

The Challenges of Model Objective Selection and Estimation for Ad-hoc Network Data Sets

Farinaz Koushanfar¹ and Davood Shamsi²

Rice University

Abstract: We introduce a new methodology for determining the difficulty of selecting the modeling objective function and estimating the parameters for an ad-hoc network data set. The method utilizes formulation of the underlying optimization problem instance that consists of an objective function and a set of constraints. The method is illustrated on real distance measurement data used for estimating the locations of wireless nodes that is the most studied and the representative problem for ad-hoc networks estimation. The properties of the data set that could affect the quality of optimization are categorized. In large optimization problems with multiple properties (characteristics) that contribute to the solution quality, it is practically impossible to analytically study the effect of each property. A number of metrics for evaluating the effectiveness of the optimization on each data set are proposed. Using the well known Plackett and Burmann fast simulation methodology, for each metric, the impact of the categorized properties of the data are determined for the specified optimization. A new approach for combining the impacts resulting from different properties on various metrics is described. We emphasize that the method is generic and has the potential to be more broadly applicable to other parameter estimation problems.

Contents

1	Introduction	334
2	Preliminaries	335
3	Metrics	337
3.1	Error Metrics	337
3.2	Objective Function (OF) Metrics	338
3.2.1	Drifting of Objective Function (OF)	338
3.2.2	Nearest Local Minimum	339
3.2.3	Measuring the Slope of OF Around the Solution	339
3.2.4	Depth of the Non-Global Local Minima	340
4	Simulation Methodology	340
5	Combining Different Ranks	341
6	Evaluation Results	342
7	Conclusion	345
	Acknowledgement	345
	References	346

¹Electrical and Computer Engineering Department, Rice University, Houston, TX 77005

²Computer Science Department, Rice University, Houston, TX 77005

1. Introduction

Wireless adhoc networks consist of multiple wireless nodes distributed in an implementation area. To be power efficient, the wireless nodes only directly communicate with the nodes in their short local range (neighbor nodes). Communication between the non-neighbor nodes is enabled by successive usage of (one or more) local neighbors as forwarding relays. Several problems in this domain include modeling and estimation of data sets that only contain pairwise exchanged data between the neighboring nodes.

Years of continuous research in building statistical models and parameter estimation has produced a multitude of readily available methods and tools that can be employed for the problems in ad-hoc networks [10]. One limitation of the available methods is that majority of the ad-hoc modeling and estimation problems concern a large body of data and do not conform with typical assumptions needed to analytically declare the known theoretical optimality criteria. In such scenarios, the quality of the modeling and estimation methods are typically evaluated by how they perform on sets of real or simulated data. For example, some statistics of the resulting prediction error and/or a defined criterion (e.g., Bayesian information criterion (BIC)) is used for experimental evaluation of the method on the adhoc network measurements. A relevant question to answer is if indeed modeling and estimation of the pertinent data set requires introduction of a new model or an estimator, or the data could have been just as well addressed by the other known methods.

Our objective is to quantify the difficulty of model selection and estimation for a given adhoc network data set. This would provide impetus for inventing newer modeling and estimation objectives and tools that can address the difficult-to-characterize data. Simultaneously, formation of new tools would depend upon finding truly challenging network data sets that need to be addressed, as opposed to building new models that have a limited practical usage. Devising sets of challenging data would also build a foundation for comparing the various modeling objective functions and estimators for the ad-hoc network data sets. The problem of finding challenging data is complicated by variations in properties of the underlying data sets collected by different sources. This includes difference in size, format, wireless ranges, hidden covariates, and the form of noise present in the collected data. Thus, it is not easy to find unique metrics that could be used for comparison of different modeling objective functions and estimation methods.

In statistics literature, sensitivity of estimation error or other discrepancy metrics to the underlying noise in data has been widely studied for a number of modeling methods [24][3]. Also, the consistency of estimators based on a number of strong assumptions on the distribution of the data has been pursued [14]. However, no generic method or tool for determining the difficulty in modeling a data set free of imposing strong assumptions – such as normality or other closed-form distributions of noise – is available for use in adhoc networks. Note that the runtime complexity of a problem is an orthogonal concept. The complexity measures the worst-case computational time for the algorithm used for addressing the problem. Analyzing the worst-case runtime complexity does not help in understanding the complexity of characterizing a specific data set.

In adhoc network scenario, after the variable selection is done and the noise models are assumed, modeling is typically done by selecting a model form (e.g., nonlinear regression) and then estimating the model parameters on the data set. For analyzing the modeling objective function and estimation performance on the data, we study the pertinent optimization problem that consists of an objective

1 function (OF) and a number of constraints. The data set is considered as the input
2 to the optimization problem. We introduce a number of metrics that measure
3 the complexity of the optimization problem based on the problem OF properties
4 and constraints. The challenge in most optimization problems is the existence of
5 nonlinearities that make the solution space coarse, causing bumpiness and multiple
6 local minimums. We propose a number of measures for the smoothness of the OF
7 and constraints space that estimate the feasibility of reaching the global minimum.

8 To enable studying the effectiveness of the optimization on an adhoc network
9 data set, one should characterize the properties of the pertinent data set. The
10 properties are specific to each data set and the problem. In this article, we focus
11 on the problem of finding the location of nodes (localization) in an adhoc wireless
12 network by using erroneous mutual distance measurements between a number of
13 node pairs. However, we emphasize that our method is generic and can be used
14 for determining the challenge in addressing many adhoc data set model objective
15 selection and estimation that includes forming an optimization problem. The lo-
16 calization problem is selected for four reasons. First, it is a very well addressed
17 problem in the literature and there are several methods that are developed for this
18 problem [6][9][19][2][20][26]. Second, there are a number of publicly available data
19 sets for the measured distance data in the networks [5][21][8]. Third, the nonlinear
20 relationship between noise in measurements data and the location of nodes makes
21 the modeling problem extremely challenging. Fourth, localization problem is an
22 NP-complete problem, i.e., in the worst case, there is no algorithm that can solve it
23 in polynomial time [25][6]. Lastly, location discovery is a precursor for a number of
24 other problems in ad hoc networks including sleeping coordination [12, 13], sensor
25 coverage [15], and sensing exposure [16].

