

Network Topology Mapping from Partial Virtual Coordinates and Graph Geodesics

Anura P. Jayasumana, Randy Paffenroth, Gunjan Mahindre, Sridhar Ramasamy, and Kelum Gajamannage

Abstract

For many important network types (e.g., sensor networks in complex harsh environments and social networks) physical coordinate systems (e.g., Cartesian), and physical distances (e.g., Euclidean), are either difficult to discern or inapplicable. Accordingly, coordinate systems and characterizations based on hop-distance measurements, such as Topology Preserving Maps (TPMs) and Virtual-Coordinate (VC) systems are attractive alternatives to Cartesian coordinates for many network algorithms. Herein, we present an approach to recover geometric and topological properties of a network with a small set of distance measurements. In particular, our approach is a combination of shortest path (often called geodesic) recovery concepts and low-rank matrix completion, generalized to the case of hop-distances in graphs. Results for sensor networks embedded in 2-D and 3-D spaces, as well as a social networks, indicates that the method can accurately capture the network connectivity with a small set of measurements. TPM generation can now also be based on various context appropriate measurements or VC systems, as long as they characterize different nodes by distances to small sets of random nodes (instead of a set of global anchors). The proposed method is a significant generalization that allows the topology to be extracted from a random set of graph shortest paths, making it applicable in contexts such as social networks where VC generation may not be possible.

Keywords: Localization, virtual coordinates, topology preserving maps, sensor networks, social networks.

I. INTRODUCTION

Large and complex networks naturally arise in communication and social networks, the Internet-of-Things (IoT), and many other systems of importance. Data embedded in such networks exhibit distributed, intricate, and dynamic patterns. Our ability to extract information from and about these networks, detect anomalies, and influence their performance can be substantially improved by a deeper understanding of their local and global structures [58], [57], [1]. Yet such operations are often impeded due to the size and complexity of these networks, and constraints such as energy, measurement cost, and accessibility that prevent *geometric and topological* structure of the entire network from being fully observed, measured, characterized, or processed. Inferences have to be made and patterns should be detected in the absence of complete information. Accordingly, in this paper, we derive and demonstrate novel techniques for detecting and representing network structures, e.g., topology, connectivity, and layout, which are scalable in their computation and communication, as well as graceful in their degradation in the presence of limited measurements. In particular, we sample a network with a small set of pair-wise hop-distances and use matrix completion techniques to capture the topology of the network. We demonstrate the technique by deriving topology preserving maps indicative of the layout of 2-D and 3-D sensor networks with far fewer measurements compared to existing schemes. Furthermore, this new approach extends applicability of hop-distance based techniques to cover complex networks such as social networks while providing a foundation for using a broader set of sampling strategies.

A. Motivating Examples

Of specific interest to us are networks of inexpensive wireless devices such as smart Radio-Frequency Identification (RFID) tags or self-powered sensor nodes embedded in complex 2-D or 3-D surfaces and volumes. Here, the network features are often characterized using the physical (Cartesian) coordinates of each node. One then relies on the Euclidean properties of the network layout (e.g., the distance between nodes) for functions such as area or volume coverage, topology control, sensing, and routing. Physical location estimation is based on measurements such as the Received Signal Strength Indicator (RSSI) and time delay [38]. In networks deployed in harsh or complex environments, such methods are hampered or made completely ineffective by issues such as multi-path interference, reflections, shadowing and clock synchronization [46], [25]. Although

A. P. Jayasumana, S. Ramasamy and G. Mahindre are with the Department of Electrical & Computer Engineering, Colorado State University, Fort Collins, CO 80525

R. Paffenroth is with the Department of Mathematical Sciences, Department of Computer Science, and the Data Science Program, Worcester Polytechnic Institute, Worcester, MA 01609

K. Gajamannage is with the Department of Mathematical Sciences, Worcester Polytechnic Institute, Worcester, MA 01609

physical coordinates provide a meaningful representation for node location, for sensor networks deployed on complex surfaces and volumes, *graph distance* based techniques that use only connectivity measurements offer a more robust, and attractive alternative [25]. For instance, in routing they overcome local minima caused by concave physical voids. Connectivity based techniques are also generalizable to networks such as social networks, for which there is no notion of a location in a physical or Euclidean space. Even if a node is associated with a location, physical proximity does not necessarily imply connectivity. In fact for many such applications, what matters is the logical topology, i.e., connectivity and the hop distances among node pairs.

B. Hop-distance based Network Sampling

We focus our attention on *undirected graphs* G defined by $G = \{V, E\}$, where V is the set of nodes or vertices and E is a set of edges corresponding to communication links, friendship status on a social network, etc. Such a graph may be represented by an *adjacency matrix* A [28], where

$$A = [a_{ij} | a_{ij} = 1 \text{ if } (i, j) \in E, 0 \text{ otherwise}]. \quad (1)$$

Note that $a_{ij} = a_{ji}$ when G is undirected. To make the current text concise, herein we focus on *unweighted* graphs where each $a_{ij} \in \{0, 1\}$. Our main object of study will be *Hop-Distance Matrices* (HDMs), $H \in \mathbb{N}_0^{N \times N}$, (where \mathbb{N}_0 denotes the non-negative integers $\mathbb{N} \cup \{0\}$), where

$$h_{ij} = \begin{array}{l} \text{the length of the shortest path} \\ \text{from node } i \text{ to node } j. \end{array} \quad (2)$$

H , by construction, is symmetric and is invariant to the situation where two nodes have multiple paths between them of the same shortest length. From a mathematical perspective, hop-distances h_{ij} may be thought of as the computation of *geodesics* (or shortest paths) in a graph [28], [53], [39], a perspective that we will make use of in the sequel. In particular, herein the term *geodesic* will refer to a shortest path between nodes, either when embedded in a particular space (e.g., the sum of straight line distances between node pairs forming a path in a Euclidean space), or when considered as a node in a graph (e.g., as computed using Dijkstra’s algorithm [43]).

Our goal is to capture the topology accurately using efficient network sampling schemes. Prominent among hop-distance based sampling methods are those relying on anchor-based Virtual Coordinate Systems (VCS), in which each node is represented by a Virtual Coordinate (VC) vector corresponding to the *minimum* hop-distances to a small set of nodes, known as *anchor nodes* [15], [42], [51]. If we have M anchor nodes, then a VCS is an $N \times M$ sub-matrix P of H where the i -th row provides an M dimensional coordinate vector for the i -th node in the graph. A VCS is then a “relative” coordinate system, as opposed to classic “absolute” coordinate systems such as Cartesian or Spherical. In particular, a VCS does not possess, or require, an absolute origin or other fixed geometric properties such as a notion of angle. However, a VCS lacks directional information as each coordinate indicates only the distance to an anchor. Thus all the information about the network layout such as shapes, voids, and boundaries are lost in a VCS. This issue leads to natural generalizations of VCSs, such as Topology Coordinates (TCs) and Topology Preserving Maps (TPMs) [25], where one performs an eigen-decomposition of P or H similar to that performed in Principal Component Analysis (PCA) [23], [26]. Such methods can recover the general shapes and boundaries of the physical network layout, thus providing an effective alternative for representing 2-D and 3-D sensor networks and carrying out operations such as routing and boundary detection, but without the need for physical distance measurements [23], [26].

C. Contribution and Significance

At a high-level, our approach is characterized by a) the choice of H to represent a network, b) network sampling schemes that allow the measurement of a small set of entries of H , e.g., to accommodate the constraints of accessibility, reduce the measurement complexity, etc., and c) filling in the incomplete sampled H to obtain the network layout and topology.

A complete H and a complete A are interchangeable, i.e., to construct A from H one merely needs to set all entries with $h_{ij} > 1$ to 0, while H can be obtained from A using procedures such as Dijkstra’s algorithm [43]. However, a *partially observed* H and a *partially observed* A are quite different. As opposed to an incomplete A , we will demonstrate that an incomplete H has quite interesting low-rank properties, both from theoretical and practical perspectives, especially for graphs arising from real-world networks. Given the properties such as low-rankness and an effective measurement scheme, the matrix H can be completed using efficient techniques for convex, low-rank matrix completion. The newly completed H can then be used to compute A or to perform any other desired analysis of the graph.

The main contribution of this paper is a technique that combines VC based techniques with low-rank matrix completion, that allows the extraction of topology and geometric features of a network from a small set of shortest path distances.

- In case of networks embedded in 2-D and 3-D physical spaces, we demonstrate that the topology preserving maps can be recovered using only a small fraction of VC values, as opposed to existing method [25] that requires the full set of VCs.
- We broaden the possible set of network sampling schemes far beyond our previous results in [24], [22], [25], and point to a theory on which new sampling schemes may be grounded, i.e., selection of samples should not hinder the ability to complete the resulting incomplete low-rank matrices. We demonstrate the reconstruction of the network topology from a small set of shortest path length measurements. Consequently, TPMs can now be generated based on a variety of geodesic measurement approaches, of which an anchor-based VCS is one instance.

