

RESEARCH ARTICLE

Clustering and visualization of single-cell RNA-seq data using path metrics

Andriana Manousidaki¹*, Anna Little²*, Yuying Xie^{1,3}*

1 Department of Statistics and Probability, Michigan State University, East Lansing, Michigan, United States of America, **2** Department of Mathematics, University of Utah, Salt Lake City, Utah, United States of America, **3** Department of Computational Mathematics, Science and Engineering, Michigan State University, East Lansing, Michigan, United States of America

* These authors contributed equally to this work.

* little@math.utah.edu (AL); xyy@msu.edu (YX)



OPEN ACCESS

Citation: Manousidaki A, Little A, Xie Y (2024) Clustering and visualization of single-cell RNA-seq data using path metrics. PLoS Comput Biol 20(5): e1012014. <https://doi.org/10.1371/journal.pcbi.1012014>

Editor: Shihua Zhang, Chinese Academy of Science, CHINA

Received: September 22, 2023

Accepted: March 21, 2024

Published: May 29, 2024

Copyright: © 2024 Manousidaki et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: The Cellmix and RNAmix data are downloaded from GEO under accession code GSE118767, and the preprocessed data are available at https://github.com/LuyiTian/CellBench_data. The PBMC4K data is available at 10x Genomics's website through <https://support.10xgenomics.com/single-cell-gene-expression/datasets/2.1.0/pbmc4k>. The Baron's pancreatic data is available in GEO with the access code GSM2230757. The simulated data were created based on the Baron's data. Simulation code is provided in <https://github.com/andrianamanousidaki/scPMP>. The mouse tissue

Abstract

Recent advances in single-cell technologies have enabled high-resolution characterization of tissue and cancer compositions. Although numerous tools for dimension reduction and clustering are available for single-cell data analyses, these methods often fail to simultaneously preserve local cluster structure and global data geometry. To address these challenges, we developed a novel analyses framework, Single-Cell Path Metrics Profiling (scPMP), using power-weighted path metrics, which measure distances between cells in a data-driven way. Unlike Euclidean distance and other commonly used distance metrics, path metrics are density sensitive and respect the underlying data geometry. By combining path metrics with multidimensional scaling, a low dimensional embedding of the data is obtained which preserves both the global data geometry and cluster structure. We evaluate the method both for clustering quality and geometric fidelity, and it outperforms current scRNAseq clustering algorithms on a wide range of benchmarking data sets.

Author summary

Advancements in single-cell technologies with the ability to measure gene expression at the cellular level have provided unprecedented opportunity to investigate the cell type (T cells, B cells, etc) and cell state diversity (active T cells and exhausted T cells) within tissues and cancers. However, analyzing this complex high-dimensional data when the noise level is high requires sophisticated tools to effectively extract useful biological information and faithfully visualize the data in a low-dimensional space (2- or 3-D). Existing computational methods such as dimension reduction and clustering (group similar cells together) for single-cell data struggle to simultaneously preserve local group structure and global data geometry (developmental relationship between cell types). To tackle this problem, we've developed a new analysis framework called scPMP (Single-Cell Path Metrics Profiling) based on a unique approach to measure distances between cells which takes into account both the density of cells (common vs rare cell types) and the overall structure of the data. We have demonstrated the ability of scPMP to better preserve the natural

scRNAseq data sets are accessible on https://figshare.com/projects/Tabula_Muris_Transcriptomic_characterization_of_20_organ_and_tissues_from_Mus_musculus_at_single_cell_resolution/27733. The code to reproduce all reported results and generate figures is available at the scPMP github repository via <https://github.com/andrianamanousidaki/scPMP>. The repository also contains a tutorial for the scPMP algorithm. Small differences observed during the reproduction of results is due to randomness introduced at the imputation step of data preprocessing.

Funding: This work was supported by the National Institutes of Health (U01DE029255, U01DE033330, R01HL166508, R01DE026728 and R03DE027399 to YX) and the National Science Foundation grants (IOS2107215 and IIS2123260 to YX; DMS1902906, DMS2131292 and DMS 2309570 to AL; DGE828149 to AM). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing interests: The authors have declared that no competing interests exist.

grouping of cells and the relationships between different groups over existing methods in numerous real and simulated data sets. This improvement could lead to more accurate identification of cell types and states.

