



DIGITAL ACCESS TO SCHOLARSHIP AT HARVARD

Tree Preserving Embedding

The Harvard community has made this article openly available.
[Please share](#) how this access benefits you. Your story matters.

Citation	Shieh, Albert D., Tatsunori B. Hashimoto, and Edoardo M. Airoidi. 2011. Tree Preserving Embedding. Proceedings of the National Academy of Sciences 108, no. 41: 16916–16921.
Published Version	doi:10.1073/pnas.1018393108
Accessed	February 16, 2015 5:23:01 PM EST
Citable Link	http://nrs.harvard.edu/urn-3:HUL.InstRepos:12724040
Terms of Use	This article was downloaded from Harvard University's DASH repository, and is made available under the terms and conditions applicable to Other Posted Material, as set forth at http://nrs.harvard.edu/urn-3:HUL.InstRepos:dash.current.terms-of-use#LAA

(Article begins on next page)

Tree preserving embedding

Albert D. Shieh *, Tatsunori B. Hashimoto *, and Edoardo M. Airolidi *

*Department of Statistics & Faculty of Arts and Sciences Center for Systems Biology, Harvard University, Cambridge, MA 02138

Submitted to Proceedings of the National Academy of Sciences of the United States of America

The goal of dimensionality reduction is to embed high-dimensional data in a low-dimensional space while preserving structure in the data relevant to exploratory data analysis such as clusters. However, existing dimensionality reduction methods often either fail to separate clusters due to the crowding problem or can only separate clusters at a single resolution. We develop a new approach to dimensionality reduction, tree preserving embedding (TPE). Our approach uses the topological notion of connectedness to separate clusters at all resolutions. We provide a formal guarantee of cluster separation for our approach that holds for finite samples. Our approach requires no parameters and can handle general types of data, making it easy to use in practice and suggesting new strategies for robust data visualization.

dimensionality reduction | multidimensional scaling | hierarchical clustering

Visualization is an important first step in the analysis of high-dimensional data [1]. High-dimensional data often has low intrinsic dimensionality, making it possible to embed the data in a low-dimensional space while preserving much of its structure [2]. However, it is rarely possible to preserve all types of structure in the embedding. Therefore, dimensionality reduction methods can only aim to preserve particular types of structure. Linear methods such as principal component analysis (PCA) [3] and classical multidimensional scaling (MDS) [4, 5, 6] preserve global distances, while non-linear methods such as manifold learning methods [7, 8, 9] preserve local distances defined by kernels or neighborhood graphs. However, most dimensionality reduction methods fail to preserve clusters [10], which are often of greatest interest.

Clusters are difficult to preserve in embeddings due to the so-called crowding problem [11]. When the intrinsic dimensionality of the data exceeds the embedding dimensionality, there is not enough space in the embedding to allow clusters to separate. Therefore, clusters are forced to collapse on top of each other in the embedding. As the embedding dimensionality increases, there is more space in the embedding for clusters to separate and the crowding problem disappears, making it possible to preserve clusters exactly [12]. However, since the embedding dimensionality is at most two or three for visualization purposes, the crowding problem is prevalent in practice. When the clusters are known, they can be used to guide the embedding to avoid the crowding problem [13]. However, the embedding is often used to help find the clusters in the first place. Therefore, it is important to solve the crowding problem without knowledge of the clusters.

Force-based methods such as stochastic neighbor embedding (SNE) [14], variants of SNE [15, 11, 16, 10], and local MDS [17], have been proposed to overcome the crowding problem. Force-based methods use attractive forces to pull together similar points and repulsive forces to push apart dissimilar points. SNE and its variants use forces based on kernels, while local MDS uses forces based on neighborhood graphs. Force-based methods have long been used in graph drawing to separate clusters [18, 19]. Although force-based methods are effective, it is difficult to balance the relative strength of attractive and repulsive forces. When repulsive forces are too weak, they will fail to separate clusters, but when repulsive forces are too strong, they will artificially create clusters. Therefore, force-based methods are sensitive to

intrinsic *resolution* parameters such as kernel bandwidths and neighborhood graph sizes that control the amount of separation between points in the embedding.

We introduce tree preserving embedding (TPE) to overcome the limitations of force-based methods. TPE aims to preserve both distances and clusters by preserving the single linkage (SL) dendrogram in the embedding. SL is a hierarchical clustering method that iteratively merges pairs of clusters with minimum nearest neighbor distance. The SL dendrogram is the associated tree with the clusters as vertices and the merge distances as vertex heights. TPE preserves the SL dendrogram in the sense that SL generates the same dendrogram from both the data and the embedding. Embeddings and dendrograms have long been used as complementary representations for dissimilarities [20]. However, there is no guarantee that embeddings and dendrograms will be consistent when used separately. In particular, clusters found by dendrograms may not be found in embeddings due to the crowding problem. TPE combines embeddings and dendrograms in a common representation.

