

Stochastic Multidimensional Scaling

Ketan Rajawat, *Member, IEEE*, and Sandeep Kumar, *Student Member, IEEE* .

Abstract

Multidimensional scaling (MDS) is a popular dimensionality reduction techniques that has been widely used for network visualization and cooperative localization. However, the traditional stress minimization formulation of MDS necessitates the use of batch optimization algorithms that are not scalable to large-sized problems. This paper considers an alternative stochastic stress minimization framework that is amenable to incremental and distributed solutions. A novel linear-complexity stochastic optimization algorithm is proposed that is provably convergent and simple to implement. The applicability of the proposed algorithm to localization and visualization tasks is also expounded. Extensive tests on synthetic and real datasets demonstrate the efficacy of the proposed algorithm.

Index Terms

Multidimensional Scaling, Stochastic SMACOF, Visualization, Localization.

I. INTRODUCTION

Multidimensional scaling addresses the problem of embedding relational data onto a low-dimensional subspace. Originally proposed in the context of psychometrics and marketing [1], MDS and its variants have since found applications in social networks [2]–[6], genomics [7], computational chemistry [8], machine learning [9], and wireless networks [10]. As an exploratory technique, MDS is often used as a first step towards uncovering the structure inherent to high-dimensional data. In the context of machine learning and data mining, the pairwise dissimilarities are calculated using high- or infinite-dimensional nodal attributes, and MDS yields a distance-preserving, low-dimensional embedding. Of particular importance are the embeddings obtained in two or three dimensional euclidean spaces, that serve as perceptual maps for visualizing relationships between objects. In the context of social networks, such representations reveal

Ketan Rajawat and Sandeep Kumar are with the Department of Electrical Engineering, Indian Institute of Technology Kanpur, Kanpur, UP 208016, India, email: {ketan, sandkr}@iitk.ac.in

interconnections between people and communities, and are often more insightful than simpler metrics such as centrality and density. Different from the classical MDS framework that utilizes principal component analysis, modern MDS formulations are based on the minimization of a non-convex stress function [1]. Since the stress function is a weighted sum of squared fitting errors, it allows for the possibility of missing and noisy dissimilarities. Consequently variants of the stress minimization problem have been developed for robust MDS [11], visualization of time-varying data [12], and cooperative localization of static [10], [13]–[15] and mobile networks [13]. Popular algorithms for solving the stress minimization problem include ‘scaling by majorizing a complicated function’ (SMACOF) [1], semidefinite programming [16], alternating directions method of multipliers [14], [15], and distributed SMACOF [10].

The attractiveness of the MDS framework has however started to diminish with the advent of the data deluge. Specifically, when embedding N objects, the per-iteration complexity and memory requirements of the aforementioned algorithms increase at least as $\mathcal{O}(N^2)$, making them impractical for large-scale problems. To this end, approximate versions of SMACOF have been proposed for large-scale visualization applications [17], [18]. Nevertheless, most approximate MDS algorithms are still too complex for large-scale data, and cannot be generalized to other applications such as cooperative localization of large networks.

Visualization or localization of time-varying data is even more challenging since the iterative majorization algorithm must converge at every time instant [10], [12], [13]. In mobile sensor networks, carrying out a large number of iterations at each time instant incurs a tremendous communication overhead, and is generally impractical. For instance, the distributed weighted MDS approach [10] still requires at least N operations per iteration per time instant, which is prohibitive for large networks. For large-scale applications, where localization or visualization is constrained by the per-iteration complexity and memory requirements, it is instead desirable to have an online algorithm. Towards this end, the goal is to obtain an adaptive algorithm that processes dissimilarity measurements in a sequential or online manner. For instance, an adaptive algorithm can allow visualization of large networks by reading and processing the pairwise dissimilarities in small batches. Similarly, the communication cost required for large-scale network localization can be reduced by processing only a few range measurements at a time.

This paper considers the stress minimization problem in a stochastic setting, where the dis-

similarity measurements and the weights are modeled as random time-varying quantities with unknown distributions. The first contribution of this paper is a novel stochastic SMACOF algorithm that processes the dissimilarities in an online fashion, and is therefore applicable to both static and time-varying scenarios (Sec. III). The proposed algorithm is not only scalable, but is also amenable to a distributed and asynchronous implementation in ad hoc networks (Sec. IV). As our second contribution, it is shown that the trajectory of the stochastic SMACOF algorithm remains close to that of an averaged algorithm, which itself converges to a stationary point of the stochastic stress minimization problem (Sec. III-B). The analysis borrows tools from spectral graph theory, stochastic approximation, and convergence analysis of the SMACOF algorithm. Finally, as the third contribution, the performance of the proposed algorithm is tested extensively on various synthetic and real-world data sets (Sec. V). The numerical tests confirm the applicability of the stochastic SMACOF algorithm to a variety of scenarios.

The notation used in this paper is as follows. Bold upper (lower) case letters denote matrices (vectors). The (m, n) -th entry of a matrix \mathbf{A} is denoted by $[\mathbf{A}]_{mn}$. \mathbf{I}_N is the $N \times N$ identity matrix, $\mathbf{0}$ denotes the all-zero matrix or vector, and $\mathbf{1}$ denotes the all-one matrix or vector, depending on the context. For a vector \mathbf{x} , $\|\mathbf{x}\|$ denotes its ℓ_2 norm. For a matrix \mathbf{A} , $\|\mathbf{A}\|$ denotes its Frobenius norm, $\|\mathbf{A}\|_2$ denotes the ℓ_2 norm, $\text{tr}(\mathbf{A})$ denotes its trace, and $\det(\mathbf{A})$ denotes its determinant.

II. BACKGROUND AND PROBLEM STATEMENT

A. Classical MDS and SMACOF

The classical MDS framework seeks P -dimensional embedding vectors $\{\mathbf{x}_n\}_{n=1}^N$, given the pairwise distances or dissimilarities $\{\delta_{mn}\}_{(m,n) \in \mathcal{E}}$, where $\mathcal{E} \subseteq \{(m, n) \mid 1 \leq m < n \leq N\}$, between N different nodes or objects, denoted by the set $\mathcal{N} := \{1, \dots, N\}$. The embedding vectors, collected into the rows of $\mathbf{X} \in \mathbb{R}^{N \times P}$, are estimated by solving the following non-convex optimization problem [1]

$$\hat{\mathbf{X}} = \arg \min_{\mathbf{X}} \sum_{1 \leq m < n \leq N} w_{mn} (\delta_{mn} - \|\mathbf{x}_m - \mathbf{x}_n\|_2)^2 \quad (1)$$

where w_{mn} is the weight associated with the measurement δ_{mn} , and is set to zero for all $(m, n) \notin \mathcal{E}$. The non-zero weights can be chosen in a number of ways, depending on the application, and are often simply set to one. The objective function in (1) is referred to as the stress function,

and is henceforth denoted by $\sigma(\mathbf{X})$. It can be seen that the optimum $\hat{\mathbf{X}}$ obtained in (1) is not unique, and exhibits translational, rotational, and reflectional ambiguity.

The stress-minimization problem in (1) is non-convex, and can be solved up to a local optimum using the well known SMACOF algorithm. Expanding the stress function, we obtain

$$\sigma(\mathbf{X}) = \sum_{m < n} w_{mn} (\delta_{mn}^2 + \|\mathbf{x}_m - \mathbf{x}_n\|^2 - 2\delta_{mn} \|\mathbf{x}_m - \mathbf{x}_n\|) \quad (2)$$

$$= \sum_{m < n} w_{mn} \delta_{mn}^2 + \text{tr}(\mathbf{X}^T \mathbf{L} \mathbf{X}) - 2\text{tr}(\mathbf{X}^T \mathbf{B}(\mathbf{X}) \mathbf{X}) \quad (3)$$

where,

$$[\mathbf{L}]_{mn} = \begin{cases} -w_{mn} & m \neq n \\ \sum_{k=1}^m w_{mk} & m = n \end{cases} \quad (4)$$

$$[\mathbf{B}(\mathbf{X})]_{mn} = \begin{cases} -\frac{w_{mn} \delta_{mn}}{\|\mathbf{x}_m - \mathbf{x}_n\|} & m \neq n, \mathbf{x}_m \neq \mathbf{x}_n \\ 0 & m \neq n, \mathbf{x}_m = \mathbf{x}_n \\ -\sum_{k=1}^m [\mathbf{B}(\mathbf{X})]_{mk} & m = n \end{cases} \quad (5)$$

The SMACOF algorithm works by iteratively majorizing the last term in (3) with a linear function and subsequently minimizing the majorized stress function with respect to \mathbf{X} . Starting with an initial $\hat{\mathbf{X}}^{(0)}$, the SMACOF update at the k -th iteration entails carrying out the following update:

$$\hat{\mathbf{X}}^{(k+1)} = \arg \min_{\mathbf{X}} \text{tr}(\mathbf{X}^T \mathbf{L} \mathbf{X}) - 2\text{tr}(\mathbf{X}^T \mathbf{B}(\hat{\mathbf{X}}^{(k)}) \hat{\mathbf{X}}^{(k)}) \quad (6)$$

$$= \mathbf{L}^\dagger \mathbf{B}(\hat{\mathbf{X}}^{(k)}) \hat{\mathbf{X}}^{(k)} \quad (7)$$

where (7) follows since $\mathbf{B}(\mathbf{X})\mathbf{X}$ lies in the range space of \mathbf{L} . Observe that since \mathbf{L} is rank-deficient, the solution to (6) is not unique. However, when the weights $\{w_{mn}\}$ specify a fully connected graph $\mathcal{G} := (\{1, \dots, N\}, \mathcal{E})$, both \mathbf{L} and $\mathbf{B}(\mathbf{X})$ have rank $N - 1$, with the null space of \mathbf{L} being $\mathbf{1}$. Therefore, any solution to (6) is of the form $\mathbf{L}^\dagger \mathbf{B}(\hat{\mathbf{X}}^{(k)}) \hat{\mathbf{X}}^{(k)} + \mathbf{1}c$ for $c \in \mathbb{R}$. Further, if the initial $\mathbf{X}^{(0)}$ is chosen such that it is centered at the origin, i.e., $\mathbf{1}^T \mathbf{X}^{(0)} = \mathbf{0}$, the updates in (7) ensure that $\mathbf{1}^T \mathbf{X}^{(k)} = \mathbf{0}$ for all $k \geq 1$.

B. Stochastic MDS

This paper considers the MDS problem in a stochastic setting, where the weights, and dissimilarities or distance measurements are random variables with unknown distributions. Specifically,

given $\{\delta_{mn}(t)\}$ and $\{w_{mn}(t)\}$, the stochastic stress minimization problem is formulated as

$$\min_{\mathbf{X}} \bar{\sigma}(\mathbf{X}) := \sum_{m < n} \mathbb{E}[w_{mn}(t)(\delta_{mn}(t) - \|\mathbf{x}_m - \mathbf{x}_n\|)^2]. \quad (8)$$

In the absence of the distribution information, the expression for $\bar{\sigma}(X)$ cannot be evaluated in closed-form, and the SMACOF algorithm cannot be applied. Instead, (8) must be solved using a stochastic optimization algorithm. Of particular interest are the so-called *online* algorithms that can process the observations $\{\delta_{mn}(t)\}, \{w_{mn}(t)\}$ in an incremental manner. Within this context, efficient implementations of the stochastic (sub-)gradient descent (SGD) method have been used to solve very large-scale problems [19]. The SGD updates utilize the subgradient of the instantaneous objective function, and for the present case, take the form:

$$\check{\mathbf{X}}_{t+1} = \check{\mathbf{X}}_t + \mu (\mathbf{B}_t(\check{\mathbf{X}}_t)\check{\mathbf{X}}_t - \mathbf{L}_t\check{\mathbf{X}}_t) \quad (9)$$

where $\mu \in (0, 1)$ is the learning rate or step size parameter. While the performance of the SGD has been well-studied for convex problems, the same is not true for non-convex problems, such as the one in (8). Indeed, the standard SGD algorithm does not necessarily converge for many non-convex problems [20]. In the present case also, the SGD method exhibits divergent behavior; see Sec. V. The general-purpose stochastic majorization-minimization method [21] is also not applicable in the present case since it requires a strongly convex surrogate function. On the other hand, problem-specific stochastic algorithms have been developed and applied with great success. Examples include the online expectation-maximization and the online matrix factorization approaches [19], [22]. Along similar lines, the next section details the stochastic version of the SMACOF algorithm, and studies its asymptotic properties.

III. ONLINE EMBEDDING VIA STOCHASTIC SMACOF

A. Algorithm outline

Given $\{\delta_{mn}(t)\}$ and $\{w_{mn}(t)\}$, and starting with an arbitrary origin-centered $\hat{\mathbf{X}}_0$, the updates for the proposed stochastic SMACOF algorithm take the form,

$$\hat{\mathbf{X}}_{t+1} = (1 - \mu)\hat{\mathbf{X}}_t + \mu \mathbf{L}_t^\dagger \mathbf{B}_t^\epsilon(\hat{\mathbf{X}}_t)\hat{\mathbf{X}}_t \quad t \geq 0 \quad (10)$$

$$\text{where, } [\mathbf{B}_t^\epsilon(\mathbf{X})]_{mn} = \begin{cases} -\frac{w_{mn}(t)\delta_{mn}(t)}{\sqrt{\|\mathbf{x}_m - \mathbf{x}_n\|^2 + \epsilon_x}} & m \neq n \\ -\sum_{k=1}^N [\mathbf{B}_t^\epsilon(\mathbf{X})]_{mk} & m = n \end{cases} \quad (11)$$

with ϵ_x being a small positive constant that ensures that the entries of $\mathbf{B}_t^\epsilon(\mathbf{X})$ stay bounded for all \mathbf{X} . The update rule can be viewed a stochastic version of the SMACOF algorithm with the following modifications (a) at each time instant, only one iteration of SMACOF is executed using the modified definition of $\mathbf{B}^\epsilon(\mathbf{X})$ in (11); (b) the estimated coordinates $\hat{\mathbf{X}}_t$ at time t are used for initialization at $t+1$; and (c) the estimated coordinates $\hat{\mathbf{X}}_{t+1}$ are constructed by taking a convex combination of $\hat{\mathbf{X}}_t$ and the SMACOF output. The last modification endows the algorithm with tracking capabilities since the parameter μ may be interpreted as the forgetting factor, and can be tuned in accordance with the rate of change of $\{\delta_{mn}(t)\}$ and $\{w_{mn}(t)\}$. For example, the embedding at time $t+1$ can be forced to be close to those at time t by setting $\mu \ll 1$. Finally, the proposed update rule subsumes the SMACOF algorithm for static scenarios, where we set $\delta_{mn}(t) = \delta_{mn}$ and $w_{mn}(t) = w_{mn}$ for all t , and $\mu = 1$.

The update rule in (10) is valid only if the graph \mathcal{G}_t defined by $\{w_{mn}(t)\}$ is connected for all $t \geq 1$. In the case when \mathcal{G}_t has more than one connected component, the coordinates within each component must be updated separately. Let \mathcal{C}_t^j be the set of nodes belonging to the j -th component and \mathbf{I}_t^j be the $|\mathcal{C}_t^j| \times N$ selection matrix containing the rows of \mathbf{I}_N corresponding to the elements in \mathcal{C}_t^j . Defining $\mathbf{L}_t^{(j)} := \mathbf{I}_t^j \mathbf{L}_t \mathbf{I}_t^{jT}$ and $\mathbf{B}_t^\epsilon(\mathbf{X}_t^{(j)}) := \mathbf{I}_t^j \mathbf{B}_t^\epsilon(\mathbf{X}_t) \mathbf{I}_t^{jT}$, the update rule for the nodes in \mathcal{C}_t^j is given by

$$\hat{\mathbf{X}}_{t+1}^{(j)} = (\mathbf{I} - \mu \mathbf{J}_t) \hat{\mathbf{X}}_t^{(j)} + \mu (\mathbf{L}_t^{(j)})^\dagger \mathbf{B}_t^\epsilon(\hat{\mathbf{X}}_t^{(j)}) \hat{\mathbf{X}}_t^{(j)} \quad (12)$$

where $\mathbf{J}_t := \mathbf{I} - \mathbf{1}\mathbf{1}^T/|\mathcal{C}_j(t)|$, is the $|\mathcal{C}_j(t)| \times |\mathcal{C}_j(t)|$ centering matrix which ensures that the coordinate center of each component does not change after the update, i.e., $\mathbf{1}^T \hat{\mathbf{X}}_{t+1}^{(j)} = \mathbf{1}^T \hat{\mathbf{X}}_t^{(j)}$.