The novelty and the significance of this work may be gauged by the following impacts it can have on network analytics:

- The theory of low-rank matrix completion can now be applied for exploiting the sparseness of complex real-world networks, and to develop communication and computation efficient techniques for large-scale networks. Although H for an arbitrary graph may not be low rank, we demonstrate that for a broad class of complex real-life networks such distance matrices are surprisingly low-rank. Our analysis is based upon ideas in low-rank matrix completion [41], [12], [13], [47], [48], which have been shown to scale to large problems with many thousands, if not millions, of entries [47], [48], and the underlying topology that we uncover is closely related to ideas in Non-linear Dimension Reduction (NDR), such as Isomap [53], [39], but generalized to the case where only hop-distances are measured.
- The applicability of hop-distance based property extraction is extended to cases where certain distances and connectivities are not observable. For large and complex networks, and especially those with access limitations, we can likely not even measure a complete set of hop-distances to a set of common anchors. Examples of such cases include sensor and communication networks with restricted access to certain nodes. As explained in [7] using examples from communication networks, it is realistic to obtain the distances between nodes in many cases, while it is difficult or impossible to obtain information about edges or absence of edges that are far away from the query node. The same argument extends to many practical problems dealing with data, e.g., pathway prediction problem in proteins [19], [54], where the distance between two nodes (proteins) can be evaluated yet finding the shortest path by some measure (e.g., folding sequence) from one node to the other is complex. Social networks such as Facebook, Instagram, and Orkut represent large scale networks. Controlled and focused crawling [16], [17] are widely used approaches for data collection in such networks for acquisition of data regarding hop distances, links, etc. As some profiles (nodes) are publicly accessible and some are not [45], we cannot gather data for all the nodes while crawling. In such cases, the hop distances acquired, or the links observed are only partial entries in the complete distance or adjacency matrix. Also, as proved in [35], missing data can have a significant impact on structural properties of a social network.
- Matrix completion algorithm finds many applications such as, finding top N recommends for users [56] and link prediction [50]. While matrix completion as applied to Euclidean distance matrices is not new [3], [44], [36], [4], [31] our approach differs in a key aspect from those currently found in the literature. Our focus is on using *hop-distances* rather than the more classic distance measures such as Euclidean distances [52]. In other words, instead of considering Euclidean Distance Matrices (EDM) where each entry of the matrix codes for the Euclidean distance between two nodes, we consider HDM where each entry of the matrix codes for the *graph shortest path distance* between two nodes.
- For a vast variety of networks, e.g., social and computer networks, *the concept of Euclidean distances is not even applicable* (e.g., what is the “Euclidean distance” between two people in a social network?). Accordingly, we leverage HDMs, which are the most faithful representations of the network to which we have access. As far as we are aware, problems of such generality have not been considered before.

Section II reviews the VC based sampling schemes. The theoretical considerations and results supporting proposed geodesic sampling schemes and low-rankness are outlined in Section III. Methodology for VC based sampling and reconstruction is presented in Section IV followed by results in Section V and conclusion.

II. RELATION TO PRIOR WORK

This section reviews hop-distance based network sampling schemes, the virtual-coordinate based representation, and their relationship to topology coordinates which recover the Euclidean properties of 2-D and 3-D networks. An anchor-based VCS is an M -dimensional abstraction of the network connectivity, where each node is represented by an M tuple, called VC vector which contains the shortest path length (in hops) from the node to each of a set of M anchors [14], [15]. In an Aligned VC system [42], each node re-evaluates its coordinates by averaging its own coordinates with its neighbors’ coordinates. Axis-based Virtual Coordinate Assignment Protocol [51] estimates a 5-tuple VC for each node corresponding to longitude, latitude, ripple, up, and down.

A key question related to anchor-based VCS is the number M and the placement of anchors. If an adequate number of anchors are not appropriately deployed, it may cause the network to suffer from identical coordinates and local minima [27], resulting in logical/virtual voids. The overhead associated with time and energy consumed for coordinate generation grows with the number of anchors. The difficulty of determining the optimal anchor set is compounded by the fact that the number of anchors and their optimal placement are dependent on each other. In Virtual Coordinate Assignment Protocol, the coordinates are defined based on three anchors [15], while all the perimeter nodes are selected in [42]. Extreme Node Search uses two initial random anchors which allows the selection of extreme nodes on internal and external boundaries as anchors [21]. An attractive alternative is to bypass anchor selection altogether, and select a set of random anchors. Note, while from a networking perspective a random selection of anchors may seem somewhat odd (i.e., a judicious placement might be viewed as more appropriate), from a matrix completion perspective such a placement is not only justified, but in many instances may be optimal [41], [12], [13], [48]. Again, such network organizational ideas have a dual perspective in the mathematical literature, and the selection of anchors is closely related to the idea of *incoherence* in low-rank matrix recovery literature [41], [12], [13], [48].

Geographical features such as boundaries and voids are missing from an anchor-based VCS. In our previous work [25], [34], recovery of geographic features from VCS is achieved by TPMs. TPMs derived from VCs are maps that are nearly homeomorphic to physical maps [25]. A TPM is a distorted version of the real physical layout (map) in such a way that the distortion accounts for connectivity information. In case of a 2D or 3D sensor network with M anchors, VCS is a mapping from the 2D or 3D network layout to an M -dimensional space. TPMs recover a 2D or 3D projection from this M -dimensional representation such that it preserves the main features such as boundaries and shapes of the network. Thus, TPMs can serve as an effective alternative for physical coordinates in many network related functions. Reconstruction of the graph from its VCs is attempted in [10] using an algorithmic approach. TC based schemes have demonstrated performance comparable, or better than the corresponding geographic coordinate based counterparts [23], [26], [34].

While the algorithms that are based on VCs or TCs [25], [18] provide a viable, competitive and robust alternative to traditional geographic coordinate based methods, these techniques have so far required the complete set of VCs in order to extract the geometric information. However, as we demonstrate here, one can recover topological properties of a network *without the need for complete knowledge of the virtual coordinates*. While complexity of generating a complete distance matrix through flooding is of order $O(N^2)$, that for VC generation from a set of $M(M \ll N)$ anchors is only $O(NM)$. Herein, we further reduce this computational cost and justify those results. In particular, we demonstrate a connection between TPMs, NDR (Nonlinear Dimension Reduction) [53], [39], and low-rank matrix completion problems [41], [12], [13], [47], [48].

Finally, we observe that the computation of a low-rank approximation of an EDM using PCA is equivalent to the Multi-Dimensional Scaling (MDS) [37], [9] algorithm. Even closer to our proposed technique, many authors consider *geodesic* adjacency matrices A [53], [39] generated by drawing short range, or *neighbor*, distances from the EDM (D). They then compute *long range* distances by way of *shortest path distances in the graph represented by A* . A low rank approximation of such a graph shortest path based distance matrix is equivalent to the Isomap [53], [39] algorithm for NDR.

III. THEORETICAL CONSIDERATIONS

Herein, we are interested in studying the low-rank structure of *graphs* that arise in real-world network applications and we begin by presenting a number of important theoretical ideas in this domain. In particular, we focus on two key questions when deriving our graph sampling and reconstruction schemes.

- First, what is the appropriate type of measurement to make? For example, it is quite classic to represent a graph as its *adjacency* matrix. However, we propose matrix completion is more effective when starting from a graph's *distance* matrix, even though the two representations are commonly viewed as being equivalent. As opposed to the adjacency matrix A , the hop distance matrix H provides *global information* about the graph. In particular, each entry of H provides constraints on *many entries* of A , with the simplest example of such global information being the Triangle-Inequality [2].
- Second, what are the appropriate structural assumptions to make for reconstruction? Given a partially observed adjacency matrix, any completion is perfectly consistent since an adjacency matrix only provides local information. Even a distance matrix can only provide bounds (e.g., the triangle inequality). Accordingly, an appropriate structural assumption must be made (e.g., low rankness of H).

A. Graph reconstruction: Adjacency Matrices

Given partial knowledge of a graph, there are several techniques for predicting the unknown information about the graph, including an array of combinatorial techniques. Examples include NP-hard combinatorial algorithms such as minimal

Hamiltonian completions [55], [30] (i.e., adding the minimum number of edges to a graph to make a Hamiltonian path that visits each vertex once) and minimal Chordal completions [49] (i.e., adding the minimum number of edges to a graph such that every sufficiently long cycle has a chord). However, with the advent of effective and scalable tools for *large scale low-rank matrix* completion, in our work we choose to focus on matrix completion techniques that make a *low-rank assumption*.

Of course, given the natural connection between adjacency matrices and graphs, one would be tempted to look at low-rank structures in matrices such as A . Much is known in such cases, including the facts that [32]:

- The only rank-0 adjacency matrix is for the graph with no edges.
- The only rank-1 adjacency matrix is the complete graph.
- The only rank-2 adjacency matrix is the complete bipartite graph (and complete tripartite is rank-3, etc.).

It is also known, for example, that for a subgraph $S \subseteq G$ $rank(S) \leq rank(G)$ [32].

Note, *Graph Laplacians* is an area in which the low-rank structure of graphs can be precisely defined and is well understood [20]. We merely observe that given an adjacency matrix A one can define the diagonal matrix D_{row} as the row sums of A . Given such a D_{row} , a Graph Laplacian for A is defined as $L = D_{row} - A$. It is well known that the Graph Laplacian is low-rank if, and only if, the underlying graph is disconnected [20]. Unfortunately, such techniques do not lead to the types of predictions that we desire.

B. Low Rank Structure of EDMs

Moving away from graphs for a moment, we observe that there is a domain in which low-rank matrix completion has been used successfully for many years, and this is in the completion of EDM [31], [36]. There is a vast literature on such matrices, and especially their low-rank structure. In particular, for a EDM D one can show that $rank(D) \leq k+2$ where k is the dimension of the embedding space of the points whose pairwise distances comprise D [36].

Unfortunately, one is then left with the task of *defining and measuring these Euclidean distances*, which can be a non-trivial task. As observed previously, the relationship between the *communication distance* and the EDM between nodes in a network can be quite complicated, especially in the presence of routing algorithms, occlusions, and anything but the most trivial spatial geometry of the sensors. *Just because the sensors can physically be embedded in a 2-D or 3-D Euclidean space does not imply that the distances implied by the communication network can be embedded in the same space.* Accordingly, a goal in the current text is to understand what elements of the theory of EDMs can be preserved while having the results be applicable to large scale networks.