1 Introduction

The advance in single-cell RNA-seq (scRNA-seq) technologies in recent years has enabled the simultaneous measurement of gene expression at the single-cell level [1–3]. This opens up new possibilities to detect previously unknown cell populations, study cellular development and dynamics, and characterize cell composition within bulk tissues. Despite its similarity with bulk RNAseq data, scRNAseq data tends to have larger variation and larger amounts of missing values due to the low abundance of initial mRNA per cell. To address these challenges, numerous computational algorithms have been proposed focusing on different aspects. Given a collection of single cell transcriptomes from scRNAseq, one of the most common applications is to identify and characterize subpopulations, e.g., cell types or cell states. Numerous clustering approaches have been developed such as *k*-means based methods SC3 [4], SIMLR [5], and RaceID [6]; hierarchical clustering based methods CIDR [7], BackSPIN [8], and pcaReduce [9]; graph based methods Rphenograph [10], SNN-Cliq [11], Seurat [12], SSNN-Louvain [13], and scanpy [14]; and deep-learning based methods scGNN [15], scVI [16], ScDeepCluster [17], DANCE [18], graph-sc [19], GraphSCC [20], scDCC [21], DESC [22], scDHA [23], scziDesk [24], scDSC [25], CELLPLM [26], scDiff [27], scMoGNN [28], scFormer [29] and scTAG [30] as summarized in [31].

To visualize and characterize relationships between cell types, it is important to represent it in a low-dimensional space. Many low-dimensional embedding methods have been proposed including UMAP [32], *t*-SNE [33], PHATE [34], and LargeVis [35]. However, a key challenge for embedding methods is to simultaneously reduce cluster variance and preserve the global geometry, including the distances between clusters and cluster shapes. For example, on a cell mixture dataset [36]: the PCA embedding preserves the global geometry but clusters have high variance; clusters are better separated in the UMAP and *t*-SNE embeddings, but the global geometric structure of the clusters is lost as shown in the result section.

When choosing a clustering algorithm, there is always an underlying tension between respecting data density and data geometry. Density based methods such as DBSCAN [37, 38] cluster data by connecting together high density regions, regardless of cluster geometry. More traditional approaches such as *k*-means require that clusters are convex and geometrically well separated. However, in many real data, clusters tend to have both nonconvex/elongated geometry and a lack of robust density separation as shown in Fig 1B which consists of three elongated Gaussian distributions and a bridge connecting two of the distributions. The data set is challenging because it exhibits elongated geometry, but methods relying only on density will fail due to the bridge. Such characteristics are commonly observed in scRNA-seq data, especially for cells sampled from a developmental process, as cell types often trace out elongated structures and frequently lack robust density separation. This elongated geometry phenomena is due to the fact that all the cell types originate from stem cells through a trajectory-like differentiation process, and the bridge structures are created by the cells in the transition states. For example, circulating monocytes in the Tabula Muris (TM) lung data set [39] have an elongated cluster structure as illustrated by the PCA plot in Fig 2A, as do the ductal cells in the TM pancreatic data set (see Fig 2C). The UMAP plots of these same data sets illustrate the lack of robust density separation: for TM lung, there is a bridge connecting the alveolar and lung cell

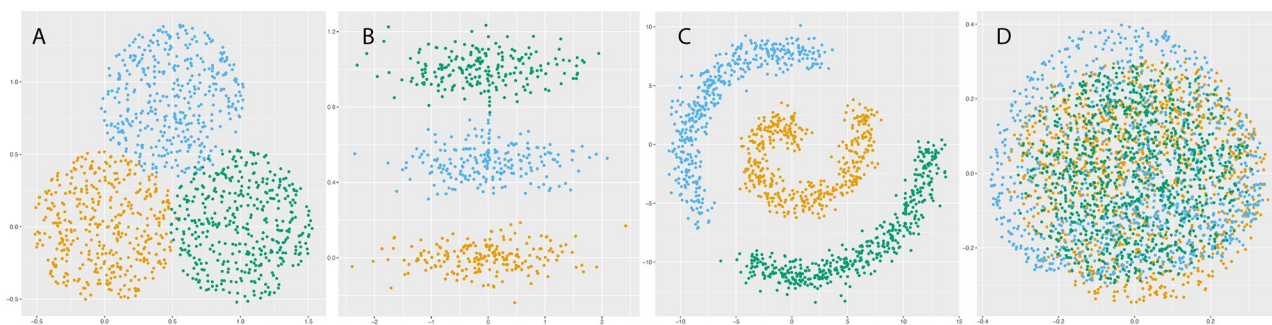


Fig 1. Toy data sets. (A) Balls; (B) elongated with bridge; (C) swiss roll; and (D) GL manifold. (A) and (B) show the 2-dimensional data sets. (C) plots the first two coordinates of the Swiss roll. (D) shows the 2-dimensional PCA plot of the SO(3) manifolds.

<https://doi.org/10.1371/journal.pcbi.1012014.g001>

types, and also an overlap/bridge between the circulating and invading monocytes (see Fig 2B); for TM pancreatic, the pancreatic A and pancreatic PP cells are not well separated. The combination of elongation and poor density separation make clustering scRNA-seq data sets a challenging task.