26 We characterize a number of properties of the measurement data set that could
27 affect the quality of location estimation. Studying the interaction between the iden-
28 tified data properties and optimization metrics requires long simulations and anal-
29 ysis. We use the well-known Plackett and Burmann [23] simulation methodology
30 to rapidly study the pairwise linear interactions of properties. A new approach for
31 combining the impacts resulting from different properties of data on various opti-
32 mization metrics is described. The sensitivity of optimization with respect to the
33 various parameter ranks are presented.

34 To the best of our knowledge, this is the first work that systematically studies
35 the impact of the adhoc network data set on the optimization employed for finding
36 the modeling objectives and estimations. Most of the previous work are devoted to
37 modeling and analysis of the worst case complexity. The results of our analysis could
38 be directly used for constructing benchmarks for the problem. The proposed work
39 aims at creating a unified framework based on real data that can help evaluation
40 and comparison of desperate efforts that address the same problem.

41 The remainder of the paper is organized as follows. In the next section, location
42 estimation problem and our notations are formally defined. In Section 3, we devise a
43 number of metrics that are used for OF evaluation. The simulation methodology is
44 described in Section 4. In Section 5, we illustrate how the results of different metrics
45 can be combined. We have applied the derived method on the measurements from
46 a real network in Section 6. We conclude in Section 7.

48 2. Preliminaries

49
50 In this section, we present the formal definition of the problem. We also describe
51 the notations that are used throughout the paper.

- 1 • P_3 – Mean squared error ($\overline{\epsilon^2}$): mean squared error of distance measurements. 1
- 2 • P_4 – Maximum allowed squared error ($\text{MAX}_{\epsilon_m^2}$): the maximum squared error 2
- 3 that can possibly exist in distance measurements. 3
- 4 • P_5 – Percent of large errors ($\text{PER}_{\epsilon_0^2}$): percentage of squared distance mea- 4
- 5 surement noises that are higher than a specific value ϵ_0^2 . 5
- 6 • P_6 – Mean degree (\overline{D}): mean degree of the nodes in the network. Degree of a 6
- 7 node v_i is define as number of nodes that have their mutual distance to v_i . 7
- 8 • P_7 – Minimum length (MINL): possible minimum length of the measured 8
- 9 distances between nodes in the network. 9
- 10 • P_8 – Maximum length (MAXL): possible maximum length of the measured 10
- 11 distances between nodes in the network. 11
- 12 • P_9 – Mean length (\overline{l}): mean length of the measured distances between nodes 12
- 13 in the network. 13
- 14 • P_{10} – Minimum degree (MIND): possible minimum degree of the nodes in the 14
- 15 network. 15

16 To study the effect of the parameters, we construct a variety of network instances 16
 17 with different properties. The networks are constructed by selecting subsets of an 17
 18 implemented network. Having specific values for parameters, we use Integer Linear 18
 19 Programming (ILP) to extract each subset such that it meets specified conditions. 19
 20 To do so, we model parameter constraints as linear equalities and inequalities. 20
 21 Some parameters such as the mean squared error, $\overline{\epsilon^2}$, can be easily stated by linear 21
 22 equalities and inequalities. But some parameters such as the mean degree of the 22
 23 nodes, \overline{D} , need a mapping to be stated in linear terms. The description of the exact 23
 24 procedure of modeling by linear constraints is beyond the scop of this paper [8]. 24
 25

26 3. Metrics 27

28 In this section, we introduce metrics for error and OF that are used for evaluating 28
 29 the importance of different parameters for location estimation. Three error metrics 29
 30 and four OF metrics are presented. Thus, a total of twelve combined metrics are 30
 31 used to evaluate the importance of parameters. 31
 32

33 3.1. Error Metrics 34

35 The three error metrics studied in this paper are: L_1 , L_2 , and the maximum 35
 36 likelihood (ML). L_1 and L_2 are the common error norms in the L_p family defined 36
 37 as: 37
 38

$$39 L_p(e_{v_n, v_m} \in E) = \left(\sum_{(v_n, v_m) \in E} |e_{v_n, v_m}|^p \right)^{1/p} \quad \text{if } 1 \leq p < \infty. \quad 39$$

40 To find the error metric corresponding to ML, we need to model the noise in dis- 40
 41 tance measurements. To model the noise, the probability density function (PDF) of 41
 42 errors, f_m , for the distance measurements should be approximated. Different meth- 42
 43 ods are developed to approximate PDF of noise, f_m [8]. We have used kernel fitting 43
 44 that is a simple and known PDF approximation method [10]. To have the maximum 44
 45 likelihood estimation for the nodes' locations, we find the nodes' coordinates such 45
 46 that they maximize 46
 47

$$48 (2) \quad \prod_{(v_n, v_m) \in E} f_m(e_{v_n, v_m}) = \exp \left\{ \sum_{(v_n, v_m) \in E} \ln(f_m(e_{v_n, v_m})) \right\} \quad 48$$