Preserving the SL dendrogram in the embedding is a natural choice for several reasons. First, the SL dendrogram is the only dendrogram consistent with the minimum spanning tree (MST) in the sense that the SL dendrograms are the same when the MSTs are the same [21, 22]. Preserving the topologies of neighborhood graphs has been shown to help overcome the crowding problem [23]. However, while the topologies of neighborhood graphs such as the MST can only be preserved *approximately* in general [24], we show that the SL dendrogram can be preserved *exactly*. Second, the SL dendrogram represents both global and local structure due to its hierarchical nature. Preserving global structure allows TPE to separate clusters, while preserving local structure prevents TPE from artificially creating clusters. Finally, TPE can separate clusters even when the SL dendrogram cannot. Although SL is often criticized as a clustering method for finding poor clusters in practice [25, 26], SL finds poor clusters due to the instability of cutting the SL dendrogram at a particular height [27]. Since TPE preserves the SL dendrogram at all heights, TPE is not sensitive to the instabilities of the SL dendrogram at any particular height.

We make cluster separation in TPE precise using the topological notion of connectedness [25]. A natural and commonly used notion of a cluster is a set of points that are connected at a particular resolution. It is well known that the SL dendrogram finds clusters of connected points at different resolutions for different heights [25]. We show that TPE preserves

Reserved for Publication Footnotes

connectedness in the sense that points in the embedding are connected if and only if they are connected in the data. Preserving connectedness guarantees that clusters separated in the data remain separated in the embedding. Therefore, TPE is guaranteed to separate clusters at all resolutions rather than just a single resolution.

Methods

In this section, we introduce TPE as an optimization problem subject to a set of constraints that preserve the SL dendrogram in the embedding. The constraints arise from a characterization of the SL dendrogram using a notion of connectedness. We introduce an algorithm similar to hierarchical clustering to implement TPE. In order to make the algorithm practical, we propose a variant based on a greedy approximation. Finally, we show that TPE preserves connectedness in a precise sense that corresponds well with separating clusters.

Algorithm. TPE is based on the framework of MDS [28]. Given a real, symmetric, non-negative, zero diagonal $n \times n$ dissimilarity matrix D for a set of n objects $S = \{1, \dots, n\}$ in a high-dimensional space, MDS finds a p -dimensional Euclidean embedding $X = \{x_1, \dots, x_n\} \subseteq \mathbb{R}^p$ of the objects that minimizes a loss function such as stress

$$\sigma(X) = \sum_{x_i, x_j \in X} (d_{i,j}(X) - D_{i,j})^2,$$

the sum of squared errors between the Euclidean distances $d_{i,j}(X) = \|x_i - x_j\|$ in the embedding and the corresponding dissimilarities $D_{i,j}$. Since loss functions such as stress emphasize approximating large dissimilarities well, minimizing them without constraints on the embedding leads to the crowding problem. TPE preserves the SL dendrogram in the embedding in order to overcome the crowding problem.

SL is a hierarchical clustering method that iteratively merges pairs of clusters $A, B \subseteq S$ with minimum nearest neighbor distance

$$\Delta(A, B) = \min_{i \in A, j \in B} D_{i,j},$$

starting with the n singleton clusters and ending with the trivial cluster. The SL dendrogram is the associated binary tree of depth $n - 1$ with singleton clusters as leaf vertices, the trivial cluster as the root vertex, merged clusters as internal vertices, and merge distances as vertex heights. For an example, see Fig. 1. There are many equivalent characterizations of the SL dendrogram [26]. We use the following notion of connectedness to express the SL dendrogram as a set of constraints on pairs of both objects and points.

Definition 1. *Objects $i, j \in S$ are ε -connected if there exists a path of objects $\alpha_1 = i, \dots, \alpha_m = j \in S$ such that $D_{\alpha_l, \alpha_{l+1}} \leq \varepsilon$ for $l = 1, \dots, m - 1$.*

Definition 2. *Points $x_i, x_j \in X$ are ε -connected if there exists a path of points $x_{\alpha_1} = x_i, \dots, x_{\alpha_m} = x_j \in X$ such that $d_{\alpha_l, \alpha_{l+1}}(X) \leq \varepsilon$ for $l = 1, \dots, m - 1$.*

Intuitively, objects are connected if there exists a path with short hops between them. The SL dendrogram contains the paths with short hops between objects. Objects are ε -connected if there exists a path of vertices with heights at most ε between their associated leaf vertices, or singleton clusters, in the SL dendrogram. Therefore, cutting the SL dendrogram at a height of ε produces clusters of ε -connected objects [25]. The relationship between the SL dendrogram and connectedness in an embedding is illustrated in Fig. 1. The merge

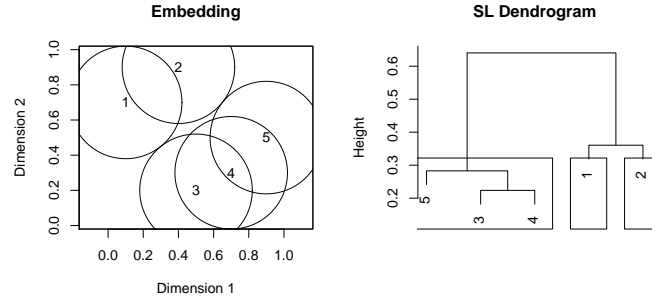


Fig. 1. Relationship between the SL dendrogram and connectedness in an embedding. Cutting the SL dendrogram at a height of $\varepsilon = 0.3$ produces three clusters of ε -connected points. Points 3 and 5 are ε -connected by point 4 because the ε -ball centered at point 4 contains points 3 and 5, while points 1 and 2 are not ε -connected to any other points because they are not contained by any ε -balls centered at other points.

distances of clusters in the SL dendrogram determine the connectedness of clusters in the embedding.