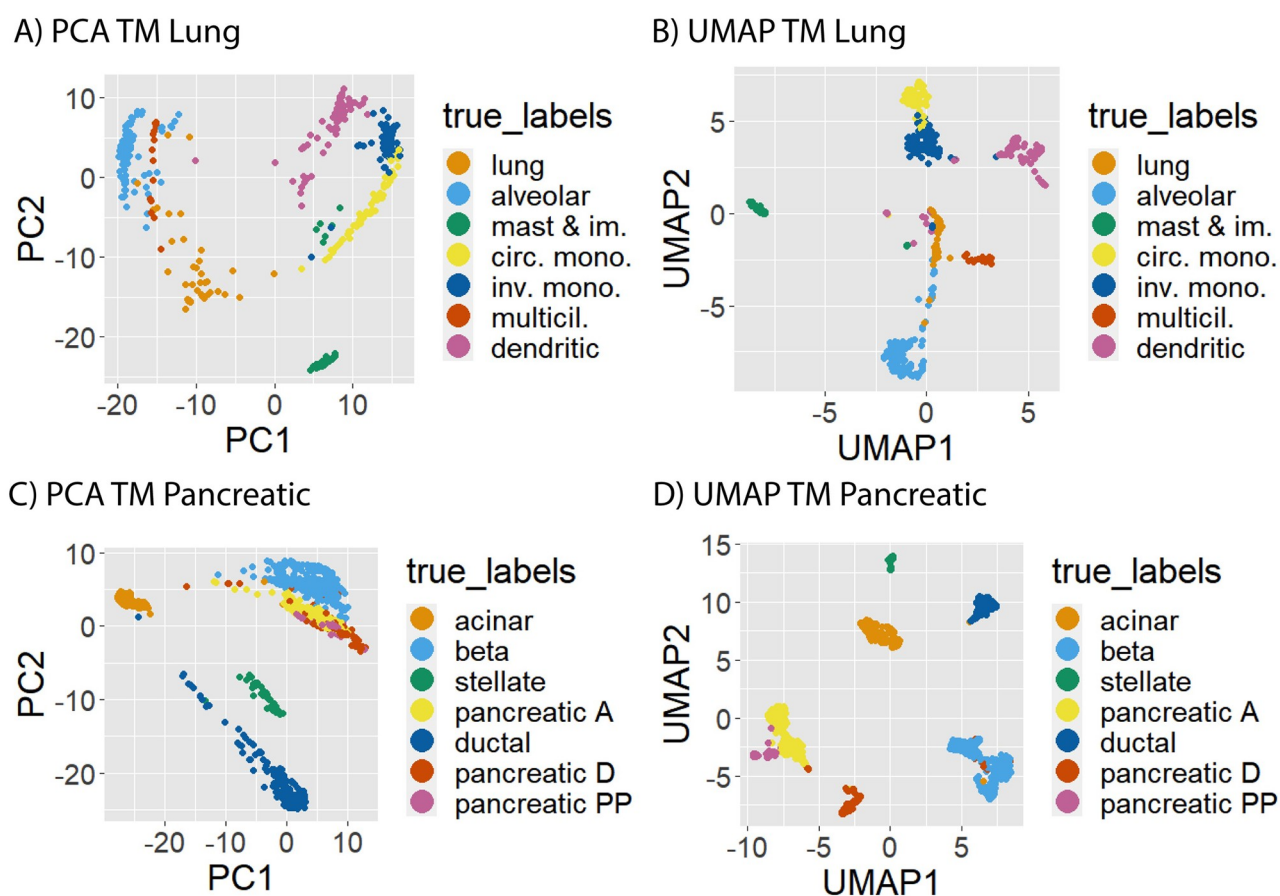


Fig 2. UMAP and PCA on Tabula Muris data sets. Tabula Muris data sets have elongated clusters in the PCA embedding and clusters connected with a bridge of points in the UMAP embedding. For both PCA and UMAP embeddings, certain clusters are not well-separated and connected by high density regions.

<https://doi.org/10.1371/journal.pcbi.1012014.g002>

We propose an embedding method based on the *power weighted path metric* which is well suited to this difficult regime. These metrics balance density and geometry considerations in the data via computation of a density-weighted geodesic distance, making them useful for many machine learning tasks such as clustering and semi-supervised learning [40–48]. They have performed well in applications such as imaging [46, 47, 49, 50], but their usefulness for the analysis of scRNAseq data remains unexplored.

Because these metrics are density-sensitive, they reduce cluster variance; in addition, these metrics also capture global distance information, and thus preserve global geometry. Using the path metric embedding to cluster the data thus yields a clustering method which balances density-based and geometric information.

2 Materials and methods

We first introduce the notations in Table 1 and our theoretical framework in Section 2.1; Section 2.2 then describes the details of the proposed scPMP algorithm, and Section 2.3 describes metrics for assessment.

2.1 Path metrics

We first define a family of power weighted path metrics parametrized by $1 \leq p < \infty$.

Definition 1 Given a discrete data set X , the discrete p -power weighted path metric between $a, b \in X$ is defined as

$$\ell_p(a, b) := \inf_{(x_0, \dots, x_s)} \left(\sum_{i=0}^{s-1} \|x_{i+1} - x_i\|_2^p \right)^{\frac{1}{p}},$$

where the infimum is taken over all sequences of points x_0, \dots, x_s in X with $x_0 = a$ and $x_s = b$.

