the MMTM statistic's magnitude to zero, while the other is slightly computationally demanding yet it extends the characterization of MMTM quality for sample sets as small as $N = 8$, by measuring the similarity between the bootstrap distribution of the MMTM statistic and the Dirac delta distribution [23]. The developed quality metrics allow the MCS system administrators to control the tradeoff between accuracy, range, and the number of MCS participants per cell, by tweaking an outlier sensitivity variable. This flexibility is significant in the design of commercial systems as it establishes a relation between cost, i.e., the number of participants, and quality.

A. MMTM Statistic and the MAD-Q Quality Metric

The selection of a suitable point in the range-accuracy tradeoff depends on the scenario, the application, as well as economic concerns. In this section, we describe the MAD-Q that allows the MCS administrator to include or exclude normal/abnormal outliers by means of a sensitivity parameter, denoted β . The MMTM statistic θ_{MMTM} is defined as

$$\begin{aligned}
\theta_{\text{MMTM}} &= \beta(\bar{x} - \bar{x}_k) + (1 - \beta)(\bar{x} - \bar{x}_{\text{MAD}}) \\
&= \bar{x} - \beta\bar{x}_k - (1 - \beta)\bar{x}_{\text{MAD}} \quad (15)
\end{aligned}$$

where $\beta \in [0, 1]$, where \bar{x}_k is the k -trimmed mean. The k -trimmed mean, discussed in Section III-C, removes the rounded $k\%$ highest and/or lowest samples of X prior the computation of the mean.

The θ_{MMTM} statistic is the difference between the mean and a weighted average of the more robust centrality estimates: the trimmed mean \bar{x}_k and the MAD-Mean \bar{x}_{MAD} . θ_{MMTM} can be rewritten in terms of the errors of the mean $\xi_{\bar{x}}$, the trimmed mean $\xi_{\bar{x}_k}$, and the MAD-mean $\xi_{\bar{x}_{\text{MAD}}}$; all errors can be derived in a manner similar to that for (13). θ_{MMTM} becomes

$$\begin{aligned}
\theta_{\text{MMTM}} &= \mu + \xi_{\bar{x}} - \beta(\mu + \xi_{\bar{x}_k}) - (1 - \beta)(\mu + \xi_{\bar{x}_{\text{MAD}}}) \\
&= \xi_{\bar{x}} - \beta\xi_{\bar{x}_k} - (1 - \beta)\xi_{\bar{x}_{\text{MAD}}} \\
&= \underbrace{\left(\hat{\mu}_{\bar{x}} - \beta\hat{\mu}_{\bar{x}_k} - (1 - \beta)\hat{\mu}_{\bar{x}_{\text{MAD}}} \right)}_{\text{estimation errors}} \\
&\quad + \underbrace{\left(\frac{N_o}{N}(\mu - \hat{\mu}_{\bar{x}}) - \beta \frac{N_{k,o}}{N_k}(\mu - \hat{\mu}_{\bar{x}_k}) - (1 - \beta) \frac{N_{\text{MAD},o}}{N_{\text{MAD}}}(\mu - \hat{\mu}_{\bar{x}_{\text{MAD}}}) \right)}_{\text{sample size errors due to the outliers' presence}} \\
&\quad + \underbrace{\sigma \left(\frac{N_o}{N} \bar{\Theta}_{\bar{x}} - \beta \frac{N_{k,o}}{N_k} \bar{\Theta}_{\bar{x}_k} - (1 - \beta) \frac{N_{\text{MAD},o}}{N_{\text{MAD}}} \bar{\Theta}_{\bar{x}_{\text{MAD}}} \right)}_{\text{deviation errors due the outliers' distance}} \quad (16)
\end{aligned}$$

Algorithm 2 Algorithm for MAD-Q Quality Assessment**Input:** Readings from N sensors: $X = \{x_1, x_2, \dots, x_n\}$ **Output:** Quality of Source: Q_{MAD} *Initialize* : $\beta, \lambda, k, \gamma_{\text{MAD}}$

- 1: $\bar{x}_{\text{MAD}} = \text{mean}(X_{\text{MAD}})$
- 2: $\bar{x} = \text{mean}(X)$
- 3: $\bar{x}_k = \text{trim_mean}(X)$
- 4: $Q_{\text{MAD}} = \log_{\gamma_{\text{MAD}}} [(\bar{x} - \beta\bar{x}_k - (1 - \beta)\bar{x}_{\text{MAD}})^{-1}]$
- 5: **if** $Q_{\text{MAD}} \rightarrow \infty$ **then**
- 6: $Q_{\text{MAD}} = 6$
- 7: **end if**
- 8: **return** Q_{MAD}

where $\hat{\mu}_k$ and $\hat{\mu}_{\text{MAD}}$ correspond to the estimation of the trimmed mean and the MAD-mean; $N_{k,o}$ and $N_{\text{MAD},o}$ correspond to the outliers present in the trimmed mean set X_k and the MAD-mean X_{MAD} ; N_k and N_{MAD} correspond to $|X_k|$ and $|X_{\text{MAD}}|$; and $\bar{\Theta}_k$ and $\bar{\Theta}_{\text{MAD}}$ correspond to the average deviation due to the presence of post-trimming outliers in X_k and X_{MAD} .

Equation (16) comprises the sum of three errors. The first being between the outlier-free estimation errors of the mean, the trimmed mean, and the MAD-mean; the second being the sample size error due to the outlier samples present in X ; and the third is the error due to the presence of outliers in the mean \bar{x} and the trimmed mean \bar{x}_k , and the \bar{x} and the MAD-mean \bar{x}_{MAD} . Furthermore, if λ , the parameter that defines the MAD-mean's outlier removal range, is less than 2, the MAD-mean becomes very aggressive in selecting which samples are nonoutliers, permitting it to be very robust with $b_{\bar{x}_{\text{MAD}}} \rightarrow 0$. Ideally, for a perfect sample, θ_{MMTM} should be zero. However, this is impossible to happen empirically. Thus, we define quality as the closeness of this value to zero, and we formulate the MAD-Q quality Q_{MAD} based on the θ_{MMTM} statistic as

$$Q_{\text{MAD}} = \log_{\gamma_{\text{MAD}}} \left(\frac{1}{\theta_{\text{MMTM}}} \right) \\ = \log_{\gamma_{\text{MAD}}} \left(\frac{1}{\bar{x} - \beta\bar{x}_k - (1 - \beta)\bar{x}_{\text{MAD}}} \right) \quad (17)$$

where γ_{MAD} is a scaling factor. Algorithm 2 shows the steps to compute the quality of an MCS sample. γ_{MAD} can be changed as needed to control the scale, however, the saturation value in line 6 in Algorithm 2 needs to be adjusted accordingly.