Cluster merges connect objects in the SL dendrogram. The ultrametric distance between objects in a dendrogram is the distance at which they are merged into the same cluster [29]. The ultrametric distance in the SL dendrogram is equivalent to the maximal sub-dominant ultrametric distance

$$L_{i,j} = \min_{\alpha_1=i, \dots, \alpha_m=j \in S} \max_{l=1}^{m-1} D_{\alpha_l, \alpha_{l+1}},$$

the maximal hop in a minimal path between objects [30], reminiscent of commute times in a graph [31]. Since the paths in the MST are minimal paths, the SL dendrogram can be constructed efficiently from the MST in practice [21]. However, it is important to emphasize that the paths in the MST are not the only possible minimal paths.

The SL dendrogram can be characterized by two constraints on each pair of objects. First, each pair of objects $i, j \in S$ must be $L_{i,j}$ -connected by their ultrametric distance $L_{i,j}$. Second, each pair of objects $i, j \in S$ cannot be ε -connected by any distance ε less than their ultrametric distance $L_{i,j}$. TPE uses the constraints on pairs of objects as constraints on corresponding pairs of points in the embedding. The first constraint guarantees that clusters in the embedding are merged at the same distances as corresponding clusters in the data. The second constraint guarantees that clusters in the embedding are merged in the same order as corresponding clusters in the data.

The following algorithm implements TPE.

1. Initialize the indices of available clusters

$$I_1 = \{1, \dots, n\},$$

the indices of the singleton clusters

$$S_1 = \{1\}, \dots, S_n = \{n\},$$

the embeddings of the singleton clusters

$$X_1 = \{x_1 = 0\}, \dots, X_n = \{x_n = 0\},$$

and the ultrametric distances for the singleton clusters

$$L_{1,1,1} = 0, \dots, L_{n,n,n} = 0$$

where $L_{c,i,j} = L_{i,j}$ denotes the ultrametric distance between objects $i, j \in S_c$ contained in cluster c .

2. For each iteration $k = 1, \dots, n - 1$:

(a) Find the next cluster merge

$$a_k, b_k = \arg \min_{a, b \in I_k: a \neq b} \Delta(S_a, S_b)$$

at a merge distance of

$$\Delta_k = \Delta(S_{a_k}, S_{b_k}).$$

(b) Merge the clusters

$$S_{n+k} = S_{a_k} \cup S_{b_k}$$

and update the indices of available clusters

$$I_{k+1} = \{i \in I_k : i \neq a_k, b_k\} \cup \{n + k\}.$$

(c) Find the ultrametric distances for the merged cluster

$$L_{n+k, i, j} = \begin{cases} L_{a_k, i, j} & \text{if } i, j \in S_{a_k} \\ L_{b_k, i, j} & \text{if } i, j \in S_{b_k} \\ \Delta_k & \text{otherwise} \end{cases} \quad \forall i, j \in S_{n+k}.$$

(d) Embed the merged cluster

$$\begin{aligned} X_{n+k} = & \arg \min_{X = \{x_i \in \mathbb{R}^p : i \in S_{n+k}\}} \sigma(X) \\ \text{s.t. } & x_i, x_j \text{ are } L_{n+k, i, j}\text{-connected} \quad \forall i, j \in S_{n+k}, \\ & d_{i, j}(X) \geq L_{n+k, i, j} \quad \forall i, j \in S_{n+k}. \end{aligned} \quad [1]$$

3. Return the embedding X_{2n-1} .

The algorithm proceeds similarly to hierarchical clustering. There are $n - 1$ iterations, one for each depth of the SL dendrogram. At each iteration, a pair of clusters is merged and the merged cluster is embedded by minimizing stress subject to the connectedness constraints. At the last iteration, the trivial cluster is embedded and returned. The number of objects being embedded changes at each iteration depending on the size of the merged cluster. Since the embeddings at each iteration are independent, only the embedding at the last iteration is needed. However, earlier embeddings can be used to help initialize later embeddings in practice.

The connectedness constraints can always be fulfilled. In fact, it is trivial to find an embedding that fulfills the connectedness constraints.

Theorem 1. *For each iteration $k = 1, \dots, n - 1$, there is a feasible solution to the optimization problem 1.*

The main difficulty in TPE is minimizing stress. The connectedness constraints may appear to be too rigid to allow TPE to find a low stress embedding since each pair of objects can be connected by an arbitrary path of objects. However, the connectedness constraints do not specify that the paths connecting pairs of points must be the same as the paths connecting corresponding pairs of objects. Preserving the paths would preserve the MST, which is not possible in general [24]. Moreover, the flexibility in choosing the paths allows points to move more freely in the embedding to lower stress. However, this flexibility comes at a cost. Due to the combinatorial nature of choosing paths, the optimization problem 1 is computationally intractable. Nevertheless, we can obtain a computationally tractable approximation by restricting the types of paths that can be chosen.

