

Capacity Preserving Mapping for High-dimensional Data Visualization

Rongrong Wang Xiaopeng Zhang

Abstract

We provide a rigorous mathematical treatment to the crowding issue in data visualization when high dimensional data sets are projected down to low dimensions for visualization. By properly adjusting the capacity of high dimensional balls, our method makes right enough room to prepare for the embedding. A key component of the proposed method is an estimation of the correlation dimension at various scales which reflects the data density variation. The proposed adjustment to the capacity applies to any distance (Euclidean, geodesic, diffusion) and can potentially be used in many existing methods to mitigate the crowding during the dimension reduction. We demonstrate the effectiveness of the new method using synthetic and real datasets.

1 Introduction

Visualizing high dimensional data via their low dimension projections is particularly useful in facilitating data analysts to understand their datasets, detect underlying data patterns and create various hypotheses about the data. To achieve this goal, we project the high dimensional dataset $X = \{X_1, \dots, X_N\} \subset \mathbb{R}^n$ down to low dimensions: $Y_i = P(X_i) \in \mathbb{R}^d$ ($i = 1, \dots, N$) with $d = 2$ or 3 and visualize the low dimensional embedding $\{Y_i\}_{i=1}^N$ via a single scatter plot. Designing the mapping P that yields a reliable visualization is the focus of this paper. A good design should take into account the special need of data visualization 1) as most human beings can only digest visual information in at most three dimensions (or four dimensions if video sequences are included), the map is required to project data of any dimensionality down to 2 or 3. 2) as revealing the data pattern is the major objective, the map should be able to preserve the geometrical structure of the data.

Most existing dimensionality reduction techniques (i.e., MDS [14], LLE [11], Isomap [13]) are only designed to reduce the data to its intrinsic dimension, which is usually higher than three. While data visualization methods SNE [6], tNSE [16], PHATE [10], UMAP [9] can reduce data of any dimension to two or three, geometric relations such as cluster radii and relative distances between clusters are usually lost. We hereby formulate the main mathematical question in data visualization: how to map datasets with a wide range of dimensionality to 2 or 3 while minimizing the geometric distortion?

The geometric distortion we care about in this paper is the relative closeness/similarity between points. It is considered successfully preserved if close points remain close and far away points remain far away. The main obstacle in preserving the geometric relation is the so-called crowding issue. Simply put, the crowding issue arises from the fact that a higher dimensional body typically has a larger capacity than a lower dimensional one, hence reducing the dimensionality causes points to be crowded.

We propose a way to adjust the capacity of the high dimensional body before the dimension reduction. The method named Capacity Preserving Mapping (CPM) is essentially a class of methods that generalize many existing methods by redefining the distance they use. Compared to the popular methods such as tSNE and UMAP, our method can better preserve geometrical structures of the dataset and does not presume existences of clusters.

Our contribution is twofold: 1. we propose a way to define a dimension-aware distance to treat the crowding issue, which greatly assists high dimensional data visualization by better preserving geometrical structures of the dataset. 2. we propose a method to compute the multi-scale intrinsic dimensionality of a dataset that is necessary for the definition of the new distance in 1.

We note that although previous methods SNE, tNSE, PHATE, UMAP also treat the crowding problem to some extent, it was done in a less rigorous manner.

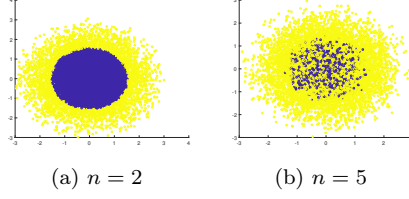


Figure 1: The crowding phenomenon: two non-overlapping classes (a) becomes overlapping after MDS embedding (b)

2 Related work

Data visualization is an important task in data mining and is closely related to dimensionality reduction, graph learning, and data clustering. A number of early works studied the so-called table data visualization problem (see the review article [4]) that visualizes N ($N > 3$) attributes of data in a table using 2D or 3D plots. This is made possible either by using multiple pixels (or attributes) in the low dimension to represent one data point in the high dimension or using multiple plots from different angles to build up the high dimensional image. The drawbacks of these methods are that the visualization is not directly meaningful, needs human effort to understand and interpret, and the relation between data points (such as whether clusters exist or how close they are) is not immediately apparent.

In a separate line of research, various dimension reduction methods [15] can be used to visualize high dimensional data in one scatter plot. The methods can be categorized as linear ones (PCA, MDS [14], ICA [7]...) and nonlinear ones (LLE [11], non-metric MDS [18], Isomap [13], Laplacian Eigenmap [1], Diffusion map [8], etc). As mentioned earlier, these methods can only reliably map the data down to its intrinsic dimension, while visualization requires the target dimension to be less than four. Hence most good visualization results are observed on artificial datasets with small intrinsic dimensions (e.g., swiss roll (2D), Helix (1D), Twin peaks (2D)).

The class of methods that are most relevant to ours includes SNE [6], tSNE [16], UMAP [9] and PHATE [10]. They are directly designed for data visualization and able to reduce the dimension from any $N \geq 2$ to two or three. However, in all these methods, a rigorous treatment of the crowding issue is missing. The purpose of this paper is to fill in this gap by proposing a distance correcting step that can work with any metric. We achieve this goal by a careful computation of the correlation dimension at various scales.

3 Motivation - the crowding phenomenon

As stated in [16], when high dimensional data is mapped to low dimensions, there is a tendency for non-overlapping groups to overlap. Theoretically, this is due to the difference in norm concentration between high and low dimensions: a ball in higher dimensions has a volume that grows faster with radius $\text{Vol}(B_2^n(r)) \sim r^n$ [3], where $B_2^n(r)$ is the ℓ_2 ball in \mathbb{R}^n with radius r . As a result, the high dimensional points are more likely to be distributed near the surface. When mapped to low dimensions, since there is less room near the surface, the points are pushed towards the center.

Let us visualize the crowding phenomenon observed during dimension reduction via Multidimensional Scaling (MDS). Assume a high dimensional ℓ_2 ball contains two classes of points. Class 1 lies inside the ball $B_2^n(1)$ and Class 2 lies in a spherical shell right outside Class 1. When $\{X_i\}_{i=1}^N \subseteq \mathbb{R}^2$ is two dimensional, a direct visualization (Figure: 1a) shows the correct relation between the two classes. When $\{X_i\}_{i=1}^N \subseteq \mathbb{R}^5$ is five dimensional and is mapped down to \mathbb{R}^2 via MDS, its geometrical structure is distorted (Figure 1b) and we see a severe crowding phenomenon: the second class is pushed towards the center. Finding a way to correct this type of distortion is our main objective.

3.1 Notation

Let $f(x) \geq 0$ and $g(x) \geq 0$ be two increasing functions of x . We use $f \simeq g$ to denote that there exist constants $0 < c < C$ such that $cg(x) \leq f(x) \leq Cg(x)$ for all x . $f(x) \lesssim g(x)$ means $\frac{f'(x)}{f(x)} < \frac{g'(x)}{g(x)}$, i.e., $\log f(x)$ grows slower than $\log g(x)$ as x increases, and $f(x) \gtrsim g(x)$ means the opposite. For a random variable z and its probability density function h , we use $P_{z \sim h}(z \in A)$ to denote the probability of the event A under h .

4 The methodology

Let us first define some mathematical terminology that helps to describe the crowding issue. Assuming each data point X_i is independently drawn from an underlying manifold \mathcal{M} according to a continuous probability distribution $f(\mathcal{M})$.

Definition 4.1 (Relative capacity) Let \mathcal{M} be a compact manifold and f be a probability distribution defined on \mathcal{M} . For any point $p \in \mathcal{M}$, we define **the relative capacity of a neighbourhood of p with radius r** as the probability that randomly drawn points according to f fall inside this neighbourhood,

$$C(p, r; \mathcal{M}, f, D) := \mathbb{P}_{z \sim f(\mathcal{M})}(D(z, p) \leq r) = \mathbb{E}_{z \sim f(\mathcal{M})} \mathbf{I}(r - D(z, p)) \quad (1)$$

where $D : \mathbb{R}^n \times \mathbb{R}^n \rightarrow [0, \infty)$ is a map that defines the pairwise distances on \mathcal{M} and \mathbf{I} is the step function used to model the neighbourhood of radius r centred on p (i.e., $\{z : D(z, p) \leq r\}$). The **relative capacity of the manifold** is the expectation of $C(p, r; \mathcal{M}, f, D)$ with respect to all $p \in \mathcal{M}$,

$$C(r; \mathcal{M}, f, D) := \mathbb{E}_{p \sim f(\mathcal{M})} C(p, r; \mathcal{M}, f, D),$$

For simplicity, we write $C(r), C(p, r)$ in short for $C(r; \mathcal{M}, f, D), C(p, r; \mathcal{M}, f, D)$, respectively.