The general update rule

$$\hat{\mathbf{X}}_{t+1} = (\mathbf{I} - \mu \mathbf{L}_t^\dagger \mathbf{L}_t) \hat{\mathbf{X}}_t + \mu \mathbf{L}_t^\dagger \mathbf{B}_t^\epsilon(\hat{\mathbf{X}}_t) \hat{\mathbf{X}}_t \quad (13)$$

subsumes the forms specified in (10) and (12), irrespective of the number of connected components in \mathcal{G}_t , since it holds that

$$[\mathbf{L}_t^\dagger \mathbf{L}_t]_{mn} = \begin{cases} 1 - 1/|\mathcal{C}_t^j| & m = n \in \mathcal{C}_t^j \\ -1/|\mathcal{C}_t^j| & m \neq n, m, n \in \mathcal{C}_t^j \\ 0 & \text{otherwise.} \end{cases} \quad (14)$$

In contrast to the classical SMACOF algorithm, the proposed algorithm is flexible enough to be used in a number of different scenarios. As already discussed, a specific choice of parameters

allows us to interpret the SMACOF algorithm as a special case of the proposed algorithm. On the other hand, the stochastic SMACOF can also be used to solve very large-scale MDS problems, where the full set of measurements $\{\delta_{mn}\}$ cannot be processed simultaneously. Instead, it is possible to apply (13) on a small subset of observations, corresponding to a subgraph \mathcal{G}_t . A special case occurs when exactly one edge is chosen per time instant and per cluster, i.e., $|\mathcal{C}_t^j| = 2$, and the updates in (13) reduce to those in encountered in the stochastic proximity embedding (SPE) algorithm [6],

$$\begin{aligned} \mathbf{x}_i(t+1) = & (1 - \mu)\mathbf{x}_i(t) + \mu \frac{\delta_{ij}(t)}{\|\mathbf{x}_i(t) - \mathbf{x}_j(t)\|} \mathbf{x}_i(t) \\ & + \mu \left(1 - \frac{\delta_{ij}(t)}{\|\mathbf{x}_i(t) - \mathbf{x}_j(t)\|} \right) \mathbf{x}_j(t) \end{aligned} \quad (15)$$

and likewise for node j . The proposed stochastic SMACOF is therefore a generalization of SPE, applied to components of arbitrary sizes. Since the updates in (13) for any two clusters \mathcal{C}_t^j and \mathcal{C}_t^k do not depend on each other, the proposed algorithm can also be implemented in a distributed and asynchronous manner. Such an implementation is particularly suited to the range-based localization problems that arise in wireless networks.

Finally, akin to the classical adaptive filtering algorithms such as LMS, the proposed algorithm can also be applied to time-varying scenarios, i.e., when $\delta_{mn}(t)$ is non-stationary. The applications of interest include localization of time-varying networks, and visualization of time-varying data. In both cases, the first term $(\mathbf{I} - \mu\mathbf{J}_t)\hat{\mathbf{X}}_t^{(j)}$ in the update (13) serves as a momentum term. That is, a small μ encourages $\hat{\mathbf{X}}_{t+1}$ to stay close to $\hat{\mathbf{X}}_t$, resulting in a smooth trajectory of $\{\hat{\mathbf{X}}_t\}$. On the other hand, a large value of μ enables tracking in highly time-varying scenarios, while making the updates sensitive to noise [23, Ch-21] [24, Ch-9]. Further implementation details pertaining to the localization and visualization problems are discussed in Sec. IV. Before proceeding with the asymptotic analysis, the following remark is due.

Remark 1. Building further on the link with adaptive algorithms, μ may be interpreted as a forgetting factor that downweights the past information. When μ is a constant that is strictly greater than zero, the algorithm forgets the old data exponentially quickly, thus offering superior tracking capability. In contrast, it is possible to have a long-memory version of the algorithm with a time-varying $\mu_t \rightarrow 0$. As $t \rightarrow \infty$, such an algorithm would no longer track the changes in $\delta_{mn}(t)$, and can be applied to a static scenarios where the algorithm can stop once the embeddings

converge. While the bounds developed here apply only to the case of constant $\mu > 0$, diminishing step size is in fact utilized in Sec. V.

B. Asymptotic Performance

In general, establishing convergence of stochastic algorithms for non-convex problems is quite challenging [20]. Here, the asymptotic performance of the proposed algorithm is established in two steps. First, it is shown that the trajectory of the stochastic SMACOF algorithm stays close to that of an averaged algorithm, in an almost sure sense. This part involves establishing a hovering theorem, and utilizes techniques from stochastic approximation [24]–[26]. Next, it is shown that the averaged algorithm converges to a stationary point of (8).

1) *Assumptions:* For the purposes of establishing convergence, a simplified setting is considered, wherein the graph \mathcal{G}_t at each t consists of $N/p \geq 1$ components of size p each. Let $j_m(t) := \{j \mid m \in \mathcal{C}_t^j\}$ be the index of the component to which node m belongs at time t , and define $\Theta_t \in \mathbb{R}^{N \times N}$ such that

$$[\Theta_t]_{mn} := \begin{cases} -1/N & j_m(t) \neq j_n(t) \\ -1/N + \mu/p & j_m(t) = j_n(t), m \neq n \\ (1 - \mu) - 1/N + \mu/p & m = n. \end{cases}$$

- (A1) The random processes $\{w_{mn}(t)\}_{t \geq 0}$ and $\{\delta_{mn}(t)\}_{t \geq 0}$ are independent identically distributed (i.i.d.).
- (A2) The random variables $\{\delta_{mn}(t)\}$ have support $(0, C_\delta]$, while the weights $\{w_{mn}(t)\}$ have support $\{0\} \cup [\epsilon_w, 1]$.
- (A3) The online algorithm is initialized such that $\|(\mathbf{I} - \mathbf{1}\mathbf{1}^T/N)\hat{\mathbf{X}}_0\| \leq C_x$.
- (A4) There exists t_0 such that for any $\mu \in (0, 1)$, there exists $\rho \in (0, 1)$ such that $\|\prod_{s=\tau+1}^t \Theta_s\|_2 < \rho^{t-\tau}$ for all $t - \tau \geq t_0$.
- (A5) For each t , the non-zero weights $\{w_{mn}(t)\}_{m,n}$ are i.i.d. with $\bar{w} := \mathbb{E}[w_{mn}(t)]$.

The i.i.d. assumption in (A1) is standard in the analysis of most stochastic approximation algorithms. For the applications at hand, the support of $\delta_{mn}(t)$ and $w_{mn}(t)$ is naturally finite. It is required from (A2) that the non-zero weights be bounded away from zero. Such a condition is required to ensure the numerical stability of the Laplacian system of equations that must be

solved at every iteration [cf. (10), (13)]. Specifically, it is shown in Appendix that **(A2)** implies the following result

Lemma 1. *Under (A2), it holds that $\|\mathbf{L}_t^\dagger\|_2 \leq \epsilon_L := (N-1)^2/2\epsilon_w$ for all $t \geq 1$.*

The proof of Lemma 1 is provided in Appendix A. The initial configuration can always be normalized to satisfy the bound in **(A3)**. Assumption **(A4)** restricts the extent to which the graphs \mathcal{G}_t can stay disconnected over time. To obtain intuition on **(A4)**, observe first the largest eigenvalue of Θ_t is $1 - \mu$ if \mathcal{G}_t has a single connected component and one otherwise. Consequently, if all $\{\mathcal{G}_s\}_{s=\tau+1}^t$ are connected, **(A4)** holds with $\varrho = 1 - \mu$. Conversely, it holds that $\|\prod_{s=\tau+1}^t \Theta_s\|_2 = 1$ if and only if (a) each $\{\mathcal{G}_s\}_{s=\tau+1}^t$ has more than one components, and (b) the components do not change over time, i.e., $j_m(t) = j_n(t)$ for all m, n , and t . Intuitively, **(A4)** allows $\{\mathcal{G}_s\}$ to have multiple connected components at each $s \geq 1$, as long as the nodes belonging to these components keep changing over time.

Finally, **(A5)** is perhaps the most restrictive, and may not always be easy to satisfy. For instance, the weights are not identically distributed in the context of dynamic network localization (cf. Sec. IV-A), since non-zero weights are often assigned to neighboring nodes only. Likewise, weights selected via Sammon mapping also result in non-identically distributed weights. The assumption however greatly simplifies the proof of convergence for the averaged algorithm. Having stated the assumptions, the averaging analysis is presented in the subsequent subsection.

2) *Hovering Theorem:* The proposed stochastic SMACOF algorithm will be related to an averaged algorithm with updates,

$$\tilde{\mathbf{X}}_{t+1} = (1 - \mu v)\tilde{\mathbf{X}}_t + \mu \mathbf{B}^a(\tilde{\mathbf{X}}_t)\tilde{\mathbf{X}}_t \quad (16)$$

where the time-invariant function $\mathbf{B}^a(\mathbf{X}) := \mathbb{E}[\mathbf{L}_t^\dagger \mathbf{B}_t^\epsilon(\mathbf{X})]$ and $v = \frac{N(p-1)}{p(N-1)}$. Assuming that both algorithms start from the same initialization, i.e., $\hat{\mathbf{X}}_0 = \tilde{\mathbf{X}}_0$, the following proposition states the main result of this section.

Proposition 1. *Under (A1)-(A5), and for $\mu < 1$, it holds for the updates generated by (13) and (16), that*

$$\max_{1 \leq t \leq 1/\mu} \|\hat{\mathbf{X}}_t - \tilde{\mathbf{X}}_t\| \leq c(\mu) \quad (17)$$

where the random variable $c(\mu) \rightarrow 0$ almost surely as $\mu \rightarrow 0$ with probability 1.

Intuitively, Proposition 1 states that the trajectory of the proposed stochastic algorithm in (13) stays close to that of the averaged algorithm in (16). Further, the stochastic "oscillations" of (13) are small if μ is also small. However, choosing too small a value of μ , which is also the step-size in (16), will generally result in a slower convergence rate for any such iterative algorithm. The parameter μ may therefore be seen as controlling the trade-off between the convergence rate and asymptotic accuracy. Further characterization of this trade-off is pursued via numerical tests in Sec. V.

Alternatively, consider the case when T updates of (13) are performed with $\mu = 1/T$. For this case, the bound in (17) becomes

$$\max_{1 \leq t \leq T} \left\| \hat{\mathbf{X}}_t - \tilde{\mathbf{X}}_t \right\| \leq c(1/T) \quad (18)$$

where $c(1/T) \rightarrow 0$ almost surely as $T \rightarrow \infty$. In other words, the stochastic oscillations can be made arbitrarily small if sufficient number of updates can be performed. It is remarked that such results are commonplace in the stochastic approximation literature [24]–[26].

Next, an outline of the proof of Proposition 1 is presented, while the details are deferred to Appendix B. The overall structure of the proof is similar to that in [24]. Significant differences exist in the details however, since workarounds are introduced in order to avoid making any assumptions on the boundedness of $\hat{\mathbf{X}}_t$. It is emphasized that such a modification is generally not possible in a vast majority of problems, and is not trivial. It is however possible here due to the special structure of the update (13) that depends only on the differences between pairs of rows of $\hat{\mathbf{X}}_t$; see (49).

Proof of Proposition 1: The difference between the iterates generated by (13) and (16) is given by

$$\begin{aligned} \Delta_{t+1} := \hat{\mathbf{X}}_{t+1} - \tilde{\mathbf{X}}_{t+1} &= \Delta_t - \mu \left(\mathbf{L}_t^\dagger \mathbf{L}_t \hat{\mathbf{X}}_t - v \tilde{\mathbf{X}}_t \right) \\ &\quad + \mu \left(\mathbf{L}_t^\dagger \mathbf{B}_t^c(\hat{\mathbf{X}}_t) \hat{\mathbf{X}}_t - \mathbf{B}^a(\tilde{\mathbf{X}}_t) \tilde{\mathbf{X}}_t \right) \end{aligned} \quad (19)$$

Assuming that both the algorithms start from the same initialization, i.e., $\hat{\mathbf{X}}_0 = \tilde{\mathbf{X}}_0$, it follows

that

$$\begin{aligned}
\Delta_{t+1} &= -\sum_{\tau=0}^t \mu \left(\mathbf{L}_\tau^\dagger \mathbf{L}_\tau \hat{\mathbf{X}}_\tau - v \tilde{\mathbf{X}}_\tau \right) \\
&\quad + \mu \sum_{\tau=0}^t \left(\mathbf{L}_\tau^\dagger \mathbf{B}_\tau^\epsilon(\hat{\mathbf{X}}_\tau) \hat{\mathbf{X}}_\tau - \mathbf{B}^a(\tilde{\mathbf{X}}_\tau) \tilde{\mathbf{X}}_\tau \right) \\
&= -\mu v \sum_{\tau=1}^t \Delta_\tau + \mu (\mathbf{K}_t^1 + \mathbf{K}_t^2 + \mathbf{K}_t^3)
\end{aligned} \tag{20}$$

where for all $t \geq 0$,

$$\mathbf{K}_t^1 = \sum_{\tau=0}^t \left(\mathbf{L}_\tau^\dagger \mathbf{B}_\tau^\epsilon(\hat{\mathbf{X}}_\tau) \hat{\mathbf{X}}_\tau - \mathbb{E}[\mathbf{L}_\tau^\dagger \mathbf{B}_\tau^\epsilon(\hat{\mathbf{X}}_\tau)] \hat{\mathbf{X}}_\tau \right) \tag{21a}$$

$$\mathbf{K}_t^2 = -\sum_{\tau=0}^t (\mathbf{L}_\tau^\dagger \mathbf{L}_\tau - v \mathbf{I}) \hat{\mathbf{X}}_\tau \tag{21b}$$

$$\mathbf{K}_t^3 = \sum_{\tau=1}^t \left(\mathbf{B}^a(\hat{\mathbf{X}}_\tau) \hat{\mathbf{X}}_\tau - \mathbf{B}^a(\tilde{\mathbf{X}}_\tau) \tilde{\mathbf{X}}_\tau \right). \tag{21c}$$

The following intermediate lemma develops bounds on the three terms in (21), and constitutes the key step in the proof.

Lemma 2. *The following bounds hold for $i = 1, 2$*

$$\|\mathbf{K}_t^i\| \leq d_t^i + C_i \mu \sum_{\tau=1}^t \pi_t^i \tag{22}$$

$$\|\mathbf{K}_t^3\| \leq C_3 \sum_{\tau=1}^t \|\Delta_\tau\| \tag{23}$$

where the constants C_1 , C_2 , and C_3 are independent of t , and the constants d_t^1 , π_t^1 , d_t^2 , and π_t^2 are such that

$$\frac{d_t^i}{t} \rightarrow 0 \qquad \frac{\pi_t^i}{t} \rightarrow 0 \tag{24a}$$

for $i = 1, 2$, almost surely as $t \rightarrow \infty$.

The proof of Lemma 2 is provided in Appendix B. The norm of Δ_{t+1} can therefore be bounded by applying triangle inequality on (19) as follows.

$$\|\Delta_{t+1}\| \leq \mu(C_3 + 1) \sum_{\tau=1}^t \|\Delta_\tau\| + f(\mu) \tag{25}$$

where we have used the fact that $\|\mathbf{J}\|_2 = 1$ and

$$f(\mu) := \max_{0 \leq t \leq 1/\mu} \mu(d_t^1 + d_t^2) + \mu^2 \sum_{\tau=1}^t C_1 \pi_t^1 + C_2 \pi_t^2. \quad (26)$$

It is further shown in Appendix B that $f(\mu) \rightarrow 0$ almost surely as $\mu \rightarrow 0$. Proposition 1 then follows from the application of the discrete Bellman-Gronwall Lemma [24] on (25), which yields

$$\begin{aligned} \|\Delta_t\| &\leq f(\mu)(1 + \mu(C_3 + 1))^t = f(\mu)e^{t \log(1 + \mu(C_3 + 1))} \\ &\leq f(\mu)e^{\mu t(C_3 + 1)} \leq f(\mu)e^{C_3 + 1} := c(\mu) \end{aligned} \quad (27)$$

■

3) *Convergence of the Averaged Algorithm:* Having established that the trajectory of the stochastic algorithm hovers around that of the averaged algorithm, we complete the proof by establishing that the averaged algorithm converges to a local minimum of (1). The challenge here is that the updates in (16) do not resemble those in other classical algorithms such as SMACOF or gradient descent. For notational brevity, let $\bar{\delta}_{mn} := \mathbb{E}[\delta_{mn}(t)] / \sqrt{\|\mathbf{x}_m - \mathbf{x}_n\|^2 + \epsilon_x}$ and $\mathbf{J} = \mathbf{I} - \mathbf{1}\mathbf{1}^T/N$, and note the following result.