Greedy Approximation. The connectedness constraints allow all points in a merged cluster to be rearranged in the embedding at each iteration. However, since each cluster being merged has already been embedded in prior iterations, it is wasteful to change the prior embeddings of the clusters. The connectedness constraints within the clusters are already fulfilled in their prior embeddings. Since connectedness is preserved under rigid transformations (rotations, translations, and reflections), if we restrict the placement of the clusters to rigid transformations, then we only need to fulfill the connectedness constraints between them. The paths that fulfill the connectedness constraints between the clusters must pass through their nearest neighbors. Therefore, placing the clusters exactly their merge distance apart fulfills the connectedness constraints between them. The remaining flexibility in placing the clusters can be used to minimize the stress between them.

In place of the optimization problem 1, the greedy approximation proceeds at iteration k as follows.

1. Find a rigid transformation that aligns the clusters while keeping them separated by exactly their merge distance

$$\begin{aligned} T^* = & \arg \min_{T \in E(p)} \sum_{x_i \in X_{a_k}, x_j \in X_{b_k}} (d_{i, j}(T) - D_{i, j})^2 \\ \text{s.t. } & \min_{x_i \in X_{a_k}, x_j \in X_{b_k}} d_{i, j}(T) = \Delta_k \end{aligned} \quad [2]$$

where $d_{i, j}(T) = \|T(x_i) - x_j\|$ is the Euclidean distance in the embedding after alignment and $E(p)$ is the set of p -dimensional rigid transformations.

2. Align the clusters

$$x_i = T^*(x_i) \quad \forall x_i \in X_{a_k}.$$

3. Return the merged cluster

$$X_{n+k} = X_{a_k} \cup X_{b_k}.$$

The alignment of the clusters in the greedy approximation is illustrated in Fig. 2. The greedy approximation is reminiscent of Procrustes analysis [32], which has been used to merge embeddings of clusters [33]. However, Procrustes analysis aligns clusters without constraints on the embedding, making it sensitive to the crowding problem. In contrast, the greedy approximation has a constraint that keeps the clusters separated in the embedding in order to preserve the SL dendrogram. The constraint makes the optimization problem

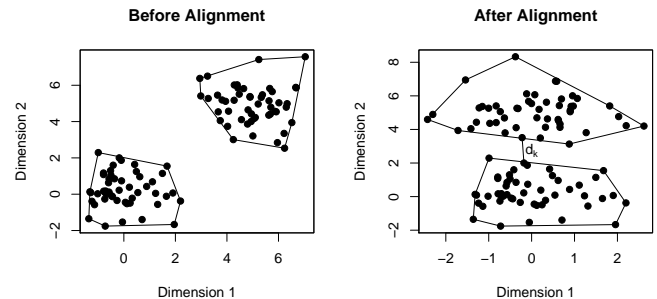


Fig. 2. Alignment of the clusters in the greedy approximation at iteration k . Before alignment, the clusters are placed far apart with high stress. During alignment, one of the clusters is moved using a rigid transformation to minimize the stress between the clusters subject to the constraint that the clusters are placed exactly their merge distance of d_k apart. After alignment, the clusters are placed close together with low stress.

2 more difficult to solve than standard Procrustes problems. Nevertheless, the optimization problem 2 can be solved efficiently in practice. Details of the greedy approximation can be found in the SI Text.

The efficiency of the greedy approximation comes at the cost of sensitivity to the merge order of the SL dendrogram. Since the greedy optimization cannot change the shapes of the clusters from prior embeddings, it cannot always align the clusters well. For small cluster merges, the prior embeddings will have little effect on the alignment. However, for large cluster merges, there may not be enough empty space in the prior embeddings to allow the clusters to be aligned well. While there are no formal guarantees for the performance of the greedy approximation, we found that good alignments of the clusters can typically be found in practice.

The greedy approximation has a time complexity of $O(n^3)$ since there are $O(n)$ iterations, each of which requires $O(1)$ evaluations of the $O(n^2)$ stress between the clusters in order to find an alignment of the clusters. The greedy approximation has a comparable time complexity to many dimensionality reduction methods, including those based on a spectral decomposition such as classical MDS and many manifold learning methods. While the cubic time complexity of TPE may be prohibitive for some applications, methods developed to improve the scalability of MDS such as landmark points [34] can be extended to TPE in principle.

Connectedness. TPE preserves clusters at all resolutions rather than just a single resolution due to the hierarchical nature of the SL dendrogram. Since the clusters found by cutting the SL dendrogram at different heights can be characterized by connectedness at different resolutions, TPE preserves connectedness in the following sense.

Theorem 2. *For any $\varepsilon > 0$, points $x_i, x_j \in X$ are ε -connected if and only if objects $i, j \in S$ are ε -connected.*

TPE preserves connectedness in the sense that points are connected in the embedding if and only if they are connected in the data. Clusters at any resolution can be neither too close together nor too far apart in the embedding without violating connectedness at some resolution. Therefore, preserving connectedness guarantees that clusters separated in the data remain separated in the embedding. Preserving connectedness is of more than just theoretical interest. Since connectedness applies to finite samples, preserving connectedness provides a formal guarantee of cluster separation for TPE in practice. To our knowledge, TPE is the first method with a formal guarantee of this kind.