C. Exactly Low-rank HDMs

While our focus is on real-world networks, and calculation of the empirical low-rank structure and predictability, theoretical considerations provide important guide-posts. For example, while the low-rankness of adjacency matrices can be used to detect structures such as complete l -partite graphs and the low-rankness of Graph Laplacians can be used to detect disconnected graphs, we wish to treat more general scenarios.

On the other hand, the theory of the low-rank structure of HDMs provides intriguing glimpses into what is possible when applying a low-rank assumption to the analysis of HDMs. For example, one has access to a vast array of theorems of the following flavor.

Theorem 1. (restated from Theorem 2.16 from [5]) *Let G be a graph with HDM H . If G has only a single even cycle of length $2k$ and a total of $2k+p$ vertices, then H is of rank k .*

Interestingly, the theory of low-rank HDMs extends well beyond such special cases. For example, it is quite common to consider k -regular graphs where each row (and column) of A sum to k [28]. A similar idea, called *transmission regular*, can be defined in the HDM case by calling a graph G k -transmission regular when each row and column of H sum to k . Such transmission regular graphs give rise to a large and interesting class of low-rank HDMs by way of the following theorem.

Theorem 2. (restated from Theorem 4.5 from [5]) *Let G_1 and G_2 be two transmission regular graphs on n_1 and n_2 vertices with transmission regularity k_1 and k_2 , with k_1 and k_2 not necessarily equal. If G is the Cartesian product of G_1 and G_2 , then the rank of H is $n_1 n_2 - (n_1 - 1)(n_2 - 2)$.*

D. Approximately Low-rank HDMs

Even though the above theorems give rise to interesting families of low-rank hop-distance matrices, real-world graphs rarely, if ever, satisfy the assumptions of these theorems, or many other similar theorems for both adjacency and hop-distance matrices. Accordingly, it is important to consider the *approximate low-rank properties* of hop-distance matrices. For example, one can consider common synthetic models that approximate the structure of real-world graphs such as *scale-free networks* or *power-law graphs*, whose degree distributions follow, at least asymptotically, a power-law. For example, one often considers graphs where the fraction of nodes N_k having k links goes like $N_k \propto k^{-\gamma}$ for some parameter γ [33], [6].

Accordingly, in Figure 1, we show a comparison between the singular values of the adjacency matrix and the hop-distance matrix for two graphs, namely a synthetic power-law graph based upon the Holme-Kim model [33] with 500 nodes and a subgraph of the real-world Gowalla social network from [40] with 2000 nodes. In addition, for the Holme-Kim synthetic data [33] we examine 100 Monte-Carlo runs to see the difference that the precise state of the random graphs makes. Note, in both cases the matrices are double-centered (as discussed in Section IV-D) and the value of all singular-values are normalized relative to the size of the largest singular value (so all curves start on the left at 1). In both cases, the singular-values of H decay much more quickly than the singular-values of A , with the normalized 100th singular value of H in both cases being very close to 0, while the same for A are approximately 0.4 and 0.2 respectively.

Finally, we consider how the placement of anchor nodes affects the low-rank structure of the HDM. In Figure 2 we compare the low-rank structure of the HDM for the Holme-Kim [33] and Gowalla social network [40] between choosing a random set of anchors of H versus choosing a set of anchors based on different centrality measures. Somewhat surprisingly, the low-rank structure of H is not strongly affected by the choice of the sampling scheme. In particular, Figure 2 suggests that the accuracy of predictions from random anchors will be similar to that from anchors with high centrality (e.g., using nodes with a large number of neighbors as anchors) [28]. Note that the difference between the singular values of the hop-distance and adjacency matrices is much larger than the difference between the various sampling schemes. Of course, this does not mean that specially chosen anchors cannot change the low-rank structure of H , but it does suggest that random anchors is a reasonable anchor selection for initial investigation.

IV. APPROACH

Proposed method for extracting network topology is as follows:

- Start with a set of geodesics. Two sampling schemes are considered: anchor-based VCs (to obtain P or a subset thereof) and random shortest paths (to obtain a subset of elements of H).
- Complete the matrix P or H using low-rank matrix completion.
- Evaluate the accuracy of the computed topology or layout. In case of 2-D and 3-D sensor networks the accuracy of resulting topology preserving maps is used as the evaluation metric, while for the social networks we evaluate the difference between the actual and completed distance matrices.

Here we describe the different elements of our approach in more detail.

A. Anchor Based VCs

We follow the notation in [34] and consider networks where M of the N nodes are designated as “anchors”. With an anchor-based VCS, each of the N nodes in the network is characterized by a VC vector of length M , i.e., each node is labeled by its shortest-path hop distance to each of the M anchors.

Let $P \in \mathbb{N}_0^{N \times M}$ be the matrix containing the VCs of all the nodes, e.g., the i -th row corresponds to the $\mathbb{N}_0^{1 \times M}$ VC vector of the i -th node, and j -th column corresponds to the M -th virtual coordinate of all the nodes in the network with respect to j -th anchor.

This matrix can be written as

$$P = \begin{bmatrix} h_{1A_1} & \dots & h_{1A_M} \\ \vdots & \ddots & \vdots \\ h_{NA_1} & \dots & h_{NA_M} \end{bmatrix} \quad (3)$$

where h_{iA_j} is the hop-distance from node i to anchor A_j . P is precisely a subset of the full hop distance matrix H derived by selecting just a few anchor nodes and constructing P from the columns corresponding to those anchor nodes. For large networks it is generally desirable to have only a small subset of nodes as anchors, i.e., $M \ll N$.

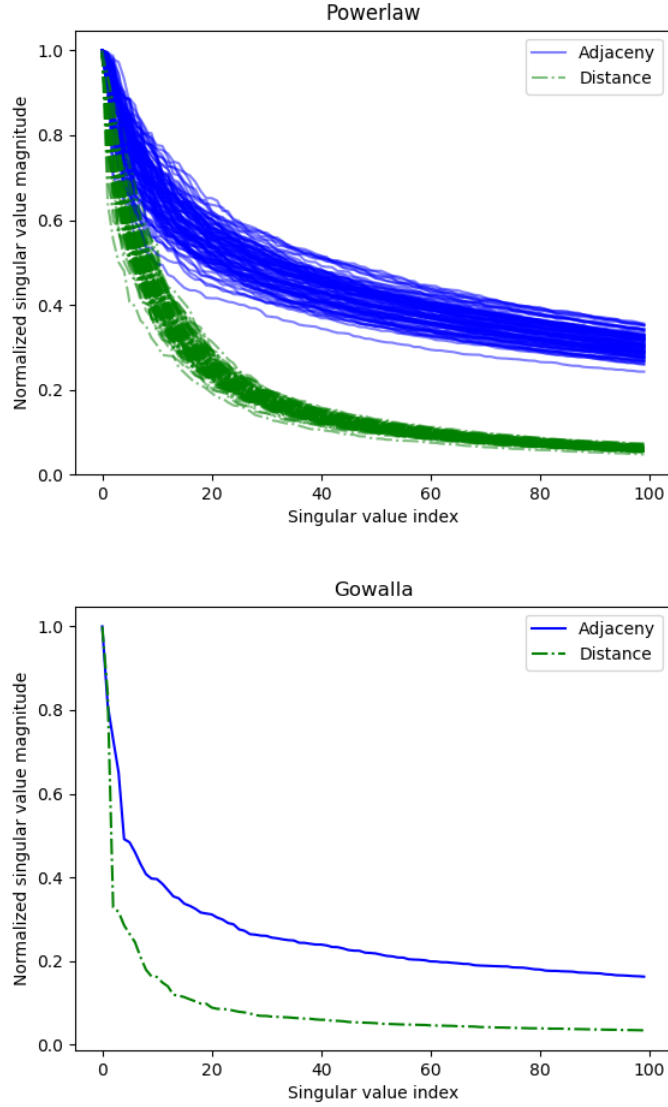


Fig. 1: Comparison of the normalized singular values of adjacency matrices versus hop-distance matrices for synthetic and real-world power-law graphs. Top figure is for a synthetic power-law graph generated using the Holme-Kim model [33] with 500 nodes and bottom figure is for 2000 node subgraph of the real-world Gowalla social network from [40].