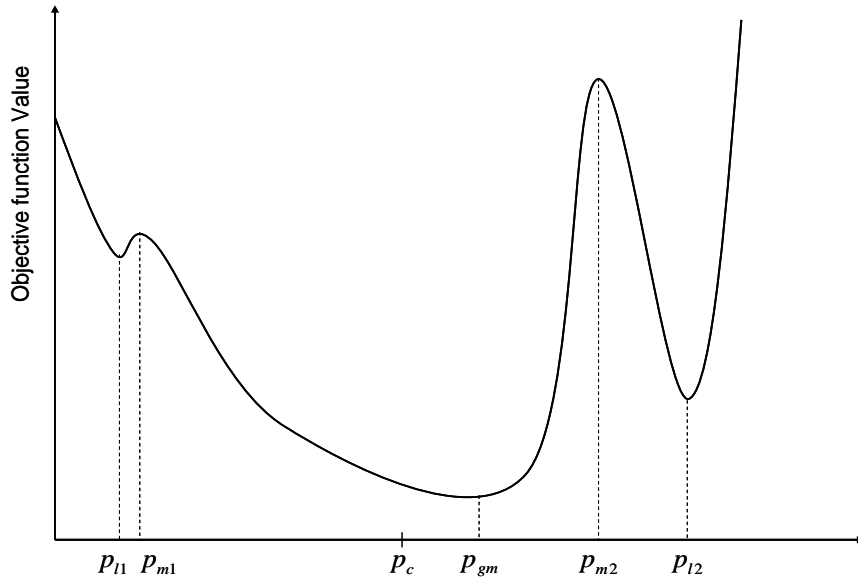


FIG 1. Metrics and objective function (OF).

or equivalently minimize

$$(3) \quad \sum_{(v_n, v_m) \in E} -\ln(f_m(e_{v_n, v_m})).$$

Note that we assume noise in distance measurements are independently identically distributed. Using the same notations as the Equation 1 and Equation 3, for the ML estimation we consider the following error metric:

$$(4) \quad L_{ML}(e_{v_n, v_m}) = -\ln(f_m(e_{v_n, v_m})).$$

3.2. Objective Function (OF) Metrics

We describe metrics that are used for evaluating OFs. The metrics are introduced based on the properties of OF that are effective in optimization. These metrics are such that they assign larger values to the more difficult-to-optimize OFs. For example, if one selects a convex OF, it may be possible to utilize convex programming depending on the form of the constraints. In defining the OF metrics, we assume that there is a fixed instance of location estimation data. Thus, for a fixed error metric, the OF would be fixed. Metrics of OF are denoted by $M : \mathcal{C} \rightarrow R^+$ where \mathcal{C} is the functional space that contains all OFs.

3.2.1. Drifting of Objective Function (OF)

Since there is noise in distance measurements, true location of the nodes is often not the global minimum of the OF. Location of the OF's global minimum is a measure of the goodness of the OF. Figure 1 illustrates the effect of noise on the OF. For the sake of presentation simplicity, an one-dimensional OF is shown. In this figure, p_c is the correct nodes' location. However, the global minimum of the

1 OF is displaced at p_{gm} because of the noise. We consider the distance between p_c 1
 2 and its displaced location p_{gm} as an OF metric and denote it by *drifting*. 2

3 To find the drifting distance, we start from the true locations as the initial point. 3
 4 Next, the steepest descent direction of the OF is followed until a local minimum is 4
 5 reached. The Euclidean distance between the true locations and this local minimum 5
 6 quantifies the drifting metric (denoted by M_1) for the pertinent OF. 6
 7

8 3.2.2. Nearest Local Minimum 8

9
 10
 11 Having a number of local minimums around the global minimum in an OF 11
 12 may cause the optimization algorithm to get trapped in one of the non-global local 12
 13 minimums. It is challenging to minimize such an OF since the global minimum 13
 14 is hard to reach. Figure 1 illustrates the phenomena. The OF has multiple local 14
 15 minima at points p_{m1} , p_{m1} and so on. The steepest decent method leads to the 15
 16 global minimum if and only if we start from a point between p_{m1} and p_{m2} . Hence, 16
 17 having a small distance between p_{m1} and p_{m2} would complicate the selection of the 17
 18 initial starting point. 18

19 We introduce a method to measure the distance of the true locations from the 19
 20 local minimums around the global minimum. Because of curse of dimensionality, 20
 21 it is not possible to find all the local minimums around the global minimum. We 21
 22 randomly sample the OF in multiple directions. The nearest local minimum is 22
 23 computed for each randomly selected direction. We statistically find the distance 23
 24 to the nearest local minimum by using multiple samples. 24

25 Assume $F : \mathbb{R}^{2K} \rightarrow \mathbb{R}^+$ is the OF. A random direction in \mathbb{R}^{2K} is a vector in this 25
 26 space. Let us denote it by $v \in \mathbb{R}^{2K}$. First, we define a new function $h : \mathbb{R}^+ \rightarrow \mathbb{R}^+$ 26
 27 such that $h(t) = F(p_c + tv)$ where p_c is a vector containing the true locations of 27
 28 nodes. Second, we find the local minimum of h with the smallest positive t and 28
 29 denote it by t_1 . We repeat this procedure for T times and find all t_i 's. T is the 29
 30 number of samples. Finally, since it is expected that the defined metric has a larger 30
 31 value for more difficult-to-optimize OF, we define the nearest local minimum metric 31
 32 to be 32

$$33$$

$$34$$

$$35 \quad (5) \quad M_2(F) = \left(\frac{1}{T} \sum_{i=1}^T t_i \right)^{-1} . \quad 35$$

$$36 \quad 36$$

$$37 \quad 37$$

38 3.2.3. Measuring the Slope of OF Around the Solution 38

39
 40
 41 The Slope of OF (*i.e.*, the norm of OF's gradient) around the global minimum is 41
 42 a very important parameter in the convergence rate of the optimization algorithm. 42
 43 OFs with a small slope around the true location converge to the global minimum 43
 44 very slowly. 44