Preserving connectedness provides implicit bounds on the Euclidean distances between pairs of points in the embedding. Let x_i, x_j be points merged into the same cluster at iteration k such that $L_{i,j} = \Delta_k$ and define

$$U_{i,j} = \sum_{k': k' \leq k, S_{n+k'} \subseteq S_{n+k}} \Delta_{k'},$$

the sum of the merge distances between clusters contained in the merged cluster up to and including iteration k . The Euclidean distance between the points is bounded as follows.

Theorem 3.

$$L_{i,j} \leq d_{i,j}(X) \leq U_{i,j}.$$

Corollary 4.

$$\sigma(X) \leq \sum_{x_i, x_j \in X} \max\{(D_{i,j} - L_{i,j})^2, (D_{i,j} - U_{i,j})^2\}.$$

Bounds on the Euclidean distances between pairs of points in the embedding provide a trivial upper bound on the stress of the embedding. Therefore, preserving connectedness prevents TPE from producing embeddings with arbitrarily high stress regardless of the quality of the optimization method such as the greedy approximation.

Results

In this section, we demonstrate the applicability of TPE by analyzing examples drawn from molecular biology, signal processing, and computer vision both qualitatively and quantitatively. Rather than being exhaustive, our goal is to highlight some of the features of TPE through each example. We compare TPE to both classical methods, PCA and non-metric MDS, and a popular force-based method, t-SNE [11], that recent studies have found separates clusters well [10].

Protein Sequences. In our first example, we analyzed 124 protein sequences of 3-phosphoglycerate kinases (3-PGKs) belonging to the domains Archaea, Bacteria, and Eukaryota collected from public databases by ref. [35]. Since protein sequences cannot be represented as real vectors, methods such as PCA that require such a representation cannot be used. Finding a good metric for protein sequences is a difficult and longstanding problem [36]. We used inverse sequence alignment scores from the basic local alignment search tool (BLAST) [37] as dissimilarities. Since BLAST scores can be highly non-metric [12], they are notoriously difficult to embed without collapsing points on top of each other.

We compared TPE to non-metric MDS, a variant of MDS for non-metric dissimilarities, and t-SNE in Fig. 3. TPE clearly separates all three domains, while non-metric MDS and t-SNE mix members of different domains together. Non-metric MDS collapses many points on top of each other, while TPE spaces the points evenly, reflecting the lack of information in the values of the dissimilarities. t-SNE separates small clusters within each domain, but does not preserve their relative locations and mixes them together, while TPE keeps each domain in a contiguous region. Since the merge order of the SL dendrogram is preserved under monotonic transformations of the dissimilarities, TPE is more sensitive to the rank order than the values of the dissimilarities. Therefore, TPE is not as sensitive to non-metric dissimilarities.

Radar Signals. In our second example, we analyzed 351 radar signals targeting free electrons in the ionosphere collected by ref. [38]. Each radar signal consisted of 34 integer and real measurements. We treated each radar signal as a 34-dimensional real vector and used Euclidean distances as dissimilarities. Good radar signals were defined as those that returned evidence of free electrons in the ionosphere, while bad radar signals were defined as those that passed through the ionosphere and returned background noise. Therefore, good radar signals are highly similar, while bad radar signals can be highly dissimilar.

We compared TPE to PCA and t-SNE in Fig. 4. TPE clearly separates good and bad radar signals, while PCA and t-SNE mix them together. PCA collapses the good and bad radar signals on top of each other with little separation. t-SNE separates small clusters of good and bad radar signals, but does not preserve their relative locations and mixes them together. TPE keeps good and bad radar signals in contiguous regions. Moreover, TPE concentrates the good radar signals and disperses the bad radar signals, reflecting the different amounts of noise in the radar signals. Therefore, TPE preserves both clusters and density.

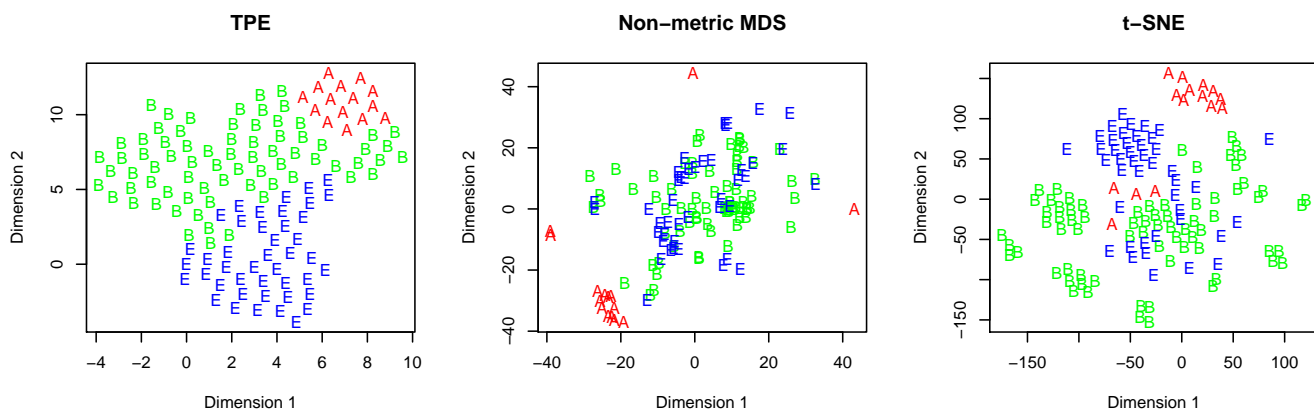


Fig. 3. Embeddings of protein sequences by TPE, non-metric MDS, and t-SNE. Each point is a protein sequence labeled by the domain it belongs to where 'A' denotes Archaea, 'B' denotes Bacteria, and 'E' denotes Eukaryota.