Intuitive, the “capacity” of a neighbourhood reflects the amount of data the neighbourhood can hold. The “relative capacity” is the normalized capacity, and hence is a probability measure. If many samples are generated, $C(p, r)$ reflects the expected number of points in a neighbourhood. If r is fixed, $C(p, r)$ as a function of p reflects how the data density varies across different locations on the manifold. If p is fixed, $C(p, r)$ as a function of r reflects how fast the capacity grows as the neighbourhood expands. The relative capacity $C(r)$ for the manifold is the percentage of pairs with a distance smaller than r . Since by definition, $C(p, r)$ and $C(r)$ are cumulative distribution functions, their partial derivatives $\frac{\partial C(p, r)}{\partial r}, \frac{\partial C(r)}{\partial r}$ are probability density functions.

Definition 4.2 (Relative density) The relative density function is defined as the derivatives of the relative capacities

$$\rho(p, r) := \frac{\partial C(p, r)}{\partial r}, \quad \rho(r) := \frac{dC(r)}{dr}$$

We also have the reversed relation $C(p, r) = \int_0^r \rho(p, t) dt$, and $C(r) = \int_0^r \rho(t) dt$.

Exmaple: Suppose $\mathcal{M} = B_2^n(1)$, $f(\mathcal{M})$ is the uniform distribution on \mathcal{M} , and $D = \|\cdot\|_2$ is the Euclidean distance. Then $C(0, r; \mathcal{M}, D) = r^n$, for any $r \leq 1$, and $C(p, r; \mathcal{M}, D) \simeq r^n$ where \simeq was defined in Sect 3.1.. Taking expectation with respect to p , one can verify that the capacity of the manifold is $C(r; \mathcal{M}, D) \simeq r^n$ and the density is about $\rho(r) \simeq nr^{n-1}$.

In what follows, we make the following two assumptions on the data.

Assumption 1: The original data $\{X_i\}_{i=1}^N$ are drawn independently from some underlying manifold \mathcal{M} according to a continuous distribution $f(\mathcal{M})$. The embedded data $Y_i = P(X_i)$, $i = 1, \dots, N$, are independent realizations of the embedded manifold $\mathcal{S} = P(\mathcal{M})$, according to the induced distribution f_I of f under P .

Assumption 2: The relative densities $\rho(r; \mathcal{M}, f, \|\cdot\|_2)$ and $\rho(r; \mathcal{S}, f_I, \|\cdot\|_2)$ of \mathcal{M} and \mathcal{S} can be fitted with the models

$$\rho(r; \mathcal{M}, f, \|\cdot\|_2) = c_m n_m(r) r^{n_m(r)-1}, \quad \rho(r; \mathcal{S}, f_I, \|\cdot\|_2) = c_s n_s(r) r^{n_s(r)}$$

where c_m and c_s are absolute constants, and $n_m(r) > 0, n_s(r) > 0$ are both slowly varying dimension functions of r .

Remark 4.3 $n(r)$ and $n(s)$ are closely related to the so-called correlation dimension [5] defined as $\dim_{\text{corr}} = \lim_{r \rightarrow 0} \frac{\partial \log C(r)}{\partial \log r}$. We can see that \dim_{corr} coincides with $n(r)$ at $r = 0$, i.e., $\dim_{\text{corr}} = \lim_{r \rightarrow 0} n_m(r)$. Therefore, $n_m(r)$ can be viewed as an extension of the correlation dimension to positive scales and be called as a dimension function.

Remark 4.4 Assumption 2 is mild as it allows the relative capacity to vary with scales.

Remark 4.5 The assumption that $n_m(r), n_s(r)$ are slowly varying functions of r ensures that they are almost constants in small intervals and therefore can be estimated by counting points in these intervals (see Sect. 5).

Remark 4.6 The slowly varying assumption can be mathematically written as: $n'_m(r), n'_s(r) \approx 0$, which implies on a small interval $[r, r + dr]$, $n_m(r)$ and $n_s(r)$ are close to constant thus we can integrate $\rho(r)$ to get the increment of $C(r)$:

$$\begin{aligned} C(r + dr; \mathcal{M}, f, \|\cdot\|_2) - C(r; \mathcal{M}, f, \|\cdot\|_2) &= c_m r^{n_m(r)} dr \\ C(r + dr; \mathcal{S}, f_I, \|\cdot\|_2) - C(r; \mathcal{S}, f_I, \|\cdot\|_2) &= c_s r^{n_s(r)} dr. \end{aligned}$$

4.1 A new distance that preserves the relative capacity

Let us explain the crowding issue with the relative capacity using the manifold $\mathcal{M} = B_2^n(1)$ endowed with the uniform distribution. Consider the ball $B_2^n(1)$ as a union of concentric shells with infinitesimal thickness and increasing radii. If these shells are to be mapped to $d = 2, 3$ without changing their inclusion order, then the relative density $\rho(0, r)$ as a function of r that reflects how many points are lying on each shell, should not change after the embedding; otherwise points will shuffle around. More explicitly, if after embedding we have

$$\rho(0, r; \mathcal{M}, f, \hat{D}) \gtrsim \rho(0, r; \mathcal{S}, f_I, \|\cdot\|_2)$$

where \gtrsim (defined in Sect 3.1) means that the left hand side grows faster with r than the right hand side, then \mathcal{M} has a larger capacity than \mathcal{S} and the crowding phenomenon will be observed. On the contrary, if

$$\rho(0, r; \mathcal{M}, f, \hat{D}) \lesssim \rho(0, r; \mathcal{S}, f_I, \|\cdot\|_2),$$

then \mathcal{S} has a larger capacity than \mathcal{M} . Points and clusters will be pulled apart to a certain degree, but is usually less harmful than the crowding.

The relative density of the high dimensional manifold $\mathcal{M} = B_2^n(1)$ obeys $\rho(0, r) \simeq nr^{n-1}$ due to Example 1. If we equip the embedded manifold \mathcal{S} with the same distribution (i.e., uniform distribution) and the same distance (i.e., the Euclidean distance), then the relative density on \mathcal{S} is $\rho(0, r) \simeq dr^{d-1}$. Since it is not growing fast enough, the crowding phenomenon will occur.

To solve this problem, we need to either 1) endow non-uniform probability distributions on the low dimensional manifold to allow points with a larger norm to be selected with a higher probability, or 2) modify the distance in the definition of $C(p, r)$ in (5) to allow the shells with larger radii to be mapped to larger shells in the low dimension. We go with the latter option in this paper and put the former as future research direction.

Explicitly, we want to design a new distance \hat{D} such that the relative densities of the original manifold and the embedded one can approximately match

$$\rho(r; \mathcal{M}, f, \hat{D}) \approx \rho(r; \mathcal{S}, f_I, \|\cdot\|_2)$$

for all $r \geq 0$. Since $\rho(r)$ is the derivative of $C(r)$, under the slowly varying assumption (Assumption 2), the above implies

$$C(r; \mathcal{M}, f, \hat{D}) \approx C(r; \mathcal{S}, f_I, \|\cdot\|_2).$$

Theorem 4.7 Let \mathcal{M} be the high dimensional data manifold and \mathcal{S} be the embedded one. Let $P : \mathcal{M} \rightarrow \mathcal{S}$ be the dimension reduction mapping. Assuming the pairwise distance in the embedded space is some function of that in the original space i.e., $\|P(X_i) - P(X_j)\|_2 = G(\|X_i - X_j\|_2)$, for any $i \neq j$, and the scaling function G is unknown. Under Assumption 1 and Assumption 2, we can define a new “distance” between any pair of points x and z as

$$\hat{D}(\|x - z\|_2) := \|x - z\|_2^{\frac{n_m(\|x - z\|_2)}{n_s(\|P(x) - P(z)\|_2)}}, \quad x, z \in \mathbb{R}^n. \quad (2)$$

This distance allows a match in the relative capacity between the original and the embedded data

$$C(r; \mathcal{M}, f, \hat{D}) = C(r; \mathcal{S}, f_I, \|\cdot\|_2),$$

where $n_m(\cdot)$ and $n_s(\cdot)$ are the same as defined in Assumption 2.

Remark 4.8 One can verify that this \hat{D} no longer defines a norm. In particular, it does not preserve the monotonicity of distances ($\|x\|_2 \leq \|y\|_2 \not\Rightarrow \hat{D}(\|x\|_2) \leq \hat{D}(\|y\|_2)$). However, the proof of Theorem 4.7 indicates that this definition of \hat{D} is both necessary and sufficient for a match of the relative capacity. Therefore we propose to bring back the monotonicity with the minimal change to \hat{D} . To that end, we start with the smallest pairwise distance in the high dimensional space and redefine \hat{D} to \tilde{D} inductively. For the smallest pairwise distance, say σ_0 , simply set $\tilde{D}(\sigma_0) = \hat{D}(\sigma_0)$.