Lemma 3. *Under (A1)-(A5), it holds that*

$$[\mathbf{B}^a(\mathbf{X})]_{mn} = \frac{v}{N} \begin{cases} -\bar{\delta}_{mn} & m \neq n \\ \sum_{k \neq m} \bar{\delta}_{mk} & m = n. \end{cases}$$

The proof of Lemma 3 is provided in the Appendix C. For the rest of the section, we will assume that $\epsilon_x \ll 1$ and thus negligible. Therefore from Lemma 3, we have that

$$\begin{aligned} \bar{\sigma}(\mathbf{X}) &= \sum_{m < n} \mathbb{E}[w_{mn}(t)\delta_{mn}(t)^2] + \text{tr}(\mathbf{X}^T \bar{\mathbf{L}} \mathbf{X}) \\ &\quad - 2\text{tr}(\mathbf{X}^T \bar{\mathbf{B}}(\mathbf{X}) \mathbf{X}) \end{aligned} \quad (28)$$

where, $\bar{\mathbf{L}} := \mathbb{E}[\mathbf{L}_t] = \bar{w}vp\mathbf{J}$ and

$$\bar{\mathbf{B}}(\mathbf{X}) := \frac{\bar{w}vp}{N} \begin{cases} -\bar{\delta}_{mn} & m \neq n \\ \sum_{k \neq m} \bar{\delta}_{mk} & m = n. \end{cases} \quad (29)$$

The main result of this subsection is stated as the following proposition.

Proposition 2. *The mean-stress values $\bar{\sigma}(\tilde{\mathbf{X}}_t)$ decrease monotonically with t and converge to a stationary point of (8).*

Proof: Without loss of generality, let $\sum_{m<n} \mathbb{E}[w_{mn}(t)\delta_{mn}(t)^2] = 1$, and define $\eta^2(\mathbf{X}) := \frac{1}{\mu\nu} \text{tr}(\mathbf{X}^T \bar{\mathbf{L}} \mathbf{X})$ and $\rho(\mathbf{X}) := \frac{1}{2}(1/\mu\nu - 1) \text{tr}(\mathbf{X}^T \bar{\mathbf{L}} \mathbf{X}) + \text{tr}(\mathbf{X}^T \bar{\mathbf{B}}(\mathbf{X}) \mathbf{X})$, and observe that $\bar{\sigma}(\mathbf{X}) = 1 + \eta^2(\mathbf{X}) - 2\rho(\mathbf{X})$. Similarly, define the mapping $\Gamma(\mathbf{X}) := (1 - \mu\nu)\mathbf{X} + \frac{\mu}{\bar{w}p} \bar{\mathbf{B}}(\mathbf{X})\mathbf{X}$, so that the updates in (16) become $\tilde{\mathbf{X}}_{t+1} = \Gamma(\tilde{\mathbf{X}}_t)$.

Given any two embeddings \mathbf{X} and \mathbf{Y} , the following bounds hold from the Cauchy-Schwarz inequality:

$$-\text{tr}(\mathbf{X}^T \bar{\mathbf{L}} \mathbf{X}) \leq -\text{tr}((2\mathbf{X} - \mathbf{Y})^T \bar{\mathbf{L}} \mathbf{Y}) \quad (30)$$

$$-\text{tr}(\mathbf{X}^T \bar{\mathbf{B}}(\mathbf{X}) \mathbf{X}) \leq -\text{tr}(\mathbf{X}^T \bar{\mathbf{B}}(\mathbf{Y}) \mathbf{Y}) \quad (31)$$

which allows us to conclude that

$$\begin{aligned} \rho(\mathbf{X}) &\geq \frac{1}{\mu\nu} \text{tr}(\mathbf{X}^T \bar{\mathbf{L}} \Gamma(\mathbf{Y})) - \frac{1 - \mu\nu}{2\mu\nu} \text{tr}(\mathbf{Y}^T \bar{\mathbf{L}} \mathbf{Y}) \\ \Rightarrow \bar{\sigma}(\mathbf{X}) &\leq 1 + \eta^2(\mathbf{X}) + \frac{1 - \mu\nu}{\mu\nu} \text{tr}(\mathbf{Y}^T \bar{\mathbf{L}} \mathbf{Y}) \\ &\quad - \frac{2}{\mu\nu} \text{tr}(\mathbf{X}^T \bar{\mathbf{L}} \Gamma(\mathbf{Y})) \\ &= 1 + (1 - \mu\nu)\eta^2(\mathbf{Y}) - \eta^2(\Gamma(\mathbf{Y})) + \eta^2(\mathbf{X} - \Gamma(\mathbf{Y})) \end{aligned} \quad (32)$$

where equalities holds for $\mathbf{X} = \mathbf{Y}$. Denote the right-hand side of (33) by $\omega_{\mathbf{Y}}(\mathbf{X})$, and observe that $\omega_{\mathbf{Y}}(\mathbf{X}) \geq \omega_{\mathbf{Y}}(\Gamma(\mathbf{Y}))$ for all \mathbf{X} . This yields the main inequality that $\bar{\sigma}(\mathbf{Y}) = \omega_{\mathbf{Y}}(\mathbf{Y}) \geq \omega_{\mathbf{Y}}(\Gamma(\mathbf{Y})) \geq \bar{\sigma}(\Gamma(\mathbf{Y}))$. In other words, we have that $\bar{\sigma}(\tilde{\mathbf{X}}_t) \geq \bar{\sigma}(\tilde{\mathbf{X}}_{t+1})$, so that the non-negative sequence $\sigma_t := \bar{\sigma}(\tilde{\mathbf{X}}_t)$ is non-increasing and therefore convergent to a limit, say $\bar{\sigma}_\infty$. By squeeze theorem for limits [27], it also holds that $\omega_{\tilde{\mathbf{X}}_t}(\tilde{\mathbf{X}}_{t+1}) \rightarrow \bar{\sigma}_\infty$, yielding the following limits

$$\lim_{t \rightarrow \infty} \eta^2(\mathbf{X}_t) = (1 - \bar{\sigma}_\infty)/\mu\nu \quad (34)$$

$$\lim_{t \rightarrow \infty} \rho(\mathbf{X}_t) = (1 - \bar{\sigma}_\infty)(1 - \mu\nu)/2 \quad (35)$$

$$\lim_{t \rightarrow \infty} \eta^2(\tilde{\mathbf{X}}_t - \tilde{\mathbf{X}}_{t+1}) = 0 \quad (36)$$

Since the matrices $\{\tilde{\mathbf{X}}_t\}_{t \geq 0}$ are origin centered, the result in (36) can equivalently be written as $\|\tilde{\mathbf{X}}_t - \tilde{\mathbf{X}}_{t+1}\| \rightarrow 0$ as $t \rightarrow \infty$. Denoting the limit point of $\tilde{\mathbf{X}}_t$ by $\tilde{\mathbf{X}}_\infty$, it can be seen that $\nabla \bar{\sigma}(\tilde{\mathbf{X}}_\infty) = 0$. ■

IV. IMPLEMENTATION ASPECTS

A. Multi-agent network localization

Multidimensional scaling has been widely used for localization, where inter-node distances are often obtained from time-of-arrival or received signal strength measurements [10], [14], [28], [29]. Wireless network localization is challenging because the pairwise distance measurements are noisy, time-varying due to mobility, fading, and synchronization errors, and often partially missing, due to the limited range of the sensors. Further, the limited battery life and resource constraints at the nodes impose restrictions on the communication and computational load that the network can tolerate [13], [14].

Towards addressing these limitations, the stochastic SMACOF algorithm for network localization works by judiciously choosing $\{w_{mn}(t)\}$ to limit the communication and computational cost at each update. The idea is to partition the network into several non-overlapping clusters (or components), chosen randomly at each time t . The coordinates within a cluster are updated as in (13). Only neighboring nodes are included within each cluster, thus eliminating the need for multihop communication between far off nodes. Finally, since the updates at different components are independent of each other, the localization algorithm is run asynchronously as follows.

- S1. At a given time t , a node j randomly declares itself as a cluster head, and solicits cluster members from among its neighbors $n \in \mathcal{N}_j$. Available neighbors respond with their current location estimates $\hat{\mathbf{x}}_n(t)$, resulting in a star shaped cluster \mathcal{C}_t^j . Once locked as cluster members, these nodes respond only to the messages from node j .
- S2. The cluster head performs distance measurements between itself and all its neighbors and collects $\delta_{jn}(t)$ for all $n \in \mathcal{C}_t^j \setminus \{j\}$.
- S3. The cluster head performs the update in (13) with appropriately chosen weights $\{w_{jn}(t)\}$, and broadcasts the new location estimates to each node in $\mathcal{C}_t^j \setminus \{j\}$.
- S4. Nodes in $\mathcal{C}_t^j \setminus \{j\}$, upon receiving the new location estimates (or upon timeout or error events), release their locks and become available.

As originally intended, the proposed algorithm can also be applied to mobile networks. The algorithm is expected to perform well as long as the node velocities are not too high. The asynchronous nature of the algorithm allows for delayed updates at nodes, balanced battery usage within the network, and communication errors. In general, it is also possible to apply multiple

updates of the form in (13) per time instant, without incurring any extra communication cost.

Nodes may declare themselves as cluster heads using a random backoff-based contention mechanism such as CSMA, and solicit neighbors by simply sending an RTS packet. An update at a cluster thus takes up at most two message exchanges. More complicated protocols that ensure recovery from collisions, and robustness or errors can also be used [30]. The online algorithm is flexible, and allows clusters of any shape or size, depending on the communication and computational resources available within the network. The non-zero weights, corresponding to available distance measurements, can be chosen according to the estimated noise variance [10], [28], following Sammon mapping [1], [31], or simply as unity.

It is remarked that the node coordinates obtained from (S1)-(S4) are relative and centered at the origin. In applications where node coordinates are required with respect to a set of GPS-enabled anchor nodes, appropriate rotation and translation operations must be applied at each node. Since the anchor nodes are generally not power constrained, it is possible for them to determine these transformations [10], and convey the result to all other nodes. As shown later in Sec. V, it is generally sufficient to calculate the transformations periodically every few time slots.

Finally, similar to the SMACOF algorithm, the stochastic SMACOF algorithm is sensitive to initialization. A random initialization may result in the algorithm getting trapped in a “poor” local minimum. In practice, superior location estimation performance is obtained if the initialization is at least roughly correct. Simple low-complexity localization algorithms can be used for initialization. For instance, nodes can roughly triangulate themselves using noisy distance estimates from the anchor nodes [32].

B. Large network visualization

It is possible to visualize N objects in a 2 or 3 dimensional euclidean space by applying MDS to the pairwise dissimilarities $\{\delta_{mn}\}$. The SMACOF algorithm is however ill-suited for large-scale visualization since it requires at least $\mathcal{O}(N^2)$ operations per iteration. Further, even processing the full measurements $\{\delta_{mn}\}$ simultaneously may not be feasible for datasets with more than a hundred thousand objects.

Visualization via stochastic embedding can be achieved by partitioning the objects into several subsets of reasonable sizes, and performing the updates in (13). The following steps are performed

for each $t \geq 1$.

- 1) Partition the N objects into random, mutually exclusive subsets \mathcal{C}_t^j with p nodes per subset.
- 2) For each subset, randomly choose a small fraction f_t of pairs and measure (calculate or fetch from memory) distances δ_{mn} for the chosen pairs. Let \mathcal{F}_t^j denote the set of chosen pairs for each cluster j and time t .
- 3) Apply the update in (13) for each subset \mathcal{C}_t^j .

Compared to the localization algorithm, in this case all pairwise distances are available a priori and without noise, but cannot be read or processed simultaneously. The aforementioned steps result in making $\{w_{mn}(t)\}$ sparse and thus reducing the per-iteration complexity. Algorithm 1 summarizes the implementation of stochastic SMACOF for large network visualization.

Algorithm 1 Stochastic SMACOF for Large Network Visualization

- 1: Initialize \mathbf{X}_0 and set μ to some value in $(0, 1)$
 - 2: **for** $t = 1, 2, \dots$ **do**
 - 3: Partition the set \mathcal{N} into C disjoint subsets $\{\mathcal{C}_t^j\}_{j=1}^C$
 - 4: **for** $j = 1, \dots, C$ **do**
 - 5: Measure or fetch from memory pairwise distances $\{\delta_{mn}(t)\}$, for a subset of object pairs $(m, n) \in \mathcal{F}_t^j$.
 - 6: Set weights $w_{mn}(t) = 1$ for all $(m, n) \in \mathcal{F}_t^j$.
 - 7: Perform the update in (13) for each subset \mathcal{C}_t^j .
 - 8: **end for**
 - 9: **end for**
-

Again, as envisioned earlier, the algorithm is also applicable to visualization of dynamic networks. The idea here is to create an animation consisting of embeddings that vary over time. By specifying a small enough value for μ in (13), it is possible to force the embeddings to change slowly over time, thus preserving the user's *mental map* [12]. Unlike existing algorithms however, the proposed algorithm can allow visualization of very large datasets.

C. Algorithm complexity

Unlike the SMACOF algorithm, whose per-iteration complexity is $\mathcal{O}(N^2)$, the stochastic SMACOF algorithm processes the data in small batches and can therefore be implemented at

near-linear complexity. This is because if \mathcal{G}_t consists of multiple components of size p each, the updates in (13) decouple and can even be carried out in parallel. Further, the weights for each cluster are chosen to be sparse, i.e., the $p \times p$ matrix $\mathbf{L}_t^{(j)}$ has at most $q \ll p^2$ non-zero elements. The problem of solving a sparse Laplacian system of equations has been well studied, and state-of-the-art solvers return a solution in time $\mathcal{O}(q \log p)$ for each component. Thus, using N/p sparse matrices $\{\mathbf{L}_t^j\}$ results in an overall complexity of $\mathcal{O}\left(\frac{Nq}{p} \log p\right)$. As we will show next, the appropriate choice of the batch size p results in a near-linear complexity. The complexity results obtained in this section are summarized in Table I.

Note that a sublinear per-iteration complexity of $\mathcal{O}(q \log(p))$ is also achievable by updating only one component per iteration. Such an implementation would however require proportionally large number of iterations. Alternatively, the per-iteration complexity of the algorithm can be calibrated using the total number of dissimilarity measurements processed per-iteration, given by $f(N) = q(N/p)$. To this end, we provide approximate rules for choosing p and q so as to minimize the per-iteration complexity, given the total number of non-zero weights $f(N)$.

First, assume that each \mathbf{L}_t^j is sparse, i.e., $q \ll p^2$, so that $f(N)/N = q/p \ll p$. In this case, since the per-iteration complexity is given by $\mathcal{O}(f(N) \log(p))$, the value of $\log(p)$ should be as small as possible. It can be seen that the choice

$$p \sim \mathcal{O}\left(\left(\frac{f(N)}{N}\right)^\beta\right) \quad q \sim \mathcal{O}\left(\left(\frac{f(N)}{N}\right)^{\beta+1}\right) \quad (37)$$

for some $\beta \gg 1$ results in the complexity $\mathcal{O}(f(N) \log(f(N)/N))$, while ensuring that \mathbf{L}_t^j is still sparse with $q \sim \mathcal{O}(p^{1+1/\beta})$. Note that it is not necessary for β to be very large, as long as the sparse Laplacian solvers can still be used. On the other hand, when \mathbf{L}_t^j is dense so that $q \sim \mathcal{O}(p^2)$, the per-iteration complexity is given by $\mathcal{O}(Nq) = \mathcal{O}(f(N)p)$. In this case, it holds that $f(N)/N = q/p \leq p$, so that one must choose $p \sim \mathcal{O}(f(N)/N)$ and $q \sim \mathcal{O}(f(N)^2/N^2)$. Consequently, the optimal iteration complexity for this case becomes $\mathcal{O}(f(N)^2/N)$.