B. MADBS-Q Bootstrap-Based Quality Metric

In the previous section, we described how the nonparametric bootstrap is useful for acquiring the sample distribution of a statistic from small-sample sizes. We combine the bootstrap with the θ_{MMTM} statistic to develop the MADBS-Q, Q_{BS} that performs in a manner similar to Q_{MAD} , but capable of assessing to even smaller sample sizes [26]. Rather than formulating MADBS-Q based on the θ_{MMTM} statistic, it is formulated based on the sample distribution of θ_{MMTM} . The bootlier, discussed in Section II, is a tool that employs human intuition in its assessment on top of being too sensitive to disregard *normal* outliers as abnormal. This “oversensitivity” to outliers

Algorithm 3 Algorithm for MADBS-Q Quality Assessment**Input:** Readings from N sensors: $X = \{x_1, x_2, \dots, x_n\}$ **Output:** Quality of Source: Q_{BS} *Initialize* : $B, \beta, k, \lambda, \gamma_{\text{BS}}$

- 1: $\theta_{\text{MMTM}}^* = \text{bootstrap}(X, B, \theta_{\text{MMTM}})$
- 2: $\mu_{\text{MMTM}} = \text{mean}(\theta_{\text{MMTM}}^*)$
- 3: $\sigma_{\text{MMTM}} = \text{variance}(\theta_{\text{MMTM}}^*)$
- 4: **return** $Q_{\text{BS}} = \frac{1}{2} [\log_{\gamma_{\text{BS}}} (\mu_{\text{MMTM}}^{-1}) + \log_{\gamma_{\text{BS}}} (\sigma_{\text{MMTM}}^{-2})]$

impacts the perception of quality, which is an important factor in reducing the MCS system's costs. To benefit from the bootstrap, we define the quality metric to relate to the closeness of the θ_{MMTM} sample distribution to the ideal impulse (Dirac delta) best captured in terms of location and spread. We define the bootstrap-based quality Q_{BS} as

$$Q_{\text{BS}} = \frac{1}{2} [\log_{\gamma_{\text{BS}}} (\mu_{\text{MMTM}}^{-1}) + \log_{\gamma_{\text{BS}}} (\sigma_{\text{MMTM}}^{-2})] \quad (18)$$

where μ_{MMTM} and σ_{MMTM}^2 are the mean and variance of the θ_{MMTM} sample distribution, respectively. γ_{BS} is a scaling factor. The mean of the resulting θ_{MMTM} distribution indicates its location, and the variance indicates its spread. However, it is best compared with the Dirac delta impulse.

The bootstrap-based quality Q_{BS} is thus defined as the average of the logs of the θ_{MMTM} 's mean and variance. The closer they are to zero, the higher the quality. This renders Q_{BS} as an absolute quality metric as it is free from any reference distributions (other than the Dirac delta) or thresholds. Algorithm 3 shows the steps required to obtain the bootstrapped sample quality Q_{BS} . In Algorithm 3, θ_{MMTM}^* refers to the vector containing the bootstrapped resamples of θ_{MMTM} .

While both metrics, MAD-Q and MADBS-Q, are interchangeable, MADBS-Q is more reliable when data sets are very small, but it comes with a computational cost, mainly from the computation of the θ_{MMTM} statistic B times. For the sample mean, the complexity of the bootstrap is #P-hard, which can be solved in polynomial time [26]. MAD-Q, on the other hand, is cheaper to compute and serves data sets which are not as small as those encountered by MADBS-Q. Ultimately, MADBS-Q is derived from, and extends, MAD-Q.

C. Impact of Control Parameters

The MMTM quality metrics, Q_{MAD} and Q_{BS} , are based on θ_{MMTM} , which is in turn based on the parameters β and λ . β , referred to as sensitivity from (15), allows the control of whether the trimmed mean \bar{x}_k or the MAD-mean \bar{x}_{MAD} is more significant in θ_{MMTM} . λ , on the other hand, is only specific to the MAD-mean, as it controls its filtering of outliers. In θ_{MMTM} , the presence of the term \bar{x}_k causes θ_{MMTM} to maintain a narrow distribution, as outliers are always filtered for $k > N_o/N$, while the presence of \bar{x}_{MAD} allows the control of the range based on the choice of λ . λ , as a parameter, allows the administrator to completely remove outliers by selecting $\lambda = 3$, or tolerate outliers where $\lambda = 4$ tolerates normal outliers and $\lambda > 4$ would tolerate abnormal outliers depending on their ODF, Θ .

As a result, θ_{MMTM} can be viewed as a sum of two distributions

$$\theta_{\text{MMTM}} = \underbrace{\theta_A}_{\bar{x}} - \underbrace{\theta_B}_{\beta\bar{x}_k + (1-\beta)\bar{x}_{\text{MAD}}} \quad (19)$$

where θ_A is the sample distribution of the mean, and θ_B represents a mixture of the sample distributions of the trimmed mean and the MAD-mean, with β as a mixing parameter.

The distribution of θ_{MMTM} , as a result, is a multimodal distribution with θ_A contributing a mode, and θ_B contributing either a mode or two based on the mixing parameter of the mixture between \bar{x}_k and \bar{x}_{MAD} . However, it is hard for θ_B to be bimodal as both estimates are robust measures of centrality, except for the case in which $\lambda > 4$ which is no longer tolerant. For $\lambda < 4$, it is safe to assume that both, the trimmed mean and the MAD-mean, follow a distribution that represents the distribution of a sample mean: the normal distribution. For a mixture of two equally weighted normal distributions with similar variability to be bimodal, the difference between their means has to be greater than the sum of both their standard errors (or deviations, if not sample distributions) [27]. For θ_B , the conditions are necessary: 1)

$$|\mu_{\bar{x}_k} - \mu_{\bar{x}_{\text{MAD}}}| > \bar{\sigma}_{\bar{x}_k} + \bar{\sigma}_{\bar{x}_{\text{MAD}}} \quad (20)$$

and 2) $\beta = 0.5$ for them to be equally weighted. Since both are robust measures of centrality, then it is always the case that $|\mu_{\bar{x}_k} - \mu_{\bar{x}_{\text{MAD}}}|$ is less than $\bar{\sigma}_{\bar{x}_k} + \bar{\sigma}_{\bar{x}_{\text{MAD}}}$, where $\bar{\sigma}_{\bar{x}_k}$ and $\bar{\sigma}_{\bar{x}_{\text{MAD}}}$ refer to the sample distribution of the trimmed mean and the MAD-mean's standard errors, respectively. Thus, bimodality is unlikely. However, for unimodality to occur in a mixed distribution, the ratio of the difference between the centers to the double of the product of the standard deviation has to be less than or equal to 1 [28], i.e.,

$$\frac{|\mu_{\bar{x}_k} - \mu_{\bar{x}_{\text{MAD}}}|}{2\sqrt{\bar{\sigma}_{\bar{x}_k}\bar{\sigma}_{\bar{x}_{\text{MAD}}}}} \leq 1. \quad (21)$$

However, this implies that there is a region between unimodality and bimodality where

$$\bar{\sigma}_{\bar{x}_k} + \bar{\sigma}_{\bar{x}_{\text{MAD}}} > |\mu_{\bar{x}_k} - \mu_{\bar{x}_{\text{MAD}}}| > 2\sqrt{\bar{\sigma}_{\bar{x}_k}\bar{\sigma}_{\bar{x}_{\text{MAD}}}}. \quad (22)$$

In this region, the distribution of θ_B is neither unimodal nor bimodal, but trimodal [29]. However, this trimodality is unlikely to happen due to the fact that both sample distributions, \bar{x}_k and \bar{x}_{MAD} are robust estimates of the true mean μ , where the difference, $|\mu_{\bar{x}_k} - \mu_{\bar{x}_{\text{MAD}}}|$ is not in the interval $(2\sqrt{\bar{\sigma}_{\bar{x}_k}\bar{\sigma}_{\bar{x}_{\text{MAD}}}}, \bar{\sigma}_{\bar{x}_k} + \bar{\sigma}_{\bar{x}_{\text{MAD}}})$.