For the second smallest distance, say σ_1 , we need to define $\tilde{D}(\sigma_1)$ as close to $\hat{D}(\sigma_1)$ as possible while having the monotonicity $\tilde{D}(\sigma_1) > \tilde{D}(\sigma_0)$. Thus we set $\tilde{D}(\sigma_1) = \max\{\hat{D}(\sigma_1), \tilde{D}(\sigma_0)\}$. Assume we have defined $\tilde{D}(\sigma_k)$ that has monotonicity up to σ_k , set $\tilde{D}(\sigma_{k+1})$ as follows guarantees $\tilde{D}(\sigma_{k+1})$ to be monotonic up to σ_{k+1}

$$\tilde{D}(\sigma_{k+1}) = \max\{\hat{D}(\sigma_{k+1}), \tilde{D}(\sigma_k)\}. \quad (3)$$

Intuitively, the maximum in the above equation makes \tilde{D} assign more room to the low dimensional space than necessary. As we discussed earlier, this is much less harmful than the crowding issue to the visualization.

Definition 4.9 We call the new “distance” defined in (2) as the Capacity Adjusted Distance (CAD), and the modified version \tilde{D} defined in (3) as the modified Capacity Adjusted Distance (modified CAD).

Remark 4.10 Theorem 2 currently assumes that the high dimensional manifold \mathcal{M} is equipped with the Euclidean distance, and provides a way to modify the Euclidean distance to allow a match of the capacity. If the equipped distance for \mathcal{M} is the geodesic distance or the diffusion distance, we can also use Theorem 2 to modify them to allow a match of the capacity, by simply replacing the quantity $\|x - z\|_2$ in (2) by $d_{\text{geodesic}}(x, z)$ and $d_{\text{diffusion}}(x, z)$. We make it an option to the user in our algorithm (Algorithm 1).

The definition of \tilde{D} requires the knowledge of the dimension functions $n_m(r)$ and $n_s(r)$ at various scales. In the next section, we introduce a dimension estimation method that allows an estimate of $n_m(r)$ from X_i . However, we cannot estimate $n_s(r)$ directly since Y_i is not yet known. To solve this problem, we make the assumption that the induced measure on the low dimension manifold is uniform, and in light of Example 1, this means $C(r, \mathcal{S}, f_I, \|\cdot\|_2) \simeq r^d$, hence $n_s(r) = d$. Once the low dimensional embedding $\{Y_i\}_{i=1}^N$ is obtained, one may use it to update n_s to make it more accurate (note that we do not update n_s in this paper).

4.2 Dimensionality reduction with the modified Capacity Adjusted Distance

Before introducing ways of estimating $n(r)$, let us first discuss the dimensionality reduction procedure based on the new distance \tilde{D} . As usual, we search for the low dimensional vectors $\{Y_i\}_{i=1}^N$ that best preserves \tilde{D} for all pairs of points. In CMP, we measure the dissimilarity using the Kullback–Leibler divergence. The low dimensional embedding $\{Y_i\}$ is the minimizer of the following optimization problem

$$\{\hat{Y}_i\}_{i=1}^N = \arg \min_{Y_i, i=1, \dots, N} \sum_{i,j} p_{i,j} \log \frac{p_{i,j}}{q_{i,j}} \quad (4)$$

where

$$p_{i,j} = \frac{(\epsilon + \tilde{D}^2(\|X_i - X_j\|))^{-1}}{\sum_{k,l,k \neq l} (\epsilon + \tilde{D}^2(\|X_k - X_l\|))^{-1}}, \quad q_{i,j} = \frac{(1 + \|Y_i - Y_j\|^2)^{-1}}{\sum_{k,l,k \neq l} (1 + \|Y_k - Y_l\|^2)^{-1}}, \quad (5)$$

where $\|\cdot\|$ is the distance measure in the high dimensional space, and $\epsilon > 0$ is some small constant used to avoid taking the reciprocal of 0. Besides the KL divergence, one can use other (dis)similarity measures and obtain different formulations of the optimization problems, which we put as future work. Putting the $p_{i,j}$ defined in (5) into a matrix, $P = [p_{i,j}]_{i,j=1}^N$. We can think of this P as a normalized network probability matrix of the random network built based on the pairwise distance matrix $\tilde{D} = [\tilde{D}(X_i, X_j)]_{i,j=1}^N$. Explicitly, let the N data points $\{X_i\}_{i=1}^N$ correspond to the N nodes of a random graph, and the random edges are constructed as follows. X_i is connected with X_j by an edge with probability $p_{i,j} \sim \tilde{D}^{-1}(X_i, X_j)$, where \sim hides a universal constant (i.e., the denominator in (5)). This is to say, closer points are more likely to be connected by an edge in this random graph. The optimization (4) is hence trying to match the KL divergence between the network probabilities of the original and the embedded graphs. The small positive constant ϵ in the expression of $p_{i,j}$ and the 1 in the expression of $q_{i,j}$ are used to avoid dividing by 0.

Although the formulation (4) looks similar to that of t-SNE (see (6) below), their performances are completely different.

$$\text{t-SNE: } \{\hat{Y}_i\}_{i=1}^N = \arg \min_{Y_i, i=1, \dots, N} \sum_{i,j} p_{i|j} \log \frac{p_{i|j}}{q_{i|j}} = \arg \min_{Y_i, i=1, \dots, N} \sum_{i,j} -p_{i|j} \log q_{i|j} \quad (6)$$

where

$$p_{i|j} = \frac{e^{-\frac{\|X_i - X_j\|_2^2}{2\sigma^2}}}{\sum_{k, k \neq j} e^{-\frac{\|X_k - X_j\|_2^2}{2\sigma^2}}}, \quad q_{i|j} = \frac{(1 + \|Y_i - Y_j\|_2^2)^{-1}}{\sum_{k, k \neq j} (1 + \|Y_k - Y_j\|_2^2)^{-1}}, \quad (7)$$

Compared to t-SNE, our formulation (4) has the following merits.

1. **Less over-stretching:** t-SNE mitigates the crowding and promotes the formation of clusters by matching Gaussian distributions with t-distributions. Intuitively, this makes close points closer and far away points further, but the degree of squeezing and stretching (i.e., how much closer and how much further) are not precisely characterized and justified. As a result, t-SNE may produce fake clusters due to the over-stretching. In contrast, our method computes and performs the amount of stretching that is necessary to avoid the crowding. In other words, the stretching in our method is much milder. This can also be seen mathematically from the definition of $p_{i,j}$ in (5) and the optimization (4), where we are fitting the t -distribution with some other heavy tailed distribution (given the way how the modified capacity adjusted distance $\tilde{D}(X_i, X_j)$ is defined), whereas t-SNE is fitting the t -distribution with the light tailed Gaussian distribution, which creates more stretching. Consequently, although our method may also tear a manifold up, it happens much less often.
2. **No tuning parameter:** the performance of t-SNE heavily depends on the choice of the bandwidth parameter σ , whereas our method does not have a key tuning parameter (the small positive constant ϵ used to avoid dividing by 0 does not affect the results much as long as it is sufficiently small).
3. **Preserving more geometry:** the probability $p_{i|j}$ in t-SNE is a conditional probability, i.e., $\sum_i p_{i|j} = 1, \forall j$, which causes a loss of the geometric information due to the following reason. In the optimization (6), $p_{i|j}$ is the weight of $\log q_{i|j}$ which contains the variables Y_i and Y_j . The total weights put on logarithms containing Y_j for a given j is then the sum $\sum_i (p_{i|j} + p_{j|i}) \approx 2$ (where we used $\sum_i p_{i|j} = 1$ and $\sum_i p_{j|i} \approx 1$). The former is due to the definition of the conditional probability. The latter is due to the symmetrisation in t-SNE on the conditional probability matrix $P = [p_{i|j}]_{i,j=1}^N$, which makes $p_{i|j} \approx p_{j|i}$ and $\sum_i p_{j|i} \approx \sum_i p_{i|j} = 1$). Since for any j the total weight on Y_j is $\sum_i (p_{i|j} + p_{j|i}) \approx 2$, all points are given equal importance in the optimization. As a result, the variability among points or clusters is lost. For instance, clusters of different sizes become of similar sizes in the embedding, and clusters of different inter-cluster distances become of similar distances after embedding¹. In contrast, the probabilities in our formulation are not conditional probabilities. Mathematically, instead of normalizing the probability matrix P row by row as in t-SNE, we normalized all entries of P by one constant. This allows variabilities to exist among rows, which carry the correct geometric information. It is worth noting that this universal normalization cannot be used in t-SNE because the fast decay of the Gaussian tail often causes the sum of some rows to be much smaller than others. Because these sums are the weights on the Y_i s in the optimization, those Y_i with too small of a weight cannot be correctly computed. In contrast, our formulation allows for both row-wise and universal normalizations. The row-wise normalization achieves a greater stability at the expense of losing more geometric information.