45 Thus, measuring the slope of the OF around the global minimum can be used 45
 46 to quantify the goodness of OF. Again, we measure slope of the OF in multiple 46
 47 random directions around the true locations, and statistically compute this metric. 47
 48 OFs with sharp slopes around the global minimum are easier to optimize. This can 48
 49 be seen in Figure 2 where the right side of the global minimum, p_{gm} , has a sharp 49
 50 slope. If the initial point of steepest descent algorithm is between p_{gm} and p_{m2} , 50
 51 it converges to the global minimum very fast. However, on the left side of global 51

1 minimum, p_{gm} , there is a gradual slope. Thus, the steepest descent algorithm would 1
 2 converge very slowly on the left side. We define the true locations' slope metric as 2

$$3 \quad (6) \quad M_3(F) = \left(\frac{1}{T} \sum_{i=1}^T \text{slope in } i\text{-th random direction} \right)^{-1} \quad 4$$

5
 6
 7 Note that the slope of the i -th random direction, v_i , is measured at $p_{gm} + \sigma v_i$ 8
 9 where σ is a small number and is a user's defined criterion. 9

10 3.2.4. Depth of the Non-Global Local Minima 10

11 Optimization problems that have an OF with deep local minimums around the 11
 12 global minimum are difficult to solve. A number of heuristic optimization methods 12
 13 take advantage of the shallow local minimums to avoid non-global local minimums, 13
 14 e.g., simulated annealing [11]. In figure 2, avoiding the local minimum at p_{l1} is much 14
 15 easier than local minimum at p_{l2} . 15
 16

17 We define the forth metric for quantifying the goodness of an OF on the data, as 17
 18 the depth of the non-global local minimums. We randomly select T local minimums 18
 19 around the true locations. Assuming that m_i is the OF value at the randomly 19
 20 selected local minimums, define 20
 21

$$22 \quad (7) \quad M_4(F) = \left(\frac{1}{T} \sum_{i=1}^T m_i \right)^{-1} \quad 23$$

24 4. Simulation Methodology 24

25 We find the linear effect of each parameter by studying all combinations of 25
 26 parameters. Assume each parameter has just two values. If we have k parameters 26
 27 then we have to study 2^k combinations that is computationally intractable. Instead, 27
 28 we use Plackett and Burman (PB) [23] fast simulation methodology that is a very 28
 29 well known method for reducing the number of simulations. Number of simulation 29
 30 in PB is proportional to the number of parameters. Although the PB method has 30
 31 not been used for the adhoc modeling and estimation problems, it was used for the 31
 32 simulations speedup in a number of other adhoc network problems [1][22][27][28]. 32
 33

34 In PB design, two values are assigned to each parameter: a normal value and 34
 35 an extreme value. The normal value is the typical value of the parameter while the 35
 36 extreme value is the value that is outside the typical range of the parameter. The 36
 37 extreme value often makes the problem either harder or easier to solve. A number 37
 38 of experiments with normal and extreme values of parameters are conducted. 38
 39

40 Experiments are arranged based on a given matrix denoted by the *design matrix*. 40
 41 Design matrix has k columns (k is the number of parameters) and s rows where s 41
 42 is the number of experiments the should be set up as follows. The elements of the 42
 43 design matrix are either 0 or 1. We set up an experiment for each row. Values of the 43
 44 parameters depend on the elements on the row: 0 indicates that the normal value 44
 45 of the parameter is used and 1 indicates that the extreme value of the parameter 45
 46 is used in the experiment corresponding to the row. 46
 47

48 Assume that we have selected an error metric, L_i , and an objective function 48
 49 metric, M_j . The OF itself denoted by F_{L_i} would be fixed. For each row of the 49
 50 51

design matrix, h , we setup an experiment based on the elements of that row and measure the goodness of the objective function $M_j(F_{L_i})$ and save it in another array element denoted by $r_{i,j,h}$. The corresponding values are summed up for computing the importance factor (IF) of each parameter. For each parameter P_t , we define

$$(8) \quad \text{IF}_{t,i,j} = \left| \sum_{h=1}^s \alpha_{h,t} r_{i,j,h} \right|$$

where s is the number of experiments (number of rows in the design matrix), and $\alpha_{h,t}$ is 1 if the extreme value of the parameter P_t is used in the h -th experiment; otherwise, $\alpha_{h,t}$ is -1 . The absolute value of IF is used to evaluate the effect of each parameter. The largest value indicates the most important parameter. For i -th error metric and j -th OF metric, $\text{IF}_{t,i,j} > \text{IF}_{u,i,j}$ means that the parameter P_t is more important than P_u . Thus, for each error metric, L_i , and for each objective function metric, M_j , we can rank parameters based on their effect on the estimated location. This ranking is denoted by $R_{i,j}$.

More precise results can be obtained by using the foldover design matrix [18]. In the foldover design matrix, all rows of the single design are repeated after its last row but 0s and 1s are exchanged in the repeated rows.

5. Combining Different Ranks

In this section, we explain how to combine the rankings of the parameters under study to obtain a global order for them. Using the ranking method in the previous section, we would have different rankings for various error metrics and OF metrics. Since there are three error metrics and four objective function metrics, there would be twelve different ranking lists for the importance of parameters; each parameter may have a different rank in each ranking list.