Moreover, the standard deviation of the sample distribution of the trimmed mean $\sigma_{\bar{x}_k}$ (not to be confused with the standard error of the trimmed mean $\bar{\sigma}_{\bar{x}_k}$) can be defined as a distribution [30]

$$\sigma_{\bar{x}_k} = \sqrt{\frac{\mathbb{E}\left[(X - \mu)^2 \mathbb{I}_{(\Phi^{-1}[k/2], \Phi^{-1}[1-k/2])}(X)\right]}{1 - k}} \quad (23)$$

where \mathbb{E} is the expectation operator, and $\mathbb{I}_{(a,b)}$ is the identity function over the interval $[a, b]$. Product with the identity function effectively trims the tails of the distribution of X .

Similarly, the standard deviation of the sample distribution of the MAD-mean $\sigma_{\bar{x}_{\text{MAD}}}$ (not to be confused with the standard error of the MAD-mean $\bar{\sigma}_{\bar{x}_{\text{MAD}}}$) can also be defined as a distribution

$$\sigma_{\bar{x}_{\text{MAD}}} = \sqrt{\frac{\mathbb{E}\left[(X - \mu)^2 \mathbb{I}_{(\Phi^{-1}[\frac{2-\lambda}{4}], \Phi^{-1}[\frac{2+\lambda}{4}])}(X)\right]}{\left(1 - \Phi^{-1}\left[\frac{2-\lambda}{4}\right]\right) - \left(1 - \Phi^{-1}\left[\frac{2+\lambda}{4}\right]\right)}}. \quad (24)$$

Since the standard errors $\bar{\sigma}_{\bar{x}_k}$ and $\bar{\sigma}_{\bar{x}_{\text{MAD}}}$ are estimates of $\sigma_{\bar{x}_k}$ and $\sigma_{\bar{x}_{\text{MAD}}}$, respectively, it can be seen that the multimodality of θ_B is mainly controlled by the administrator's choice of λ rather than the choice of k , as a choice of $k \neq 0$ will always trim at least a single outlier.

Consequently, θ_{MMTM} have either a bimodal or a trimodal distribution (with three modes) as θ_A contributes an extra mode due to the nonrobust mean, while θ_B could be unimodal or trimodal, depending on the choice of λ . There are five possible scenarios for the sample distribution of θ_{MMTM} based on the ratio of outliers present, N_o/N .

- 1) No outliers are present ($N_o/N = 0$) and all measures of centrality are close to each other and θ_{MMTM} is unimodal.
- 2) The outlier ratio is less than the breakdown point of the MAD-mean term ($b_{\bar{x}_{\text{MAD}}} > N_o/N > 0$) and $\lambda \leq 2$, which will cause θ_B to be unimodal, thus resulting in a bimodal θ_{MMTM} due to the nonrobustness of θ_A .
- 3) The outlier ratio is less than the breakdown point of the MAD-mean term and $4 > \lambda > 2$, which will cause θ_B to be bimodal due to the presence of normal outlier, however that bimodality will not be apparent (i.e., it will not cause a local minima between the two modes), and therefore θ_{MMTM} will be trimodal. However, this trimodality will be a superposition of leaking impulses as previously mentioned in Section II-C, where the leakage covers the mode due to the trimmed mean term.
- 4) The outlier ratio is less than the breakdown point of the MAD-mean term and $\lambda > 4$, which will cause θ_B to be bimodal due to the presence of abnormal outliers, and thus θ_{MMTM} will be trimodal.
- 5) The outlier ratio is more than the breakdown point. In this case, the outliers are no longer outliers and the modality of θ_{MMTM} is unknown as the variance would be large due to the spread of the values.

IV. MCS-SPECIFIC COVERAGE METRIC

While MCS systems can function without full coverage, it is important for MCS administrators to be able to enhance coverage by identifying which spatiotemporal cells are below a desired coverage threshold. In this section, we provide a simple scheme that builds upon the quality metrics, discussed in Section III, to evaluate the quality of the area as a whole during a specific sensing cycle. This quality evaluation comes after the cell-specific quality has been evaluated. The method described in this section provides the MCS administrators with two quantities, MAD-Q and MADBS-Q, that allow the characterization of MCS coverage quality. The first definition is an indicator of how the average quality is over the whole space,

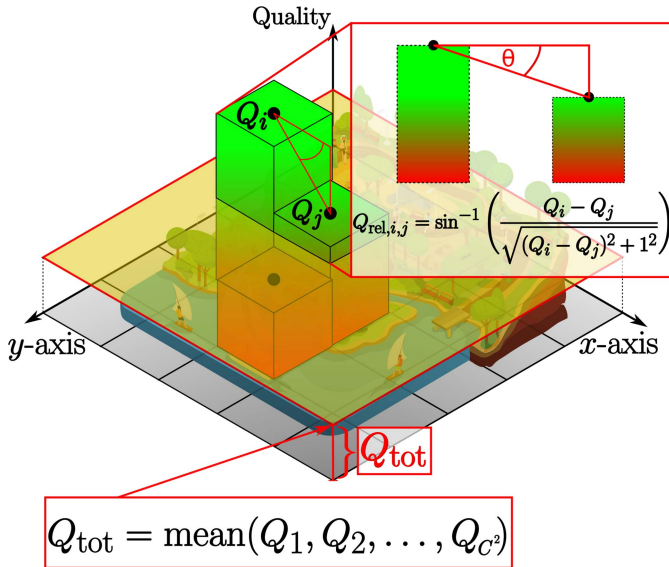


Fig. 3. Description of MCS-specific coverage metric over a $C \times C$ square grid.

while the second definition is a detailed map of which cells have a higher quality relative to their neighboring cells. By using the second definition, the MCS system can recruit participants from cells with oversatisfied quality and ask them to move to cells with lower quality and perform a required MCS task there. By having a metric that indicates where participants need to be moved, the MCS system becomes capable of achieving uniform MCS coverage as a part of participatory sensing, unlocking a scheme in which participants could be asked to (or incentivized to) voluntarily move from a place to another. Coverage, in that sense, extends the definition of the MAD-Q and MADBS-Q quality metrics, or any other cell-specific metric, over the area of interest, while uniformity refers to having good quality overall. While MCS coverage is independent of both MAD-Q and MADBS-Q, it is capable of augmenting both techniques for more control over participatory sensing. Fig. 3 provides a description of the developed coverage metric.

A. Overall Coverage Quality

For algorithmic convenience, the cells in the area of interest are assumed to be equally spaced resulting in a $C \times C$ square grid. For each cell throughout the grid, the algorithms covered in Section III return the quality of each cell. This allows the definition of an overall quality metric, Q_{tot}

$$Q_{\text{tot}} = \text{mean}(\mathbf{Q}_{\text{map}}) = \text{mean} \begin{bmatrix} Q_{1,1} & Q_{1,2} & \dots & Q_{1,C} \\ Q_{2,1} & Q_{2,2} & \dots & Q_{2,C} \\ \vdots & \vdots & \ddots & \vdots \\ Q_{C,1} & Q_{C,2} & \dots & Q_{C,C} \end{bmatrix} \quad (25)$$

where \mathbf{Q}_{map} is the resulting matrix of all quality evaluations for all cells, and $Q_{i,j}$ corresponds to the cell's coordinates. Q_{tot} can be seen in Fig. 3 where it corresponds to a plane elevated to the average value of all cells' MMTM quality.

It is important that the MCS administrator sets a minimum threshold for quality Q_{min} , which can be useful in two ways: 1) to evaluate whether the overall quality is satisfied, i.e., is $Q_{\text{tot}} \geq Q_{\text{min}}$ or 2) to identify specific cells that are below the threshold. Combining these together allows the characterization of overall quality, however, it does not indicate *how* the quality can be adjusted to achieve uniformity.