In particular, one can equivalently think of P as a (non-principal) *sub-matrix* of the full hop-distance matrix H . If we decompose H into blocks by writing

$$H = \begin{bmatrix} A & B^T \\ B & C \end{bmatrix},$$

then $A \in \mathbb{N}_0^{M \times M}$ contains the hop-distances between the M anchors and themselves, $B \in \mathbb{N}_0^{(N-M) \times M}$ contains the hop-distances between the M anchors and the $N - M$ non-anchor nodes, and $C \in \mathbb{N}_0^{(N-M) \times (N-M)}$ contains the hop-distances between the $N - M$ non-anchor nodes and themselves, which in our case are missing entries. In this way, P can equivalently be written as

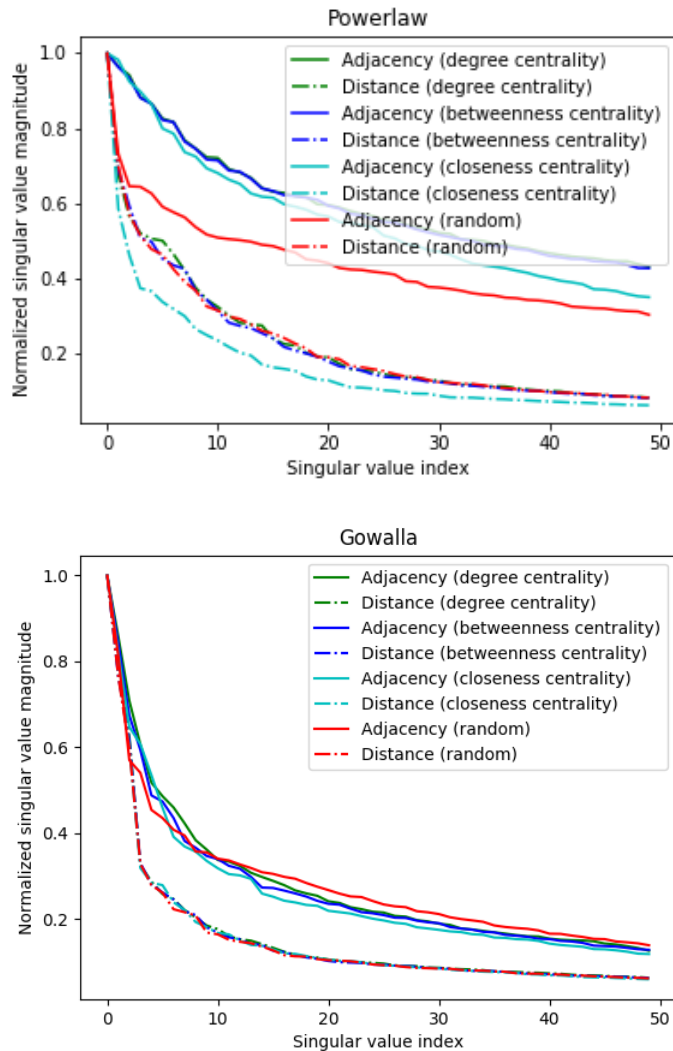


Fig. 2: Comparison of the normalized singular values of adjacency matrices versus hop-distance matrices for a variety of anchor selection strategies for 100 anchor nodes. Top figure is for a synthetic power-law graph generated using the Holme-Kim model [33] with 500 nodes. Bottom figure is for a subgraph of the real-world Gowalla social network from [40] with 2000 nodes.

$$P = \begin{bmatrix} A \\ B \end{bmatrix}.$$

Note, the prediction of both P and H from partial observations are of interest in the current context. Accordingly, in Section V, we will provide numerical results for both problems.

B. Low Rank Matrices

As can be inferred from our prior work [25], and as we further demonstrated above, hop distance matrices H for many interesting and realistic networks are, somewhat surprisingly, *approximately low-rank*. It is this empirical observation that inspires our work.

A widely used tool for analyzing low-rank structure in matrices is the Singular Value Decomposition (SVD) [29] and the closely related idea of Principal Component Analysis [8]. In particular, given our VC matrix $P \in \mathbb{R}^{N \times M}$ one can write P as

$$P = U \Sigma V^T$$

where, assuming $M < N$ as $U \in \mathbb{R}^{N \times M}$, $\Sigma \in \mathbb{R}^{M \times M}$, and $V \in \mathbb{R}^{M \times M}$. In addition, the columns of U and V are orthonormal (i.e., U and V are each sub-matrices of a *unitary* matrix) and Σ is diagonal. Finally, the diagonal entries of Σ are called the singular values of P and the rank of P is precisely the number of non-zero singular values.

Accordingly, one can compute an approximation of P by setting the “small” entries of Σ to 0, using an appropriate threshold, to generate an approximation $\Sigma \approx \hat{\Sigma}$. P can then be approximated similarly by setting $P \approx \hat{P} = U \hat{\Sigma} V^T$. Such ideas have a long history, with an important milestone being the seminal work of Eckart and Young in 1936 [29].

C. Topology Coordinates and Topology Preserving Maps

The mathematical foundation of our previous work in TPM generation from a VCS follows from the above formulation [25]. Consider the principle components of P given by,

$$P_{SVD} = U \Sigma$$

In the TC domain, each node in a 3-D network is characterized by a triple of Cartesian coordinates $(x_T(i), y_T(i), z_T(i))$. Let $[X_T, Y_T, Z_T]$ be the matrix of TCs for the entire set of nodes, i.e., the i -th row is the TCs of node i . Then from [25],

$$[X_T, Y_T, Z_T] = [P_{SVD}^{(2)}, P_{SVD}^{(3)}, P_{SVD}^{(4)}] \quad (4)$$

where, $P_{SVD}^{(j)}$ is the j -th column of P_{SVD} . Note, in the derivation of TCs as presented in [25] the first singular vector $P_{SVD}^{(1)}$ is *not used* in the representation. In Section IV-D, we will discuss the relationship between generating TCs without $P_{SVD}^{(1)}$ and the idea of “double centering” [39].

The importance of TCs is that they capture the geometric features such as the shape and boundaries in spite of the fact that no Euclidean distance measurements are used. However, if some physical locations are known, then the TCs can be transferred to approximate physical coordinates as well [11].

D. Connections to Non-linear Dimension Reduction (NDR)

The Topology Preserving Map (TPM) generation above is closely related to several algorithms in NDR [39]. In particular, given a *squared* EDM D , one can compute a “double centering” of D by writing

$$S = -\frac{1}{2} \left(D - \frac{1}{N} \mathbb{1} \mathbb{1}^T D - \frac{1}{N} D \mathbb{1} \mathbb{1}^T + \frac{1}{N^2} \mathbb{1} \mathbb{1}^T D \mathbb{1} \mathbb{1}^T \right) \quad (5)$$

where $\mathbb{1} \in \mathbb{R}^{n \times 1}$ is the vector all of whose entries are 1 [39]. In effect, $\frac{1}{N} \mathbb{1} \mathbb{1}^T D$ is the matrix which contains all of the column averages of D , $\frac{1}{N} D \mathbb{1} \mathbb{1}^T$ is the matrix containing all of the row averages of D , and $\frac{1}{N^2} \mathbb{1} \mathbb{1}^T D \mathbb{1} \mathbb{1}^T$ is the matrix containing the average of all the entries of D .

Note that a TPM is a low-rank approximation of a double centered S computed from H rather than a squared EDM D . In this sense, a TPM is analogous to the Multi-Dimensional Scaling (MDS) algorithm [37], [9], [39].

Even closer to our proposed technique, others have considered *geodesic* distance matrices D_G [53], [39] generated by drawing short range distances from D , say by using a proscribed number of neighbors or only considering distances below a certain threshold, but computing the rest of the distances by other means. In particular, one computes D_G by selecting some number of neighboring points for each point x (e.g., all of the points laying in some ϵ -ball around x) and then completing D_G by computing shortest-paths in the resulting weighted graph. A low-rank approximation of such a geodesic based distance matrix D (after double centering as in (5)) is equivalent to the *Isomap* [53], [39] algorithm for NDR.

Intuitively, Isomap can be thought of as a relaxation of MDS to the case where Euclidean distances are “trusted” for short range interactions, but not “trusted” for long range interactions. The long range interactions are instead approximated by geodesic distances, which are thought to be more faithful to the true geometry and topology of the network. Our method generalizes this argument by assuming that not even short range distances are to be “trusted” and instead our HDM is computed from unweighted connectivity information. Accordingly, our method uses geodesic distance matrix H analogous to D_G in Isomap. *However, all of our short range distances are presumed to be 1*, i.e., we use the number of hops.

E. Matrix Completion

Prior work on TPM generation is based on the case where entire columns are taken from H and used to construct P . However, in the current work, we consider the more interesting, and practically important case, where each anchor node only has a *partial* set of measurements to the rest of the network. Accordingly, some entries in P are *not observed* and the matrix P is therefore incomplete. Predicting the unobserved entries in P can be phrased as a low-rank *matrix completion* problem. In particular, we have leveraged modern ideas in low-rank *matrix completion* [41], [12], [13], [47], [48]. Space does not afford a fulsome treatment of the theory and implementation details for these algorithms. Accordingly, we merely endeavor to provide the reader with the intuition for such approaches in the context of predicting unobserved entries in HDMs.

The key idea of such methods can be phrased as the following optimization problem

$$L = \underset{L_0}{\operatorname{argmin}} \rho(L_0), \text{ s.t. } \mathcal{P}_\Omega(M) = \mathcal{P}_\Omega(L_0) \quad (6)$$

where M is an arbitrary matrix, ρ is the rank operator and \mathcal{P}_Ω is an operator that extracts from M the set of observed entries designated by Ω (i.e., the constraint in (6) is only enforced at the observed points). In other words, we seek to find a matrix L_0 such that the rank of L_0 (denoted $\rho(L_0)$) is minimized while enforcing the constraint that the matrix we construct matches our observed entries $\mathcal{P}_\Omega(M)$. Since, we enforce the constraint that $\mathcal{P}_\Omega(M) = \mathcal{P}_\Omega(L_0)$ the returned matrix L_0 will be faithful to our measured hop-distances but L_0 is free to take on any values it likes outside of Ω to minimize its rank.

Unfortunately, as stated, (6) is an NP-hard optimization problem, and can only be solved for small networks. Recent results [41], [12], [13], [47], [48] allow, under mild assumptions, for the NP-hard optimization in (6) to be recast as a convex optimization problem

$$L = \underset{L_0}{\operatorname{argmin}} \|L_0\|_*, \text{ s.t. } \mathcal{P}_\Omega(M) = \mathcal{P}_\Omega(L_0) \quad (7)$$

where $\|L_0\|_*$ sum of the singular values of L_0 , often called the *nuclear-norm* of L_0 . The optimization problem in (7) is *convex* and can easily be solved for millions of nodes on commodity computing hardware using splitting techniques and iterative matrix decomposition algorithms [41], [48], [47].