Each rank is obtained based on a specific property of the optimization problem. As it is explained in Section 3, for each error and objective function metric, the parameters are ranked based on the importance factor obtained from PB-design. IFs with large discrepancies lead to a stronger ranking compared to IFs with small discrepancies. Simply summing up the rankings would not necessarily determine which of the importance factors were better differentiating among the parameters.

For each ranking, $R_{i,j}$, and for each pair of parameters, P_s , P_t , we find the probability that P_s is more important than P_t . Based on the probabilities, we construct the global ranking.

Consider a specific error metric, L_i , and a specific objective function metric, M_j . Assume that the importance factor of the parameter P_t , $\text{IF}_{t,i,j}$, is normally distributed $\mathcal{N}(\lambda_{t,i,j}, \sigma^2)$. The observed value of $\text{IF}_{t,i,j}$ in a specific experiment is denoted by $if_{t,i,j}$. We normalize the importance factors to have a maximum value W . The mean of IFs are assumed to be uniformly distributed in $[0, W]$.

For each two parameters, P_s and P_t , given the BP-design experiment importance values $if_{s,i,j}$, and $if_{t,i,j}$, we find the probability: $Pr(\lambda_{s,i,j} \geq \lambda_{t,i,j} | IF_{s,i,j} = if_{s,i,j}, IF_{t,i,j} = if_{t,i,j})$. The conditional probability can be written in the Bayesian format as

$$(9) \quad \beta_{s,t,i,j} = \frac{Pr(\lambda_{s,i,j} \geq \lambda_{t,i,j} | IF_{s,i,j} = if_{s,i,j}, IF_{t,i,j} = if_{t,i,j}) \cdot Pr(IF_{s,i,j} = if_{s,i,j}, IF_{t,i,j} = if_{t,i,j} | \lambda_{s,i,j} \geq \lambda_{t,i,j}) \cdot Pr(\lambda_{s,i,j} \geq \lambda_{t,i,j})}{Pr(IF_{s,i,j} = if_{s,i,j}, IF_{t,i,j} = if_{t,i,j})}$$

Parameter	N_S	B_S	$\overline{e_S^2}$	$\text{MAX}_{\epsilon_2}^2$	$\text{PER}_{\epsilon_0}^2$	$\overline{D_S}$	MIN_{L_S}	MAX_{L_S}	$\overline{l_S}$	MIND_S
Normal Value	55	12	10 (m^2)	200 (m^2)	50	10	5 (m)	40 (m)	20 (m)	4
Extreme Value	80	3	50 (m^2)	500 (m^2)	20	6	10 (m)	60 (m)	30 (m)	3

TABLE 1
Normal and extreme values for the parameters.

Since there is no prior information about the distributions of $\lambda_{s,i,j}$ and $\lambda_{t,i,j}$, we assume that $\Pr(\lambda_{s,i,j} \geq \lambda_{t,i,j}) = \frac{1}{2}$. Furthermore,

$$\begin{aligned}
 & \Pr(IF_{s,i,j} = if_{s,i,j}, IF_{t,i,j} = v_{t,i,j} | \lambda_{s,i,j} \geq \lambda_{t,i,j}) = \\
 & \int_{x=0}^W \int_{y=x}^W \Pr(IF_{s,i,j} = if_{s,i,j}, IF_{t,i,j} = if_{t,i,j} | \lambda_{s,i,j} = y, \lambda_{t,i,j} = x) \frac{dy}{W} \frac{dx}{W} = \\
 (10) \quad & \frac{1}{W^2} \int_{x=0}^W \int_{y=x}^W \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-if_{s,i,j})^2}{2\sigma^2}} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-if_{t,i,j})^2}{2\sigma^2}} dy dx.
 \end{aligned}$$

Similarly, one can find

$$\begin{aligned}
 & \Pr(IF_{s,i,j} = if_{s,i,j}, IF_{t,i,j} = if_{t,i,j}) \\
 & = \Pr(IF_{s,i,j} = if_{s,i,j}, IF_{t,i,j} = if_{t,i,j} | \lambda_{s,i,j} \geq \lambda_{t,i,j}) \Pr(\lambda_{s,i,j} \geq \lambda_{t,i,j}) \\
 & + \Pr(IF_{s,i,j} = if_{s,i,j}, IF_{t,i,j} = if_{t,i,j} | \lambda_{s,i,j} < \lambda_{t,i,j}) \Pr(\lambda_{s,i,j} < \lambda_{t,i,j}).
 \end{aligned}$$

Now, for each parameter, P_t , we define the global importance factor, if_t ,

$$(11) \quad if_t = \sum_{i=1}^{N_{em}} \sum_{j=1}^{N_{om}} \sum_{s=1, s \neq t}^{N_p} \beta_{s,t,i,j}.$$

Parameters with a larger if_t have a higher probability of being important compared to the other parameters. We sort the parameters based on their corresponding if_t values.

6. Evaluation Results

We have applied the developed method to real distance measurement data for location estimation problem. Parameters that were described in Section 2 are ranked using our methodology. We illustrate how the various ranking lists differ. Then, we combine the rankings to obtain a global ranking.

The distance measurements data from the CENS lab [4] is used to evaluate the effect of each parameter. This database is based on the real distance measurements for SH4 nodes [8]. 91 nodes are located in fixed locations. Distance measurement is done multiple times and in different days. The distance measurements are based on the time of flight (ToF)[17] of the signals. In this method, the time of flight of an acoustic signal is used to determine the distance between two nodes. It was previously shown that the noise in the measurements is strongly non-static [7]. Therefore, parametric methods based on optimizing the results according to a fixed noise distribution do not yield good location estimations.