B. Relative Coverage Quality

In order to characterize *how* quality can be improved over a space, we developed a method that measures the *angle* between each two adjacent quality points and constructs a $C^2 \times C^2$ matrix that maps this relation, denoted Q_{rel} whose elements are defined as per Pythagoras's theorem as

$$Q_{\text{rel}}(Q_{i,j}, Q_{a,b}) = \sin^{-1} \left[\frac{Q_{i,j} - Q_{a,b}}{\sqrt{[Q_{i,j} - Q_{a,b}]^2 + 1}} \right] \quad (26)$$

where (i,j) are the coordinates of the current cell, and (a,b) are the coordinates of the adjacent cell, and \sin^{-1} comes from trigonometric ratios of a right-angled triangle on the midpoints of quality values $Q_{i,j}$ and $Q_{a,b}$. Fig. 3 provides an illustration of (26) with $Q_i \equiv Q_{i,j}$ and $Q_j \equiv Q_{a,b}$.

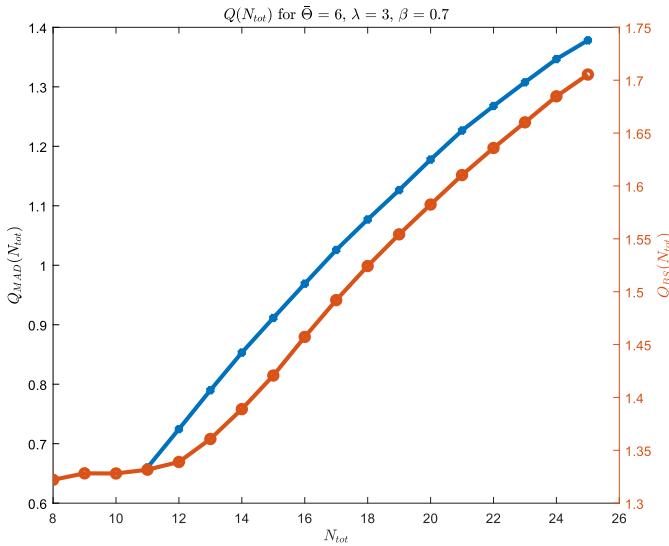
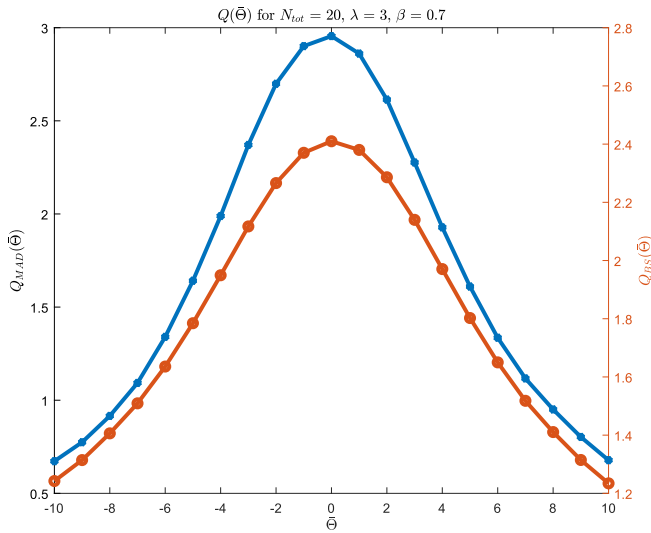
In order to implement this in an automated algorithm, the *2-tuple map* of a $C \times C$ square grid, \mathbf{Q}_{map} maps to a *singleton map*, $\hat{\mathbf{Q}}_{\text{map}}$ such that

$$\underbrace{\begin{bmatrix} Q_{1,1} & Q_{1,2} & \dots & Q_{1,C} \\ Q_{2,1} & Q_{2,2} & \dots & Q_{2,C} \\ \vdots & \vdots & \ddots & \vdots \\ Q_{C,1} & Q_{C,2} & \dots & Q_{C,C} \end{bmatrix}}_{\mathbf{Q}_{\text{map}}} \mapsto \underbrace{\begin{bmatrix} Q_1 & Q_2 & \dots & Q_C \\ Q_{C+1} & Q_{C+2} & \dots & Q_{2C} \\ \vdots & \vdots & \ddots & \vdots \\ Q_{\Delta+1} & Q_{\Delta+2} & \dots & Q_{C^2} \end{bmatrix}}_{\hat{\mathbf{Q}}_{\text{map}}} \quad (27)$$

where $\Delta = C(C - 1)$.

Within the singleton map, cells adjacent to the k th cell are corresponding to the cardinal directions: north ($k - C$), west ($k - 1$), east ($k + 1$), and south ($k + C$), for any cell that is not an edge or a corner. As a result, Q_{rel} can then be described by substituting (26) for the corresponding elements.

The relative coverage quality Q_{rel} looks at the quality, described by \mathbf{Q}_{map} , as a surface and measures the angles between adjacent points on it. In that sense, the matrix \mathbf{Q}_{rel} is defined as an antisymmetric matrix in which transposing elements is a pair of alternate interior angles, and thus the sign change. As a result, \mathbf{Q}_{rel} is a sparse matrix, which can be visualized using color-coded tables or spy plots, as illustrated in Fig. 11 in Section V. If the corresponding value for $(k \rightarrow l)$ is positive, then the uniformity of coverage can be improved by moving participants from cell k to cell l , while negative implies that participants need to be moved from l to k and/or incentive new participants in k .

Fig. 4. $Q_{MAD}(N)$ and $Q_{BS}(N)$.Fig. 5. $Q_{MAD}(\bar{\Theta})$ and $Q_{BS}(\bar{\Theta})$.

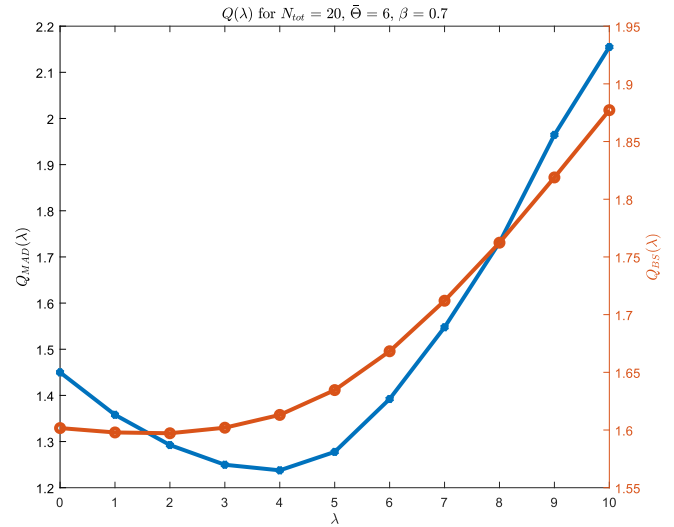
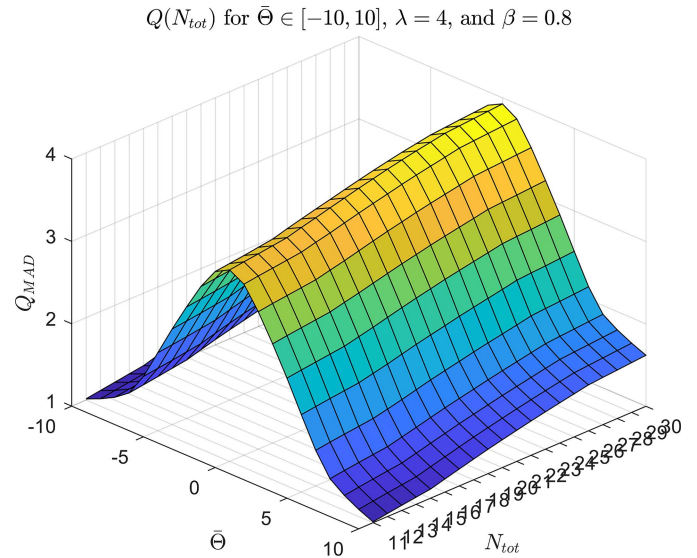
V. SIMULATION AND RESULTS

In order to test the proposed algorithms, we consider a temperature evaluation scenario in MATLAB, where temperature measurements are acquired from $N(24, 0.5)$, and two outliers are obtained from an abnormal outlier distribution located $N(24 + \bar{\Theta}, 0.3)$. $\gamma = 3.1$ is chosen for the MAD-Q, while $\gamma = 10$ is selected for MADBS-Q to ensure comparable scales. $k = 10\%$ is decided for the trimmed mean. The cell-specific quality metrics are evaluated over an ensemble of 30 randomly generated samples for each data point. For the MADBS-Q, $B = 250$ is employed for the purpose of this simulation.