Table 1. Notations.

Notation	Definition
ARI	Adjusted Rand Index
ECP	Entropy of Cluster Purity
ECA	Entropy of Cluster Accuracy
UMAP	Uniform Manifold Approximation and Projection
t -SNE	t -distributed Stochastic Neighbor Embedding
PCA	Principal Component Analysis
MDS	Mutlidimensional Scaling
scPMP	Single-Cell Path Metrics Profiling
PM	Path Metrics
p	Parameter of power weighted path metrics
PM_p	scPMP clustering with path metric parameter p
NN	Nearest Neighbors
K_1	Number of NN in local averaging
K_2	Number of NN for path metric distance
k	Number of clusters
d	Number of features of data set
n	Number of samples
f	Density function of samples
π	Geometric pertrubations

<https://doi.org/10.1371/journal.pcbi.1012014.t001>

Note as $p \rightarrow \infty$, ℓ_p converges to the “bottleneck edge” distance

$$\ell_\infty(a, b) := \inf_{(x_0, \dots, x_s)} \max_i \|x_{i+1} - x_i\|_2,$$

which is well studied in the computer science literature [51–54]. Two points are close in ℓ_∞ if they are connected by a high-density path through the data, regardless of how far apart the points are. On the other hand, when $p = 1$, ℓ_1 reduces to Euclidean distance. If path edges are furthermore restricted to lie in a nearest neighbor graph, ℓ_1 approximates the geodesic distance between the points, i.e. the length of the shortest path lying on the underlying data structure, which is a highly useful metric for manifold learning [55]. The parameter p governs a trade-off between these two extremes, i.e. it determines how to balance density and geometry considerations when determining which data points should be considered close.

The relationship between ℓ_p and density can be made precise. Assume n independent samples from a continuous, nonzero density function f supported on a d -dimensional, compact Riemannian manifold \mathcal{M} (a manifold is a smooth, locally linear surface; see [56]). Then for $p > 1$, $\ell_p(a, b)$ converges (after appropriate normalization) to

$$\mathcal{L}_p(a, b) := \inf_\gamma \left(\int_\gamma f(\gamma(t))^{-\frac{(p-1)}{d}} |\gamma'(t)| dt \right)^{\frac{1}{p}}, \quad (1)$$

as $n \rightarrow \infty$, where the infimum is taken over all smooth curves $\gamma : [0, 1] \rightarrow \mathcal{M}$ connecting a, b [57–59]. Note $|\gamma'(t)|$ is simply the arclength element on \mathcal{M} , so \mathcal{L}_1 reduces to the standard geodesic distance. When $p \neq 1$, one obtains a density-weighted geodesic distance.

The optimal \mathcal{L}_p path is not necessarily the most direct: a detour may be worth it if it allows the path to stay in a high-density region; see Fig 3. Thus the metric is *density-sensitive*, in that distances across high-density regions are smaller than distances across low-density regions; this is a desirable property for many machine learning tasks [60], including trajectory estimation for developmental cells and cancer cells. However, the metric is also *geometry preserving*, since it is computed by path integrals on \mathcal{M} . The parameter p controls the balance of these