Table I shows a few example choices of $f(N)$ and the corresponding per-iteration complexity values. It can be observed that when $f(N)$ is almost linear in N , so is the per-iteration complexity, regardless of the sparsity of \mathbf{L}_t^j . On the other hand, using a sparse \mathbf{L}_t^j becomes important when $f(N)$ is large.

Non-zero weights $f(N)$	sparse \mathbf{L}_t^j	dense \mathbf{L}_t^j
$\mathcal{O}(N^{1+\omega}), 0 < \omega \ll 1$	$\mathcal{O}(N^{1+\omega} \log(N))$	$\mathcal{O}(N^{1+2\omega})$
$\mathcal{O}(N \log N)$	$\mathcal{O}(N \log N \log \log N)$	$\mathcal{O}(N \log^2 N)$
$\mathcal{O}(N^{3/2})$	$\mathcal{O}(N^{3/2} \log N)$	$\mathcal{O}(N^2)$

TABLE I: Algorithm complexity for different choices of $f(N)$

V. SIMULATION RESULTS

This section provides simulation results evaluating the performance of the proposed algorithm. The general properties of the stochastic SMACOF algorithm are first characterized using numerical tests. Next, simulation results are provided for the online localization algorithm, evaluating its performance in various mobile network scenarios. Finally, applicability to large-scale visualization is demonstrated by running the algorithm on two different datasets. Before proceeding, it is remarked that the proposed stochastic SMACOF is better suited to applications where the size of the dataset is large, preferably $N > 50$. Indeed, if the problem at hand is small (say $N < 20$), conventional SMACOF would likely be faster, since the proposed algorithm generally requires more iterations to converge. The computational advantage arising from processing only a few distance measurements per time instant becomes significant only when N is sufficiently large.

A. Algorithm Behavior

This section provides several numerical tests that allow us to study various properties of the stochastic SMACOF algorithm. Towards this end, consider a network with 100 nodes, distributed uniformly over a 10×10 planar area. The measured distances between nodes m and n are given by $\delta_{mn}(t) = \|\mathbf{x}_m - \mathbf{x}_n\| + v_{mn}(t)$, where $v_{mn}(t) \sim \mathcal{N}(0, 0.01)$. Negative distance measurements were discarded by setting the corresponding $w_{mn}(t) = 0$. The algorithm is run for different values of μ , with $p = 25$ and about 35% density of non-zeros¹. All non-zero weights are chosen to be unity.

1) *Transient performance*: Fig 1 (Top) shows the sequence of normalized stress values obtained from an example run of the algorithm [cf. (13)]. For comparison, the stress values

¹Non-zero locations are generated randomly, and the number of non-zeros vary between different instantiations.

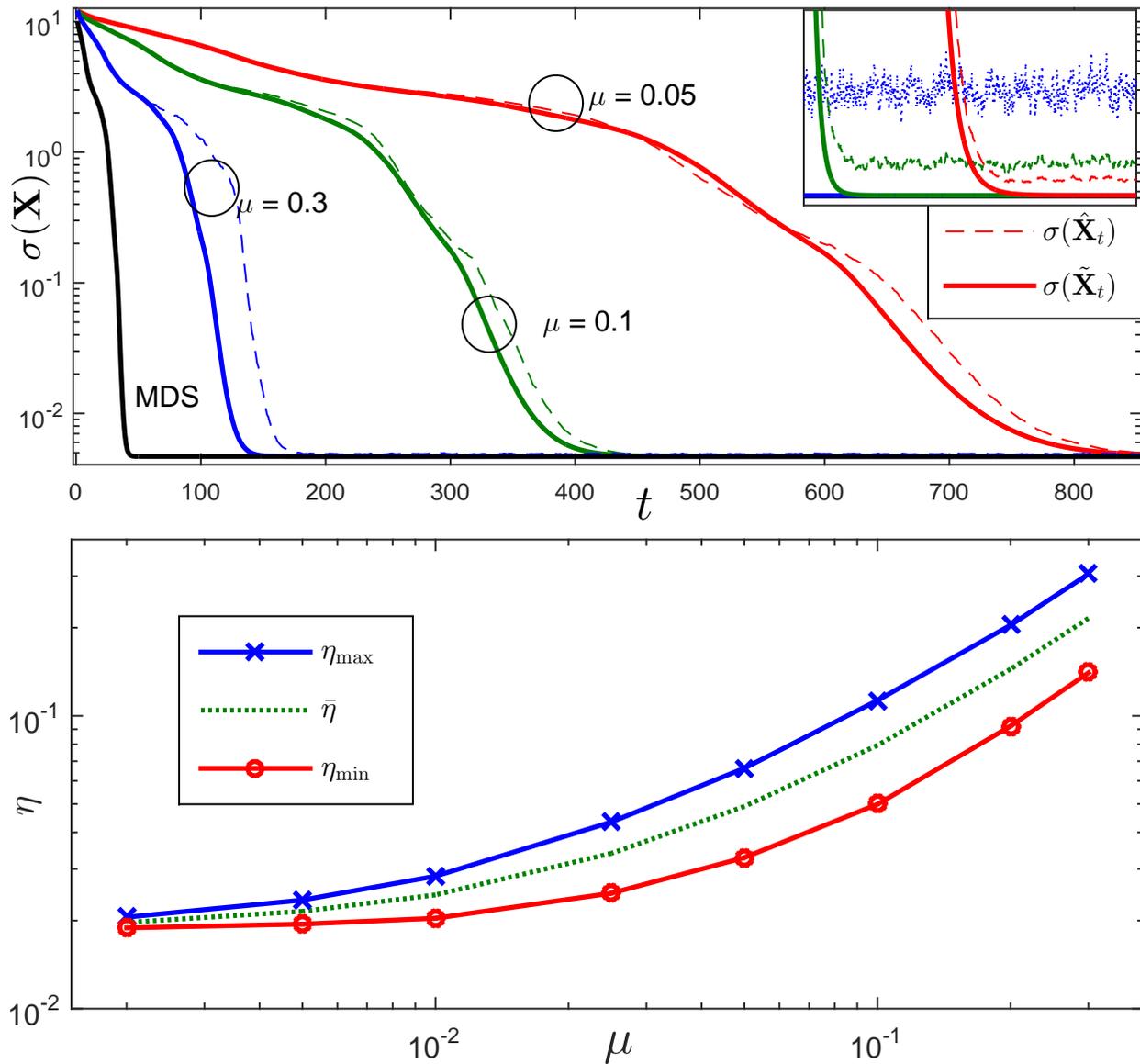


Fig. 1: (Top) Performance of the stochastic SMACOF algorithm, the averaged algorithm, and the SMACOF algorithm; (Bottom) Steady state fluctuations in the stress.

obtained from running the averaged algorithm (cf. (16)) and the SMACOF algorithm for weighted MDS (cf. (7)) are also plotted. All algorithms are initialized with the same randomly chosen configuration. The MDS algorithm runs with all-one weights, while the updates for the averaged algorithm are obtained via empirical averaging.

As expected, the convergence speed of the algorithm varies monotonically with μ . Consistent with Proposition 1, the trajectory of the proposed algorithm follows that of the averaged algorithm. As expected, the steady-state stress value achieved by the averaged algorithm is very close to that of SMACOF. Further, as shown in the inset, the proposed algorithm hovers above the averaged algorithm, with steady-state deviation decreasing with μ .

It is remarked that the SGD algorithm, with updates specified in (9), tended to diverge in the presence of noisy distance measurements, different weight choices, and poor initializations. For instance, when using Sammon mapping, i.e., $w_{mn} = 1/\delta_{mn}$, the noisy measurement model specified earlier, and $\mu = 0.05$, the SGD algorithm converged for only 19 out of 100 test runs. In contrast, no divergent behavior was ever observed for the proposed algorithm even with measurement noise $v_{ij} \sim \mathcal{N}(0, 10)$.

2) *Steady state performance:* The algorithm is allowed to run for 5000 time instants with different values of μ , and the minimum, mean, and maximum steady-state stress values are evaluated. We set $\mathcal{T}_{ss} = [4801, \dots, 5000]$ and evaluate

$$\eta_{\min} = \min_{t \in \mathcal{T}_{ss}} \sigma(\hat{\mathbf{X}}_t) \quad \bar{\eta} = \sum_{t \in \mathcal{T}_{ss}} \frac{\sigma(\hat{\mathbf{X}}_t)}{|\mathcal{T}_{ss}|} \quad \eta_{\max} = \max_{t \in \mathcal{T}_{ss}} \sigma(\hat{\mathbf{X}}_t).$$

Starting with the same initialization, the entire experiment is repeated for 100 Monte Carlo iterations. Fig. 1 (Bottom) shows the minimum, mean, and maximum steady state errors plotted against μ . As expected, the stress values converge to a small non-zero value that decreases with μ .

B. Dynamic Network Localization

The localization performance of the proposed algorithm is studied on a mobile network. Video² shows an example run of the algorithm on a mobile network with $N = 8$ and $\mu = 0.3$. The performance of the algorithm is further analyzed by carrying out simulations over networks with different sizes and node velocities. For a mobile network with N nodes, nodes are deployed randomly with an average density of one node per unit area. Nodes can measure distances and communicate within a radius of $\sqrt{N}/2$. For all values of N , five nodes are randomly chosen to be anchors. The node velocities are initialized randomly and updated according to the following

²<https://www.youtube.com/watch?v=-MQFR3yiv7U>

model $\mathbf{v}_{mn}(t+1) = \alpha\mathbf{v}_{mn}(t) + \sqrt{1-\alpha^2}\mathbf{n}_v(t)$, where $\mathbf{v}_{mn}(0), \mathbf{n}_v(t) \sim \mathcal{N}(0, \sigma_v^2\mathbf{I})$. The mobility parameter σ_v is directly proportional to the average speed of the nodes, and influences the tracking performance of the algorithms used.

The performance of the proposed algorithm is compared with the weighted MDS solution obtained by running the SMACOF algorithm till convergence. The non-zero weights, corresponding to node pairs within the communication radius of each other, are all set to one. Note however that a direct comparison between the SMACOF solution and the proposed algorithm is unfair, since SMACOF is too complex to be directly implemented in a mobile network. Even among cooperative localization techniques that focus on efficient implementation (see e.g. [10], [14], [28], [29]), localization requires several iterations per time instant. In contrast, the proposed algorithm is asynchronous, and incurs linear or sublinear complexity, but is inaccurate for the first few time instants.

In order to perform a fair comparison between algorithms, the following modifications are adopted. First, a time-slotted version of the stochastic SMACOF algorithm is considered. Within each time slot, the network forms several clusters, and performs steps (S1)-(S4). In order to reduce the overhead associated with cluster formation, nodes with fewer than 5 neighbors do not form clusters. Similarly, to limit the computational complexity at each node, cluster heads respond to at most 10 nearest neighbors. With these settings, the computational and communication complexity incurred by the network at every time slot is approximately $N/5$. The computational and communication complexity of the SMACOF variants in [10], [14], [33] is also normalized appropriately. As a first order approximation, it is assumed that these algorithms require $\mathcal{O}(N)$ message exchanges per iteration. Equivalently, if we allow $N/5$ message exchanges per iteration, and assume that 10 iterations are required for convergence, SMACOF requires about 50 time slots for convergence. For obtaining the plots however, SMACOF is run till convergence, and the number of iterations incurred was often more than 50. Both algorithms start with an initial estimate of the node locations. Approximate node estimates can be quickly obtained using simple techniques such as those in [32]. For the purpose of simulations, the initial locations are chosen as $\hat{\mathbf{x}}_m(0) = \mathbf{x}_m(0) + \mathbf{v}_m$ where $v_m \sim \mathcal{N}(0, N/100)$. Warm starts are utilized at subsequent time slots by initializing SMACOF with the previously estimated node locations.

Fig. 2(top) shows an example run of the two algorithms with $\sigma_v = 0.01$, $N = 50$, and $\mu = 0.5$. The best possible estimation error obtained by solving the MDS problem is also

shown for comparison. Observe that the proposed algorithm is inaccurate initially, and gradually approaches its steady state value. Interestingly, the transient period required by the proposed algorithm is small, especially when compared to the 50 time slots required by the SMACOF implementation.

Next, the steady-state localization error of the two algorithms is compared. Both algorithms are run for 700 iterations, and the maximum localization error incurred in the last 200 iterations is evaluated as $e_\ell = \max_{t \in \mathcal{T}_{ss}} \frac{1}{N} \left\| \hat{\mathbf{X}}_t - \mathbf{X}_t \right\|$ where $\mathcal{T}_{ss} = [501, \dots, 700]$. The entire process is repeated for 100 Monte-Carlo repetitions. For the proposed algorithm, the value of μ is tuned a priori to minimize the localization error. Fig. 2(bottom) shows the steady-state localization error incurred by the online and SMACOF algorithms, plotted for different values of N and σ_v . It is evident that the proposed algorithm performs significantly better than the complexity-normalized SMACOF. In particular, while the performance of the two algorithms deteriorates with increasing node mobility, the gap between their performance also increases. This is because at higher node speeds, the node locations change significantly within the 50 time slots required by SMACOF to run. Observe that for a given average node velocity, the performance of all algorithms appears to improve with N . However, this is simply because the average node distances increase with N , thereby reducing the relative average node speeds.

Fig. 3 shows an example run of the algorithm on a mobile network with $N = 8$ and $\mu = 0.3$. The network has four static anchors placed at the four corners of the 1×1 region, that provide the necessary translation and rotation information to all other nodes. For simplicity, only one 8-node cluster is formed at each time instant by a randomly selected node. The actual and estimated node locations are shown as circles and squares respectively, with markers drawn every 10 time instants. The nodes move in the direction indicated by decreasing marker sizes. As evident from the figure, the trajectory of the estimated node locations converges to the actual trajectory within 30-40 time instants, and follows it thereafter.

C. Large-scale Visualization

This section demonstrates the use of the stochastic SMACOF algorithm for large-scale visualization. Given the plethora of highly sophisticated visualization algorithms a full-fledged comparison is beyond the scope of the present work. Instead, we only present the visualizations obtained from running the proposed algorithm for both static and dynamic datasets. The proposed

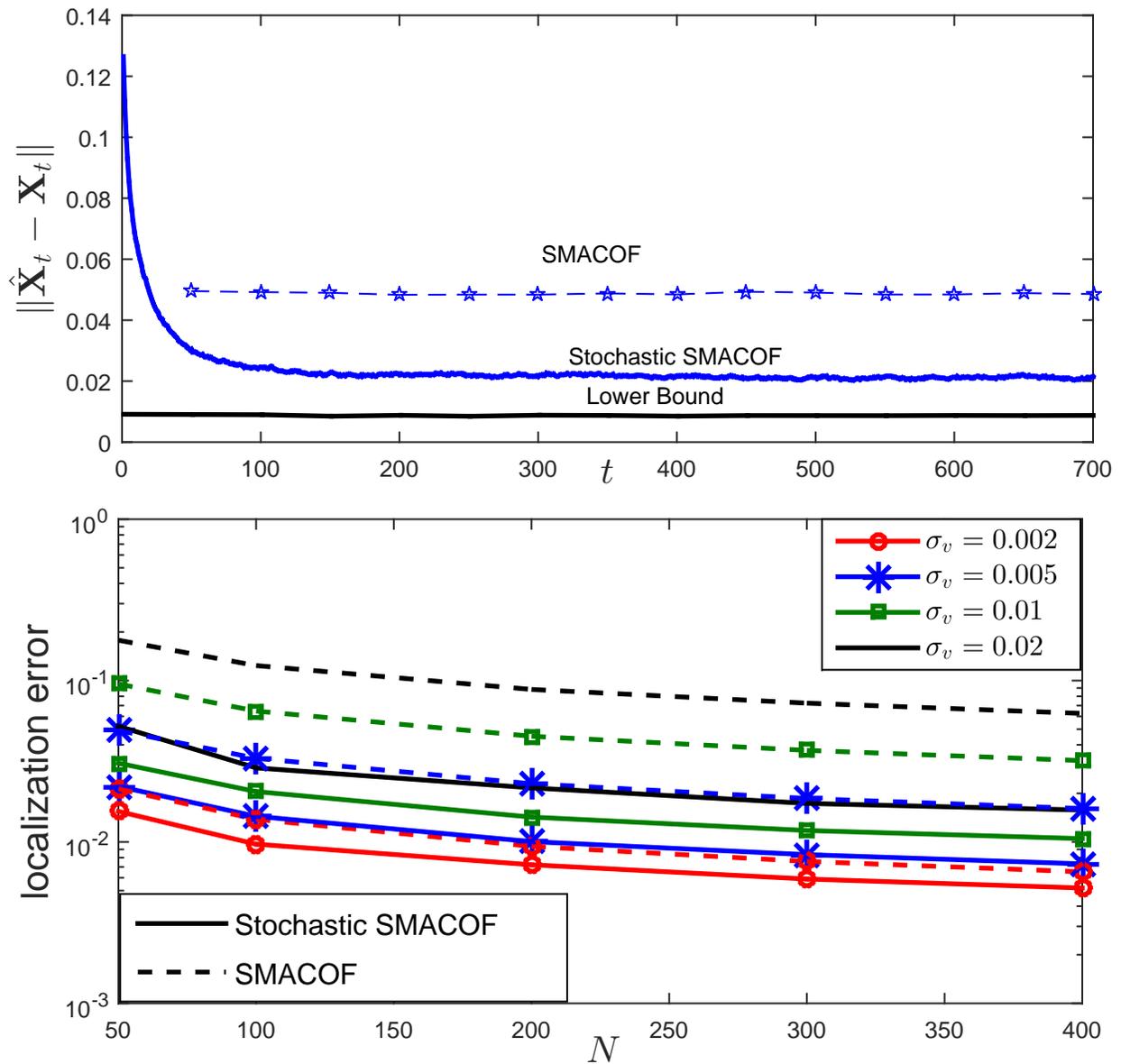


Fig. 2: (Top) Estimation error for an example run of the Stochastic SMACOF and SMACOF algorithms; (Bottom) Localization error for different network sizes and average node velocities.

algorithms are implemented in MATLAB and run on an Intel Core i7 CPU. This is in contrast to the state-of-the-art visualization algorithms that require large compute clusters with hundreds of processors for similar-sized datasets [17].