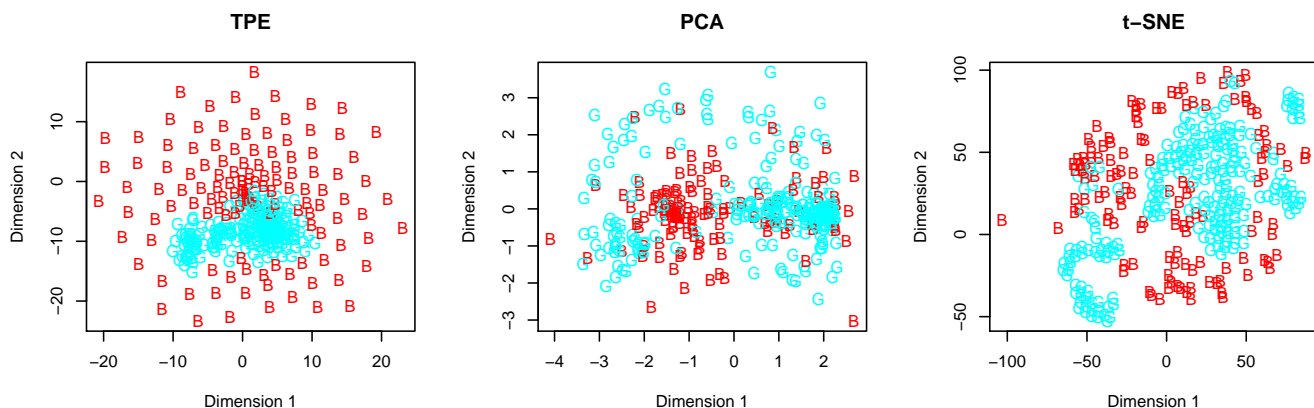


Fig. 4. Embeddings of radar signals by TPE, PCA, and t-SNE. Each point is a radar signal labeled by its quality where 'G' denotes a good radar signal and 'B' denotes a bad radar signal.

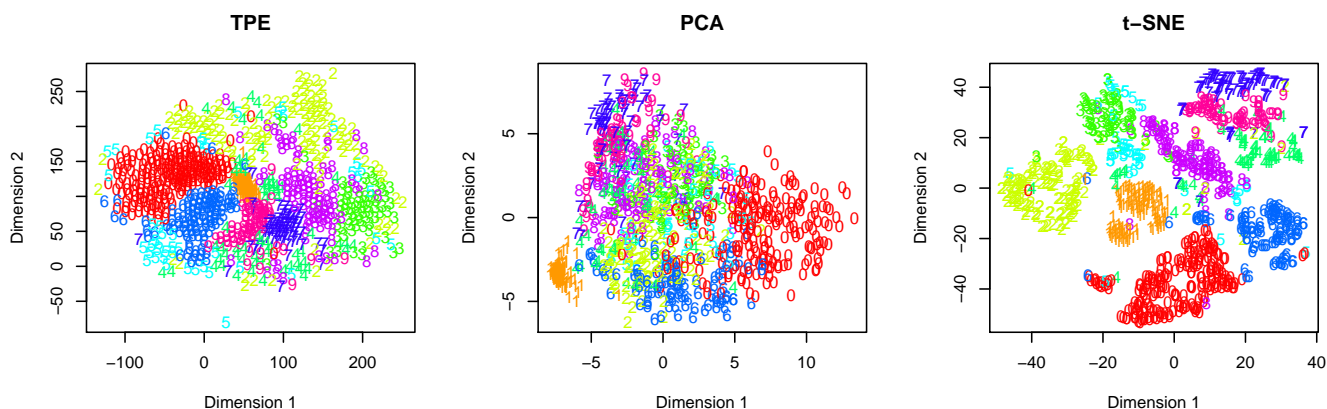


Fig. 5. Embeddings of handwritten digits by TPE, PCA, and t-SNE. Each point is an image labeled by the digit it represents.

Handwritten Digits. In our final example, we analyzed 1000 images of handwritten digits collected by the United States

Postal Service in ref. [39]. Each image was 16×16 pixels and greyscale color. We treated each image as a 256-

dimensional real vector and used Euclidean distances as dissimilarities. Since the intrinsic dimensionality of handwritten digits is thought to be much higher than two or three [40], it is notoriously difficult to separate all ten digits in an embedding due to the crowding problem.

We compared TPE to PCA and t-SNE in Fig. 5. TPE and t-SNE separate all ten digits, while PCA can only separate some of them. t-SNE most clearly separates all ten digits by creating empty space between them. Since TPE cannot create empty space between clusters without violating connectedness, it is more sensitive than t-SNE to the crowding problem, particularly when there are many clusters. However, t-SNE does not preserve density well. For example, t-SNE separates small clusters within the digit one, while TPE and PCA show that the digit one is the densest cluster, reflecting the minimal amount of variation in how it is written. Although TPE does not separate clusters as well as t-SNE, TPE preserves density better than t-SNE. Therefore, TPE strikes a balance between preserving clusters and density.