Algorithm 1: Capacity Preserving Mapping (CPM)

Output: low dimension embedding $Y_i, i = 1, \dots, N$

Input: original data $X_i, i = 1, \dots, N$; target dimension $d = 2$ or 3 ; choices of the high dimensional distance D : Euclidean, geodesic or diffusion; number of scales: M

```

1 for  $r = 1: M$  do
2   | estimate the dimension  $n_m(r)$  from  $\{X_i\}_{i=1}^N$  via solving (17) ;
3 end
4 for  $i=1:N$  do
5   | for  $j=i+1:N$  do
6     | use the estimated  $n_m(r)$  at  $r = D(X_i, X_j)$  to defined the modified distance  $\tilde{D}(X_i, X_j)$ 
7     | as in (3) ;
5   | end
4   | end
3 end
9 Solve the optimization (4) to construct  $\{Y_i\}_{i=1}^N$ .
```

¹We refer the readers to the website <https://distill.pub/2016/misread-tsne/> for more such examples.

4.3 Comparisons to other methods

Besides t-SNE, our method is also related to the non-metric MDS (NMDS) [18] and the multi-scale SNE [17]. Similar to our method, the non-metric MDS also aims at preserving the pairwise dissimilarity as closely as possible. It approaches this goal by minimizing a scaled distances between points. Let $d_{i,j}$ be the high dimensional distances, the NMDS embedding $\{Y_i\}_{i=1}^N$ is obtained by solving the optimization problem

$$\min_{\{Y_i\}_{i=1}^N, f \in \mathcal{F}} \sum_{i,j} \|\|Y_i - Y_j\|_2 - f(d_{i,j})\| \quad (8)$$

where \mathcal{F} is the set of positive monotonically increasing functions. The scaling function f plays the role of mitigating the crowding. Indeed, the scaled distance $f(d_{i,j})$ essentially corresponds to our capacity adjusted distance $\tilde{D}(d_{i,j})$. From this perspective, our formulation provides an explicit way to compute f which avoids the trouble of solving it from an optimization. Compared to the proposed method, the NMDS in its original form (8) fails to render the correct small-scale information of data (see Figure 3). While replacing the absolute value in (8) with a KL-divergence type of dissimilarity measure can help preserving the small-scale structure, the resulting optimization is difficult to solve due to the existence of the function variable $f(x)$.

Multi-scale SNE [17] is similar to our approach in the sense that both methods assume a scale-varying dimension of the data manifold. However, since it inherits the structure of SNE, the third drawback of t-SNE mentioned in Section 4.2 applies.

5 Dimensional estimation at various scales

We propose a way to estimate the dimension function $n(r)$ from the data. To the best of our knowledge, no existing method calculates the multi-scale correlation dimension for all scales. The multi-scale dimension estimation method in [17] calculates the *average* dimension from scale 0 to scale r instead of the instantaneous dimension at scale r . Recall that in Assumption 2, $n(r)$ is defined as the growth rate of the relative density $\rho(r)$

$$\rho(r) \equiv \rho(r; \mathcal{M}, f, \|\cdot\|_2) = cn(r)r^{n(r)-1} \quad (9)$$

and c is some unknown absolute constant. Therefore, it is possible to find $n(r)$ by fitting the slope of $\log \rho$

$$n(r) \approx \frac{\log(\rho(r + \delta r)) - \log(\rho(r))}{\log(r + \delta r) - \log(r)}. \quad (10)$$

However, this estimate is highly unstable especially when the data is insufficient, and may even output negative dimensions.

Therefore, we seek a way to estimated $n(r)$ from $C(r)$, which can be more stably approximated from data by counting the number of points falling inside a neighbourhood of radius r ,

$$C(r) \leftarrow \hat{C}(r) = \frac{1}{N} \cdot \frac{1}{N-1} \sum_{i=1}^N \sum_{j=1, j \neq i}^N \chi_{\{\hat{D}(X_i, X_j) \leq r\}}. \quad (11)$$

At scale zero, as mentioned in Remark 4.3, $n(r)|_{r=0}$ corresponds to the well-known correlation dimension, and can be estimated by counting points in the ball $B(0, \Delta r)$ with various radius $\Delta r \rightarrow 0$ and then be solved from $C(\Delta r) \sim \Delta r^{n(0)}$ via a linear fitting procedure between $\log C$ and $\log \Delta r$ [3, 5, 12]. However, at non-zero scales, $n(r)|_{r=r_0 \neq 0}$ can only be estimated through counting points in a shell $B(r_0 + \Delta r) \setminus B(r_0)$ that is thin enough on which one can assume $n(r)$ to be a constant. Unlike the counting number of a ball, the counting number of a shell, denoted as $C_S(r)$ does not obey the simple relation $C \sim r^n$, but the following relation as obtained in Remark 4.6.

$$C_S(r) \equiv C(r + \Delta r) - C(r) \approx c(r + \Delta r)^{n(r)} - cr^{n(r)}. \quad (12)$$

Assumption 3 For some fixed $r > 0$ and $\Delta r_0 \ll r_0$, $n(r)$ is a constant on the interval $[r_0, r_0 + \Delta r_0]$.

For illustration purposes, we introduce our method under Assumption 3, but the proposed method works for slowly varying dimensions as well.

Under Assumption 3, (12) becomes an equality, ,

$$C_S(r) \equiv C(r + \Delta r) - C(r) = c(r + \Delta r)^n - cr^n \quad \text{for } [r, r + \Delta r] \subseteq [r_0, \Delta r_0] \quad (13)$$

where $n = n(r_0)$ is the constant dimension on the small interval. The counting number $C_S(r)$ of the shell $B(r_0 + \Delta r) \setminus B(r_0)$ on the left hand side of (13) can be estimated by its empirical version $\hat{C}(r + \Delta r) - \hat{C}(r)$, where \hat{C} was defined in (11). Now the problem boils down to: given $r > 0$, solving for n from (13) when Δr is some small constant chosen in advance and $C_S(r)$ has been estimated from the data (note that the constant c is still unknown). As mentioned above, this new problem is easier to solve when $r = 0$ than when $r > 0$, because when $r = 0$, we have $C(0) = 0$ which helps to reduce (13) to

$$C(\Delta r) = c(\Delta r)^n.$$

Then n equals to $\frac{\partial \log C(r)}{\partial \log r} \Big|_{r=0}$, which can be computed numerically via the finite difference [5]. However, this approach does not generalize to $r > 0$, as the presence of the second term cr^n in the right hand side of Equation (13) prevents taking logarithm from simplifying the equation. Without simplification, solving for n from (13) is difficult as c is unknown and this equation is highly non-linear. We propose to a way to reduce the nonlinearity to a degree that allows a meaningful solution to be found by fixed-point iterations. The following key observation (14) to be proven in Theorem 5.2 states that, with the help of the data, the term cr^n can be approximated by a linear function of $1/n$,

$$cr^n \approx \psi(r + \Delta r, r) \left(\frac{1}{n} - \Phi(r + \Delta r, r) \right) =: p(n, \Delta r), \quad (14)$$

where

$$\Phi(r', r) = \frac{1}{\hat{C}(r') - \hat{C}(r)} \int_r^{r'} \frac{\hat{C}(s) - \hat{C}(r)}{s} ds, \quad \psi(r', r) = \frac{\hat{C}(r') - \hat{C}(r)}{\log(r') - \log(r)},$$

are both empirical quantities that can be directly computed from the data. Here \approx means asymptotically equal as $N \rightarrow \infty$, where N is the number of samples. As a result, $p(n, \Delta r)$ is a linear function in $1/n$. Intuitively, this lower order $p(n, \Delta r)$ can approximate the exponential quantity cr^n because the empirical quantities Φ and ψ have taken into account all the high order information. Plugging (14) into (13), we obtain

$$\hat{C}(r + \Delta r) - \hat{C}(r) + p(n, \Delta r) \approx c(r + \Delta r)^n \quad (15)$$

or equivalently,

$$n \approx \frac{\partial_{\Delta r} \log [\hat{C}(r + \Delta r) - \hat{C}(r) + p(n, \Delta r)]}{\partial_{\Delta r} \log(r + \Delta r)} \Big|_{\Delta r=0}. \quad (16)$$

where $\partial_{\Delta r}$ denotes the derivative with respect to Δr . Let $q(n, \Delta r)$ be the finite difference approximation of the right hand side

$$q(n, \Delta r) := \frac{\log [\hat{C}(r + \Delta r) - \hat{C}(r) + p(n, \Delta r)] - \log [\hat{C}(r + \Delta r/2) - \hat{C}(r) + p(n, \Delta r/2)]}{\log(r + \Delta r) - \log(r + \Delta r/2)}$$

The following Proposition states that this resulting equation

$$n = q(n, \Delta r) \quad (17)$$

is close to quadratic provided that Δr is much smaller than r .

Proposition 5.1 *When $N \rightarrow \infty$ and $\Delta r \ll r$, the equation $n = q(n, \Delta r)$ defined in (17) is approximately linear in n .*

Therefore, we initialize n using the solution to the linear approximation and solve it via fixed point iterations $n^{k+1} = q(n^k)$.

Theorem 5.2 *Assume c and $n(r)$ are the c_m and $n_m(r)$ defined in Assumption 2, and $n(r) = n$ is a constant on the small interval $[r, r + \Delta r]$. Let p be the same as defined in (14), then for any $0 < \delta r \leq \Delta r$, we have*

$$cr^n = \lim_{N \rightarrow \infty} p(n, \delta r),$$

where N is the number of samples.