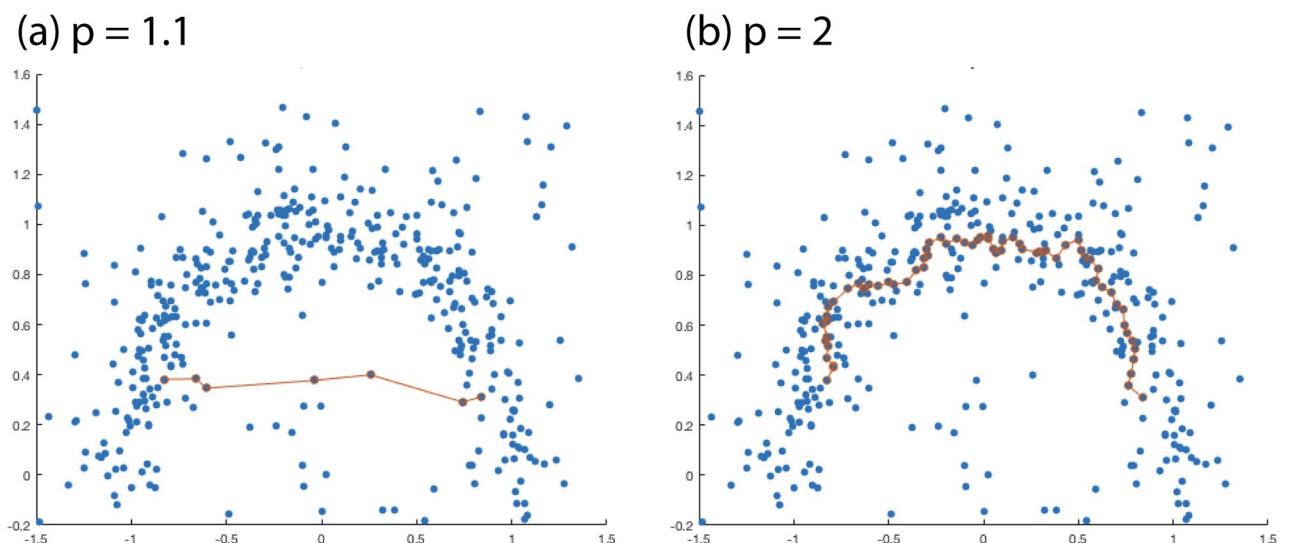


Fig 3. Optimal ℓ_p path between two points in a moon data set.

<https://doi.org/10.1371/journal.pcbi.1012014.g003>

two properties: when p is small, \mathcal{L}_p depends mainly on the geometry of the data, while for large p , \mathcal{L}_p is primarily determined by data density.

Although path metrics are defined in a complete graph, i.e. Definition 1 considers *every* path in the data connecting a, b , recent work [46, 61–63] has established that it is sufficient to only consider paths in a K -nearest neighbors (KNN) graph, as long as $K \geq C \log n$ for a constant C depending on p, d, f , and the geometry of the data. By restricting to a KNN graph, all pairwise path distances can be computed in $O(Kn^2)$ with Dijkstra's algorithm [64].

2.2 Algorithm

Algorithm 1 scPMP

```

1: Input: noisy data  $\tilde{X} \in \mathbb{R}^{n \times d}$ , parameter  $p$ , number of clusters  $k$ 
2: Optional input:  $K_1, K_2, r_{\min}, r_{\max}, \tau$ 
3:   (Defaults: 12,  $n \wedge 500$ , 3, 39, 0.01)
4: Output: scPMP embedding  $Y \in \mathbb{R}^{n \times r}$ , label vector  $\hat{\ell} \in [k]^n$ 
5:
6: % Denoise data:
7:  $x_i \leftarrow \frac{1}{K_1} \sum_{j \in \mathcal{N}_{i,K_1}} \tilde{x}_j$ 
8:
9: % Compute path metrics:
10:  $\mathcal{G}_{K_2}^p \leftarrow K_2\text{NN graph on } X \text{ with edge weights } \|x_i - x_j\|^p$ 
11:  $D_{ij}^p \leftarrow \text{length of shortest path connecting } x_i, x_j \text{ in } \mathcal{G}_{K_2}^p$ 
12:  $(D_{\text{PM}})_{ij} \leftarrow (D_{ij}^p)^{\frac{1}{p}}$ 
13:
14: % Compute MDS embedding of path metrics:
15:  $B \leftarrow -\frac{1}{2} \mathcal{J} D_{\text{PM}}^{(2)} \mathcal{J}$ 
16:  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n) \leftarrow \text{eigenvalues of } B \text{ in descending order}$ 
17:  $V = (v_1, \dots, v_n) \leftarrow \text{corresponding eigenvectors of } B$ 
18:  $r \leftarrow \text{index maximizing } \lambda_i / \lambda_{i+1} \text{ for } i \text{ satisfying } r_{\min} \leq i \leq r_{\max}, \lambda_i / \lambda_1 \geq \tau$ 
19:  $Y \leftarrow (\sqrt{\lambda_1} v_1, \dots, \sqrt{\lambda_r} v_r) \in \mathbb{R}^{n \times r}$ 
20:
21: % Cluster the data:
22:  $\hat{\ell} \leftarrow \text{constrained } k\text{-means}(Y, k)$ 

```

We consider a noisy data set of n data points $\tilde{x}_1, \dots, \tilde{x}_n \in \mathbb{R}^d$, which form the rows of noisy data matrix $\tilde{X} \in \mathbb{R}^{n \times d}$. We first denoise the data with a local averaging procedure, which has been shown to be advantageous for manifold plus noise data models [65] and contributes to the improvement of clustering performance on scRNAseq data sets as explored in [S1 Text](#). More specifically, we replace \tilde{x}_i with its local average:

$$x_i := \frac{1}{K_1} \sum_{j \in \mathcal{N}_{i,K_1}} \tilde{x}_j, \quad \mathcal{N}_{i,K_1} = \{j : \tilde{x}_j \text{ is a } K_1\text{NN of } \tilde{x}_i\},$$

and let $X \in \mathbb{R}^{n \times d}$ denote the denoised data matrix.