F. Completion of Partially Observed Hop-Distance Matrices

Simply stated, our proposed method for computing VCs from partially observed HDMs revolves around combining the NDR ideas from Section IV-D with the matrix completion ideas from Section IV-E. However, one impediment remains, namely, the double centering operation in (5), *prima facie*, would seem to require a fully observed matrix H , negating our ability to analyze partially observed HDM matrices.

However, this difficulty in computing a “double centering” of a partially observed P can be overcome by way of the following equation, similar to Equation (5),

$$S_{i,j} = -\frac{1}{2} (P_{i,j}^2 - \mu_j(P^2) - \mu_i(P^2) + \mu_{i,j}(P^2)) \quad (8)$$

where $\mu_j(P)$ is the mean of the observed entries in the j -th column of P , $\mu_i(P)$ is the mean of the observed entries in the i -th row of P , and $\mu_{i,j}(P)$ is the mean of all of the entries in H . In effect, each entry of the double-centered matrix $S_{i,j}$ only depends on the square of the single entry $P_{i,j}$, along with *mean values* of the *rows* and *columns* of P . Accordingly, estimates of these mean values can be computed even for a partially observed matrix such as $\mathcal{P}_\Omega(P)$, by performing the required mean over just the observed entries of the appropriate column, row, or the entire matrix. Of course, if a particular node has no measurements such a node cannot be predicted. Also, classically, (8) is defined for EDMs, and one might wonder if it is applicable to HDMs? In fact, that is a salient point of our work, to find a space in which the Euclidean-distances and the hop-distances coincide.

In some sense, such restrictions on acceptable sampling schemes for H should not be surprising. For example, if a particular row of P (or H) contains no measurements then the distance from this node to any anchor is not known. In effect, nothing is known about this node and therefore no predictions can be made.

Identifying classes of acceptable and unacceptable sampling schemes is a topic for future research. However, such ideas are a close cousin of the incoherence requirements that arise in matrix completion problems [41], [12], [13], [47]. Accordingly, drawing inspiration from that literature and for simplicity, we choose *random nodes as anchors* and *random nodes* whose distance we measure from each anchor.

Our algorithm for recovery of complete P from partial entries and generating TCS from a partially observed HDM is described in Procedure 1 below. Topology coordinates in [25] are given by 2nd and 3rd singular vectors in case of 2-D

networks and 2nd, 3rd and 4th in case of 3-D networks as given in Equation 4. Thus as a comparison we use Procedure 2, which carries out matrix completion directly on P , and follows the approach in [25] to generate TPMs. Procedure 1, based on double centering followed by the completion of the Gramian matrix S , however, follows the approaches such as MDS for NDR in Euclidean spaces more closely.

Procedure 1 Computing TPMs from a partially observed HDM matrix P via completion of Gramian matrix.

Input: Partially complete P of a graph $G = \{V, E\}$

Input: A target dimension k

- 1: **procedure** COMPLETE S FROM PARTIAL P
 - 2: Compute $\mathcal{P}_\Omega(S)$ from $\mathcal{P}_\Omega(P)$ using (8)
 - 3: Compute approximate Gramian matrix S from $\mathcal{P}_\Omega(S)$ using (7)
 - 4: Compute the SVD $S = U\Sigma V^T$
 - 5: **return** The first k columns of $U\Sigma$ are the TCs
 - 6: **end procedure**
-

Procedure 2 Computing TPMs from a partially observed HDM matrix P via matrix completion.

Input: Partially complete P of a graph $G = \{V, E\}$

Input: A target dimension k

- 1: **procedure** COMPLETE P FROM PARTIAL P
 - 2: Compute approximate distance vector matrix P from $\mathcal{P}_\Omega(P)$ using (7)
 - 3: Compute the SVD $P = U\Sigma V^T$
 - 4: **return** Columns 2 through $k+1$ (i.e., the first column is excluded) are the TCs
 - 5: **end procedure**
-

V. RESULTS

Here, we demonstrate the effectiveness of the proposed HDM based approach described in Procedure 1 in constructing accurate topology maps from a small set of hop-distances among node pairs. In particular, we demonstrate how P or H can be recovered from a set of partial observations, and how topological coordinates that arise from the eigen-decomposition of P provide accurate recovery of relationships among network nodes. The evaluation is carried out for a set of 2-D and 3-D sensor networks, and in Appendix also for a social network representing classes of networks where physical coordinates play no role. All the results have been averaged over 100 experimental iterations. As these results demonstrate, our methods provide surprisingly accurate predictions, even when only a tiny fraction of the network has been measured.

A. Recovery of Networks Embedded in 2-D/3-D Spaces

Four networks representative of 2-D and 3-D sensor network deployments covering a range of shapes and sizes are used for the evaluation. They contain complex features such as convex and concave boundaries and voids.

- A concave 2-D network with 550 nodes, the physical layout of which is shown in Figure 4(a) [25].
- A 2-D circular network with multiple circular voids of 496 nodes as shown in Figure 5(a) [25].
- A 3-D network, shown in Figure 6(a), consisting of 1640 nodes, which occupies a cube shaped volume with a hollow region in the shape of an hourglass devoid of nodes [34].
- A 3-D surface network, shown in Figure 7(a), consisting of 1245 nodes, which is comprised of two hollow cylinders joined in a “T” configuration [25].

As an initial exploration of the applicability of the techniques we propose, we first examine the rank of the VC matrices (P) for each of our networks. Twenty random anchors were selected in each case, i.e., $M = 20$, which corresponds to approximately 3.6%, 4%, 1.2%, and 1.6% of the nodes, respectively, for each of our four networks. The singular values of the full VC matrices of the four networks are shown in Figure 3. If we were considering EDMs, then the rank of the first two networks would be 4 (since they are embedded in \mathbb{R}^2) and the rank of the second two networks would be 5 (since they are embedded in \mathbb{R}^3) [31], [36]. As seen in Figure 3, the rank of the HDMs is certainly higher than their embedding dimensions would indicate, if they were

EDMs. However, somewhat surprisingly, even though the four networks are quite different, all of their ranks are substantially smaller than 20, for the chosen random anchors. Our interest is in the recovery of topological information and geometric relationships such as the general shapes of boundaries, voids in the networks, and node neighborhood preservation. Thus, the question is whether such information is preserved and can be extracted from small numbers of anchors and partial observations of P .

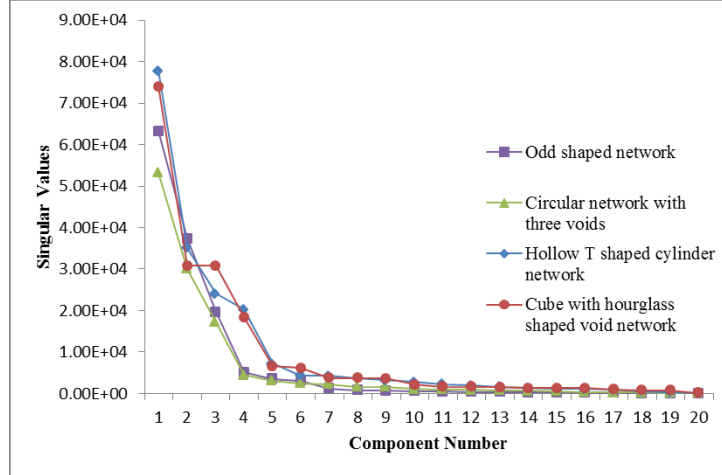


Fig. 3: Singular values of the VC matrix for Circular Network, Odd-Shaped Network, Hollow T cylinder network, and the 3-D network with void indicating the low-rankness of VCS data.

Two-dimensional TPMs extracted using the full set of VCs following [25] are shown in Figure 4(b), Figure 5(b), Figure 6(b), and Figure 7(b), respectively, for the four networks. It is important to note that even the full set of VCs corresponding to 20 anchors, which corresponds to 20 random columns of H , contains only approximately 3.6%, 4%, 1.2%, and 1.6% elements of the corresponding HDM.

Next, we randomly discard 10%, 20%, 40% and 60% respectively of this already small sample of the elements of H . The TPMs recovered using low-rank matrix completion followed by TPM extraction are shown in Figure 4(c-e), Figure 5(c-e), Figure 6(c-e), and Figure 7(c-e) for these various sub-samplings. The results indicate that accurate TPMs of networks are obtainable with only a fraction of virtual coordinates. It is important to recognize that the goal of this work is not to necessarily recover the maps in subfigures (a) of Figures 4-7. Rather, we wish to recover subfigures (b) (the fully observed topology map) from a sparse set of hop-distance observations.

To precisely quantify the error introduced to the TPM due to missing VCs, we define the mean error E as follows:

$$E = \frac{\left[\sum_{i=1}^N \sum_{j=1}^M |d_{ij}(f) - d_{ij}(0)| \right]}{\left[\sum_{i=1}^N \sum_{j=1}^M d_{ij}(0) \right]}, \quad (9)$$

where, $d_{ij}(f)$ refers to the Euclidean distance between nodes i and j on the TPM when f fraction of random anchor coordinates are missing. The percentage mean error with percentage of missing VCs for the four networks are shown in Figure 8. It is important to note that even when mean error is high, much of the local neighborhood and shape information is preserved.