Parameter	DOF			NLM			SMAS			DNGLM		
	L_1	L_2	ML	L_1	L_2	ML	L_1	L_2	ML	L_1	L_2	ML
N_S	4	4	2	6	5	6	2	2	2	3	2	1
B_S	2	1	1	4	2	3	4	9	4	1	4	3
ϵ_S^2	1	2	3	2	3	4	3	3	3	5	9	6
$MAX_{\epsilon_m^2}$	6	8	7	9	10	10	6	4	5	7	7	8
$PER_{\epsilon_0^2}$	7	9	10	10	8	5	7	8	7	9	10	9
$\overline{D_S}$	3	3	4	1	1	1	1	1	1	2	1	2
$MINL_S$	8	6	8	7	6	9	8	5	10	4	3	10
$MAXL_S$	10	10	9	5	9	8	10	10	8	10	5	4
$\overline{I_S}$	9	5	5	8	7	7	9	7	9	6	6	5
$MIND_S$	5	7	6	3	4	2	5	6	6	8	8	7

TABLE 2
Importance of different parameters for different objective functions and metrics.

We have used Integer Linear Programming (ILP) to sample the database for drawing instances with specific properties. In each experiment, the PB-design matrix implies a specific value for each parameter. Extreme and normal values for parameters are shown in Table 1. The values are determined based on the real measurements' error. In all experiments, ϵ_0^2 is equal to $20(m^2)$.

The following abbreviations are used in this section.

- ML : Maximum Likelihood
- DOF : Drifting of the Objective Function (M_1)
- NLM : Nearest Local Minimum (M_2)
- SMAS : Slope Measurement Around the Solution (M_3)
- DNGLM: Depth of Non-Global Local Minimum (M_4)

Table 2 shows the result of PB-based evaluations. Each parameter is ranked based on the specific error metric and the specific OF metric. It can be seen that a specific parameter has different rankings under various error metrics and OF metrics. For example, the total number of nodes, N_S , is ranked 1, 2, 3, 4, 5, and 6 in different cases. Thus, a specific parameter does not have the same importance under various metrics. It can be seen that the number of nodes, N_S , and the number of beacons, B_S , are the two important parameters in most evaluations; $PER_{\epsilon_0^2}$ and $MAXL_S$ have overall low rankings.

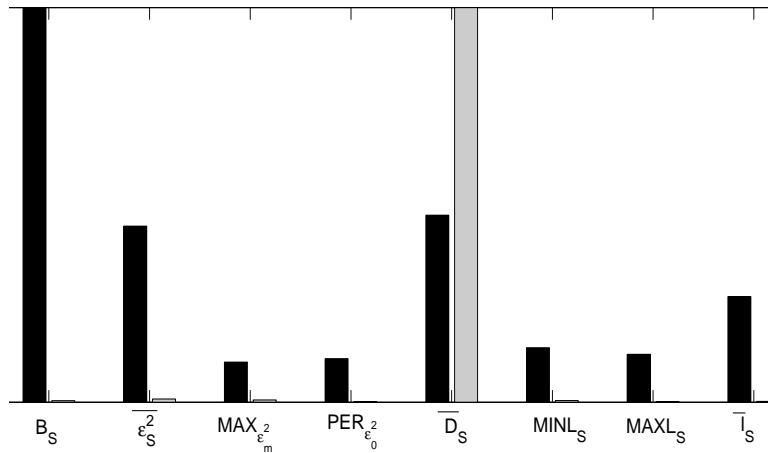


FIG 2. Importance of different parameters for different objective functions and metrics.

Parameter	N_S	B_S	$\overline{\epsilon_S^2}$	$\text{MAX}_{\epsilon_m^2}$	$\text{PER}_{\epsilon_0^2}$	$\overline{D_S}$	MINL_S	MAXL_S	$\overline{l_S}$	MIND_S
N_S	0	0.071	0.417	0.725	0.716	0.403	0.708	0.691	0.598	0.748
B_S	0.929	0	0.899	0.993	0.984	0.884	0.977	0.981	0.939	0.984
$\overline{\epsilon_S^2}$	0.583	0.101	0	0.787	0.786	0.476	0.756	0.754	0.660	0.798
$\text{MAX}_{\epsilon_m^2}$	0.275	0.007	0.213	0	0.490	0.193	0.464	0.477	0.354	0.515
$\text{PER}_{\epsilon_0^2}$	0.284	0.016	0.214	0.510	0	0.202	0.469	0.499	0.357	0.528
$\overline{D_S}$	0.597	0.116	0.524	0.807	0.798	0	0.785	0.795	0.678	0.821
MINL_S	0.292	0.023	0.244	0.536	0.531	0.215	0	0.519	0.392	0.545
MAXL_S	0.309	0.019	0.246	0.523	0.501	0.205	0.481	0	0.371	0.537
$\overline{l_S}$	0.402	0.061	0.340	0.646	0.643	0.322	0.608	0.629	0	0.671
MIND_S	0.252	0.016	0.202	0.485	0.472	0.179	0.455	0.463	0.329	0

TABLE 3

Drifting of objective function and L_2 metric: $\Pr(\lambda_{i,j,s} \geq \lambda_{i,j,t} | V_{i,j,s} = v_{i,j,s}, V_{i,j,t} = v_{i,j,t})$
 where the first column is P_s and the first row is P_t .

The comparative ranks of parameter pairs tend to vary as well. Figure 2 shows the normalized importance factor (IF) for two cases: DOF and SMAS with L_2 error metric. For DOF, the number of beacons B_S is strongly more important than the mean squared error $\overline{\epsilon_S^2}$. The mean degree of nodes, $\overline{D_S}$, is weakly more important than the mean squared error $\overline{\epsilon_S^2}$. The same behavior can be seen in SMAS. From our visual inspections, the number of nodes N_S and the mean degree of nodes $\overline{D_S}$ are the most important while others almost have the same importance factor (IF). The ranks of the mean squared error $\overline{\epsilon_S^2}$ and maximum edge length MAXL_S are 3 and 10 respectively. However, their importance factors are very close.