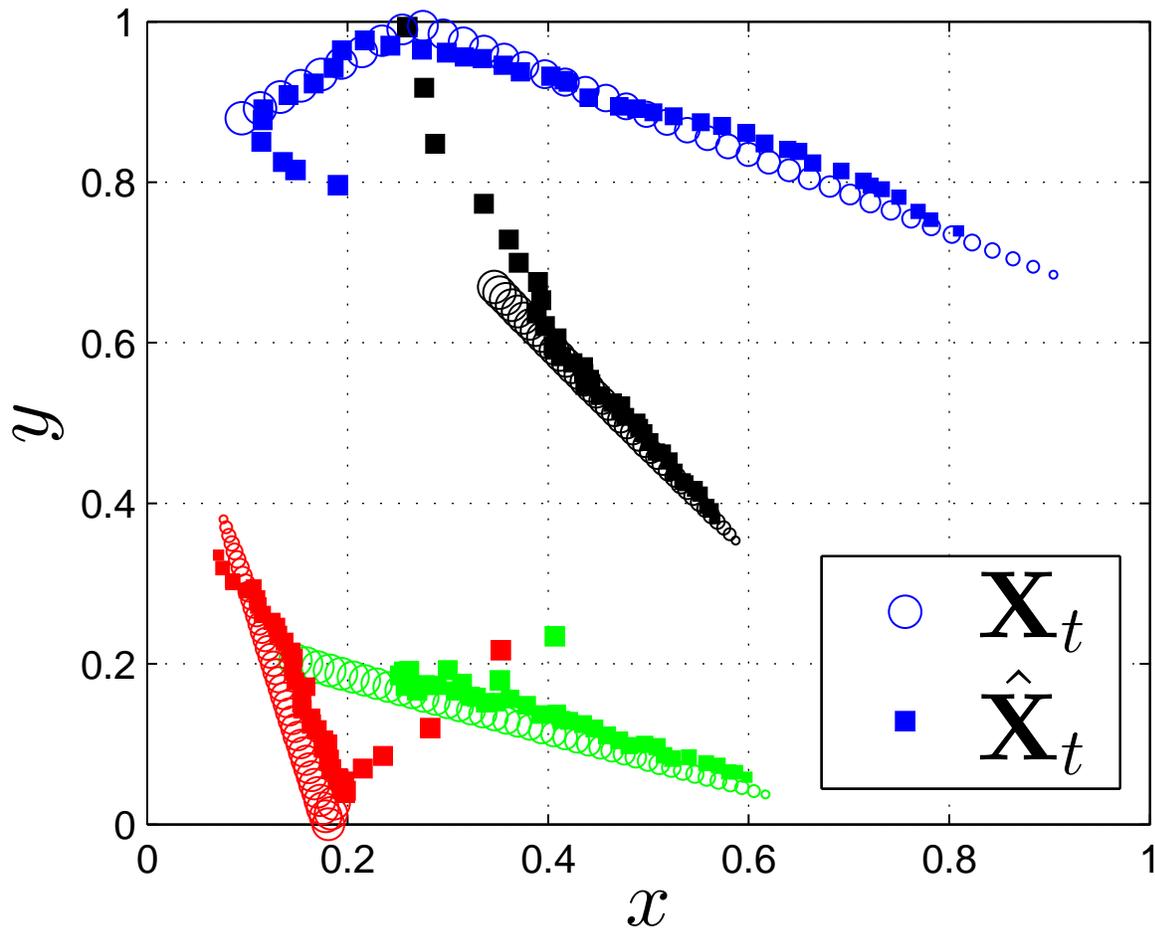


Fig. 3: Example run of the dynamic network localization algorithm. Marker size decreases with time to indicate the direction of motion.

1) *PubChem Dataset*: We consider a subset of 800,000 unique chemical compounds taken from the pubchem compound database [34], [35]. The structural information about each compound is represented by its 166 bit MACCS fingerprint. Dissimilarities between two compounds is calculated using the Tanimoto score. Dissimilarities between two compounds with binary fingerprints \mathbf{h} and \mathbf{g} is calculated using the Tanimoto score [36, Ch-8], given by

$$\gamma = 1 - \frac{\sum_i h_i \cap g_i}{\sum_i h_i \cup g_i} \quad (38)$$

where \cap and \cup denote the logical AND and OR operators respectively. It is remarked that for this case, it is no longer possible to load an $N \times N$ matrix in the memory. Following the

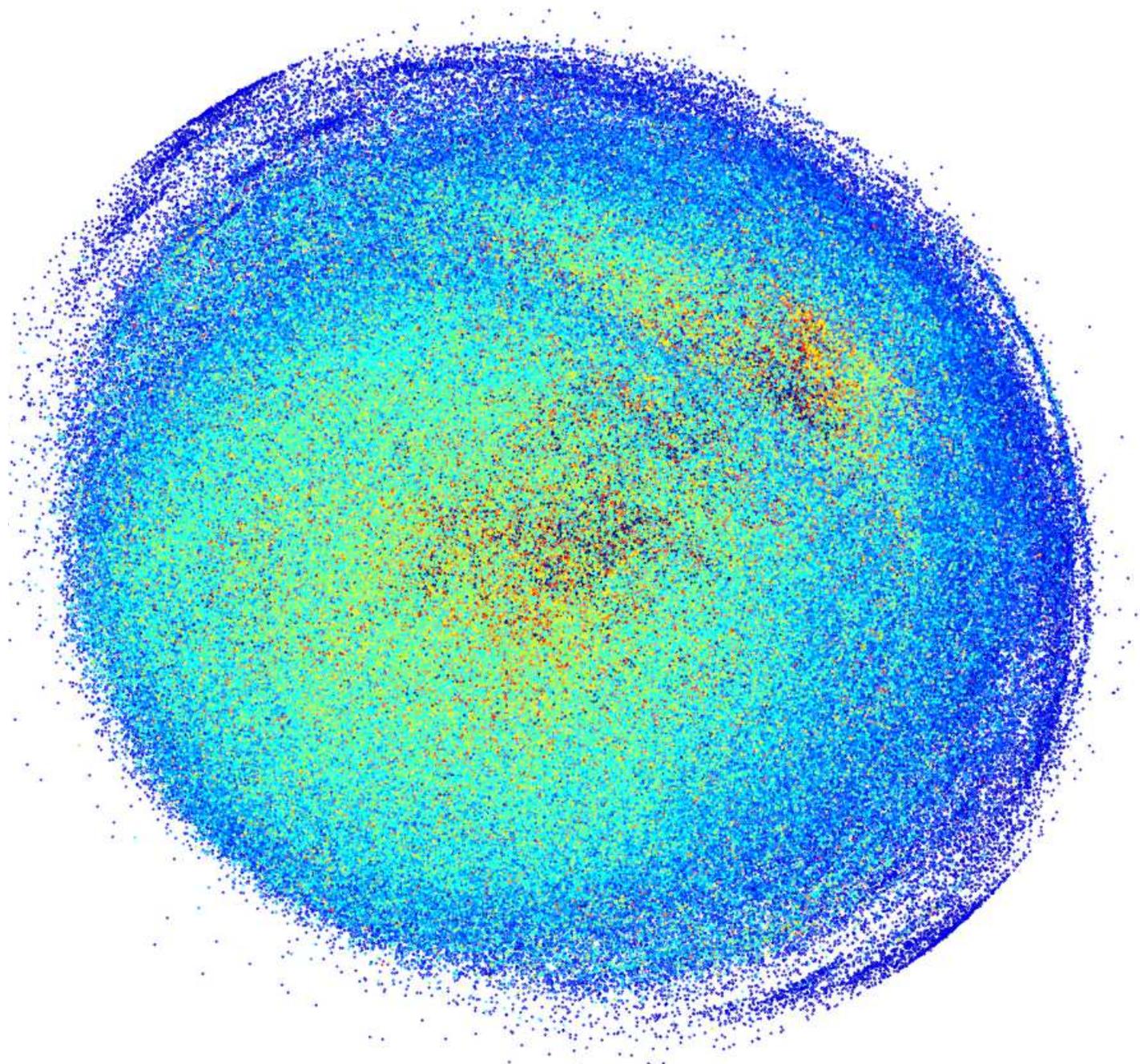


Fig. 4: Visualization of PubChem Datasets.

discussion in Sec. IV-B, we use $p = 100$ and $q = 50$, so as to obtain linear complexity per iteration. The simulation is run for 5000 iterations, and the value of μ is reduced every 1000 iterations from 0.2 to 0.001. Figure. 4 shows the visualization obtained from the stochastic

SMACOF algorithm. Each dot represents a compound, and is colored according to its *molecular complexity*, a measure available from the PubChem dataset. Specifically, the blue dots represent simpler (lower complexity) molecules, while green, yellow, and red colored dots represent progressively more complex molecules. It is observed that MDS yields two distinct clusters of compounds, while the lower complexity compounds are scattered towards the edges. The visualization obtained here is comparable to those obtained in [8], [17].

2) *MovieLens Dataset*: The proposed algorithm is used to perform dynamic visualization of the 27,000 movies on the MovieLens database [37]. To this end, the time-stamp associated with each movie rating is utilized to generate a dynamic network \mathcal{G}_t that only contains the movies released and rated till the week number t . The distance between two movies is estimated from their cosine similarities. Video shows a visualization of the evolution of the movie-space over the duration 1995-2015. Each movie is colored in accordance with its popularity, and the newly released movies start at the origin. From the video, it can be seen that the popular movies move quickly (within few weeks) towards the edge of the graph, while the less popular ones tend to remain near the center. See the video at the link ³.

3) *Newcomb Fraternity's Dataset*: The dynamic visualization of the Newcomb Fraternity dataset [38] is considered. Since the dataset consists of only 16 nodes, and yields only 14 snapshots overall, computational complexity is not an issue. Nevertheless, the dynamic visualization is obtained so that it may be compared with the regularized MDS technique of [12]. Video⁴ shows the dynamic visualization obtained from running the stochastic SMACOF algorithm for 50 iterations per time slot with $\mu = 0.2$. The video is generated following the procedure similar to that in [12]. The resulting video is quite similar to the one obtained via the graph-regularized framework of [12]. Intuitively, the momentum term in the updates in (13) plays the role of the regularization term here, and keeps the embeddings from changing too quickly.

VI. CONCLUSION

The multidimensional scaling (MDS) problem is considered within a stochastic setting, and a novel stochastic scaling by majorizing a complicated function (SMACOF) is proposed. The

³<https://www.youtube.com/watch?v=iJbY3HPHAUM>

⁴<https://www.youtube.com/watch?v=G9geUI3U7Tw&feature=youtu.be>

proposed algorithm is highly scalable, and is applicable to visualization and localization problems of very large sizes. Asymptotic analysis of the stochastic SMACOF algorithm shows that it stays close to the trajectory of an averaged algorithm, which itself converges to a stationary point of the stochastic stress minimization problem. Implementation details, as well as the computational complexity analysis of the proposed algorithms are also provided. The performance of the proposed algorithm is discussed for large-scale localization and visualization examples. The efficacy of the proposed algorithm is demonstrated for localization of mobile networks, and visualization of both, static and dynamic networks.

APPENDIX A

LOWER BOUND ON THE ALGEBRAIC CONNECTIVITY

In order to obtain intuition on **(A3)**, consider the undirected graph \mathcal{G}_t whose edges have weights $\{w_{mn}(t)\}$, and recall that \mathbf{L}_t is the graph Laplacian of \mathcal{G}_t . The eigenvalues of \mathbf{L}_t constitute the spectrum of the graph \mathcal{G}_t [39]. If \mathcal{G}_t is connected, the smallest eigenvalue of \mathbf{L}_t is zero, while the second-smallest eigenvalue $a(\mathcal{G}_t) = 1 / \left\| \mathbf{L}_t^\dagger \right\|_2$ is always non-zero and is referred to as the algebraic connectivity of \mathcal{G}_t . As the name suggests, $a(\mathcal{G})$ captures the overall connectivity of the graph. On the other hand, if \mathcal{G}_t has $K \geq 2$ connected components $\{\mathcal{G}_t^k\}_{k=1}^K$, the K smallest eigenvalues of \mathbf{L}_t are zero, so the smallest non-zero eigenvalue is simply $a(\mathcal{G}_t) = \min_k a(\mathcal{G}_t^k)$. Next, we establish a lower bound on the algebraic connectivity of the weighted graph \mathcal{G}_t .

Proof of Lemma 1: If \mathcal{G}_t is connected, the second smallest eigenvalue is given by

$$a(\mathcal{G}_t) = N \min_{\mathbf{1}^T \mathbf{y} = 0, \mathbf{y} \neq \mathbf{0}} \frac{\sum_{m < n} w_{mn} (y_m - y_n)^2}{\sum_{m < n} (y_m - y_n)^2}. \quad (39)$$

Here, the minimum is attained by the corresponding eigenvector $\check{\mathbf{y}}$, that satisfies $\mathbf{L}_t \check{\mathbf{y}} = a(\mathcal{G}_t) \check{\mathbf{y}}$. Recall that $\mathcal{E} := \{(m, n) \mid w_{mn} \in [\epsilon_w, 1]\}$, and observe that since \mathcal{G}_t is connected, there exists a path \mathcal{P} between any two nodes m and n , such that

$$(\check{y}_m - \check{y}_n)^2 = \left[\sum_{(i,j) \in \mathcal{P}} \check{y}_i - \check{y}_j \right]^2 \leq (N-1) \sum_{(i,j) \in \mathcal{P}} (\check{y}_i - \check{y}_j)^2 \quad (40)$$

$$\leq (N-1) \sum_{(i,j) \in \mathcal{E}} (\check{y}_i - \check{y}_j)^2 \quad (41)$$

where, (40) holds since \mathcal{P} may contain at most $N-1$ edges. Summing both sides over all edges in the graph, we have that

$$\sum_{m < n} (\check{y}_m - \check{y}_n)^2 \leq \frac{N(N-1)^2}{2} \sum_{(m,n) \in \mathcal{E}} (\check{y}_m - \check{y}_n)^2 \quad (42)$$

Substituting (42) into (39) for $\mathbf{y} = \check{\mathbf{y}}$, we have that

$$a(\mathcal{G}) = N \frac{\sum_{m < n} w_{mn} (\check{y}_m - \check{y}_n)^2}{\sum_{m < n} (\check{y}_m - \check{y}_n)^2} \quad (43)$$

$$\geq \frac{2}{(N-1)^2} \frac{\sum_{(m,n) \in \mathcal{E}} w_{mn} (\check{y}_m - \check{y}_n)^2}{\sum_{(m,n) \in \mathcal{E}} (\check{y}_m - \check{y}_n)^2} \geq \frac{2\epsilon_w}{(N-1)^2} \quad (44)$$

which is the required bound. If \mathcal{G}_t is not connected, it holds for a component \mathcal{G}_t^k with p nodes that $a(\mathcal{G}_t^k) \geq 2\epsilon_w/(p-1)^2 \geq 2\epsilon_w/(N-1)^2$, so that we again have $a(\mathcal{G}_t) = \min a(\mathcal{G}_t^k) \geq 2\epsilon_w/(N-1)^2$, which is the desired result. \blacksquare

APPENDIX B

PROOF OF LEMMA 2

Before proceeding with the proof, we state some basic results, and introduced necessary notation. In the subsequent analysis, we will repeatedly use the following inequalities [40]

$$\|\mathbf{AB}\| \leq \|\mathbf{A}\|_2 \|\mathbf{B}\| \leq \|\mathbf{A}\| \|\mathbf{B}\| \quad (45)$$

where \mathbf{A} and \mathbf{B} matrices of compatible sizes. For notational brevity, $d_{mn} := \sqrt{\|\mathbf{x}_m - \mathbf{x}_n\|^2 + \epsilon_x}$ and $\check{d}_{mn} := \sqrt{\|\check{\mathbf{x}}_m - \check{\mathbf{x}}_n\|^2 + \epsilon_x}$, and note that $d_{mn}, \check{d}_{mn} \geq \sqrt{\epsilon}$.