A. Cell-Specific Quality

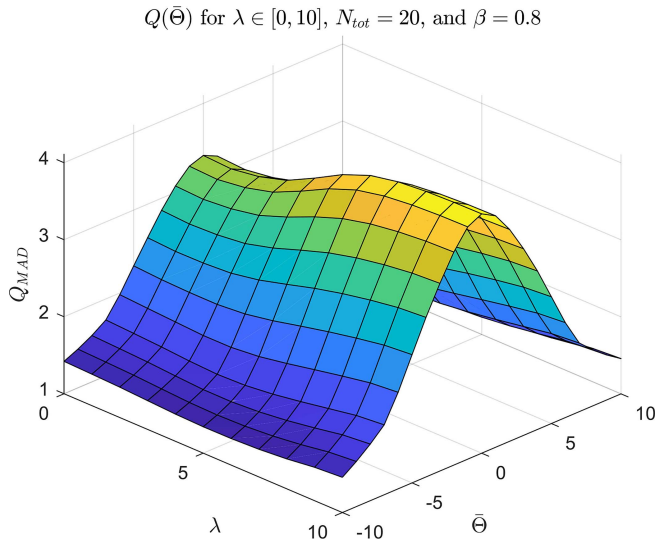
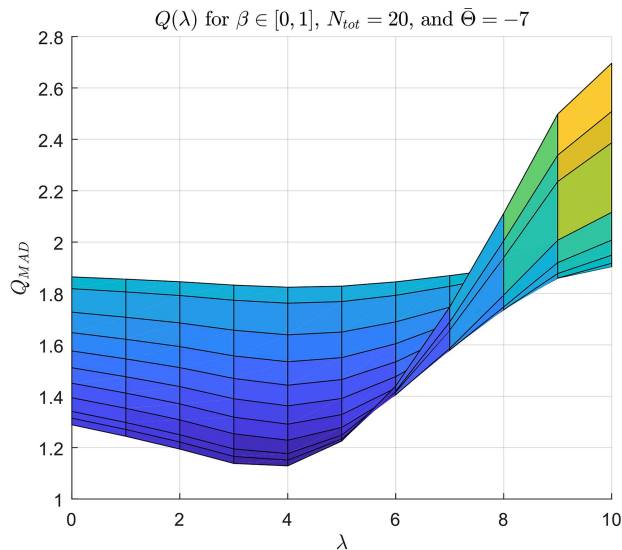
We tackle the case of variables N_{tot} , $\bar{\Theta}$, λ , and β , in order to illustrate the impact of these variables on the quality.

Fig. 4 illustrates Q_{MAD} and Q_{BS} for different sample sizes. Both Q_{MAD} and Q_{BS} show that with larger sample size, the quality increases. Moreover, this figure illustrates the sample

Fig. 6. $Q_{MAD}(\lambda)$ and $Q_{BS}(\lambda)$.Fig. 7. $Q_{MAD}(N, \bar{\Theta})$.

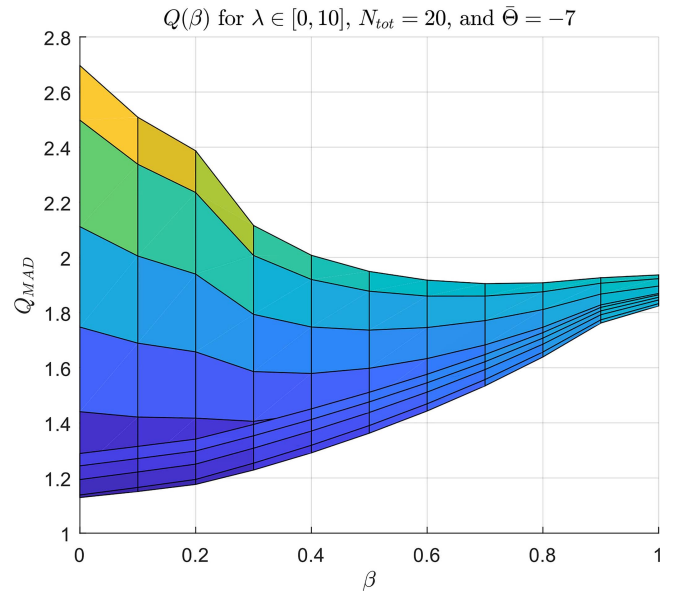
size required to achieve a certain quality threshold, which is useful in selecting a suitable number of participants for recruitment in a specific MCS application. This can be further seen in Fig. 7, for which the sample size is evaluated for Q_{MAD} at different AODFs, $\bar{\Theta}$. Fig. 5 represents Q_{MAD} and Q_{BS} for different $\bar{\Theta}$. It is observable that the curves are seemingly symmetric as they tend to follow the shape of a bell curve. This is due to the fact that the closer the outliers are to the center of the true value's population, the less of an outlier they are. At $\bar{\Theta} = 0$, the outliers constitute the sample itself, i.e., no outlier is present. At the extremes, $\bar{\Theta} = \pm 10$, the quality worsens significantly. Figs. 4 and 5 are slices of Fig. 7 at $\bar{\Theta} = 6$ and $N_{tot} = 20$, respectively.

In Fig. 6, Q_{MAD} and Q_{BS} for different λ are shown. Fig. 8 provides a view of how $Q(\lambda)$ varies with $\bar{\Theta}$, raising the general quality for high λ . For low λ , $\lambda < 2$, the quality is consistent as no outliers are tolerated. However, for larger values of $\lambda > 4$, the metric becomes more tolerant toward outliers, and thus raises the quality perceived as λ increases.


 Fig. 8. $Q_{MAD}(\bar{\theta}, \lambda)$.

 Fig. 9. $Q_{MAD}(\lambda)$ for different β .

Figs. 9 and 10 provide a view of Q_{MAD} for different λ and β . Since β controls the presence of the MAD-mean term in the θ_{MMTM} statistic, at $\beta = 0$, the trimmed mean term is entirely absent, thus causing a larger variation in quality. For $\lambda < 2$, the quality is not very high nor very low. This is due to the fact that $Q_{MAD,BS}$ measures the closeness of a sample set's values. For $\lambda = 4$, the quality can be seen to have started increasing as outliers gradually become tolerable. For $\lambda > 4$, the quality increases as λ increases, as it allows more outliers to be tolerated. However, for $\beta = 1$, when the trimmed mean is fully present, λ has no effect at all. Similarly for β , the quality increases from $\beta = 0$ to $\beta = 1$ except when $\lambda \geq 4$.