We then fix p and compute the p -power weighted path distance between all points in X to obtain pairwise distance matrix $D_{\text{PM}} \in \mathbb{R}^{n \times n}$. More precisely, we let $\mathcal{G}_{K_2}^p = (X, E)$ be the graph on X where x_i, x_j are connected with edge weight $E_{ij} = \|x_i - x_j\|_2^p$ if x_i is a K_2 NN of x_j or x_j is a K_2 NN of x_i . We then compute D_{ij}^p as the total length of the shortest path connecting x_i, x_j in $\mathcal{G}_{K_2}^p$, and define D_{PM} by $(D_{\text{PM}})_{ij} = (D_{ij}^p)^{\frac{1}{p}}$.

We next apply classical multidimensional scaling [66] to obtain a low-dimensional embedding which preserves the path metrics. Specifically, we define the path metric MDS matrix $B = -\frac{1}{2}JD_{\text{PM}}^{(2)}J$ where $J = I_n - \frac{1}{n}\mathbf{1}\mathbf{1}^T$ is the centering matrix, $\mathbf{1} \in \mathbb{R}^n$ is a vector of all 1's, and $D_{\text{PM}}^{(2)}$ is obtained from D_{PM} by squaring all entries. We let the spectral decomposition of B be denoted by $B = V\Lambda V^T$, where $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$, $V = (v_1, \dots, v_n) \in \mathbb{R}^{n \times n}$ contain the eigenvalues and eigenvectors of B in descending order. The embedding dimension r is then chosen as the index i which maximizes the eigenratio λ_i/λ_{i+1} [67], with the following restrictions: we constrain $3 \leq i \leq 39$ and only consider ratios λ_{i+1}/λ_i between “large” eigenvalues, i.e. we require $\lambda_i/\lambda_1 \geq 0.01$. The scPMP embedding is then defined by $Y = (\sqrt{\lambda_1}v_1, \dots, \sqrt{\lambda_r}v_r) \in \mathbb{R}^{n \times r}$.

Finally, we apply k -means to the scPMP embedding to obtain cluster labels. Specifically, we let $\hat{\ell}_i \in [k] = \{1, \dots, k\}$ be the cluster label of x_i returned by running k -means on Y with k clusters and 20 replicates. Since k -means may return highly imbalanced clusters, cluster sample sizes were constrained to be at least $\sqrt{n}/2$. Specifically, if k -means returned a tiny cluster, k was increased to $k + 1$, and the tiny cluster merged with the closest non-trivial cluster. This entire procedure is summarized in the pseudocode in Algorithm 1.

We note that the computational bottleneck for scPMP is the computation and storage of all pairwise path distances, which has complexity $O(n^2 \log n)$ when $K_2 = O(\log n)$. However this quadratic cost can be avoided by utilizing a low rank approximation of the squared distance matrix via the Nystrom method [68–72]. For example, [73] propose a fast, quasi-linear implementation of MDS which only requires the computation of path distances from a set of q landmarks, so that the complexity of computing path distances is reduced to $O(qn \log n)$. Our implementation of scPMP includes the option to use this landmark-based approximation and is thus highly scalable.

We also note that an important consideration in the fully unsupervised setting is how to select the number of clusters k . This is a rather ill-posed question with multiple reasonable answers due to hierarchical cluster structure. We do not focus on this in the current article, and scPMP assumes the number of clusters is given. However we emphasize that when k is unknown, the scPMP embedding offers a useful tool for selecting a reasonable number of clusters. For example, Line 21 of Algorithm 1 can be repeated for a range of candidate k values to obtain candidate clusterings $\hat{\ell}_k$; \hat{k} can then be chosen so that $\hat{\ell}_k$ optimizes a cluster validity criterion such as the silhouette criterion [74, 75]. Alternatively, one could build a graph with distances computed in the scPMP embedding, and estimate k as the number of small eigenvalues of a corresponding graph Laplacian [47, 76].

2.3 Assessment

We evaluate the performance of scPMP with respect to (1) cluster quality and (2) geometric fidelity on a collection of labeled benchmarking data sets with ground truth labels ℓ . There are many helpful metrics for the quality of the estimated cluster labels $\hat{\ell}$, and we compute the adjusted rand index (ARI), entropy of cluster accuracy (ECA), and entropy of cluster purity (ECP). Definitions of ECA and ECP can be found in S2 Text. We compare our clustering results with the output of k -means, DBSCAN [37, 38], k -means on t -SNE embedding [33], DBSCAN on UMAP embedding [32] and for scRNAseq data sets additionally with the following scRNAseq clustering methods: SC3 [4], scanpy [14], RaceID3 [77], SIMLR [5] and Seurat [12].