B. Accuracy of 2-D Topology Preservation

The accuracy of neighborhood preservation of reconstructed topology maps of 2-D networks is evaluated below using the topology preservation error E_{TP} , which captures the degree to which the neighborhood relationships are altered. We provide only a brief explanation of E_{TP} and refer the reader to [25] for its precise definition. Consider the network in Figure 4(a). The network is scanned along the set of lines in the horizontal direction (\vec{H}) and the vertical direction (\vec{V}). Let α and β denote the sets of lines in \vec{H} and \vec{V} directions respectively. Consider one such line which contains the ordered set of some m nodes $\{n_1, \dots, n_m\}$, which we call the original placement. Now consider the projection of this set of nodes in a TPM (e.g.,

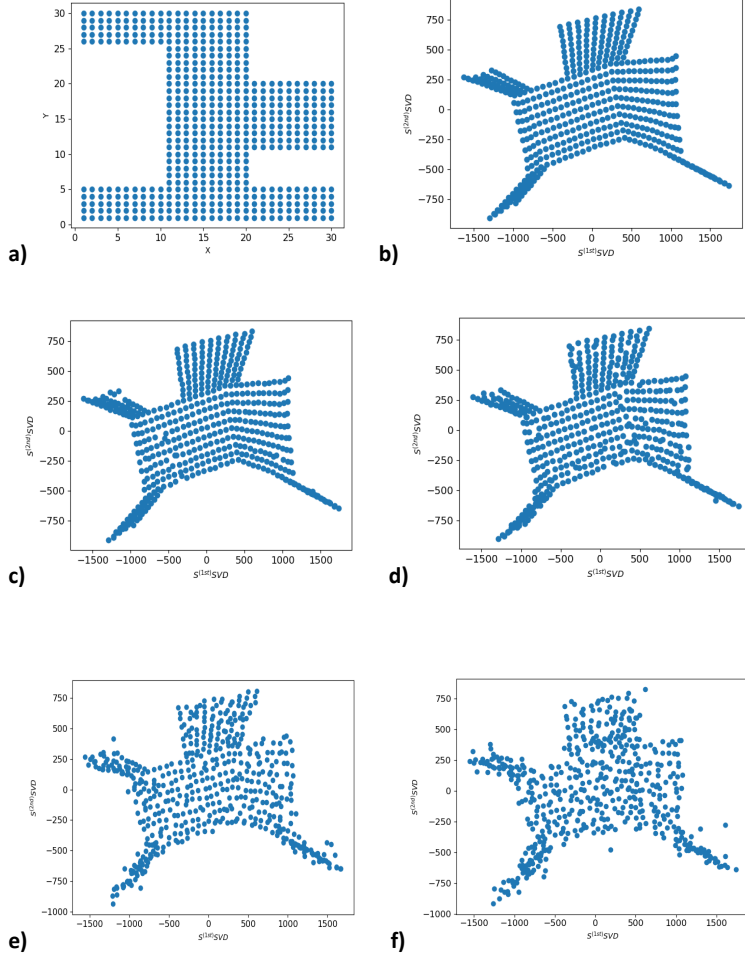


Fig. 4: Concave network: (a) Original layout, and (b) TPM recovered from full set of VCs with 20 random anchors; Recovered TPM with (c) 10%, (d) 20%, (e) 40%, and (f) 60% of sampled coordinates randomly discarded.

Figure 4(b)) on the corresponding axis. An error indicator function $I_{i,j}$ is defined by comparing the order of projected nodes with the original placement as:

$$I_{i,j} = \begin{cases} 1, & \text{nodes } i \text{ and } j \text{ are out of order,} \\ 0, & \text{nodes } i \text{ and } j \text{ are in the same order or } i=j. \end{cases} \quad (10)$$

For the line under consideration, the neighborhood preservation error is quantified by $\sum_{\forall i,j} (I_{i,j})/{}^m P_2$, where ${}^m P_2$ is permutation of m objects taken 2 at a time. However, in case of TPM, we are interested in the overall neighborhood preservation error E_{TP} over the set of lines in \vec{H} and \vec{V} directions, which is given by [25]:

$$E_{TP} = \left[\sum_{\alpha} \sum_{\forall i,j} (I_{i,j}) + \sum_{\beta} \sum_{\forall i,j} (I_{i,j}) \right] / \left[\sum_{\alpha} {}^m P_2 + \sum_{\beta} {}^m P_2 \right] \quad (11)$$

Next we evaluate the effectiveness of matrix completion with HDMs for TPM generation using the approaches in Procedure 1 and Procedure 2. E_{TP} (as a percentage) for the 2-D networks for these two cases provided in Table I show that with merely 20 random anchors, which corresponds to only 4% of the nodes, networks can be recovered with an error less than 5% even after deleting 80% of the entries.

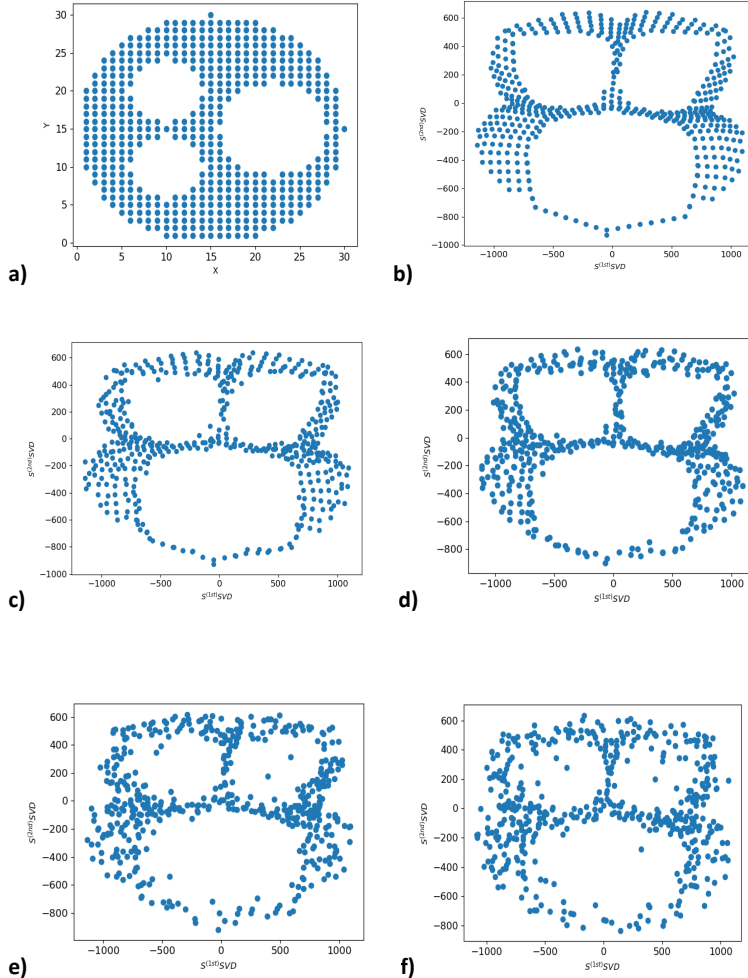


Fig. 5: Circular network: (a) Original layout, and (b) TPM recovered from full set of VCs with 20 random anchors; Recovered TPM with (c) 10%, (d) 20%, (e) 40%, and (f) 60% of sampled coordinates randomly discarded.

Deletions(%)	10%	20%	40%	60%	80%
$Circular_S$	3.31	3.55	3.91	4.25	4.29
$Circular_P$	2.47	2.92	3.73	4.04	4.30
$Concave_S$	3.09	3.47	3.48	4.43	4.50
$Concave_P$	2.94	3.23	3.65	3.82	3.94

TABLE I: Topology preservation errors ($E_{TP}\%$) for circular and concave networks with 10% to 80% of sampled coordinates randomly discarded. Subscript S denotes results with completion of Grammian matrix and taking 1st and 2nd singular vectors (Procedure 1), while subscript P denotes completion of the P matrix and taking 2nd and 3rd singular vectors (Procedure 2).

VI. CONCLUSION

This paper addresses the problem of recovering network features from a small set of hop-based graph geodesics. For networks deployed on 2-D surfaces and 3-D spaces the geometric and physical layout features are of importance. The approach starts with anchor-based VCs but, unlike prior techniques that required the entire VC set, the proposed approach requires only a fraction of the measured VCs to recover accurate topology preserving maps. Our technique is based on the theory of low-rank matrix completion that reconstructs missing VCs, the result of which is used to recover layout maps using topology preserving map generation techniques. The results presented here not only allow the reduction of cost (communication, power,

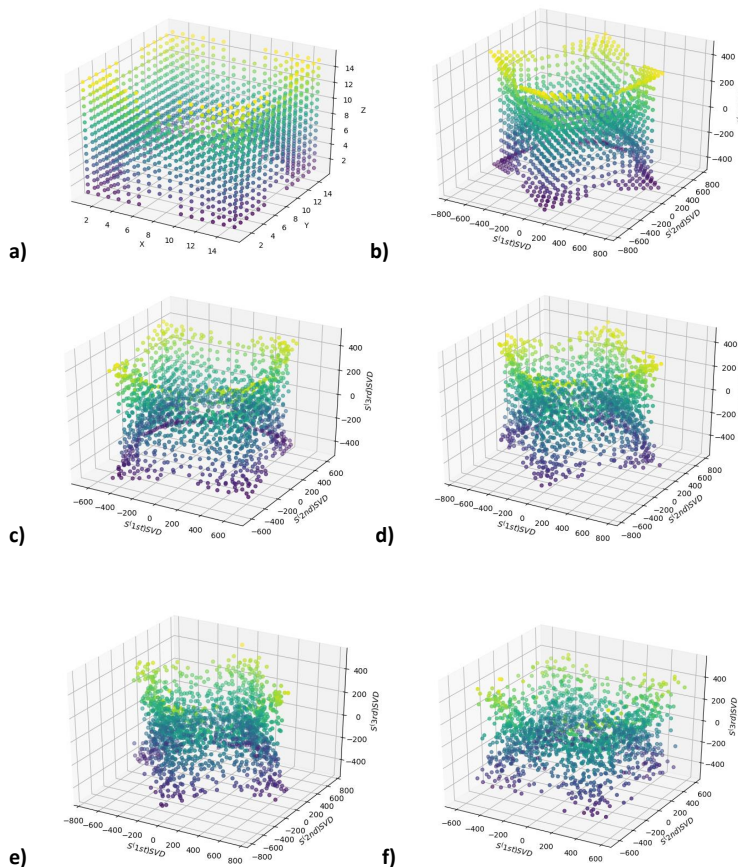


Fig. 6: Cube with hourglass shaped void: (a) Original layout, and (b) TPM recovered from full set of VCs with 20 random anchors; Recovered TPM with (c) 10%, (d) 20%, (e) 40%, and (f) 60% of sampled coordinates randomly discarded.

etc.) of VC generation but, more importantly, open the possibility of using topology coordinate based techniques for large networks and even those involving soft-state systems, where some coordinate values may be allowed to expire, thus allowing for more resilient network operations.