The difference between both methods lies in the fact that the MADBS-Q can quantify, consistently, samples as small as $N = 8$ at a computational cost without falling into degeneracies, unlike the MAD-Q which often falls in the log function's negative domain for $N < 11$, being sometimes undefined


 Fig. 10. $Q_{MAD}(\beta)$ for different λ .

(returning a NaN for a negative θ_{MMTM}). However, the behavior for Q_{BS} , as a surface, is similar to that illustrated in Figs. 7–10.

B. Coverage Quality

To test our coverage quality metric, we generate a 3×3 map with values obtained from a distribution $N(2.2, 0.5)$, to simulate diverse values obtained from the metrics discussed in the previous section, where \mathbf{Q}_{map} is

$$\mathbf{Q}_{map} = \begin{bmatrix} 3.09 & 1.04 & 2.23 \\ 2.81 & 2.65 & 2.22 \\ 1.56 & 1.28 & 3.32 \end{bmatrix}. \quad (28)$$

The overall quality metric Q_{tot} is found to be 2.24, and \mathbf{Q}_{rel} , illustrated in Fig. 11, is obtained. We use a color code to identify quality, where green denotes that the corresponding cells are near each other, i.e., \mathbf{Q}_{rel} is closer to zero at these two cells, and red denotes that there is a significant contrast between the values, i.e., \mathbf{Q}_{rel} far from zero. The resultant matrix, illustrated as a sparse adjacency matrix in Fig. 11, is antisymmetric as the relations between cells are bidirectional. This result is very important as it shows the *status-quo* to the administrator and sheds the light on where more participants are needed in order to balance the distribution among the cells.

C. Discussion

The proposed algorithm has two main parameters, λ and β , both which impact the θ_{MMTM} parameter, as previously discussed in Section III-C. Another parameter that is also considered is k , the amount of trimming done by the trimmed mean term \bar{x}_k . Increasing k increases the intolerance of the trimmed mean to outliers, however, in most cases, 10% is more than enough in SD scenarios where $N \in [8, 30]$, as increasing it would be of little benefit as removing 20%, for example, of 20 samples reduces it to 16 samples; which is unlikely for

	1	2	3	4	5	6	7	8	9
1	Q(1,1) 0.00	Q(1,2) 64.04	Q(1,3) 0.00	Q(1,4) 15.51	Q(1,5) 0.00	Q(1,6) 0.00	Q(1,7) 0.00	Q(1,8) 0.00	Q(1,9) 0.00
2	Q(2,1) -64.04	Q(2,2) 0.00	Q(2,3) -50.14	Q(2,4) 0.00	Q(2,5) -58.24	Q(2,6) 0.00	Q(2,7) 0.00	Q(2,8) 0.00	Q(2,9) 0.00
3	Q(3,1) 0.00	Q(3,2) 50.14	Q(3,3) 0.00	Q(3,4) 0.00	Q(3,5) 0.00	Q(3,6) 0.90	Q(3,7) 0.00	Q(3,8) 0.00	Q(3,9) 0.00
4	Q(4,1) -15.51	Q(4,2) 0.00	Q(4,3) 0.00	Q(4,4) 0.00	Q(4,5) 9.12	Q(4,6) 0.00	Q(4,7) 51.41	Q(4,8) 0.00	Q(4,9) 0.00
5	Q(5,1) 0.00	Q(5,2) 58.24	Q(5,3) 0.00	Q(5,4) -9.12	Q(5,5) 0.00	Q(5,6) 23.42	Q(5,7) 0.00	Q(5,8) 53.85	Q(5,9) 0.00
6	Q(6,1) 0.00	Q(6,2) 0.00	Q(6,3) -9.90	Q(6,4) 0.00	Q(6,5) -23.42	Q(6,6) 0.00	Q(6,7) 0.00	Q(6,8) 0.00	Q(6,9) -17.62
7	Q(7,1) 0.00	Q(7,2) 0.00	Q(7,3) 0.00	Q(7,4) -51.41	Q(7,5) 0.00	Q(7,6) 0.00	Q(7,7) 0.00	Q(7,8) 15.44	Q(7,9) 0.00
8	Q(8,1) 0.00	Q(8,2) 0.00	Q(8,3) 0.00	Q(8,4) 0.00	Q(8,5) -53.85	Q(8,6) 0.00	Q(8,7) -15.44	Q(8,8) 0.00	Q(8,9) -63.79
9	Q(9,1) 0.00	Q(9,2) 0.00	Q(9,3) 0.00	Q(9,4) 0.00	Q(9,5) 0.00	Q(9,6) 17.62	Q(9,7) 0.00	Q(9,8) 63.79	Q(9,9) 0.00

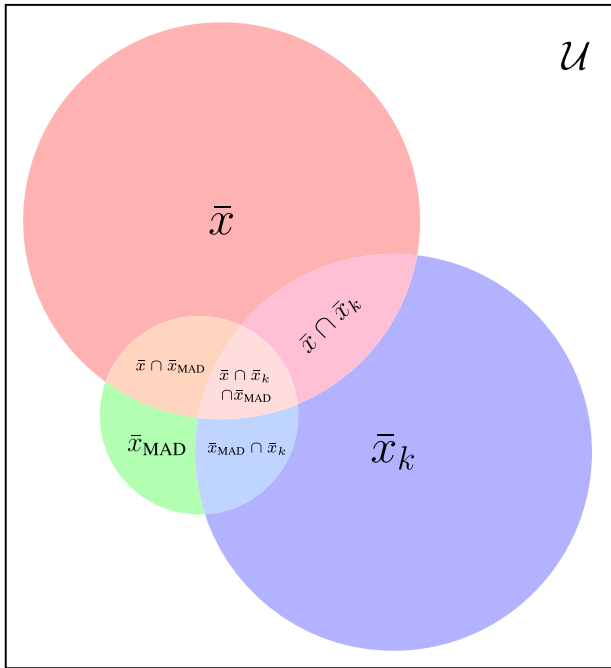
Fig. 11. Spy plot of the adjacency matrix, Q_{rel} .

Fig. 12. Venn diagram of the bootstrapped sample space and its subsets.

abnormal outliers. The choice of β and λ relies on the application on a question. If the application requires being more tolerant toward outliers, a small value of β would increase the significance of the \bar{x}_{MAD} term and thus λ would be the only parameter impacting the filtration of outliers. However, the application requires intolerance to outliers, then choosing a high β increases the significance of the trimmed mean term, \bar{x}_k thus filters outliers aggressively. However, changing k , as previously mentioned, is of little benefit.

VI. CONCLUSION

We have presented a new direction for looking at scenarios in which MCS participants present are not enough or the data reported by them is not entirely reliable due to the presence of abnormalities, and thus errors. Sometimes, the available data may not be enough to fully characterize the MCS area of interest. To enhance the MCS system's robustness, we introduced SD quality techniques to allow MCS systems to evaluate a group of participants' reports in a specific spatiotemporal cell based on the data they reported in scenarios where the

scale of data is insufficient for proper inference. This evaluation was performed without the knowledge of the true value, ascertaining as much as possible from the least amount of data.

The proposed quality metrics, MAD-Q and MADBS-Q, allow the characterization of a sample's quality—typically an MCS sample being a set of physical readings whose readings come from a symmetric distribution—when the number of contributors is minimal. The MAD-Q technique allows the MCS system to evaluate the quality of samples as small as 11, while MADBS-Q performs the evaluation of samples as small as eight with a little computational cost by employing the non-parametric bootstrap. Both metrics are based on the MMTM statistic, which is designed to allow the formulation of a flexible quality metric, that takes into consideration tolerating specific outliers, by tweaking the sensitivity parameter β that controls the weights of the trimmed mean ($\beta = 1$) and the MAD-mean ($\beta = 0$), and a range parameter λ that allows the control of the MAD-mean's outlier tolerance. This control, while allowing the tolerance of certain values, allows MCS administrators to define, subjectively, what a usable MCS data set is in a domain-agnostic manner. Furthermore, the importance of the quality metrics lies in their relationship to the sample size, as illustrated in the MAD-Q technique, allowing the administrators control in the quality-cost tradeoff.