Assessing the geometric fidelity of the low-dimensional embedding Y is more delicate; we want to assess whether the embedding procedure preserves the global relative distances between clusters. We first compute the mean of each cluster as in [33] using the ground truth

labels, i.e. $\mu_j(X) = \frac{1}{|\mathcal{I}_j|} \sum_{i \in \mathcal{I}_j} x_i$, where $\mathcal{I}_j = \{i : \ell_i = j\}$; we then define $D_{\mu,X}(i, j) = \|\mu_{\ell_i}(X) - \mu_{\ell_j}(X)\|_2$. Similarly, we compute the means $\mu_j(Y)$ in the scPMP embedding, and define $D_{\mu,Y}(i, j) = \|\mu_{\ell_i}(Y) - \mu_{\ell_j}(Y)\|_2$; we then compare $D_{\mu,X}$ and $D_{\mu,Y}$. Specifically, we define the geometric perturbation π by:

$$\pi(X, Y, \ell) = \min_c \frac{\|D_{\mu,X} - cD_{\mu,Y}\|_F^2}{\|D_{\mu,X}\|_F^2},$$

where $\|\cdot\|_F$ is the Frobenius norm. The c achieving the minimum is easy to compute, and one obtains

$$\pi(X, Y, \ell) = \frac{\|D_{\mu,X} - c^*D_{\mu,Y}\|_F^2}{\|D_{\mu,X}\|_F^2}, \quad c^* = \frac{\langle D_{\mu,X}, D_{\mu,Y} \rangle}{\|D_{\mu,Y}\|_F^2}.$$

We compare $\pi(X, Y, \ell)$ with the geometric perturbation of other embedding schemes for X , i.e. with $\pi(X, U, \ell)$ for U equal to the UMAP [32] and t -SNE [33] embeddings. Note that π is not always a useful measure: for example if X consisted of concentric spheres sharing the same center, the metric would be meaningless, as the distance between cluster means would be zero. Nevertheless, in most cases π is a helpful metric for quantifying the preservation of global cluster geometry.

3 Results

We apply scPMP to both a collection of toy manifold data sets and a collection of scRNAseq data sets. Results are reported in Sections 3.1 and 3.2 respectively. The default parameter values reported in scPMP were used on all data sets.

3.1 Manifold data

We apply scPMP for $p = 1.5, 2, 4$ to the following four manifold data sets:

- **Balls** ($n = 1200, d = 2, k = 3$): Clusters were created by uniform sampling of 3 overlapping balls in \mathbb{R}^2 ; see Fig 1A.
- **Elongated with bridge** (denoted EWB, $n = 620, d = 2, k = 3$): Clusters were created by sampling from 3 elongated Gaussian distributions. A bridge was added connecting two of the Gaussians; see Fig 1B.
- **Swiss roll** ($n = 1275, d = 3, k = 3$): Clusters were created by uniform sampling from three distinct regions of a Swiss roll; 3-dimensional isotropic Gaussian noise ($\sigma = 0.75$) was then added to the data. Fig 1C shows the first two data coordinates.
- **SO(3) manifolds** ($n = 3000, d = 1000, k = 3$): For $1 \leq i \leq 3$, the 3-dimensional manifold $\mathcal{M}_i \subseteq \mathbb{R}^9$ is defined by fixing three eigenvalues $D_i = \text{diag}(\lambda_1, \lambda_2, \lambda_3)$ and then defining $\mathcal{M}_i = \cup_{V \in \text{SO}(3)} VD_iV^T$, where $\text{SO}(3)$ is the special orthogonal group. After fixing D_i , we randomly sample from \mathcal{M}_i by taking random orthonormal bases V of \mathbb{R}^3 . A noisy, high-dimensional embedding was then obtained by adding uniform random noise with standard deviation $\sigma = 0.0075$ in 1000 dimensions. Fig 1D shows the first two principal components of the data, which exhibits no cluster separation.

The data sets were chosen to illustrate various cluster separability characteristics. For the balls, the clusters have good geometric separation but are not separable by density. For the Swiss roll and SO(3), the clusters have a complex and inter-twined geometry but are well

Table 2. The results of clustering accuracy (ARI) for manifold data.

Method	Balls	EWB	Swiss	SO(3)
<i>k</i> -means	0.955	-0.001	0.373	0.010
DBSCAN	0.055	0.550	1	1
UMAP+DBSCAN	0.600	0.645	1	1
<i>t</i> -SNE+ <i>k</i> -means	0.895	0.359	1	0.532
Seurat	0.777	0.837	1	1
PM _{1.5}	0.921	0.489	1	0.501
PM ₂	0.907	0.990	1	1
PM ₄	0.781	0.584	1	1

<https://doi.org/10.1371/journal.pcbi.1012014.t002>

separated in terms of density. For EWB, clusters are both elongated and lack robust density separability due to the bridge, and one expects that methods which rely too heavily on either geometry or density will fail. The ARIs achieved by scPMP, *k*-means based methods, DBSCAN based methods, and Seurat are reported in Table 2. See Table A in S2 Text and Table B in S2 Text for ECP and ECA. As expected, *k*-means outperforms all methods on the balls but performs very poorly on all other data sets. DBSCAN and Seurat achieve perfect accuracy on the Swiss roll and SO(3) but perform rather poorly on the balls and EWB, although Seurat does noticeably better than DBSCAN. scPMP with $p = 2$ (PM₂), is the only method which achieves a high ARI (> 90%) and a low ECP and ECA (< 0.15) on all data sets.