We begin by defining the *total deviation* functions corresponding to \mathbf{K}_t^1 and \mathbf{K}_t^2 as

$$\mathbf{D}_t^1(\mathbf{X}) := \sum_{\tau=1}^t (\mathbf{L}_\tau^\dagger \mathbf{B}_\tau^\epsilon(\mathbf{X}) \mathbf{X} - \mathbb{E}[\mathbf{L}_\tau^\dagger \mathbf{B}_\tau^\epsilon(\mathbf{X}) \mathbf{X}]) \quad (46)$$

$$\mathbf{D}_t^2(\mathbf{X}) := \sum_{\tau=1}^t (\mathbf{L}_\tau^\dagger \mathbf{L}_\tau - \mathbb{E}[\mathbf{L}_\tau^\dagger \mathbf{L}_\tau]) \mathbf{X} \quad (47)$$

The following lemma lists several preliminary results required in deriving the bounds in Lemma 2.

Lemma 4. *There exists $t_0 < \infty$, such that for all $t \geq t_0$, it holds that*

$$\left\| \mathbf{L}_t^\dagger \mathbf{B}_t^\epsilon(\mathbf{X})\mathbf{X} - \mathbf{L}_t^\dagger \mathbf{B}_t^\epsilon(\check{\mathbf{X}})\check{\mathbf{X}} \right\| \leq C_3 \|\mathbf{X} - \check{\mathbf{X}}\| \quad (48a)$$

$$\left\| \mathbf{L}_t^\dagger \mathbf{B}_t^\epsilon(\mathbf{X})\mathbf{X} \right\| \leq C_4 \quad (48b)$$

$$\left\| \mathbf{J}\hat{\mathbf{X}}_t \right\| \leq C_5 \quad (48c)$$

$$\left\| \mathbf{D}_t^1(\mathbf{X}) \right\| \leq d_t^1 \quad (48d)$$

$$\left\| \mathbf{D}_t^1(\mathbf{X}) - \mathbf{D}_t^1(\check{\mathbf{X}}) \right\| \leq \pi_t^1 \|\mathbf{X} - \check{\mathbf{X}}\| \quad (48e)$$

$$\left\| \mathbf{D}_t^2(\mathbf{X}) \right\| \leq d_t^2 \quad (48f)$$

$$\left\| \mathbf{D}_t^2(\mathbf{X}) - \mathbf{D}_t^2(\check{\mathbf{X}}) \right\| \leq \pi_t^2 \|\mathbf{X} - \check{\mathbf{X}}\| \quad (48g)$$

where $\mathbf{J} = \mathbf{I} - \mathbf{1}\mathbf{1}^T/N$, C_3 and C_4 are constants, while the random variables d_t^1 , d_t^2 , π_t^1 , and π_t^2 follow (24). Results in (48f) and (48g) also require \mathbf{X} to be such that $\|\mathbf{J}\mathbf{X}\| \leq C_5$.

The proof organized into four steps, each considering one or more inequalities.

Proof of (48a) and (48b): Observe that the m -th row of $\mathbf{B}_t^\epsilon(\mathbf{X})\mathbf{X}$ for each $t \geq 0$ can be written as

$$[\mathbf{B}_t^\epsilon(\mathbf{X})\mathbf{X}]_{m,:} = \sum_{n \neq m} \frac{w_{mn}(t)\delta_{mn}(t)}{d_{mn}} (\mathbf{x}_m - \mathbf{x}_n)$$

which implies that

$$\left\| [\mathbf{B}_t^\epsilon(\mathbf{X})\mathbf{X}]_{m,:} \right\| \leq \sum_{n \neq m} |w_{mn}(t)\delta_{mn}(t)| \leq NC_\delta. \quad (49)$$

The bound in (48b) therefore follows from the use of (45),

$$\left\| \mathbf{L}_t^\dagger \mathbf{B}_t^\epsilon(\mathbf{X})\mathbf{X} \right\| \leq \left\| \mathbf{L}_t^\dagger \right\|_2 \|\mathbf{B}_t^\epsilon(\mathbf{X})\mathbf{X}\| \leq \frac{N^2 C_\delta}{\epsilon_{\mathbf{L}}}. \quad (50)$$

which yields $C_4 = N^2 C_\delta / \epsilon_{\mathbf{L}}$. Likewise, the m -th row of $\mathbf{B}_t^\epsilon(\mathbf{X})\mathbf{X} - \mathbf{B}_t^\epsilon(\check{\mathbf{X}})\check{\mathbf{X}}$ becomes

$$\begin{aligned} & [\mathbf{B}_t^\epsilon(\mathbf{X})\mathbf{X} - \mathbf{B}_t^\epsilon(\check{\mathbf{X}})\check{\mathbf{X}}]_{m,:} \\ &= \sum_{n \neq m} w_{mn}(t)\delta_{mn}(t) \left(\frac{\mathbf{x}_m - \mathbf{x}_n}{d_{mn}} - \frac{\check{\mathbf{x}}_m - \check{\mathbf{x}}_n}{\check{d}_{mn}} \right). \end{aligned} \quad (51)$$

Adding and subtracting the term $(\check{\mathbf{x}}_m - \check{\mathbf{x}}_n)/d_{mn}$ to each term within the summation in (51), it can be seen that

$$\begin{aligned}
& \frac{\mathbf{x}_m - \mathbf{x}_n}{d_{mn}} - \frac{\check{\mathbf{x}}_m - \check{\mathbf{x}}_n}{\check{d}_{mn}} \\
&= \frac{\mathbf{x}_m - \check{\mathbf{x}}_m}{d_{mn}} - \frac{\mathbf{x}_n - \check{\mathbf{x}}_n}{d_{mn}} + (\check{\mathbf{x}}_m - \check{\mathbf{x}}_n) \left(\frac{1}{d_{mn}} - \frac{1}{\check{d}_{mn}} \right) \\
&= \frac{\mathbf{x}_m - \check{\mathbf{x}}_m}{d_{mn}} - \frac{\mathbf{x}_n - \check{\mathbf{x}}_n}{d_{mn}} + \frac{\check{\mathbf{x}}_m - \check{\mathbf{x}}_n}{\check{d}_{mn}} \frac{\check{d}_{mn}^2 - d_{mn}^2}{d_{mn}(\check{d}_{mn} + d_{mn})}.
\end{aligned} \tag{52}$$

Further, the term $\check{d}_{mn}^2 - d_{mn}^2$ can be written compactly as

$$\begin{aligned}
\check{d}_{mn}^2 - d_{mn}^2 &= \check{\mathbf{x}}_m^T \check{\mathbf{x}}_m + \check{\mathbf{x}}_n^T \check{\mathbf{x}}_n - 2\check{\mathbf{x}}_m^T \check{\mathbf{x}}_n - \mathbf{x}_m^T \mathbf{x}_m - \mathbf{x}_n^T \mathbf{x}_n + 2\mathbf{x}_m^T \mathbf{x}_n \\
&= (\mathbf{x}_m - \mathbf{x}_n + \check{\mathbf{x}}_m - \check{\mathbf{x}}_n)^T (\mathbf{x}_m - \check{\mathbf{x}}_m + \check{\mathbf{x}}_n - \mathbf{x}_n)
\end{aligned} \tag{53}$$

Consequently, it is possible to write (51) as,

$$\begin{aligned}
& [\mathbf{B}_t^\epsilon(\mathbf{X})\mathbf{X} - \mathbf{B}_t^\epsilon(\check{\mathbf{X}})\check{\mathbf{X}}]_{m,:} \\
&= \sum_{n \neq m} w_{mn}(t) \delta_{mn}(t) \mathbf{A}_{mn} ((\mathbf{x}_m - \check{\mathbf{x}}_m) - (\mathbf{x}_n - \check{\mathbf{x}}_n))
\end{aligned}$$

where the matrix \mathbf{A}_{mn} is given by

$$\mathbf{A}_{mn} = \frac{1}{d_{mn}} \mathbf{I} + \frac{(\check{\mathbf{x}}_m - \check{\mathbf{x}}_n)(\mathbf{x}_m - \mathbf{x}_n + \check{\mathbf{x}}_m - \check{\mathbf{x}}_n)^T}{d_{mn}\check{d}_{mn}(\check{d}_{mn} + d_{mn})}. \tag{54}$$

Thus, the full difference becomes

$$\mathbf{B}_t^\epsilon(\mathbf{X})\mathbf{X} - \mathbf{B}_t^\epsilon(\check{\mathbf{X}})\check{\mathbf{X}} = \mathbf{A}_t(\mathbf{X}, \check{\mathbf{X}}) \text{vec}(\mathbf{X} - \check{\mathbf{X}}) \tag{55}$$

where the (m, n) -th $p \times p$ block of $\mathbf{A}_t(\mathbf{X}, \check{\mathbf{X}})$ is given by

$$[\mathbf{A}_t(\mathbf{X}, \check{\mathbf{X}})] := \begin{cases} -\mathbf{A}_{mn} w_{mn}(t) \delta_{mn}(t) & m \neq n \\ \sum_{n \neq m} \mathbf{A}_{mn} w_{mn}(t) \delta_{mn}(t) & m = n \end{cases} \tag{56}$$

Next, repeated use of the triangle inequality yields

$$\begin{aligned}
& \|\mathbf{A}_{mn}\|^2 \\
& \leq \frac{2}{d_{mn}^2} \left(\|\mathbf{I}\|^2 + \frac{\|\check{\mathbf{x}}_m - \check{\mathbf{x}}_n\|^2}{\check{d}_{mn}^2} \frac{\|\check{\mathbf{x}}_m - \check{\mathbf{x}}_n + \mathbf{x}_m - \mathbf{x}_n\|^2}{(\check{d}_{mn} + d_{mn})^2} \right)
\end{aligned}$$

Here, it holds from the definition of \check{d}_{mn} that $\|\check{\mathbf{x}}_m - \check{\mathbf{x}}_n\| / \check{d}_{mn} \leq 1$. Similarly, it holds that

$$\|\check{\mathbf{x}}_m - \check{\mathbf{x}}_n + \mathbf{x}_m - \mathbf{x}_n\|^2 \quad (57)$$

$$\begin{aligned} &\leq \|\check{\mathbf{x}}_m - \check{\mathbf{x}}_n\|^2 + \|\mathbf{x}_m - \mathbf{x}_n\|^2 + 2\|\check{\mathbf{x}}_m - \check{\mathbf{x}}_n\| \|\mathbf{x}_m - \mathbf{x}_n\| \\ &\leq \check{d}_{mn}^2 + d_{mn}^2 + 2\check{d}_{mn}d_{mn} = (\check{d}_{mn} + d_{mn})^2 \end{aligned} \quad (58)$$

Therefore, the bound on $\|\mathbf{A}_{mn}\|^2$ becomes

$$\|\mathbf{A}_{mn}\|^2 \leq \frac{2(N+1)}{\epsilon} \quad (59)$$

Similarly, it holds for $\|\mathbf{A}_t(\mathbf{X}, \check{\mathbf{X}})\|$ that

$$\begin{aligned} \|\mathbf{A}_t(\mathbf{X}, \check{\mathbf{X}})\|^2 &\leq C_\delta^2 \sum_m \sum_{n \neq m} \|\mathbf{A}_{mn}\|^2 + \left(\sum_{n \neq m} \|\mathbf{A}_{mn}\| \right)^2 \\ &\leq 3C_\delta^2 \sum_m \sum_{n \neq m} \|\mathbf{A}_{mn}\|^2 \end{aligned} \quad (60)$$

$$\leq 3C_\delta^2 \frac{N(N-1)(N+1)}{\epsilon} < C_\delta^2 \frac{6N^3}{\epsilon_x} \quad (61)$$

which in turn, yields the bound

$$\|\mathbf{A}_t(\mathbf{X}, \check{\mathbf{X}})\|^2 \leq 6N^3 \frac{C_\delta^2}{\epsilon_x}. \quad (62)$$

The Lipschitz continuity of $\mathbf{L}_t^\dagger \mathbf{B}_t(\mathbf{X})\mathbf{X}$ thus follows as

$$\begin{aligned} \left\| \mathbf{L}_t^\dagger \mathbf{B}_t^\epsilon(\mathbf{X})\mathbf{X} - \mathbf{L}_t^\dagger \mathbf{B}_t^\epsilon(\check{\mathbf{X}})\check{\mathbf{X}} \right\| &\leq \left\| \mathbf{L}_t^\dagger \right\|_2 \left\| \mathbf{B}_t^\epsilon(\mathbf{X})\mathbf{X} - \mathbf{B}_t^\epsilon(\check{\mathbf{X}})\check{\mathbf{X}} \right\| \\ &\leq \frac{NC_\delta}{\epsilon_L} \sqrt{\frac{6N}{\epsilon_x}} \|\mathbf{X} - \check{\mathbf{X}}\|, \end{aligned} \quad (63)$$

so that $C_3 = \frac{NC_\delta}{\epsilon_L} \sqrt{\frac{6N}{\epsilon_x}}$. ■

Proof of (48c): Observe that $\mathbf{L}_t \mathbf{J} = \mathbf{L}_t$ and $\mathbf{J} \mathbf{L}_t^\dagger = \mathbf{L}_t^\dagger$. Right multiplying both sides of (13) by \mathbf{J} , it follows that

$$\mathbf{J} \hat{\mathbf{X}}_{t+1} = \mathbf{J}(\mathbf{I} - \mu \mathbf{L}_t^\dagger \mathbf{L}_t) \hat{\mathbf{X}}_t + \mu \mathbf{J} \mathbf{L}_t^\dagger \mathbf{B}_t^\epsilon(\hat{\mathbf{X}}_t) \hat{\mathbf{X}}_t \quad (64)$$

$$= (\mathbf{J} \mathbf{J} - \mu \mathbf{J} \mathbf{L}_t^\dagger \mathbf{L}_t \mathbf{J}) \hat{\mathbf{X}}_t + \mu \mathbf{L}_t^\dagger \mathbf{B}_t^\epsilon(\hat{\mathbf{X}}_t) \hat{\mathbf{X}}_t \quad (65)$$

$$= (\mathbf{J} - \mu \mathbf{L}_t^\dagger \mathbf{L}_t) \mathbf{J} \hat{\mathbf{X}}_t + \mu \mathbf{L}_t^\dagger \mathbf{B}_t^\epsilon(\hat{\mathbf{X}}_t) \hat{\mathbf{X}}_t \quad (66)$$

$$\begin{aligned} &= (\mathbf{J} - \mu \mathbf{L}_t^\dagger \mathbf{L}_t) (\mathbf{J} - \mu \mathbf{L}_{t-1}^\dagger \mathbf{L}_{t-1}) \mathbf{J} \hat{\mathbf{X}}_{t-1} + \mu \mathbf{L}_t^\dagger \mathbf{B}_t^\epsilon(\hat{\mathbf{X}}_t) \hat{\mathbf{X}}_t \\ &\quad + \mu (\mathbf{J} - \mu \mathbf{L}_t^\dagger \mathbf{L}_t) \mathbf{L}_{t-1}^\dagger \mathbf{B}_{t-1}(\hat{\mathbf{X}}_{t-1}) \hat{\mathbf{X}}_{t-1} \end{aligned} \quad (67)$$