Note that the multi-scale *correlation* dimension is the object of study here as it is directly linked to the relative capacity. Our method cannot estimate the actual intrinsic dimension of the manifold [2].

The Assumption 3 used in Theorem 5.2 can be relaxed from constant dimension to slowly varying ones (Assumption 2).

Theorem 5.3 Assume c and $n(r)$ are the c_m and $n_m(r)$ defined in Assumption 2, and $n(r)$ is slowly varying on the small interval $[r, r + \Delta r_0]$, i.e., $n'(r) \ll 1$. Let p be the same as defined in (14), then

$$cr^{n(r)} = \lim_{\Delta r \rightarrow 0} \lim_{N \rightarrow \infty} p(n(r), \Delta r),$$

where N is the number of data points. In addition, $n(r) = q(n(r))$.

6 Numerical simulation

To evaluate the performance of the proposed Capacity Preserving Mapping (CPM) method, we first compare it with the landmark methods non-metric MDS, Isomap, and t-SNE on three datasets: 1) the motivating example introduced in Sect. 3, 2) the augmented Swiss roll to be defined shortly, and 3) the MNIST dataset.

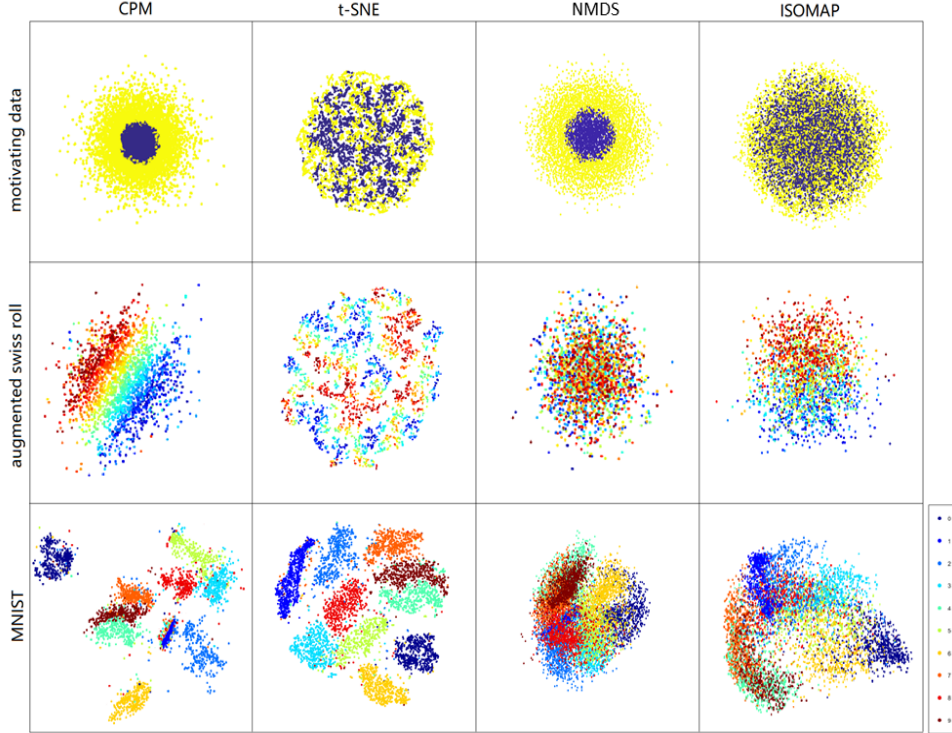


Figure 2: A comparison of CPM, t-SNE, NMDS and ISOMAP

The motivating example: in Sect. 3, we observed the crowding phenomenon when using MDS to reduce the dimension of two non-overlapping objects from 5 to 2, with one object being an ℓ_2 ball and the other being a shell lying right outside the ball. Figure 2 shows that the same crowding problem arises in t-SNE (perplexity = 30), and Isomap. In contrast, CPM and non-metric MDS are able to mitigate the crowding and reveal the correct relation between the classes. To further compare CPM and NMDS, we plot the Shepard diagrams, which show the goodness-of-fit after the dimension reduction. Figure 3 displays the Shepard diagrams of NMDS, CPM and t-SNE (perplexity 30, other perplexity values produce similar or worse results) when they map 5 dimensional (bottom row) and 20 dimensional (top row) Gaussian point clouds (each containing 1000 i.i.d. sampled points according to $N(0, I_d)$, with $d = 5, 20$, respectively) down to 2D. We see that t-SNE is only good at preserving small scale distances. NMDS is good at preserving medium and large distances, but is worse than t-SNE at small scales (as discussed in Sect. 4.3). CPM has a similar performance to t-SNE at small scales and a similar performance to NMDS at medium and large scales and its advantage becomes more visible as the dimension gets higher.

Augmented Swiss roll: the Swiss roll is a popular synthetic test dataset which can evaluate a method’s ability in preserving the geometric structure. However, the fact that its intrinsic dimension is only 2 makes it less representative. In order to test the ability of our algorithm to map the data to below its intrinsic dimension, we construct the augmented Swiss roll dataset $X = [x_1, \dots, x_p]$, where the first three coordinates $[x_1, x_2, x_3]$ are exactly the same as those in the original Swiss

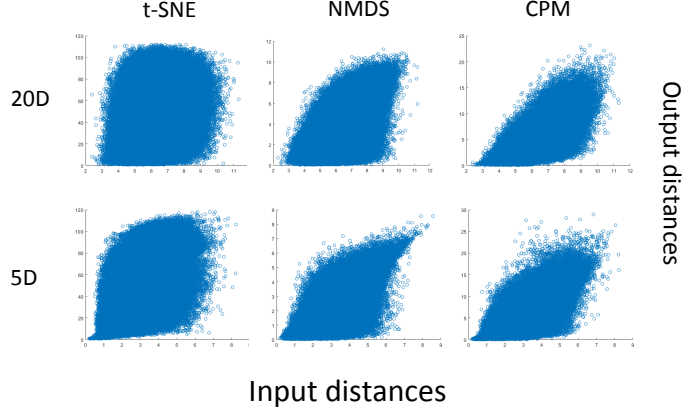


Figure 3: Shepard diagrams of t-SNE, NMDS and CPM applied to Gaussian point clouds

roll, and the rest of the coordinates x_4, \dots, x_p are filled with i.i.d. Gaussian entries. Therefore the intrinsic dimension of X is $p - 1$. Here we set $p = 6$ and map the data to 2D. Explicitly, we set

$$\begin{aligned} x_1(t) &= (t + 1) \cos(t), \\ x_2(t) &= (t + 1) \sin(t), \\ x_j(t) &= g_j(t), \quad j = 3, \dots, 6 \end{aligned}$$

where t takes values in a finite grid of the interval $[0, 1]$, $g_j(t) \sim 5N(0, 1)$ for any j and t , and $g_j(t)$ are independent among various values of j and t .

We apply to this dataset Isomap, NMDS with the geodesic distance, and CPM with the geodesic distance, all with the same number of neighbours ($=10$). From the second row of Figure 2, we see that in terms of unfolding the manifold, CPM did the best job among all. The result of t-SNE (perplexity 30) is also included for comparison.

MNIST dataset We repeat the previous experiment on the MNIST data set to test the method’s ability in preserving and revealing clusters. The MNIST dataset contains 60000 training images of handwritten digits. We apply t-SNE, Isomap, NMDS with geodesic distance and CPM with geodesic distance to a subset of 6000 randomly selected images from the training set. Figure 2 (last row) shows that CPM separates the clusters equally well as t-SNE, while revealing some new structures. For example, it shows that the cluster of the digit 1 has the smallest variance, which is consistent with our intuition that the handwritten digit 1 has the least variation among different writing styles. To confirm this observation, we computed the variance of each cluster based on the original data and obtained Table 1. We can see that the cluster of the digit 1 indeed has a much smaller variance than all other clusters. The second and third smallest clusters digit 7 and digit 9 also appear smaller than others in the visualization. Another piece of information conveyed by the CPM visualization is that the cluster of digit 1 is close to a lot of other clusters, which is aligned with the intuition that the digit 1 look similar to 2, 3, 7, and maybe 9.

Table 1: Variances of clusters in MNIST dataset (normalized)

Digit	0	1	2	3	4	5	6	7	8	9
Variance	0.991	0.448	1.000	0.880	0.824	0.952	0.871	0.750	0.891	0.756

6.1 Coil 20

The Coil 20 dataset contains images of 20 objects captured from different angles while they rotate. Previous methods are only aiming at separating the images into 20 clusters. While we also care about the formation of clusters, we are more interested in examining how well the geometric features are correctly displayed. In particular, we check the following properties of the embedding,

1. objects with large variances during rotations should correspond to clusters with large sizes in the visualization;

2. objects similar to each other should correspond to clusters close to each other in the visualization;
3. if an object is symmetric with respect to its center, then the corresponding data should form a trajectory similar to a folded circle in the visualization;
4. if an object is nearly isometric (looks similar from all angles), then its corresponding cluster should have a small size.