Quantitative Evaluation. Since qualitative evaluation of the quality of the embeddings can be subjective, quantitative evaluation is important. We used several popular performance metrics to evaluate the extent to which the embeddings preserve the dissimilarities, the local neighborhoods, and the known clusters. Although TPE sacrifices preserving the dissimilarities by preserving connectedness, the loss is relatively small compared to the gain made in preserving the local neighborhoods and the known clusters, which are widely believed to be more important for visualization purposes [10]. Details of the quantitative evaluation can be found in the SI Text.

Discussion

Revealing clusters is one of the main goals of visualization. However, most dimensionality reduction methods have difficulty preserving clusters due to the crowding problem. In three difficult examples, TPE was able to separate clusters of interest well compared to other dimensionality reduction methods. It is important to emphasize that the success of TPE is not a mere consequence of the ability of the SL dendrogram to separate clusters. In all three examples, cutting the SL dendrogram produced clusters at a single resolution that were no better than random clusters in terms of accuracy with respect to the known clusters. TPE succeeds by preserving clusters at all resolutions rather than just a single resolution.

Dimensionality reduction methods that separate clusters often have issues with artificially creating clusters. It is well known that force-based methods such as t-SNE can find clusters in data where there are none [10]. TPE is not as susceptible to this problem. While preserving connectedness prevents clusters from being too close together, it also prevents clusters from being too far apart. Therefore, it is difficult for TPE

to artificially create clusters without violating connectedness. In order to empirically test whether TPE artificially creates clusters, we simulated an example of a difficult convex manifold, the Swiss roll. TPE and Isomap [7], a popular manifold learning method, were able to preserve the continuity of the manifold well, while t-SNE artificially created clusters. Details of the experiment can be found in the SI Text.

Dimensionality reduction methods often have issues with robustness to noise. The dependence of TPE on the SL dendrogram may raise concerns about the sensitivity of the SL dendrogram and TPE to sampling variability. However, the SL dendrogram has been shown to be stable in the sense that small perturbations of the data do not change the structure of the SL dendrogram significantly [41, 30]. Therefore, we expect that TPE will also be stable. In order to empirically test the stability of TPE, we simulated 100 samples from a difficult non-convex manifold, the barbell, and computed the average sample variance of the coordinates of the points in the embeddings. TPE and Isomap had comparable stability to the exact embedding, while t-SNE was two orders of magnitude less stable. Details of the experiment can be found in the SI Text.

Preserving connectedness allows TPE to preserve different types of structure. However, preserving connectedness is a strong constraint that may not be effective for certain types of structure. Therefore, TPE will not always be able to perform as well as other dimensionality reduction methods in specific applications. For example, manifold learning methods may preserve certain manifolds such as the Swiss roll better and force-based methods may preserve certain clusters such as the handwritten digits better. Nevertheless, we believe that the robustness of TPE is what makes it useful in practice.

TPE is a promising approach to visualization because it has a formal guarantee of cluster separation, requires no parameters, and can handle general types of data. However, there are a few issues with TPE that may limit its applicability. First, TPE has a cubic time complexity, which can be prohibitively slow for large data sets. Second, since TPE only provides an embedding rather than a mapping, it cannot be applied to out-of-sample data. Finally, although we have found that the greedy approximation works well in practice, better optimization methods may significantly improve the performance of TPE. We hope that these issues will be addressed by future research.

ACKNOWLEDGMENTS. An earlier version of this work appeared in [42]. This work was partially supported by the National Science Foundation under grants no. DMS-0907009 and no. IIS-1017967, by the National Institute of Health under grant no. R01 GM096193, and by the Army Research Office Multidisciplinary University Research Initiative under grant no. 58153-MA-MUR all to Harvard University. Additional funding was provided by the Harvard Medical School's Milton Fund. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the National Institute of Health, the Army Research Office, the National Science Foundation, or the U.S. government.

1. R M Shiffrin and K Börner. Mapping knowledge domains. *Proc Natl Acad Sci USA*, 101:5183–5185, 2004.
2. G E Hinton and R R Salakhutdinov. Reducing the dimensionality of data with neural networks. *Science*, 313:504–507, 2006.
3. I T Joliffe. *Principal component analysis*. Springer-Verlag, New York, 2002.
4. J B Kruskal and M Wish. *Multidimensional Scaling*. Sage University Press, Newbury Park, 1978.
5. T F Cox and M A A Cox. *Multidimensional Scaling*. Chapman & Hall/CRC, Boca Raton, 2001.
6. I Borg and P Groenen. *Modern Multidimensional Scaling: Theory and Applications*. Springer-Verlag, New York, 2005.
7. J B Tenenbaum, V de Silva, and J C Langford. A global geometric framework for nonlinear dimensionality reduction. *Science*, 290:2319–2323, 2000.
8. S T Roweis and L K Saul. Nonlinear dimensionality reduction by locally linear embedding. *Science*, 290:2323–2326, 2000.
9. M Belkin and P Niyogi. Laplacian eigenmaps for dimensionality reduction and data representation. *Neural Comput*, 15:1373–1396, 2003.
10. J Venna, J Peltonen, K Nybo, H Aidos, and K Samuel. Information retrieval perspective to nonlinear dimensionality reduction for data visualization. *J Mach Learn Res*, 11:451–490, 2010.
11. L van der Maaten and G E Hinton. Visualizing data using t-sne. *J Mach Learn Res*, 9:2579–2605, 2008.
12. V Roth, J Laub, M Kawanabe, and J M Buhmann. Optimal cluster preserving embedding of nonmetric proximity data. *IEEE T Pattern Anal*, 25:1540–1551, 2003.