Continuing in a similar manner, taking norm on both sides of (67), applying triangle inequality, and using (48b) yields

$$\left\| \mathbf{J}\hat{\mathbf{X}}_{t+1} \right\| \leq \left\| \mathbf{Q}_t^0 \right\|_2 \left\| \mathbf{J}\hat{\mathbf{X}}_0 \right\| + \mu \left(1 + \sum_{\tau=1}^t \left\| \mathbf{Q}_t^\tau \right\|_2 \right) C_4 \quad (68)$$

where $\mathbf{Q}_t^\tau := \prod_{k=\tau}^t (\mathbf{J} - \mu \mathbf{L}_k^\dagger \mathbf{L}_k)$. Next, from (A4), there exists some $t_0 < \infty$ and $\varrho < 1$ such that $\left\| \mathbf{Q}_t^\tau \right\| \leq \varrho^{t-\tau+1}$ for all $t - \tau + 1 \geq t_0$. Since $\left\| \mathbf{Q}_t^\tau \right\| \leq 1$ for all $t \geq \tau + 1$, bound in (68) becomes

$$\begin{aligned} \left\| \mathbf{J}\hat{\mathbf{X}}_{t+1} \right\| &\leq C_x \varrho^t + \mu C_4 \left(1 + t_0 + \frac{\varrho^t}{1 - \varrho} \right) \\ &= C_x + \mu C_4 \left(1 + t_0 + \frac{1}{1 - \varrho} \right) =: C_5 \end{aligned} \quad (69)$$

for all $t \geq t_0$. ■

Proof of (48d) and (48e): Observe that each term of $\mathbf{D}_t^1(\mathbf{X})$ in (46) is zero mean, and bounded as

$$\left\| (\mathbf{L}_\tau^\dagger \mathbf{B}_\tau^\epsilon(\mathbf{X}) \mathbf{X} - \mathbb{E}[\mathbf{L}_\tau^\dagger \mathbf{B}_\tau^\epsilon(\mathbf{X}) \mathbf{X}]) \right\| \quad (70)$$

$$\leq \left\| (\mathbf{L}_\tau^\dagger \mathbf{B}_\tau^\epsilon(\mathbf{X}) \mathbf{X}) \right\| + \left\| \mathbb{E}[\mathbf{L}_\tau^\dagger \mathbf{B}_\tau^\epsilon(\mathbf{X}) \mathbf{X}] \right\| \quad (71)$$

$$\leq \left\| (\mathbf{L}_\tau^\dagger \mathbf{B}_\tau^\epsilon(\mathbf{X}) \mathbf{X}) \right\| + \mathbb{E}[\left\| \mathbf{L}_\tau^\dagger \mathbf{B}_\tau^\epsilon(\mathbf{X}) \mathbf{X} \right\|] \leq 2C_4 \quad (72)$$

The law of large numbers therefore implies that $\mathbf{D}_t^1(\mathbf{X})/t \rightarrow 0$ almost surely as $t \rightarrow \infty$. This also implies that there exists d_t^1 such that $\left\| \mathbf{D}_t^1(\mathbf{X}) \right\| \leq d_t^1$ and $d_t^1/t \rightarrow 0$ as $t \rightarrow \infty$.

The Lipschitz continuity of $\mathbf{D}_t^1(\mathbf{X})$ can similarly be shown using (48a). Towards this end, observe that

$$\begin{aligned} \mathbf{D}_t^1(\mathbf{X}) - \mathbf{D}_t^1(\check{\mathbf{X}}) &= \sum_{\tau=1}^{t-1} (\mathbf{L}_\tau^\dagger \mathbf{B}_\tau^\epsilon(\mathbf{X}) \mathbf{X} - \mathbb{E}[\mathbf{L}_\tau^\dagger \mathbf{B}_\tau^\epsilon(\mathbf{X}) \mathbf{X}]) \\ &\quad - \sum_{\tau=1}^{t-1} (\mathbf{L}_\tau^\dagger \mathbf{B}_\tau^\epsilon(\check{\mathbf{X}}) \check{\mathbf{X}} - \mathbb{E}[\mathbf{L}_\tau^\dagger \mathbf{B}_\tau^\epsilon(\check{\mathbf{X}}) \check{\mathbf{X}}]) \\ &= \sum_{\tau=0}^{t-1} \mathbf{L}_\tau^\dagger (\mathbf{B}_\tau^\epsilon(\mathbf{X}) \mathbf{X} - \mathbf{B}_\tau^\epsilon(\check{\mathbf{X}}) \check{\mathbf{X}}) \\ &\quad - \mathbb{E}[\mathbf{L}_\tau^\dagger (\mathbf{B}_\tau^\epsilon(\mathbf{X}) \mathbf{X} - \mathbf{B}_\tau^\epsilon(\check{\mathbf{X}}) \check{\mathbf{X}})] \end{aligned} \quad (73)$$

The vectorized version of the first term can be written as

$$\begin{aligned} & \text{vec} (\mathbf{L}_\tau^\dagger \mathbf{B}_\tau^\epsilon(\mathbf{X})\mathbf{X} - \mathbf{L}_\tau^\dagger \mathbf{B}_\tau^\epsilon(\check{\mathbf{X}})\check{\mathbf{X}}) \\ &= (\mathbf{I} \otimes \mathbf{L}_\tau^\dagger) \text{vec} (\mathbf{B}_\tau^\epsilon(\mathbf{X})\mathbf{X} - \mathbf{B}_\tau^\epsilon(\check{\mathbf{X}})\check{\mathbf{X}}) \end{aligned} \quad (74)$$

$$= (\mathbf{I} \otimes \mathbf{L}_\tau^\dagger) \mathbf{A}_\tau(\mathbf{X}, \check{\mathbf{X}}) \text{vec} (\mathbf{X} - \check{\mathbf{X}}) \quad (75)$$

Using a similar transformation on the second term of (73), the vectorized version of the right-hand side can be written as

$$\begin{aligned} & \text{vec} (\mathbf{D}_t^1(\mathbf{X}) - \mathbf{D}_t^1(\check{\mathbf{X}})) \\ &= \left(\sum_{\tau=0}^{t-1} \mathbf{C}_\tau(\mathbf{X}, \check{\mathbf{X}}) - \mathbb{E}[\mathbf{C}_\tau(\mathbf{X}, \check{\mathbf{X}})] \right) \text{vec} (\mathbf{X} - \check{\mathbf{X}}) \end{aligned} \quad (76)$$

where $\mathbf{C}_\tau(\mathbf{X}, \check{\mathbf{X}}) = (\mathbf{I} \otimes \mathbf{L}_\tau^\dagger) \mathbf{A}_\tau(\mathbf{X}, \check{\mathbf{X}})$ is bounded as $\|\mathbf{C}_\tau(\mathbf{X}, \check{\mathbf{X}})\| \leq \|\mathbf{L}_\tau^\dagger\|_2 \|\mathbf{A}_\tau(\mathbf{X}, \check{\mathbf{X}})\| \leq C_3$. It is therefore possible to write

$$\|\mathbf{D}_t^1(\mathbf{X}) - \mathbf{D}_t^1(\check{\mathbf{X}})\| \leq \pi_t \|\mathbf{X} - \check{\mathbf{X}}\| \quad (77)$$

$$\text{where, } \pi_t = \left\| \sum_{\tau=0}^{t-1} \mathbf{C}_\tau(\mathbf{X}, \check{\mathbf{X}}) - \mathbb{E}[\mathbf{C}_\tau(\mathbf{X}, \check{\mathbf{X}})] \right\| \quad (78)$$

Since the term within the norm is a bounded zero-mean random variable, it follows from law of large numbers that

$$\frac{1}{t} \sum_{\tau=0}^{t-1} \mathbf{C}_\tau(\mathbf{X}, \check{\mathbf{X}}) - \mathbb{E}[\mathbf{C}_\tau(\mathbf{X}, \check{\mathbf{X}})] \rightarrow \mathbf{0} \quad (79)$$

with probability 1 as $t \rightarrow \infty$. This also implies that $\pi_t/t \rightarrow 0$ almost surely as $t \rightarrow \infty$.

1) *Proof of (48f) and (48g):* Observe that the zero mean random variable $\mathbf{D}_t^2(\mathbf{X})$ can be written as

$$\mathbf{D}_t^2(\mathbf{X}) = \sum_{\tau=0}^t (\mathbf{L}_\tau^\dagger \mathbf{L}_\tau - \mathbb{E}[\mathbf{L}_\tau^\dagger \mathbf{L}_\tau]) \mathbf{J} \mathbf{X} \quad (80)$$

so that it follows from (48c) that $\|(\mathbf{L}_\tau^\dagger \mathbf{L}_\tau - \mathbb{E}[\mathbf{L}_\tau^\dagger \mathbf{L}_\tau]) \mathbf{J} \mathbf{X}\| \leq 2C_5$ for all \mathbf{X} such that $\|\mathbf{J} \mathbf{X}\| \leq C_5$. Invoking the law of large numbers as before, $\mathbf{D}_t^2(\mathbf{X})/t \rightarrow 0$ almost surely as $t \rightarrow \infty$. Consequently, there exists d_t^2 such that $\|\mathbf{D}_t^2(\mathbf{X})\| \leq d_t^2$ and $d_t^2/t \rightarrow 0$ almost surely as $t \rightarrow \infty$.

In order to establish the Lipschitz continuity of $\mathbf{D}_t^2(\mathbf{X})$, observe that $\mathbf{D}_t^2(\mathbf{X}) - \mathbf{D}_t^2(\check{\mathbf{X}}) = \mathbf{C}'_t(\mathbf{X} - \check{\mathbf{X}})$, where

$$\mathbf{C}'_t := \sum_{\tau=0}^t \mathbf{L}_\tau^\dagger \mathbf{L}_\tau - \mathbb{E}[\mathbf{L}_\tau^\dagger \mathbf{L}_\tau] \quad (81)$$

Since each summand in (81) is zero mean and bounded, it holds from law of large numbers that $\mathbf{C}'_t/t \rightarrow 0$ almost surely as $t \rightarrow \infty$. Consequently, there exists π_t^2 such that $\|\mathbf{D}_t^2(\mathbf{X}) - \mathbf{D}_t^2(\check{\mathbf{X}})\| \leq \pi_t^2 \|\mathbf{X} - \check{\mathbf{X}}\|$, and $\pi_t^2/t \rightarrow 0$ almost surely as $t \rightarrow \infty$. \blacksquare

Proof of Lemma 2: Bounds in (22) can be derived by observing that for $1 \leq \tau \leq t$ and $\iota = 1, 2$, it holds that

$$\begin{aligned} \mathbf{D}_\tau^\iota(\hat{\mathbf{X}}_\tau) - \mathbf{D}_{\tau-1}^\iota(\hat{\mathbf{X}}_{\tau-1}) &= \mathbf{K}_\tau^\iota - \mathbf{K}_{\tau-1}^\iota + \\ &\quad \mathbf{D}_{\tau-1}^\iota(\hat{\mathbf{X}}_\tau) - \mathbf{D}_{\tau-1}^\iota(\hat{\mathbf{X}}_{\tau-1}). \end{aligned} \quad (82)$$

Summing (82) over $\tau = 1, \dots, t$, it follows that

$$\mathbf{D}_t^\iota(\hat{\mathbf{X}}_t) - \mathbf{D}_0^\iota(\hat{\mathbf{X}}_0) = \mathbf{K}_t^\iota - \mathbf{K}_0^\iota + \sum_{\tau=1}^t \left(\mathbf{D}_\tau^\iota(\hat{\mathbf{X}}_{\tau+1}) - \mathbf{D}_\tau^\iota(\hat{\mathbf{X}}_\tau) \right)$$

Observing that $\mathbf{K}_0^\iota = \mathbf{D}_0^\iota(\hat{\mathbf{X}}_0)$, a bound on \mathbf{K}_t^ι can be derived by using (48d) and (48e) as follows:

$$\|\mathbf{K}_t^\iota\| \leq \left\| \mathbf{D}_t^\iota(\hat{\mathbf{X}}_t) \right\| + \sum_{\tau=1}^t \left\| \mathbf{D}_\tau^\iota(\hat{\mathbf{X}}_{\tau+1}) - \mathbf{D}_\tau^\iota(\hat{\mathbf{X}}_\tau) \right\| \quad (83)$$

$$\leq d_t^\iota + \sum_{\tau=1}^t \pi_\tau^\iota \left\| \hat{\mathbf{X}}_{\tau+1} - \hat{\mathbf{X}}_\tau \right\| \quad (84)$$

$$= d_t^\iota + \mu \sum_{\tau=1}^t \pi_\tau^\iota \left\| \mathbf{L}_\tau^\dagger \mathbf{B}_\tau^c(\hat{\mathbf{X}}_\tau) \hat{\mathbf{X}}_\tau - \mathbf{L}_\tau^\dagger \mathbf{L}_\tau \hat{\mathbf{X}}_\tau \right\| \quad (85)$$

$$\leq d_t^\iota + \mu \sum_{\tau=1}^t \pi_\tau^\iota \left(C_4 + \left\| \mathbf{L}_\tau^\dagger \mathbf{L}_\tau \mathbf{J} \hat{\mathbf{X}}_\tau \right\| \right) \quad (86)$$

$$\leq d_t^\iota + \mu(C_4 + C_5) \sum_{\tau=1}^t \pi_\tau^\iota \quad (87)$$

so that $C_1 = C_2 = (C_4 + C_5)$ for $\iota = 1, 2$.

The bound on $\|\mathbf{K}_t^3\|$ follows from applying triangle inequality on (21c), and using (48a) as follows:

$$\|\mathbf{K}_t^3\| \leq \sum_{\tau=1}^{t-1} \left\| \mathbb{E}[\mathbf{L}_\tau^\dagger \mathbf{B}_\tau^\epsilon(\hat{\mathbf{X}}_\tau) \hat{\mathbf{X}}_\tau - \mathbf{L}_\tau^\dagger \mathbf{B}_\tau^\epsilon(\tilde{\mathbf{X}}_\tau) \tilde{\mathbf{X}}_\tau] \right\| \quad (88)$$

$$\leq \sum_{\tau=1}^{t-1} \mathbb{E} \left[\left\| \mathbf{L}_\tau^\dagger \mathbf{B}_\tau^\epsilon(\hat{\mathbf{X}}_\tau) \hat{\mathbf{X}}_\tau - \mathbf{L}_\tau^\dagger \mathbf{B}_\tau^\epsilon(\tilde{\mathbf{X}}_\tau) \tilde{\mathbf{X}}_\tau \right\| \right] \quad (89)$$

$$\leq \sum_{\tau=1}^{t-1} C_3 \left\| \hat{\mathbf{X}}_\tau - \tilde{\mathbf{X}}_\tau \right\| = C_3 \sum_{\tau=1}^{t-1} \|\Delta_\tau\| \quad (90)$$

Finally, to show that $f_t(\mu) \leq f_T(\mu) \rightarrow 0$ for the interval $0 \leq t \leq T/\mu$, observe that for $\iota = 1, 2$, it holds that $\mu d_t^\iota \leq T d_t^\iota / t$. From (24), it is known that given any ε , there exists $t_0(\varepsilon)$ and C_d such that

$$\mathbb{P}[d_t^\iota / t \leq C_d] = 1 \quad \forall t, \quad (91)$$

$$\text{and} \quad \mathbb{P}[d_t^\iota / t \leq \varepsilon] = 1 \quad \forall t > t_0(\varepsilon). \quad (92)$$

Such a $t_0(\varepsilon)$ exists within $[0, T/\mu]$ for all $\mu \leq T/t_0(\varepsilon)$. Therefore, given ε , if $t \leq t_0(\varepsilon)$, it holds that

$$\mathbb{P}[\mu d_t^\iota \leq \varepsilon] = 1 \quad (93)$$

for all $\mu \leq \varepsilon / t_0(\varepsilon) C_d$. On the other hand, if $t > t_0(\varepsilon)$, (93) holds for all $\mu \leq T/t_0(\varepsilon/T)$. Combining the two cases, it holds that $\max_{0 \leq t \leq T/\mu} \mu d_t^\iota \rightarrow 0$, with probability one as $\mu \rightarrow 0$.