The discrepancy in the rank and comparative ranks confirms our postulation that averaging the parameter ranks is not the best way for combining them. Thus, we use the combining method that was introduced in Section 5. The probability comparisons for the values in Figure 2 are shown in Tables 3 and 4. The tables compare the importance of parameters. For example, for the DOF- L_2 , Figure 2 states that B_S is strongly more important than $\text{PER}_{\epsilon_0^2}$. Table 3 shows that the probability that the mean of B_S is larger than the mean of $\text{PER}_{\epsilon_0^2}$ is 0.984. Similarly, $\text{MAX}_{\epsilon_m^2}$ and $\text{PER}_{\epsilon_0^2}$ have approximately the same importance. The probability that the mean of $\text{MAX}_{\epsilon_m^2}$ is larger than the mean of $\text{PER}_{\epsilon_0^2}$ is 0.49. This probability value is close to 0.5, meaning that there is not enough information to compare the values.

Table 4 compares the importance factors of SMAS for the L_2 error metric. Table 4 confirms the result. The rows corresponding to N_S , and $\overline{D_S}$ have values close to 1 confirming the high importance of the two parameters. When comparing other parameters, the probability that one parameter is greater than the other is about 0.5. It confirms our previous postulation that simple rankings are not sufficient for concluding the global parameter ordering and the importance factors are significant as well.

The global ranking based on the introduced combining method in Section 5 is shown in Table 5. The table indicates that the mean degree of nodes $\overline{D_S}$ is the most important parameter. This result is consistent with Table 2 where the mean degree of nodes $\overline{D_S}$ is the most important parameter in the seven scenarios.

The global ranking results could be used to improve the goodness of location estimations in ad-hoc networks. To deploy a network or on an already deployed network, one could exploit the results by considering the analyzed effect of each parameter on the estimated location's accuracy. Based on the constraints of the problem, the best parameters for improving the estimated locations could be determined. For example, when there are limitations for the mean degree of the graph,

Parameter	N_S	B_S	$\overline{\epsilon_S^2}$	$\text{MAX}_{\epsilon_m^2}$	$\text{PER}_{\epsilon_0^2}$	$\overline{D_S}$	MINL_S	MAXL_S	$\overline{l_S}$	MIND_S
N_S	0	0.947	0.947	0.936	0.944	0.285	0.939	0.931	0.937	0.958
B_S	0.053	0	0.493	0.506	0.504	0.017	0.501	0.496	0.500	0.504
$\overline{\epsilon_S^2}$	0.053	0.507	0	0.509	0.505	0.018	0.504	0.516	0.507	0.511
$\text{MAX}_{\epsilon_m^2}$	0.064	0.494	0.491	0	0.505	0.017	0.504	0.506	0.499	0.499
$\text{PER}_{\epsilon_0^2}$	0.056	0.496	0.495	0.495	0	0.017	0.492	0.496	0.502	0.500
$\overline{D_S}$	0.715	0.983	0.982	0.983	0.983	0	0.975	0.984	0.974	0.980
MINL_S	0.061	0.499	0.496	0.496	0.508	0.025	0	0.506	0.501	0.488
MAXL_S	0.069	0.504	0.484	0.494	0.504	0.016	0.494	0	0.502	0.494
$\overline{l_S}$	0.063	0.500	0.493	0.501	0.498	0.026	0.499	0.498	0	0.505
MIND_S	0.042	0.496	0.489	0.501	0.500	0.020	0.512	0.506	0.495	0

TABLE 4
SMAS and L₂ metric: Pr(λ_{i,j,s} ≥ λ_{i,j,t} | V_{i,j,s} = v_{i,j,s}, V_{i,j,t} = v_{i,j,t}) where the first column is P_s and the first row is P_t.

Parameter	N_S	B_S	$\overline{\epsilon_S^2}$	$\text{MAX}_{\epsilon_m^2}$	$\text{PER}_{\epsilon_0^2}$	$\overline{D_S}$	MINL_S	MAXL_S	$\overline{l_S}$	MIND_S
Rank	2	3	4	8	10	1	6	9	7	5

TABLE 5
Global ranks.

one can increase the number of nodes in the network to increase the accuracy of the estimated location. Note that, changing one parameter typically only improves the accuracy up to a certain point; further changing the parameter would not yield an improvement in the estimation accuracy.

7. Conclusion

We introduce a systematic methodology for determining the challenge of modeling a pertinent adhoc network data set. The complex modeling problem is studied as an instance of a nonlinear optimization problem that consists of an objective function (OF) and a set of constraints. The data set is the optimization input and the estimated model is the output. We characterize the input by a set of its characteristic parameters. We define four new metrics that can be used to evaluate the goodness of an input for being optimized by a specific OF. The introduced metrics are: (1) drifting of the OF, (2) distance to the nearest local minimum, (3) the slope of the OF around the solution, and (4) the depth of the non-global local minima. We employ Plackett and Burmann simulation methodology to systematically evaluate the linear impact of various input parameters under each metric. Finally, we present a method for combining the effect of parameters under different metrics to determine the global impact of each parameter. We utilize the new methodology for estimating the locations of the nodes in an ad-hoc network where the distance measurement data is available. Three common forms of OF are considered: L_1 , L_2 and L_∞ . Our evaluations show that the mean degree on the nodes and the number of nodes in the network are the two most important parameters for estimating the locations.