13. E P Xing, A Y Ng, M I Jordan, and S Russell. Distance metric learning, with application to clustering with side-information. In *Advances in Neural Information Processing Systems*, volume 14, 2002.
14. G E Hinton and S T Roweis. Stochastic neighbor embedding. In *Advances in Neural Information Processing Systems*, volume 15, 2003.
15. J A Cook, I Sutskever, A Mnih, and G E Hinton. Visualizing similarity data with a mixture of maps. In *Proceedings of the 11th International Conference on Artificial Intelligence and Statistics*, 2007.
16. M A Carreira-Perpiñán. The elastic embedding algorithm for dimensionality reduction. In *Proceedings of the 27th International Conference on Machine Learning*, 2010.
17. L Chen and A Buja. Local multidimensional scaling for nonlinear dimension reduction, graph drawing, and proximity analysis. *J Am Stat Assoc*, 104:209–219, 2009.
18. G Di Battista, P Eades, R Tamassia, and I G Tollis. *Graph Drawing: Algorithms for the Visualization of Graphs*. Prentice Hall, New York, 1998.
19. M Kaufmann and D Wagner. *Drawing Graphs: Methods and Models*. Springer-Verlag, New York, 2001.
20. R N Shepard. Multidimensional scaling, tree-fitting, and clustering. *Science*, 210:390–398, 1980.
21. J C Gower and G J S Ross. Minimum spanning trees and single linkage cluster analysis. *Appl Stat*, 18:54–64, 1969.
22. R B Zadeh and S Ben-David. A uniqueness theorem for clustering. In *Proceedings of the 25th Annual Conference on Uncertainty in Artificial Intelligence*, 2009.
23. B Shaw and T Jebara. Structure preserving embedding. In *Proceedings of the 26th International Conference on Machine Learning*, 2009.
24. P Eades. The realization problem for euclidean minimum spanning trees is np-hard. *Algorithmica*, 16:60–82, 1996.
25. J A Hartigan. *Clustering Algorithms*. Wiley, New York, 1975.
26. J A Hartigan. Statistical theory in clustering. *J Classif*, pages 63–76, 1985.
27. W Stuetzle. Estimating the cluster tree of a density by analyzing the minimal spanning tree of a sample. *J Classif*, 20:25–47, 2003.
28. A Buja, D F Swayne, M L Littman, N Dean, H Hofmann, and L Chen. Data visualization with multidimensional scaling. *J Comput Graph Stat*, 17:444–472, 2008.
29. S C Johnson. Hierarchical clustering schemes. *Psychometrika*, 32:241–254, 1967.
30. G Carlsson and F Mémoli. Characterization, stability and convergence of hierarchical clustering methods. *J Mach Learn Res*, 11:1425–1470, 2010.
31. H Qiu and E R Hancock. Clustering and embedding using commute times. *IEEE T Pattern Anal*, 29:1873–1890, 2007.
32. J C Gower and G B Dijkstra. *Procrustes Problems*. Oxford University Press, Oxford, 2004.
33. M Quist and G Yona. Distributional scaling: an algorithm for structure-preserving embedding of metric and non-metric spaces. *J Mach Learn Res*, 5:399–420, 2004.
34. V de Silva and J B Tenenbaum. Global versus local methods for nonlinear dimensionality reduction. In *Advances in Neural Information Processing Systems*, volume 15, 2003.
35. J D Pollack, Q Li, and Pearl D K. Taxonomic utility of a phylogenetic analysis of phosphoglycerate kinase proteins of archaea, bacteria, and eukaryota: Insights by bayesian analyses. *Mol Phylogenet Evol*, 35:420–430, 2005.
36. W R Atchley, J Zhao, A D Fernandes, and T Driike. Solving the protein sequence metric problem. *Proc Natl Acad Sci USA*, 102:6395–6400, 2005.
37. S F Althscul, W Gish, W Miller, E W Myers, and D J Lipman. Basic local alignment search tool. *J Mol Biol*, 215:403–410, 1990.
38. V G Sigillito, S P Wing, L V Hutton, and K B Baker. Classification of radar returns from the ionosphere using neural networks. *J Hopkins Apl Tech D*, 10:262–266, 1989.
39. J J Hull. A database for handwritten text recognition research. *IEEE T Pattern Anal*, 16:550–554, 1994.
40. L K Saul and S T Roweis. Think globally, fit locally: Unsupervised learning of low dimensional manifolds. *J Mach Learn Res*, 4:119–155, 2003.
41. J A Hartigan. Consistency of single linkage for high-density clusters. *J Am Stat Assoc*, 76:388–394, 1981.
42. A D Shieh, T B Hashimoto, and E M Airoldi. Tree preserving embedding. In *Proceedings of the 28th International Conference on Machine Learning*, 2011.