For the other two terms, observe similarly that given ε , there exists T_ε and C_π such that

$$\mathbb{P}[\pi_t^\iota / t \leq C_\pi] = 1 \quad \forall t, \quad (94)$$

$$\text{and} \quad \mathbb{P}[\pi_t^\iota / t \leq \varepsilon] = 1 \quad \forall t > T_\varepsilon. \quad (95)$$

Thus, given ε , if $t \leq T_\varepsilon$, it holds that

$$\mathbb{P} \left[\mu^2 \sum_{\tau=2}^t \pi_\tau^\iota \leq \varepsilon \right] = 1, \quad \forall \mu, \quad \text{s.t.} \quad \mu \leq \frac{1}{T_\varepsilon} \sqrt{\frac{\varepsilon}{C_\pi}}. \quad (96)$$

Similarly, the result in (96) holds for $t > T_\varepsilon$ for all $\mu \leq \frac{T}{T_\varepsilon/T^2}$. ■

APPENDIX C

PROOF OF LEMMA 3

For notational convenience, let $\check{\delta}_{mn}(t) := \frac{\delta_{mn}(t)}{\sqrt{\|\mathbf{x}_m - \mathbf{x}_n\|^2 + \epsilon_x}}$ and recall that $\bar{\delta}_{mn} = \mathbb{E}[\check{\delta}_{mn}(t)]$. The proof is divided into two parts. In the first part, we consider the case when \mathcal{G}_t is connected, so that $p = N$. In this case, the goal is to show that

$$N \left[\mathbb{E}[\mathbf{L}_t^\dagger \mathbf{B}_t^\epsilon(\mathbf{X})] \right]_{mn} = \begin{cases} -\bar{\delta}_{mn} & m \neq n \\ \sum_{n \neq m} \bar{\delta}_{mn} & m = n. \end{cases} \quad (97)$$

Since the graph is connected, it holds that $\mathbf{L}_t^\dagger = (\mathbf{L}_t + \mathbf{1}\mathbf{1}^T/N)^{-1} - \mathbf{1}\mathbf{1}^T/N$. Let ψ_{mn} denote the (m, n) -th co-factor of $\mathbf{L}_t + \mathbf{1}\mathbf{1}^T/N$ and $\Psi := \det(\mathbf{L}_t + \mathbf{1}\mathbf{1}^T/N)$, so that $[\mathbf{L}_t^\dagger]_{mn} = \psi_{mn}/\Psi - 1/N$. Since \mathbf{L}_t^\dagger has zero row and column sums, we also have that $\sum_{n=1}^M \psi_{mn} = \Psi$. Therefore, expanding along the m -th row, the expression for Ψ becomes

$$\Psi = \sum_{n \neq m} w_{mn}(t)(\psi_{mm} - \psi_{mn}) + \frac{1}{N} \sum_{n=1}^M \psi_{mn} \quad (98)$$

$$= \frac{N}{N-1} \sum_{n \neq m} w_{mn}(t)(\psi_{mm} - \psi_{mn}) \quad (99)$$

for each $1 \leq m \leq N$. Straightforward manipulations allow us to conclude that

$$\left[\mathbf{L}_t^\dagger \mathbf{B}_t^\epsilon(\mathbf{X}) \right]_{mn} = \frac{1}{\Psi} \begin{cases} -\check{\delta}_{mn}(t)w_{mn}(t)(\psi_{mm} - \psi_{mn}) & m \neq n \\ -\sum_{k \neq m, n} w_{nk}(t)\check{\delta}_{nk}(t)(\psi_{mn} - \psi_{mk}) \\ \sum_{k \neq m} w_{mk}(t)\check{\delta}_{mk}(t)(\psi_{mm} - \psi_{mk}) & m = n. \end{cases}$$

Next, we show that the random variables ψ_{mn} and ψ_{mk} are identically distributed for $n \neq k \neq m$. Without loss of generality, let $m = 1$. Also, let \mathbf{L}_i^{nk} denote the $(N-2) \times (N-2)$ submatrix of $\mathbf{L}_t + \mathbf{1}\mathbf{1}^T/N$ after the removal of rows $(1, i)$ and columns (n, k) . The Laplace expansion of

ψ_{1n} along the k -th column yields

$$\begin{aligned}
\psi_{1n} &= - \sum_{i \neq 1, n, k} \left(\frac{1}{N} - w_{ki}(t) \right) (-1)^{n+i+k} |\mathbf{L}_i^{nk}| \\
&\quad - \left(\frac{1}{N} - w_{kn}(t) \right) (-1)^k |\mathbf{L}_n^{nk}| - \left(\frac{1}{N} + \sum_{i \neq k} w_{ki}(t) \right) (-1)^n |\mathbf{L}_k^{nk}| \\
&= - \sum_{i \neq 1} \left(\frac{1}{N} - w_{ki}(t) \right) (-1)^{n+i+k} |\mathbf{L}_i^{nk}| \\
&\quad - \left(\sum_{i \neq k, n} w_{ki}(t) + 2w_{kn}(t) \right) (-1)^n |\mathbf{L}_k^{nk}| \tag{100}
\end{aligned}$$

Likewise, the expansion of ψ_{1k} along the n -th column yields

$$\begin{aligned}
\psi_{1k} &= - \sum_{i \neq 1} \left(\frac{1}{N} - w_{ni}(t) \right) (-1)^{n+i+k} |\mathbf{L}_i^{nk}| \\
&\quad - \left(\sum_{i \neq k, n} w_{ki}(t) + 2w_{kn}(t) \right) (-1)^k |\mathbf{L}_n^{nk}| \tag{101}
\end{aligned}$$

It can be seen that the first terms in (100) and (101) are identically distributed since $w_{ni}(t)$ and $w_{ki}(t)$ are identical (cf. **(A5)**). Further, performing $n - k$ row exchanges on \mathbf{L}_n^{nk} , it is possible to obtain $\tilde{\mathbf{L}}_n^{nk}$ which only differs from \mathbf{L}_k^{nk} in the k -th row. Indeed, the elements of the k -th row of $\tilde{\mathbf{L}}_n^{nk}$ are $\{(1/N - w_{ki}(t))\}_{i \neq k, n}$, while the elements of the k -th row of \mathbf{L}_n^{nk} are $\{(1/N - w_{ni}(t))\}_{i \neq k, n}$. Since the determinant is linear in its rows, it follows that $|\mathbf{L}_n^{nk}|$ and $|\tilde{\mathbf{L}}_n^{nk}| = (-1)^{n+k} |\mathbf{L}_k^{nk}|$ are identically distributed. In summary, we have that the distributions of ψ_{mn} and ψ_{mk} are identical for all $k \neq n \neq m$.

Next, define identical random variables $\chi_{mn} := w_{mn}(t)(\psi_{mm} - \psi_{mn})$ for each $n \neq m$, so that $\Psi = \frac{N}{N-1} \sum_{n \neq m} \chi_{mn}$. Since \mathcal{G}_t is connected, it holds that $\Psi > 0$. Therefore from symmetry, we have that

$$\mathbb{E} \left[\frac{\chi_{mn}}{\Psi} \right] = \frac{N-1}{N} \mathbb{E} \left[\frac{\chi_{mn}}{\sum_{n \neq m} \chi_{mn}} \right] = \frac{1}{N} \tag{102}$$

Further, using the fact that $\mathbb{E}[\chi_{mn}] = \mathbb{E}[\chi_{mk}]$ for each $k \neq n$, it can be seen that

$$\mathbb{E} \left[\mathbf{L}_t^\dagger \mathbf{B}_t^\epsilon(\mathbf{X}) \right]_{mn} = \frac{1}{N} \begin{cases} -\bar{\delta}_{mn} & m \neq n \\ \sum_{k \neq m} \bar{\delta}_{mk} & m = n. \end{cases}$$

which is the required result.

Finally, if \mathcal{G}_t consists of multiple connected components, the quantity $\mathbf{L}_t^\dagger \mathbf{B}_t^\epsilon(\mathbf{X})$ is a permuted version of the block-diagonal matrix with N/p block matrices of size $p \times p$ each. Let Ψ^j denote the determinant of j -th block, and the random variables χ_{mn}^j be similarly defined block-wise. Proceeding along similar lines, it can be seen that

$$\mathbb{E}\left[\frac{\chi_{mn}^j}{\Psi^j}\right] = \frac{p-1}{p} \mathbb{E}\left[\frac{\chi_{mn}^j}{\sum_{n \neq m} \chi_{mn}^j}\right] = \frac{1}{p}. \quad (103)$$

Consequently, $[\mathbf{L}_t^\dagger \mathbf{B}_t^\epsilon(\mathbf{X})]_{mn}$ is non-zero if and only if the node pair (m, n) belong to the same component, and is zero otherwise. From (A5), we have that the probability that a given pair of nodes (m, n) belongs to the same connected component is given by $(p-1)/(N-1)$, yielding the required expression

$$\mathbb{E}\left[\mathbf{L}_t^\dagger \mathbf{B}_t^\epsilon(\mathbf{X})\right]_{mn} = \frac{p-1}{p(N-1)} \begin{cases} -\bar{\delta}_{mn} & m \neq n \\ \sum_{k \neq m} \bar{\delta}_{mk} & m = n. \end{cases}$$

REFERENCES

- [1] I. Borg and P. J. Groenen, *Modern multidimensional scaling: Theory and applications*. Springer Science & Business Media, 2005.
- [2] L. V. D. Maaten and G. Hinton, “Visualizing data using t-SNE,” *Journal of Machine Learning Research*, vol. 9, no. Nov, pp. 2579–2605, 2008.
- [3] A. Platzter, “Visualization of SNPs with t-SNE,” *PloS One*, vol. 8, no. 2, 2013.
- [4] J. Tang, M. Qu, M. Wang, M. Zhang, J. Yan, and Q. Mei, “Line: Large-scale information network embedding,” in *Proc. of the Intl. Conf. on World Wide Web*, 2015, pp. 1067–77.
- [5] L. van der Maaten and G. Hinton, “Visualizing non-metric similarities in multiple maps,” *Machine Learning*, vol. 87, no. 1, pp. 33–55, 2012.
- [6] D. K. Agrafiotis, “Stochastic proximity embedding,” *Journal of computational chemistry*, vol. 24, no. 10, pp. 1215–21, 2003.
- [7] J. Tzeng, H. H. Lu, and W.-H. Li, “Multidimensional scaling for large genomic data sets,” *BMC bioinformatics*, vol. 9, no. 1, p. 179, 2008.
- [8] J. Y. Choi, S.-H. Bae, J. Qiu, B. Chen, and D. Wild, “Browsing large-scale cheminformatics data with dimension reduction,” *Concurrency and Computation: Practice and Experience*, vol. 23, no. 17, pp. 2315–25, 2011.
- [9] M. Beatty and B. Manjunath, “Dimensionality reduction using multi-dimensional scaling for content-based retrieval,” in *Proc. of the ICIP*, vol. 2, 1997, pp. 835–838.
- [10] J. A. Costa, N. Patwari, and A. O. Hero III, “Distributed weighted-multidimensional scaling for node localization in sensor networks,” *ACM Trans. on Sensor Networks*, vol. 2, no. 1, pp. 39–64, 2006.
- [11] P. A. Forero and G. B. Giannakis, “Sparsity-exploiting robust multidimensional scaling,” *IEEE Trans. on Signal Proc.*, vol. 60, no. 8, pp. 4118–34, 2012.

- [12] K. S. Xu, M. Klinger, and A. O. Hero III, "A regularized graph layout framework for dynamic network visualization," *Data Mining and Knowledge Discovery*, vol. 27, no. 1, pp. 84–116, 2013.
- [13] S. Kumar, R. Kumar, and K. Rajawat, "Cooperative localization of mobile networks via velocity-assisted multidimensional scaling," *IEEE Trans. on Signal Process.*, vol. 64, no. 7, pp. 1744–1758, 2016.
- [14] A. Simonetto and G. Leus, "Distributed maximum likelihood sensor network localization," *IEEE Trans. on Signal Proc.*, vol. 62, no. 6, pp. 1424–1437, 2014.
- [15] S. Kumar, R. Jain, and K. Rajawat, "Asynchronous optimization over heterogeneous networks via consensus admm," *IEEE Trans. on Signal and Inf. Proc. over Networks*, 2016 (to be published).
- [16] P. Biswas, T.-C. Lian, T.-C. Wang, and Y. Ye, "Semidefinite programming based algorithms for sensor network localization," *ACM Transactions on Sensor Networks*, vol. 2, no. 2, pp. 188–220, 2006.
- [17] S.-H. Bae, J. Qiu, and G. Fox, "Adaptive interpolation of multidimensional scaling," *Procedia Computer Science*, vol. 9, pp. 393 – 402, 2012.
- [18] S. Ingram, T. Munzner, and M. Olano, "Glimmer: Multilevel MDS on the GPU," *IEEE Trans. on Visualization and Computer Graphics*, vol. 15, no. 2, pp. 249–261, 2009.
- [19] L. Bottou, *On-line Learning in Neural Networks*, D. Saad, Ed. New York, NY, USA: Cambridge University Press, 1998.
- [20] C. D. Sa, C. Re, and K. Olukotun, "Global convergence of stochastic gradient descent for some non-convex matrix problems," in *Proc. of the Intl. Conf. on Machine Learning*, 2015, pp. 2332–41.
- [21] J. Mairal, "Stochastic majorization-minimization algorithms for large-scale optimization," in *Advances in Neural Information Processing Systems*, 2013, pp. 2283–2291.
- [22] O. Cappé and E. Moulines, "On-line expectation–maximization algorithm for latent data models," *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, vol. 71, no. 3, pp. 593–613, 2009.
- [23] A. H. Sayed, *Adaptive filters*. John Wiley & Sons, 2011.
- [24] V. Solo and X. Kong, *Adaptive signal processing algorithms: stability and performance*. Prentice-Hall, Inc., 1994.
- [25] H. Kushner and G. G. Yin, *Stochastic approximation and recursive algorithms and applications*. Springer, 2003, vol. 35.
- [26] V. S. Borkar *et al.*, "Stochastic approximation," *Cambridge Books*, 2008.
- [27] W. Rudin, *Principles of mathematical analysis*. McGraw-Hill New York, 1964, vol. 3.
- [28] N. Patwari, J. N. Ash, S. Kyperountas, A. O. Hero, R. L. Moses, and N. S. Correal, "Locating the nodes: cooperative localization in wireless sensor networks," *IEEE Signal Process. Mag.*, vol. 22, no. 4.
- [29] H. Wymeersch, J. Lien, and M. Z. Win, "Cooperative localization in wireless networks," *Proceedings of the IEEE*, vol. 97, no. 2, pp. 427–450, 2009.
- [30] I. Demirkol, C. Ersoy, F. Alagoz *et al.*, "Mac protocols for wireless sensor networks: a survey," *IEEE Commun. Mag.*, vol. 44, no. 4.
- [31] W. S. Torgerson, "Multidimensional scaling of similarity," *Psychometrika*, vol. 30, no. 4, pp. 379–393, 1965.
- [32] C. Savarese, J. M. Rabaey, and J. Beutel, "Location in distributed ad-hoc wireless sensor networks," in *Proc. of the IEEE ICASSP*, vol. 4, 2001, pp. 2037–2040.
- [33] L. Dong, "Cooperative localization and tracking of mobile ad hoc networks," *IEEE Trans. on Signal Proc.*, vol. 60, no. 7.
- [34] S. Kim, P. A. Thiessen, E. E. Bolton, J. Chen, G. Fu, A. Gindulyte, L. Han, J. He, S. He, B. A. Shoemaker *et al.*, "Pubchem substance and compound databases," *Nucleic acids research*, p. gkv951, 2015.
- [35] E. E. Bolton, Y. Wang, P. A. Thiessen, and S. H. Bryant, "Pubchem: integrated platform of small molecules and biological activities," *Annual reports in computational chemistry*, vol. 4, pp. 217–241, 2008.

- [36] P.-N. Tan, M. Steinbach, and V. Kumar, *Introduction to data mining*. Pearson Education India, 2006.
- [37] F. M. Harper and J. A. Konstan, “The movielens datasets: History and context,” *ACM Trans. on Interactive Intelligent Systems (TiiS)*, vol. 5, no. 4, p. 19, 2016.
- [38] T. M. Newcomb, *The acquaintance process*. Holt, Rinehart & Winston, 1961.
- [39] B. Mohar, “Laplace eigenvalues of graphs: a survey,” *Discrete mathematics*, vol. 109, no. 1-3, pp. 171–183, 1992.
- [40] G. H. Golub and C. F. Van Loan, *Matrix computations*. JHU Press, 2012, vol. 3.