We also proposed an MCS-specific coverage metric to employ the notion of a cell-specific quality, no matter how it is defined, in the definition of an MCS system's coverage. The coverage metric aims to achieve uniformity by describing quality as a whole, and also in a relative manner between the MCS system's spatiotemporal cells. This gives the MCS system a way of measuring which cells are overstaffed or understaffed by MCS participants as well as a metric that indicates how to manipulate the current participant topology (i.e., the participants' distribution over the grid) to achieve coverage uniformity, either by recruiting more participants or by requesting their movement around the area of interest. The developed quality metrics, MAD-Q and MADBS-Q, along with the MCS-specific coverage metric, provide a clear relationship for the MCS system between sensing quality and sensing cost.

APPENDIX

BOOTSTRAP PROBABILITY ANALYSIS

The purpose of this appendix is to introduce the reader to the MMTM statistic defined in (15) under the statistical bootstrap. The reader will find a derivation of the sampling distribution of the sample median within the bootstrap. This result would be helpful for further derivations under the bootstrap.

Within the bootstrapped sample space, the sample distributions of the centrality estimates (mean, trimmed mean, and MAD mean) are different, but they all share an overlap. Equation (15) is a function of three random variables, this can be illustrated as a sample space denoted \mathcal{U} , containing subsets corresponding to each of the random variables as shown in Fig. 12.

Inferences can be made, at this point, regarding each of the random variables. The sample distribution \bar{x} and the \bar{x}_k both

$$P(\tilde{x}_b = x_n) = \frac{\sum_{r_U=0}^{\frac{N+1}{2}} \sum_{r_L=0}^{\frac{N+1}{2}} \frac{N!}{r_L!(N-r_U-r_L)!r_U!} \binom{n-1}{N}^{r_L} \left(\frac{1}{N}\right)^{(N-r_U-r_L)} \binom{N-n}{N}^{r_U}}{\sum_{n=1}^N \sum_{r_U=0}^{\frac{N+1}{2}} \sum_{r_L=0}^{\frac{N+1}{2}} \frac{N!}{r_L!(N-r_U-r_L)!r_U!} \binom{n-1}{N}^{r_L} \left(\frac{1}{N}\right)^{(N-r_U-r_L)} \binom{N-n}{N}^{r_U}} \quad (29)$$

$$P(\tilde{x}_b = \tilde{x}_j) = \frac{\sum_{r_{M1}=1}^{N/2} \sum_{r_{M2}=1}^{N/2} \frac{N!}{r_L!r_{M1}!r_{M2}!r_U!} \binom{N-f_{jl}}{N}^{r_L} \left(\frac{N-f_{ju}}{N}\right)^{r_U} \left(\frac{1}{N}\right)^{r_{M1}+r_{M2}}}{\sum_{j_l=1}^{\sum_{i=0}^{N-1} (N-i)} \sum_{j_u=1}^{\sum_{i=0}^{N-1} (N-i)} \sum_{r_{M1}=1}^{N/2} \sum_{r_{M2}=1}^{N/2} \frac{N!}{r_L!r_{M1}!r_{M2}!r_U!} \binom{N-f_{jl}}{N}^{r_L} \left(\frac{N-f_{ju}}{N}\right)^{r_U} \left(\frac{1}{N}\right)^{r_{M1}+r_{M2}}} \quad (30)$$

Algorithm 4 Algorithm for Computing Median Probability for Even Samples

Input: Sample $X = \{x_1, x_2, \dots, x_n\}$

Output: \tilde{x}_j

```

1: for  $a = 1$  to  $N$  do
2:   for  $b = a$  to  $N$  do
3:     store  $(x_a + x_b)/2$  in  $\tilde{x}_j[i, 1]$ 
4:     store  $x_a$  in  $\tilde{x}_j[i, 2]$ 
5:     store  $x_b$  in  $\tilde{x}_j[i, 3]$ 
6:   end for
7: end for
8:  $[\sim, \text{idx}] = \text{sort}(\tilde{x}_j[:, 1])$ 
9:  $\tilde{x}_j = \tilde{x}_j[\text{idx}, :]$ 
10: return  $\tilde{x}_j$ 

```

follow normal distributions, thanks to the CLT. The distribution of the MAD-mean \bar{x}_{MAD} , however, is different as it is a result of the distribution of the median. The median, however, thanks to its discrete nature has various cases. If the sample size was odd, the probability that a specific bootstrap sample, x_n is selected as the median \tilde{x}_b , can be expressed by (29), shown at the top of the page, where r_U and r_L correspond to the number of values greater than the n th, and values lower than the n th, respectively, and N is the sample size. The numerator is the probability of a specific event, whereas the denominator is the total, which is an extension of the multinomial distribution that treats the possibilities of values greater or smaller than the median coexisting with it, in a bootstrap resample such that the bootstrap median would reflect the n th sample.

For an even sample size, the expression becomes more tedious since the possibilities of the bootstrap median increase, as the averages among samples become median candidates in addition to the samples themselves. Equation (30), shown at the top of the page, is an expression we have derived for that, where r_{M1} and r_{M2} represent the two upper and lower values which would contribute to the median, while r_U and r_L represent the rest of the values within the sample that are remote to the center.

It is possible to see from both equations that the probability distribution tends to have a central limit. However, it is not a smooth distribution since the median is, by nature, discrete. Thus, the distribution of the medians attempts to approach that of the normal distribution—due to the central limit—however, the shape is not always normal. As a result, the MAD-mean’s

distribution is the discrete impulses representing that of the median convolved with the normal distribution. While it does have a clear central tendency to be normal, it is not necessarily normal if the medians are distant from each other, which would result in a bumpy distribution. Ultimately, to obtain the closed-form approximation of the bootstrap, the compound probabilities formed from the simple ones described in the Venn diagram in Fig. 12, and have to be computed, as well as their expectations.

Finally, f_{jl} and f_{ju} are the number of occurrences (frequency) necessary for the j th median out of the range of all possible medians, for the even sample size. The expression for the median probability for an even sample can be obtained using Algorithm 4.