We now evaluate the performance of CPM based on these four criteria.

1. Table 2 summarizes the variance of each object as it rotates. The variance V_i for the i th object is computed using the formula

$$V_i = \frac{1}{n_i} \sum_{j=1}^{n_i} \|X_j - \bar{X}_i\|_2^2$$

where n_i is the number of points in the i th cluster, $X_j, j = 1, \dots, n_i$ are vectorized images of the i th object, and $\bar{X}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} X_j$ is the mean. Table 3 and Figure 4 together confirm that objects with large variances during rotations corresponding to clusters with large sizes in the visualization.

2. To characterize the preservation of the inter-cluster distances of the 20 classes/objects in Coil 20, we propose to use the distance index defined as follows. For a given cluster say Cluster i , we compute its distances to all other 19 clusters (the distance between two clusters are defined as the averaged pairwise distances between one cluster and the other), rank the 19 distances and find the p th percentile for some given p . Those clusters within the p th percentile are said to be close to Cluster i , otherwise is said to be far to Cluster i . After doing this for each $i = 1, \dots, 20$, we build a 20×20 proximity matrix C . If Cluster j is close to Cluster i , then $C_{i,j} = 1$, otherwise $C_{i,j} = 0$. We compute the proximity matrix of the original data (obtain C^{orig}) and that of the embedded data (obtain C^{new}), and then compute the dissimilarity between C^{orig} and C^{new} :

$$E(p) = \frac{\sum_{i,j} 1_{\{C_{i,j}^{orig}=1, C_{i,j}^{new}=0\}}}{\sum_{i,j} 1_{\{C_{i,j}^{orig}=1\}}}$$

That is to say, we compute the percentage of close clusters that are no longer close after the embedding. For any predefined percentile p , such a dissimilarity $E(p)$ can be computed. Let p range from 0%-50%, we derive the percentile versus dissimilarity plot in Figure 5 for the three embedding methods CPM with Euclidean distance, t-SNE (perplexity 30) and MDS. Clearly, MDS is good at preserving large distances, our method is better at preserving smaller distances, and t-SNE is not good at preserving inter-cluster distances.

3. Table 5 shows that the CPM embedding indeed produces folded circles for symmetric objects.

4. Table 4 shows that the CPM embedding indeed produce small clusters for the isotropic objects. In addition, Table 6 also suggested the trajectory of these objects are (unfolded) circles.

Table 2: Variances of clusters in the 2D visualization of the COIL 20 dataset

Label	1	2	3	4	5	6	7	8	9	10
Variance	51.19	69.92	38.42	32.83	42.55	48.58	25.39	12.34	45.69	27.34
Label	11	12	13	14	15	16	17	18	19	20
Variance	28.31	3.45	31.72	21.10	2.03	1.90	2.97	10.72	43.18	8.37

Table 3: Objects with large variances


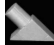






Objects				
Labels	19	2	5	13
Variance	43.18	69.92	42.55	31.72

Table 4: Objects with small variances

Objects				
Labels	17	15	12	16
Variance	2.97	2.03	3.45	1.90

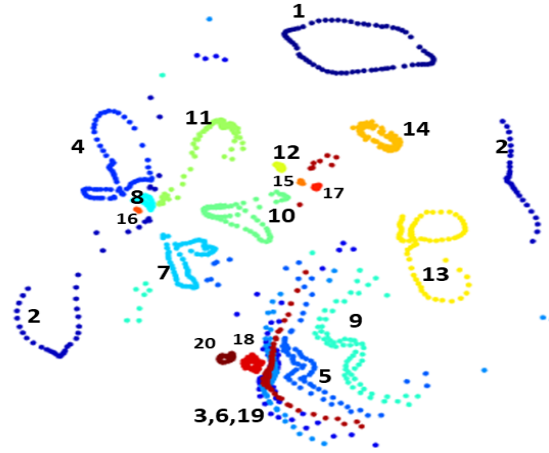


Figure 4: Visualization by CPM.

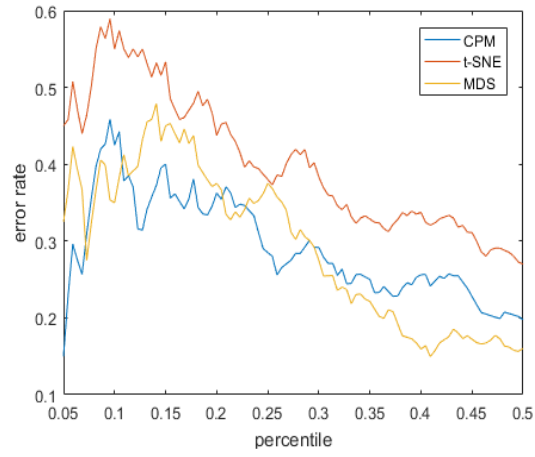


Figure 5: Error rates in preserving relative cluster proximity by CPM,t-SNE and MDS

Table 5: Visualizations of symmetric objects








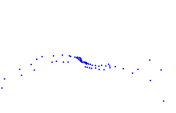
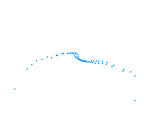
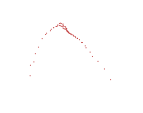





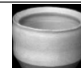

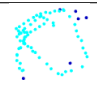






Objects					
Visualizations					

Table 6: Visualizations of isotropic objects

Objects							
Visualizations							

7 Conclusion and future directions

In this paper, we asked the question of how to rigorously characterize and treat the intrinsic crowding issue in data visualization. After giving a mathematical notation to the capacity, we discussed two possible directions to mitigate the crowding: altering the density or altering the distance. We chose the latter for simplicity, but it will be interesting to explore the former as well.

After the Capacity Adjusted Distance was defined, we proposed to find the low dimensional embedding by matching the dissimilarity measured by the KL divergence. There are many other dissimilarity measures in the literature. The performance of the combination of the Capacity Adjusted Distance with other measures is yet to be explored.

The definition of \tilde{D} requires the knowledge of the dimension functions $n_s(r)$ which cannot be estimated directly from data since Y_i is not readily known. We made the assumption that the induced measure on the low dimension manifold is uniform and hence $n_s(r) = d$. A better way to obtain $n_s(r)$ is to initialize it with $n_s(r) = d$, and update it once the low dimensional embedding $\{Y_i\}_{i=1}^N$ is obtained. These two steps (updating Y_i and updating n_s) can be performed iteratively.

In the numerical experiments, we only tested CPM under Euclidean and Geodesic distances. It would be interesting to see how it works with other distances (e.g., the diffusion distance).

From the discussion in Sect. 4.2. The preservation of geometry essentially corresponds to allowing different row sums of the probability matrix. Our method achieves it by the universal normalization, but there are other ways to achieve it. That is to say, although t-SNE does not admit the universal normalization, there might be other remedies to bring back the geometric information.

8 References

- [1] M. Belkin and P. Niyogi. Laplacian Eigenmaps and spectral techniques for embedding and clustering. In *Advances in Neural Information Processing Systems*, volume 14, pages 585–591, Cambridge, MA, USA, 2002. The MIT Press.
- [2] F. Camastra and A. Vinciarelli. Estimating the intrinsic dimension of data with a fractal-based method. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 24, no. 10, pp. 1404–1407, 2002.
- [3] K. Falconer. *Fractal Geometry - Mathematical Foundations and Applications*, John Wiley and Sons, 3rd edition, 2014.
- [4] M.C. Ferreira de Oliveira and H. Levkowitz. From visual data exploration to visual data mining: A survey. *IEEE Transactions on Visualization and Computer Graphics*, 9(3): 378–394, 2003.
- [5] P. Grassberger and I. Procaccia. Measuring the strangeness of strange attractors. *Physica D. Nonlinear Phenomena*, vol. 9, no. 1-2, pp. 189–208, 1983.
- [6] G.E. Hinton and S.T. Roweis. Stochastic Neighbor Embedding. In *Advances in Neural Information Processing Systems*, volume 15, pages 833–840, Cambridge, MA, USA, 2002. The MIT Press.
- [7] A. Hyvarinen and E. Oja. Independent Component Analysis: Algorithms and Applications. *Neural Networks*, 13(4-5):411–430, 2000.
- [8] S. Lafon and A.B Lee. Diffusion maps and Coarse-Graining: A Unified Framework for Dimensionality Reduction, Graph Partitioning and Data Set Parametrization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28(9):1393–1403, 2006.
- [9] L. McInnes and J. Healy. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. *Arxiv e-prints*, February 2018.
- [10] Moon K.R. et al. (2017) PHATE: A Dimensionality Reduction Method for Visualizing Trajectory Structures in High-Dimensional Biological Data. *bioRxiv* Published online March 24, 2017. <http://dx.doi.org/10.1101/120378>
- [11] S.T. Roweis and L.K. Saul. Nonlinear dimensionality reduction by Locally Linear Embedding. *Science*, 290(5500):2323–2326, 2000.
- [12] N. Tatti, T. Mielikäinen, A. Gionis, and H. Mannila. What is the dimension of your binary data? In *Proceedings of the 6th International Conference on Data Mining (ICDM)*, pp. 603–612, December, 2016.
- [13] J.B. Tenenbaum, V. de Silva, and J.C. Langford. A global geometric framework for nonlinear dimensionality reduction. *Science*, 290(5500):2319–2323, 2000.
- [14] W.S. Torgerson. Multidimensional scaling I: Theory and method. *Psychometrika*, 17:401–419, 1952.
- [15] L. van der Maaten, E. Postma, and J. van der Herik. Dimensionality Reduction: A Comparative Review. *J Mach Learn Res*, 10:66–71, 2009.
- [16] L. Van der Maaten and G. Hinton. Visualizing data using t-SNE. *Journal of Machine Learning Research*, 9(2579–2605): 85, 2008 .
- [17] J. A. Lee, D. H. Peluffo-Ordóñez, and M. Verleysen . Multi-scale similarities in stochastic neighbour embedding: Reducing dimensionality while preserving both local and global structure. *Neurocomputing*, 169, 246–261, 2015.
- [18] G. B. Rabinowitz. An introduction to nonmetric multidimensional scaling. *American Journal of Political Science*, 343–390, 1975.