Acknowledgement

This work is partly supported by the National Science Foundation (NSF) CAREER Award under grant number 0644289.

References

- [1] BARRETT, C. and MARATHE, A. and MARATHE, M.V. and DROZDA, M. (2002). Characterizing the interaction between routing and MAC protocols in ad-hoc networks. *MobiHoc*, 92-103.
- [2] BISWAS, P. and YE, Y. (2004). Semidefinite programming for ad hoc sensor network localization. *Information Processing in Sensor Networks (IPSN)*, 2673-2684.
- [3] BULLARD, C. and SEBALD, A. (1988). Monte carlo sensitivity analysis of input-output models. *The Review of Economics and Statistics*, **22.4** 708-712.
- [4] *CENS: Center for Embedded Networked Sensing at UCLA*. <http://research.cens.ucla.edu/>.
- [5] *The Cricket indoor location system at MIT*. <http://cricket.csail.mit.edu>.
- [6] EREN, T. and GOLDENBERG, D. and WHITELEY, W. and YANG, Y.R. and MORSE, A. and ANDERSON, B. and BELHUMEUR, P. (2004). Rigidity, computation, and randomization in network localization. *INFOCOM*, 2673-2684.
- [7] FENG, J. and GIROD, L. and POTKONJAK, M. (2006). Consistency-based on-line localization in sensor networks. *Distributed Computing in Sensor Systems (DCOSS)*, 529-545.
- [8] FENG, J. and GIROD, L. and POTKONJAK, M. (2006). Location discovery using data-driven statistical error modeling. *INFOCOM*, 1-14.
- [9] GOLDENBERG, D. and KRISHNAMURTHY, A. and MANESS, W. and YANG, Y. and YOUNG, A. and MORSE, A. and SAVVIDES, A. (2006). Network localization in partially localizable networks. *INFOCOM*, 313-326.
- [10] HASTIE, T. and TIBSHIRANI, R. and FRIEDMAN, J. (2001). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer-Verlag, Germany.
- [11] KIRKPATRICK, S. and GELATT, C.D. JR. and VECCHI, M.P. (1983). Optimization by simulated annealing. *Science*, **220.4598**, 671-680.
- [12] KOUSHANFAR, F. and DAVARE, A. and NGUYEN, D. and SANGIOVANNI-VINCENTELLI, A. and POTKONJAK, M. (2007). Techniques for maintaining connectivity in wireless ad-hoc networks under energy constraints. *ACM Trans. on Embedded Computing Systems*, **6.3**, 16-37.
- [13] KOUSHANFAR, F. and TAFT, N. and POTKONJAK, M. (2006). Sleeping coordination for comprehensive sensing using isotonic regression and domatic partitions. *INFOCOM*, 1-13.
- [14] LEHMANN, E. (2006). *Nonparametrics: Statistical Methods Based on Ranks*. Springer-Verlag, Germany.
- [15] MEGERIAN, S. and KOUSHANFAR, F. and QU, G. and VELTRI, G. and POTKONJAK, M. (2002). Exposure in wireless sensor networks: Theory and practical solutions. *ACM Journal of Wireless Networks*, **8.5**, 443454.
- [16] MEGERIAN, S. and KOUSHANFAR, F. and POTKONJAK, M. and SRIVASTAVA, M. (2005). Worst- and best-case coverage in sensor networks. *IEEE Transactions on Mobile Computing*, **4.1**, 84-92.
- [17] LANZISERA, S. and LIN, D. and PISTER, K. (2006). Rf time of flight ranging for wireless sensor network localization. *WISES*, 1-12.
- [18] MONTGOMERY, D. (2001). *Design and Analysis of Experiments*, 5th ed. Wiley, New York.
- [19] PATWARI, N. and ASH, J. and KYPEROUNTAS, S. and HERO, A. III and MOSES, R. and CORREAL, R. (2005). Locating the nodes: cooperative localization in wireless sensor networks. *IEEE Signal Processing Magazine*, **22(4)** 54-69.
- [20] PRIYANTHA, N. and CHAKRABORTY, A. and BALAKRISHNAN, H. (2000). The cricket location-support system. *MOBICOM*, 32-43.
- [21] *OSL: Open Systems Laboratory at UIUC*. <http://www-osl.cs.uiuc.edu/research?action=topic&topic=Sensor+Networks>.
- [22] PERKINS, D. and HUGHES, H. and OWEN, C. (2002). Factors affecting the performance of ad hoc networks. *ICC*, 2048-2052.
- [23] PLACKETT, R.L. and BURMAN, J.P. (1946). The design of optimum multifactorial experiments. *Biometrika*, **33.4**, 305-325.
- [24] SALTELLI, A. and TARANTOLA, S. and CAMPOLONGO, F. and RATTO, M. (2004). *Sensitivity Analysis in Practice: A Guide to Assessing Scientific Models*. Wiley, England.
- [25] SAXE, J. (1979) Embeddability of weighted graphs in k-space is strongly np-hard. *Allerton Conf. in Communications, Control, and Computing*, 480-489.
- [26] SLIJEPEVIC, S. and MEGERIAN, S. and POTKONJAK, M. (2002). Location errors in wireless embedded sensor networks: sources, models, and effects on applications. *Mobile Computing and Communications Review*, **6.3**, 67-78.
- [27] TOTARO, M. and PERKINS, D. (2005). Using statistical design of experiments for analyzing mobile ad hoc networks. *MSWiM*, 159-168.
- [28] VADDE, K. and SYROTIUK, V. (2004). Factor interaction on service delivery in mobile ad hoc networks. *IEEE Journal on Selected Areas in Communications (JSAC)*, **22.7**, 13351346.