With random anchor selection, the VC matrix (P) is equivalent to a small set ($< 10\%$) of random columns of an HDM. Thus, the partially complete P matrix where a large set of random VC entries are missing is equivalent to an incomplete HDM with only a very small number of entries. Beyond results for physically embedded networks, we also demonstrated the ability to make accurate predictions in a large social network also from a small set of random geodesic measurements. We also demonstrated that the HDMs of many real-world networks are low-rank. Therefore, the approach presented here provides a foundation for designing novel graph sampling techniques that allows the capture of complex real-world networks with a small number of measurements.

REFERENCES

- [1] Network Science (NRC). *Report by Committee on Network Science for Future Army Applications*, 2005.
- [2] M. Abramowitz and I. A. Stegun. *Handbook of Mathematical Functions with Formulas, Graphs and Mathematical Tables*, page 300. 1964.
- [3] S. Al-Homidan and H. Wolkowicz. Approximate and exact completion problems for Euclidean distance matrices using semidefinite programming. *Linear Algebra and its Applications*, vol. 406, pages 109–141, Sep 2005.
- [4] A. Y. Alfakih. Graph rigidity via Euclidean distance matrices. *Linear Algebra and Its Applications*, vol. 310, pages 149–165, 2000.
- [5] M. Aouchiche and P. Hansen. Distance spectra of graphs: A survey, vol. 458, pages 301–386, Oct 2014.
- [6] A. Barabasi and R. Albert. Emergence of scaling in random networks. *Science (New York, N.Y.)*, vol. 286, pages 509–12, Oct 1999.
- [7] Z. Beerliova, F. Eberhard, T. Erlebach, A. Hall, M. Hoffmann, M. Mihal'ak, and L. S. Ram. Network discovery and verification. *IEEE Journal on Selected Areas in Communications*, vol. 24, no. 12, pages 2168–2181, Dec 2006.
- [8] C. M. Bishop. *Pattern Recognition and Machine Learning*. Springer New York, vol. 4, 2006.

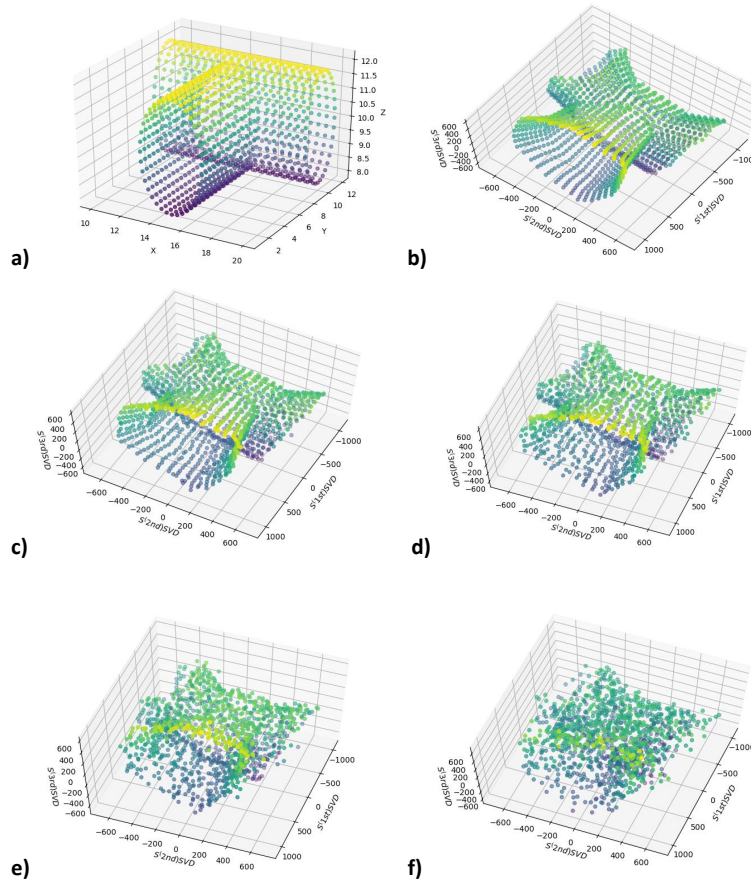


Fig. 7: Hollow T shaped Cylinder: (a) Original layout, and (b) TPM recovered from full set of VCs with 20 random anchors; Recovered TPM with (c) 10%, (d) 20%, (e) 40%, and (f) 60% of sampled coordinates randomly discarded.

- [9] I. Borg and P. Groenen. *Modern Multidimensional Scaling - Theory and Applications*. Springer New York, 2005.
- [10] T. Bouchoucha, C. N. Chuah, and Z. Ding. Finding link topology of large scale networks from anchored hop count reports. In *Proc. GLOBECOM - IEEE Global Communications Conference*, pages 1–6, Dec 2017.
- [11] A. F. Buoud and A. P. Jayasumana. Topology preserving map to physical map - a thin-plate spline based transform. In *Proc. IEEE 41st Conference on Local Computer Networks (LCN)*, pages 262–270, Nov 2016.
- [12] E. J. Candes and Y. Plan. Matrix completion with noise. *Proc. IEEE*, vol. 98, 2009.
- [13] E. J. Candes and B. Recht. Exact matrix completion via convex optimization. *Foundations of Computational Mathematics*, Dec 2009.
- [14] Q. Cao and T. Abdelzaher. Scalable logical coordinates framework for routing in wireless sensor networks. In *Proc. ACM Transactions on Sensor Networks*, vol. 2, pages 557–593, Nov 2006.
- [15] A. Caruso, S. Chessa, S. De, and A. Urpi. GPS free coordinate assignment and routing in wireless sensor networks. In *Proc. INFOCOM*, 2005.
- [16] S. A. Catanese, P. De Meo, E. Ferrara, G. Fiumara, and A. Proveti. Crawling Facebook for social network analysis purposes. In *Proc. International Conference on Web Intelligence, Mining and Semantics, WIMS '11*, pages 52:1–52:8. ACM, 2011.
- [17] S. Chakrabarti, M. van den Berg, and B. Dom. Focused crawling: a new approach to topic-specific web resource discovery. *Computer Networks*, vol. 31, pages 1623–1640, 1999.
- [18] S. Chatterjee. Matrix estimation by universal singular value thresholding. *Ann. Statist.*, vol. 43, pages 177–214, 2015.
- [19] K. Chou. Applications of graph theory to enzyme kinetics and protein folding kinetics: Steady and non-steady-state systems. *Biophysical Chemistry*, vol. 35, pages 1–24, 1990.
- [20] F. R. K. Chung. *Spectral Graph Theory (CBMS Regional Conference Series in Mathematics)*. American Mathematical Society, 1996.
- [21] D. C. Dhanapala and A. P. Jayasumana. Anchor selection and topology preserving maps in WSNs - A directional virtual coordinate based approach. In *Proc. IEEE Conference on Local Computer Networks (LCN)*, pages 571–579, 2011.
- [22] D. C. Dhanapala and A. P. Jayasumana. Directional virtual coordinate systems for wireless sensor networks. In *Proc. IEEE International Conference on Communications (ICC)*, pages 1–6, Jun 2011.
- [23] D. C. Dhanapala and A. P. Jayasumana. Geo-logical routing in wireless sensor networks. *Sensor, Mesh and Ad Hoc*, 2011.
- [24] D. C. Dhanapala and A. P. Jayasumana. Clueless nodes to network-cognizant smart nodes: Achieving network awareness in wireless sensor networks. In *Proc. IEEE Consumer Communications and Networking Conference (CCNC)*, pages 174–179, Jan 2012.

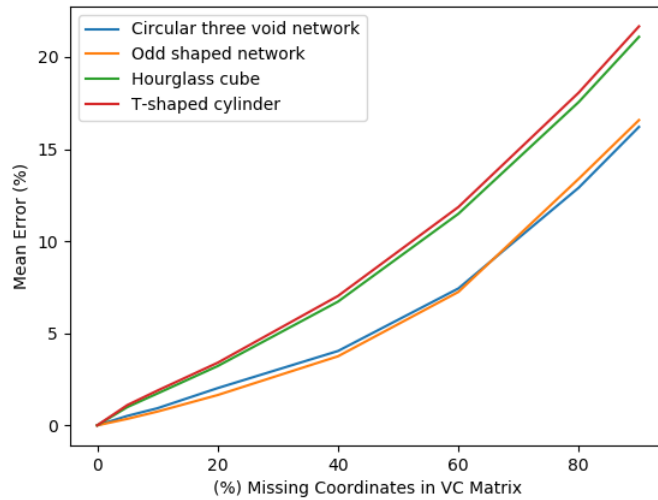


Fig. 8: Mean error [as defined in (9)] versus the percentage of missing virtual coordinates for sensor networks.