9 Appendix

9.1 Proof of Theorem 4.7

Proof: The assumption that the embedded pairwise distances can be written as a function of the pairwise distances in the original space, i.e., $\|P(x) - P(z)\|_2 = G(\|x - z\|_2)$ (for some G), allows us to define the new distance $\hat{D}(x, z)$ to be a function of $\|x - z\|_2$ only. We will show that if \hat{D} is defined as in the theorem, i.e., $\hat{D}(\|x - z\|_2) = \|x - z\|_2^{\frac{n_m(\|x - z\|_2)}{n_s(\|x - z\|_2)}}$, then the capacity is preserved. Without any ambiguity from the context, we write $n_m(\|x - z\|_2)$, $n_s(\|x - z\|_2)$ in short as n_m and n_s . By the definition of the relative capacity, for any $\tilde{r} > 0$, and $d\tilde{r} \ll \tilde{r}$ (i.e., the interval $[\tilde{r}, \tilde{r} + d\tilde{r}]$ is small), n_m and n_s can be treated as constants and

$$\begin{aligned} C(\tilde{r} + d\tilde{r}; \mathcal{M}, f, \hat{D}) - C(\tilde{r}; \mathcal{M}, f, \hat{D}) &= \mathbb{P}_{x, z \sim f(\mathcal{M})}(\tilde{r} \leq \hat{D}(\|x - z\|_2) \leq \tilde{r} + d\tilde{r}) \\ &= p_{x, z \sim f(\mathcal{M})}(\hat{D}(\|x - z\|_2) = \tilde{r})d\tilde{r} \\ &= \rho(\tilde{r}; \mathcal{M}, f, \hat{D})d\tilde{r} \end{aligned}$$

where $p_{x, z}$ denotes the probability density function with respect to the random variables x and z . The first and last equalities used the definitions of $C(r; \mathcal{M}, f, \hat{D})$ and $\rho(\tilde{r}; \mathcal{M}, f, \hat{D})$. Hence the relative capacities match if the relative densities match. We will show that relative densities match by showing

$$\rho(\tilde{r}; \mathcal{M}, f, \hat{D})d\tilde{r} = cn_s \tilde{r}^{n_s - 1} d\tilde{r}, \quad (18)$$

because the right hand side is exactly the relative density for the low dimensional embedding $\rho(\tilde{r}; \mathcal{S}, f_I, \|\cdot\|_2)$ by Assumption 2. To compute the left hand side, we define $r = \tilde{r}^{\frac{n_s}{n_m}}$, then $\hat{D}(\|x - z\|_2) = \tilde{r}$ is equivalent to $\|x - z\|_2 = r$. Hence

$$\rho(\tilde{r}; \mathcal{M}, f, \hat{D})d\tilde{r} = p_{x, z \sim f(\mathcal{M})}(\hat{D} = \tilde{r})d\tilde{r} = p_{x, z \sim f(\mathcal{M})}(\|x - z\|_2 = r)dr \quad (19)$$

This last term can be easily calculated using Assumption 2,

$$p_{x, y \sim f(\mathcal{M})}(\|x - z\|_2 = r) = \rho(r; \mathcal{M}, f, \|\cdot\|_2) = cn_m r^{n_m - 1}$$

Inserting this into (19) we obtain

$$\rho(\tilde{r}; \mathcal{M}, f, \hat{D}) = f_{x, y \sim f(\mathcal{M})}(\|x - z\|_2 = r) \cdot \frac{dr}{d\tilde{r}} = cn_m r^{n_m - 1} \cdot \frac{dr}{d\tilde{r}} = cn_m r^{n_m - 1} \cdot \frac{n_s}{n_m} \tilde{r}^{\frac{n_s}{n_m} - 1} = cn_s \tilde{r}^{n_s - 1}$$

Hence we proved (18). \square

9.2 Proof of Theorem 5.2

Proof: It is immediate to verify that the empirical counting number $\hat{C}(r)$ defined by

$$\hat{C}(r) = \frac{1}{N} \cdot \frac{1}{N-1} \sum_{i=1}^N \sum_{j=1, j \neq i}^N \chi_{\{\hat{D}(X_i, X_j) \leq r\}}$$

is a consistent estimate of $C(r)$, i.e., $C(r) = \lim_{N \rightarrow \infty} \hat{C}(r)$. Therefore, by letting $N \rightarrow \infty$, we can replace all $\hat{C}(r)$ in the definition of p by $C(r)$ and therefore we only need to verify

$$cr^n = \tilde{p}(n, \delta r) \quad (20)$$

where $\tilde{p} = \lim_{N \rightarrow \infty} p$. In terms of the limits $\tilde{\Phi} = \lim_{N \rightarrow \infty} \Phi$ and $\tilde{\psi} = \lim_{N \rightarrow \infty} \psi$, \tilde{p} can be written as

$$\tilde{p}(n, \Delta r) = \tilde{\psi}(r + \Delta r, r) \left(\frac{1}{n} - \tilde{\Phi}(r + \Delta r, r) \right) \quad (21)$$

where

$$\tilde{\Phi}(r', r) = \frac{1}{C(r') - C(r)} \int_r^{r'} \frac{C(s) - C(r)}{s} ds, \quad \tilde{\psi}(r', r) = \frac{C(r') - C(r)}{\log(r') - \log(r)}. \quad (22)$$

To verify (20), let us first compute $\tilde{\Phi}(r', r)$

$$\begin{aligned}
\tilde{\Phi}(r', r) &= \frac{1}{C(r') - C(r)} \int_r^{r'} \frac{C(s) - C(r)}{s} ds \\
&= \frac{1}{C(r') - C(r)} \left[(C(r') - C(r)) \log(r') - \int_r^{r'} \log(s) \frac{dC(s)}{ds} ds \right] \\
&= \frac{1}{C(r') - C(r)} \left[(C(r') - C(r)) \log(r') - c \int_r^{r'} \log(s) n s^{n-1} ds \right] \\
&= \frac{1}{C(r') - C(r)} \left[(C(r') - C(r)) \log(r') - c(r')^n \log(r') + cr^n \log(r) + c \int_r^{r'} s^{n-1} ds \right] \\
&= -cr^n \frac{\log(r') - \log(r)}{C(r') - C(r)} + \frac{c \int_r^{r'} s^{n-1} ds}{C(r') - C(r)} \\
&= -cr^n \frac{\log(r') - \log(r)}{C(r') - C(r)} + \frac{c \int_r^{r'} s^{n-1} ds}{c \int_r^{r'} n s^{n-1} ds} \\
&= -cr^n \tilde{\psi}^{-1}(r', r) + \frac{1}{n}.
\end{aligned}$$

In the above calculation, we only used integration by part and the fact that $C(r') - C(r) = c \int_r^{r'} n s^{n-1} ds$. Rearranging the equation, we obtain

$$\tilde{p}(n, \Delta r) \equiv \tilde{\psi}(r', r) \left(\frac{1}{n} - \tilde{\Phi}(r', r) \right) = cr^n.$$

□

9.3 Proof of Proposition 5.1

Proof: We only need to examine the limiting behaviour of the equation $n = q(n, \Delta r)$. Recall that $q(n)$ is

$$q(n, \Delta r) := \frac{\log [\hat{C}(r + \Delta r) - \hat{C}(r) + p(n, \Delta r)] - \log [\hat{C}(r + \Delta r/2) - \hat{C}(r) + p(n, \Delta r/2)]}{\log(r + \Delta r) - \log(r + \Delta r/2)} \quad (23)$$

and $p(n, \Delta r)$ is

$$p(n, \Delta r) := \psi(r + \Delta r, r) \left(\frac{1}{n} - \Phi(r + \Delta r, r) \right). \quad (24)$$