REFERENCES

- [1] V. Freschi, S. Delpriori, E. Lattanzi, and A. Bogliolo, “Bootstrap based uncertainty propagation for data quality estimation in crowdsensing systems,” *IEEE Access*, vol. 5, pp. 1146–1155, 2017.
- [2] S. B. Azmy, N. Zorba, and H. S. Hassanein, “Robust quality metric for scarce mobile crowd-sensing scenarios,” in *Proc. IEEE Int. Conf. Commun. Workshops (ICC Workshops)*, Kansas City, MO, USA, May 2018, pp. 1–5.
- [3] S. B. Azmy, N. Zorba, and H. S. Hassanein, “Bootstrap-based quality metric for scarce sensing systems,” in *Proc. IEEE Global Commun. Conf. (GLOBECOM)*, Abu Dhabi, UAE, Dec. 2018, pp. 1–6.
- [4] B. Guo *et al.*, “Mobile crowd sensing and computing: The review of an emerging human-powered sensing paradigm,” *ACM Comput. Surveys*, vol. 48, no. 1, pp. 1–31, Aug. 2015.
- [5] M. Zappatore, C. Loglisci, A. Longo, M. A. Bochicchio, L. Vaira, and D. Malerba, “Trustworthiness of context-aware urban pollution data in mobile crowd sensing,” *IEEE Access*, vol. 7, pp. 154141–154156, 2019.
- [6] G. Cardone, A. Corradi, L. Foschini, and R. Ianniello, “ParticipAct: A large-scale crowdsensing platform,” *IEEE Trans. Emerg. Topics Comput.*, vol. 4, no. 1, pp. 21–32, Jan./Mar. 2016.
- [7] R. F. El Khatib, N. Zorba, and H. S. Hassanein, “Multi-tasking for cost-efficient mobile crowdsensing under uniformity constraints,” in *Proc. IEEE Int. Conf. Commun. (ICC)*, Kansas City, MO, USA, May 2018, pp. 1–6.
- [8] D. Zhang, L. Wang, H. Xiong, and B. Guo, “4W1H in mobile crowd sensing,” *IEEE Commun. Mag.*, vol. 52, no. 8, pp. 42–48, Aug. 2014.
- [9] L. Wang, D. Zhang, Y. Wang, C. Chen, X. Han, and A. M’hamed, “Sparse mobile crowdsensing: Challenges and opportunities,” *IEEE Commun. Mag.*, vol. 54, no. 7, pp. 161–167, Jul. 2016.
- [10] X. Xu, X. Liu, Z. Xu, F. Dai, X. Zhang, and L. Qi, “Trust-oriented IoT service placement for smart cities in edge computing,” *IEEE Internet Things J.*, vol. 7, no. 5, pp. 4084–4091, May 2020.
- [11] J. Franklin, “Global and local,” *Math. Intell.*, vol. 36, no. 4, pp. 4–9, 2014.
- [12] L. Wang, D. Zhang, Z. Yan, H. Xiong, and B. Xie, “effSense: A novel mobile crowd-sensing framework for energy-efficient and cost-effective data uploading,” *IEEE Trans. Syst., Man, Cybern., Syst.*, vol. 45, no. 12, pp. 1549–1563, Dec. 2015.
- [13] J. Ni, K. Zhang, Q. Xia, X. Lin, and X. S. Shen, “Enabling strong privacy preservation and accurate task allocation for mobile crowdsensing,” *IEEE Trans. Mobile Comput.*, vol. 19, no. 6, pp. 1317–1331, Jun. 2020.

- [14] R. Figliola and D. Beasley, *Systematic and Random Errors*. Hoboken, NJ, USA: Wiley, 2011.
- [15] Y. Qu *et al.*, "Posted pricing for chance constrained robust crowdsensing," *IEEE Trans. Mobile Comput.*, vol. 19, no. 1, pp. 188–199, Jan. 2020.
- [16] T. Luo, J. Huang, S. S. Kanhere, J. Zhang, and S. K. Das, "Improving IoT data quality in mobile crowd sensing: A cross validation approach," *IEEE Internet Things J.*, vol. 6, no. 3, pp. 5651–5664, Jun. 2019.
- [17] P. J. Rousseeuw and C. Croux, "Alternatives to the median absolute deviation," *J. Amer. Stat. Assoc.*, vol. 88, no. 424, pp. 1273–1283, 1993.
- [18] E. Grafarend and J. Awange, *Applications of Linear and Nonlinear Models* (Springer Geophysics). New York, NY, USA: Springer, 2012.
- [19] C. Leys, O. Klein, P. Bernard, and L. Licata, "Detecting outliers: Do not use standard deviation around the mean, use absolute deviation around the median," *J. Exp. Soc. Psychol.*, vol. 49, no. 4, pp. 764–766, 2013.
- [20] G. Buzzi-Ferraris and F. Manenti, "Outlier detection in large data sets," *Comput. Chem. Eng.*, vol. 35, no. 2, pp. 388–390, 2011.
- [21] B. Efron, "Bootstrap methods: Another look at the jackknife," *Ann. Stat.*, vol. 7, no. 1, pp. 1–26, Jan. 1979.
- [22] K. Singh and M. Xie, "Bootlier-plot: Bootstrap based outlier detection plot," *Sankhya Indian J. Stat. (2003-2007)*, vol. 65, no. 3, pp. 532–559, 2003.
- [23] G. Arfken and H. Weber, *Mathematical Methods for Physicists International Student Edition*. Amsterdam, The Netherlands: Elsevier Sci., 2005.
- [24] N. Stubbs and S. Park, "Optimal sensor placement for mode shapes via Shannon's sampling theorem," *Comput.-Aided Civil Infrastruct. Eng.*, vol. 11, no. 6, pp. 411–419, 1996.
- [25] R. Hogg and E. Tanis, *Probability and Statistical Inference*. Upper Saddle River, NJ, USA: Prentice-Hall, 2006.
- [26] D. Bertsimas and B. Sturt, "Computation of exact bootstrap confidence intervals: Complexity and deterministic algorithms," in *Eprints for the Optimization Community*. Ann Arbor, MI, USA: Univ. Michigan, 2017.
- [27] M. F. Schilling, A. E. Watkins, and W. Watkins, "Is human height bimodal?" *Amer. Stat.*, vol. 56, no. 3, pp. 223–229, 2002.
- [28] H. Holzmann and S. Vollmer, "A likelihood ratio test for bimodality in two-component mixtures with application to regional income distribution in the EU," *Adv. Stat. Anal.*, vol. 92, no. 1, pp. 57–69, Feb. 2008.
- [29] S. Doslá, "Conditions for bimodality and multimodality of a mixture of two unimodal densities," *Kybernetika*, vol. 45, no. 2, pp. 279–292, 2009.
- [30] J. Högel, W. Schmid, and W. Gaus, "Robustness of the standard deviation and other measures of dispersion," *Biometrical J.*, vol. 36, no. 4, pp. 411–427, 1994.



Sherif B. Azmy (Student Member, IEEE) received the B.Sc. degree in electrical engineering from Qatar University, Doha, Qatar, in 2017. He is currently pursuing the M.Sc. degree in electrical and computer engineering with Queen's University, Kingston, ON, Canada.

His current research interests include operations research, mobile crowdsensing, and Internet of Things.



Nizar Zorba (Senior Member, IEEE) received the B.Sc. degree in electrical engineering from JUST University, Irbid, Jordan, in 2002, and the Ph.D. degree in signal processing for communications from UPC, Barcelona, Spain, in 2007.

He is a Professor with the Electrical Engineering Department, Qatar University, Doha, Qatar. He has authored five international patents and coauthored over 120 papers in peer-reviewed journals and international conferences.

Prof. Zorba is an Associate/Guest Editor of IEEE COMMUNICATIONS LETTERS, IEEE ACCESS, *IEEE Communications Magazine*, and IEEE NETWORK. He is currently the Vice-Chair of the IEEE ComSoc Communication Systems Integration and Modeling Technical Committee.



Hossam S. Hassanein (Fellow, IEEE) received the Ph.D. degree in computing science from the University of Alberta, Edmonton, AB, Canada, in 1990.

He is a leading authority in the areas of broadband, wireless and mobile networks architecture, protocols, control, and performance evaluation. His record spans more than 500 publications in journals, conferences, and book chapters, in addition to numerous keynotes and plenary talks in flagship venues.

Dr. Hassanein has received several recognition and best papers awards at top international conferences. He is also the Founder and the Director of the Telecommunications Research Lab, Queen's University School of Computing, with extensive international academic and industrial collaborations. He is a Former Chair of the IEEE Communication Society Technical Committee on Ad hoc and Sensor Networks.