- [25] D. C. Dhanapala and A. P. Jayasumana. Topology preserving maps—Extracting layout maps of wireless sensor networks from virtual coordinates. *IEEE/ACM Transaction on Networking*, vol. 22, pages 784–797, 2014.
- [26] D. C. Dhanapala, A. P. Jayasumana, and S. Mehta. On boundary detection of 2-D and 3-D wireless sensor networks. *Proc. IEEE GLOBECOM*, pages 1–5, 2011.
- [27] D.C. Dhanapala and A.P. Jayasumana. CSR: Convex subspace routing protocol for wsns. In *Proc. IEEE Conference on Local Computer Networks (LCN)*, 2009.
- [28] R. Diestel. *Graph Theory*. Springer Verlag, vol. 173, 2005.
- [29] C. Eckart and G. Young. The approximation of one matrix by another of lower rank. *Psychometrika*, vol. 1, pages 211–218, 1936.
- [30] D. Gamarnik and M. Sviridenko. Hamiltonian completions of sparse random graphs. *Discrete Applied Mathematics*, vol. 152, pages 139–158, 2005.
- [31] J. C. Gower. Euclidean distance geometry. *Mathematical Scientist*, vol. 7, pages 1–14, 1982.
- [32] I. Gutman and B. Borovicanin. Nullity of graphs: an updated survey. *Selected topics on applications of graph spectra, Math. Inst., Belgrade*, pages 137–154, 2011.
- [33] P. Holme and B. J. Kim. Growing scale-free networks with tunable clustering. *Physical Review E - Statistical, Nonlinear, and Soft Matter Physics*, 65, 2002.
- [34] A. P. Jayasumana, R. P. Paffenroth, and S. Ramasamy. Topology maps and distance-free localization from partial virtual coordinates for IoT networks. In *Proc. IEEE International Conference on Communications (ICC)*, pages 1–6, May 2016.
- [35] G. Kossinets. Effects of missing data in social networks. *Social Networks*, vol. 28, pages 247–268, 2006.
- [36] N. Krislock and H. Wolkowicz. *Euclidean Distance Matrices and Applications*. Springer, Sep 2011.
- [37] J. B. Kruskal. Multidimensional scaling by optimizing goodness of fit to a nonmetric hypothesis. *Psychometrika*, vol. 29, pages 1–27, 1964.
- [38] T. K. Le and N. Ono. Closed-form and near closed-form solutions for TDOA-based joint source and sensor localization. *IEEE Transactions on Signal Processing*, vol. 65, pages 1207–1221, Mar 2017.
- [39] J. A. Lee and M. Verleysen. *Nonlinear Dimensionality Reduction*. Springer, 2007.
- [40] J Leskovec and A Krevl. SNAP Datasets: Stanford Large Network Dataset Collection. Jun 2014.
- [41] Z. Lin, M. Chen, and Y. Ma. The augmented Lagrange multiplier method for exact recovery of corrupted low-rank matrices. *arXiv:1009.5055*, page 23, 2013.
- [42] K. Liu and N. Abu-Ghazaleh. Aligned virtual coordinates for greedy routing in WSNs. In *Proc. IEEE International Conference on Mobile Ad Hoc and Sensor Systems (MASS)*, pages 377–386, 2007.
- [43] K. Mehlhorn and P. Sanders. *Algorithms and Data Structures: The Basic Toolbox*. 2008.
- [44] B. Mishra, G. Meyer, and R. Sepulchre. Low-rank optimization for distance matrix completion. *Proc. IEEE Conference on Decision and Control*, pages 4455–4460, 2011.
- [45] A. Mislove, M. Marcon, K. P. Gummadi, P. Druschel, and B. Bhattacharjee. Measurement and analysis of online social networks. In *Proc. 7th ACM SIGCOMM Conference on Internet Measurement, IMC '07*, pages 29–42. ACM, 2007.
- [46] H. Oliveira, E. F. Nakamura and A.F. Loureiro, and A. Boukerche. Error analysis of localization systems for sensor networks. In *Proc. 13th annual ACM International Workshop on Geographic Information Systems*, pages 71–78. ACM, 2005.
- [47] R. C. Paffenroth, R. Nong, and P. C. Du Toit. On covariance structure in noisy, big data. *Proc. SPIE*, vol. 8857, 2013.
- [48] R. C. Paffenroth, P. Du Toit, R. Nong, L. L. Scharf, A. P. Jayasumana, and V. Bandara. Space-time signal processing for distributed pattern detection in sensor networks. *IEEE Journal of Selected Topics in Signal Processing*, vol. 7, pages 38–49, 2013.
- [49] A. Parra and P. Scheffler. Characterizations and algorithmic applications of chordal graph embeddings. *Discrete Applied Mathematics*, vol. 79, pages 171–188, 1997.
- [50] R. Pech, L. Pan D. Hao, and T. Zhou H. Cheng. Link prediction via matrix completion. *EPL (Europhysics Letters)*, vol. 117, no. 2, 2017.

- [51] A. Rao, S. Ratnasamy, C. Papadimitriou, S. Shenker, and I. Stoica. Geographic routing without location information. In *Proc. 9th annual international conference on Mobile computing and networking*, pages 96–108. ACM, 2003.
- [52] Y. Shin T. L. N. Nguyen. Matrix completion optimization for localization in wireless sensor networks for intelligent IoT. *Sensors*, vol. 16, 2016.
- [53] J. B. Tenenbaum, V. De Silva, and J. C. Langford. A global geometric framework for nonlinear dimensionality reduction. *Science*, vol. 290, pages 2319–2323, 2000.
- [54] J. T. L. Wang, M. J. Zaki, H. T. T. Toivonen, and D. Shasha. *Introduction to Data Mining in Bioinformatics*, pages 3–8. Springer London, 2005.
- [55] Q. Wu, C. L. Lu, and R. C. T. Lee. The approximability of the weighted Hamiltonian path completion problem on a tree. *Theoretical Computer Science*, vol. 341, pages 385–397, 2005.
- [56] Q. Cheng Z. Kang, C. Peng. Top-N recommender system via matrix completion. *Proc. 13th AAAI Conference on Artificial Intelligence*, 2016.
- [57] J. Zhang and I. C. Paschalidis. An improved composite hypothesis test for Markov models with applications in network anomaly detection. In *Proc. 54th IEEE Conference on Decision and Control (CDC)*, pages 3810–3815, Dec 2015.
- [58] J. Zhang and I. C. Paschalidis. Statistical anomaly detection via composite hypothesis testing for Markov models. *IEEE Transactions on Signal Processing*, vol. 66, pages 589–602, Feb 2018.

APPENDIX

Note: The material in this Appendix is included for reviewer information only. It will be included in the supplementary material in the final version.

A. Table of Acronyms

Acronym	Expansion
EDM	Euclidean Distance Matrix
HDM	Hop Distance Matrix
MDS	Multi-Dimensional Scaling
NDR	Nonlinear Dimension Reduction
PCA	Principal Component Analysis
SVD	Singular Value Decomposition
TC	Topology Coordinates
TPM	Topology Preserving Map
VC	Virtual Coordinates
VCS	Virtual Coordinate System

B. Recovery of a Social Network

Gowalla is a mobile web based social media application for which the connectivity dataset is available from the Stanford Large Network Dataset Collection [40]. A subnetwork of 10,000 nodes with 140,866 edges among them is used for the evaluation. As opposed to the results in Section V-A, for this experiment, we assume that a random set of elements of H is known. This is equivalent to making every node of H an anchor node. *However, each anchor node only measures its distance to a small set of other nodes.*

As a Euclidean distance metric, such as in (9), is not applicable in the case of a social network, we use mean error E_m corresponding to the percentage error in the prediction of H , and the absolute hop-error in geodesic lengths E_a is defined as follows:

$$E_m = \left[\sum_{i,j=1}^{N,N} |\hat{h}_{ij}(f) - h_{ij}| \right] / \left[\sum_{i,j=1}^{N,N} h_{ij} \right], \quad (12)$$

$$E_a = \left[\sum_{i,j=1}^{N,N} |\hat{h}_{ij}(f) - h_{ij}| \right] / N^2, \quad (13)$$

where, $\hat{h}_{ij}(f)$ refers to the element ij of estimated HDM (\hat{H}) when f fraction of random elements are missing. Note that $\hat{h}_{ij}(0) = h_{ij}$.

The mean error and absolute hop distance error for different percentages of missing elements are shown in Figure 9. Even with 20% of the elements of H , i.e., 80% of coordinates missing, the network can be recovered with an error of approximately 6%, while the absolute hop error is less than 1.

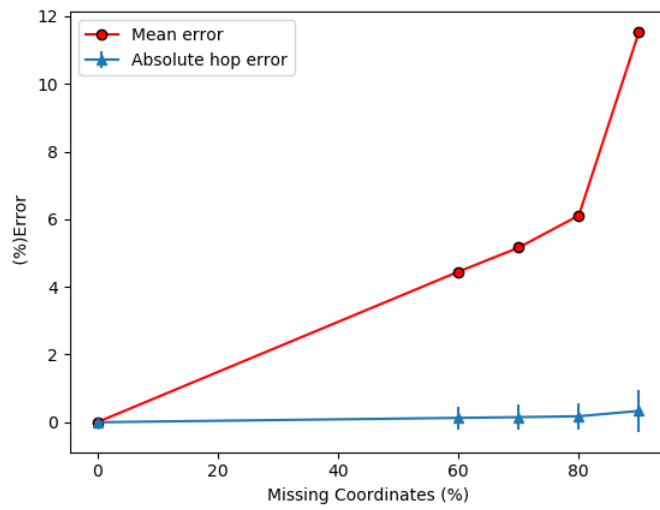


Fig. 9: Mean error and absolute hop distance error as defined in equations (12) (in red) and (13) (in blue) versus the percentage of missing elements in H for Gowalla social network.