Let \tilde{p} and \tilde{q} be the limit of p and q as $N \rightarrow \infty$, i.e,

$$\tilde{p}(n, \Delta r) = \tilde{\psi}(r + \Delta r, r) \left(\frac{1}{n} - \tilde{\Phi}(r + \Delta r, r) \right)$$

and

$$\tilde{q}(n, \Delta r) := \frac{\log [C(r + \Delta r) - C(r) + \tilde{p}(n, \Delta r)] - \log [C(r + \Delta r/2) - C(r) + \tilde{p}(n, \Delta r/2)]}{\log(r + \Delta r) - \log(r + \Delta r/2)}$$

where $\tilde{\psi}$ and $\tilde{\Phi}$ are the same as in (22). Since we proved in Theorem 5.2 that

$$\tilde{p}(n, \delta r) = cr^n$$

this immediately implies that $\lim_{\Delta r \rightarrow 0} \tilde{q}(n, \Delta r) = n$. Therefore, for sufficiently small Δr , it makes sense to solve n from $\tilde{q}(n, \Delta r) = n$. To prove the equation $\tilde{q}(n, \Delta r) = n$ is approximately linear in n , we truncate the Taylor expansion of $\tilde{q}(n, \Delta r)$ after the first term,

$$\begin{aligned}
\tilde{q}(n, \Delta r) &= \frac{\log [C(r + \Delta r) - C(r) + \tilde{p}(n, \Delta r)] - \log [C(r + \Delta r/2) - C(r) + \tilde{p}(n, \Delta r/2)]}{\log(r + \Delta r) - \log(r + \Delta r/2)} \\
&= \frac{1}{\gamma_1} \log \left[1 + \frac{C(r + \Delta r) - C(r + \Delta r/2)}{C(r + \Delta r/2) - C(r) + \tilde{p}(n, \Delta r/2)} \right] \quad \left(\gamma_1 := \log(1 + \frac{\Delta r}{2r + \Delta r}) \right) \\
&= \frac{1}{\gamma_1} \log \left[1 + \alpha \cdot \frac{\gamma_2}{\gamma_2 + \frac{1}{n} - \tilde{\Phi}} \right] \quad \left(\alpha = \frac{C(r + \Delta r) - C(r + \Delta r/2)}{C(r + \Delta r/2) - C(r)}, \gamma_2 = \log(1 + \frac{\Delta r}{2r}) \right) \\
&= \frac{\beta}{\gamma_2 + \frac{1}{n} - \tilde{\Phi}} + o((\Delta r)^2), \quad (\beta = \alpha \gamma_2 / \gamma_1) \quad (25)
\end{aligned}$$

where the second line used the fact that $\tilde{p}(n, \Delta r) = \tilde{p}(n, \Delta r/2)$ (Theorem 5.2). We obtained the third line by inserting the definition of \tilde{p} , and the last line is obtained by keeping the first term in the Taylor expansion

of the logarithm $\log(1+x) \approx x$, where $x = \frac{C(r+\Delta r) - C(r+\Delta r/2)}{C(r+\Delta r/2) - C(r) + \tilde{p}(n, \Delta r/2)} = \frac{\alpha\gamma_2}{\gamma_2 + \frac{1}{n} - \tilde{\Phi}}$. The approximation has a small error if $x \ll 1$, this is indeed the case here as the numerator of x can be written as

$$C(r + \Delta r) - C(r + \Delta r/2) = c(r + \Delta r)^n - c(r + \Delta r/2)^n$$

due to (13), the denominator is greater than

$$\tilde{p}(r, \Delta r) = c(r)^n$$

due to (13) and (15). Clearly the denominator is much larger than the numerator as long as $\Delta r \ll r$ (We comment here that ideally, Δr can be taken as small as we want. However, when the data is not dense enough, taking Δr too small might result in no point falling inside the interval $[r, r + \Delta r]$. Hence in the numerical simulation, Δr is taken based on some small percentile of the sorted pairwise distances of the data). Inserting (25) into $\tilde{q}(n, \Delta r) = n$, we have

$$n = \frac{\beta}{\gamma_2 + \frac{1}{n} - \tilde{\Phi}} + O(\Delta r).$$

By the definition of $\tilde{\Phi}$ in (22), we know that $\tilde{\Phi}(r + \Delta r/2, r) = \frac{\Delta r}{4r} + o(\Delta r)$. It is also easy to verify that $\alpha = 1 + O(\Delta r)$, $\beta = 1 + O(\Delta r)$, $\gamma_2 = \frac{\Delta r}{2r} + o(\Delta r)$. Rearranging the above equation, we get the following quadratic equation of n

$$n = \frac{(\beta - 1)}{\gamma_2 - \tilde{\Phi}} + O(\Delta r),$$

For sufficiently small Δr , so the equation has one positive root. As $\Delta r \rightarrow 0$, the equation becomes linear. \square

9.4 Proof of Theorem 5.3

Proof: The proof is similar to that of Theorem 5.2. Let \hat{C} is the counting number defined by

$$\hat{C}(r) = \frac{1}{N} \cdot \frac{1}{N-1} \sum_{i=1}^N \sum_{j=1, j \neq i}^N \chi_{\{\hat{D}(X_i, X_j) \leq r\}}.$$

Then $C(r) = \lim_{N \rightarrow \infty} \hat{C}(r)$. Therefore, by letting $N \rightarrow \infty$, the limit \tilde{p} of p is

$$\tilde{p}(n(r), \Delta r) = \lim_{N \rightarrow 0} p(n(r), \Delta r) = \tilde{\psi}(r + \Delta r, r) \left(\frac{1}{n(r)} - \tilde{\Phi}(r + \Delta r, r) \right) \quad (26)$$

where

$$\tilde{\Phi}(r', r) = \frac{1}{C(r') - C(r)} \int_r^{r'} \frac{C(s) - C(r)}{s} ds, \quad \tilde{\psi}(r', r) = \frac{C(r') - C(r)}{\log(r') - \log(r)}.$$

In order to prove the theorem, we need to verify

$$cr^{n(r)} = \lim_{\Delta r \rightarrow 0} \tilde{p}(n(r), \Delta r).$$

To that end, we compute $\tilde{\Phi}(r', r)$

$$\begin{aligned} \tilde{\Phi}(r', r) &= \frac{1}{C(r') - C(r)} \int_r^{r'} \frac{C(s) - C(r)}{s} ds \\ &= \frac{1}{C(r') - C(r)} \left[(C(r') - C(r)) \log(r') - \int_r^{r'} \log(s) \frac{dC(s)}{ds} ds \right] \\ &= \frac{1}{C(r') - C(r)} \left[(C(r') - C(r)) \log(r') - c \int_r^{r'} \log(s) n(s) s^{n(s)-1} ds \right] \\ &= \frac{1}{C(r') - C(r)} \left[(C(r') - C(r)) \log(r') - c \int_r^{r'} \log(s) n(r) s^{n(r)-1} ds \right] + R \end{aligned}$$

where

$$R = \frac{c}{C(r') - C(r)} \int_r^{r'} \log(s) \left(n(r) s^{n(r)-1} - n(s) s^{n(s)-1} \right) ds = n'(r) R_1 + R_2$$

with $R_1 = O(\Delta r)$ and $R_2 = o(\Delta r)$. Continuing the above calculation

$$\begin{aligned}
\tilde{\Phi}(r', r) &= \frac{1}{C(r') - C(r)} \left[(C(r') - C(r)) \log(r') - c \int_r^{r'} \log(s) n(r) s^{n(r)-1} ds \right] + R \\
&= \frac{1}{C(r') - C(r)} \left[(C(r') - C(r)) \log(r') - c(r')^{n(r)} \log(r') + cr^{n(r)} \log(r) + c \int_r^{r'} s^{n(r)-1} ds \right] + n'(r)R_1 + R_2 \\
&= -cr^{n(r)} \frac{\log(r') - \log(r)}{C(r') - C(r)} + \frac{c \int_r^{r'} s^{n(r)-1} ds}{C(r') - C(r)} + n'(r)R_1 + R_2 \\
&= -cr^{n(r)} \frac{\log(r') - \log(r)}{C(r') - C(r)} + \frac{c \int_r^{r'} s^{n(r)-1} ds}{c \int_r^{r'} n(r) s^{n(r)-1} ds} + n'(r)R_1 + o(\Delta r) \\
&= -cr^{n(r)} \tilde{\psi}^{-1}(r', r) + \frac{1}{n(r)} + n'(r)R_1 + o(\Delta r).
\end{aligned}$$

Rearranging the equation and letting $r' \rightarrow r$ we obtain

$$\lim_{\Delta r \rightarrow 0} p(n(r), \Delta r) \equiv \lim_{r' \rightarrow r} \tilde{\psi}(r', r) \left(\frac{1}{n(r)} - \tilde{\Phi}(r', r) \right) = cr^{n(r)}.$$

When $n'(r) \ll 1$, one can also verify that Proposition 5.1 approximately holds using the same prove as in Section 9.3. \square