Table 3 reports the geometric perturbation of the embedding produced by scPMP and compares with UMAP and *t*-SNE. Since scPMP generally selects an embedding dimension $r > 2$, to ensure a fair comparison the geometric perturbation was computed in both the 2d and r -dimensional (rd) embeddings for all methods, where for UMAP r is the dimension selected by Algorithm 1 and for *t*-SNE $r = 3$ (note $r \leq 3$ was required in Rtsne implementation). Overall PM_{1.5} achieved the lowest geometric perturbation, although all methods had small perturbation on the Balls data set and *t*-SNE had the lowest perturbation on EWB. We point out however that for both the Swiss roll and SO(3), the metric may not be meaningful due to the complicated cluster geometry.

Table 3. Geometric perturbation for manifold data.

Method	Balls	EWB	Swiss	SO(3)
2d UMAP	0.001	0.006	0.305	0.071
rd UMAP	0	0.033	0.339	0.054
2d <i>t</i> -SNE	0	0.004	0.187	0.171
rd <i>t</i> -SNE	0	0.042	0.074	0.157
2d PM _{1.5}	0	0.033	0.002	0.103
rd PM _{1.5}	0	0.023	0.011	0.154
2d PM ₂	0	0.146	0.025	0.156
rd PM ₂	0	0.068	0.025	0.179
2d PM ₄	0.003	0.191	0.056	0.194
rd PM ₄	0.004	0.157	0.056	0.194

The rd UMAP embeddings were computed with an embedding dimension of $r = 5$ for the balls, EWB, Swiss roll and $r = 7$ for SO(3), which corresponded to the estimated dimension for both PM_{1.5} and PM₂. For *t*-SNE, $r = 3$ for all data sets.

<https://doi.org/10.1371/journal.pcbi.1012014.t003>

3.2 scRNAseq data

We apply scPMP for $p = 1.5, 2, 4$ to the following synthetic scRNAseq data sets:

- **RNA mixture:** Benchmarking scRNAseq data set from [36]. RNAmix1 was processed with CEL-seq2 and has $n = 296$ cells and $d = 14687$ genes. RNAmix2 was processed with Sort-seq and has $n = 340$ cells and $d = 14224$ genes. For the creation of the two data sets, RNA was extracted in bulk for each of the following cell lines: H2228, H1975, HCC827. Then the RNA was mixed in $k = 7$ different proportions (each defining a ground truth cluster label), diluted to single cell equivalent amounts ranging from 3.75pg to 30pg, and processed using CEL-seq2 and SORT-seq.
- **Simulated beta:** Simulated data set of $n = 473$ beta cells and $d = 2279$ genes, created based on SAVER [78] and scImpute [79]. First, we subset the Baron's Pancreatic data set [80] to include only Beta cells. As in [79], we randomly choose 10% of the genes to operate as marker genes. Then, we split the cells to $k = 3$ clusters and each cluster is assigned a different group of marker genes. For each cluster we scale up the mean expression of its marker genes. Lastly, to simulate the drop out effect, as in [78], we multiply each cell by an efficiency loss constant drawn by Gamma(10, 100). Using S to refer to the data matrix resulting from the above steps, the final simulated data X is obtained by letting X_{ij} be drawn from Poisson(S_{ij}).

In addition to the synthetic data, we evaluate the performance of scPMP on the following real scRNAseq data sets:

- **Cell mixture data set:** Another benchmarking data set from [36] consisting of a mixture of $k = 5$ cell lines created with 10x sequencing platform. The cell line identity of a cell is also its true cluster label. The data set consists of $n = 3822$ cells and $d = 11786$ genes; we removed multiplets, based on the provided metadata file and kept 3000 most variable genes after SCT transformation [81, 82].
- **Baron's pancreatic:** Human pancreatic data set generated by [80]. After quality control and SAVER imputation, there are $d = 14738$ genes and $n = 1844$ cells. For analysis purposes cells that belong in a group with less than 70 members were filtered out to reduce to $k = 8$ cell types. Also, we kept only the 3000 most variable genes after SCT transformation [81, 82]. The cell types associated with each cell were obtained by an iterative hierarchical clustering method that restricts genes enriched in one cell type from being used to separate other cell types. The enriched markers in every cluster defined the cell type of the cells that belong in that cluster.
- **Tabula Muris data sets:** Mouse scRNAseq data for different tissues and organs [39]. We select the pancreatic data (TM Panc) with $n = 1444$ cells and $d = 23433$ genes and the lung data (TM Lung) with $n = 453$ cells and $d = 23433$ genes. Both data sets have $k = 7$ different cell types which were characterized by an FACS-based full length transcript analysis.
- **PBMC4k data set:** This data set includes the gene expression of Peripheral Blood Mononuclear Cells. The raw data are available from 10X Genomics. After quality control, saver imputation, and removing the two smallest cell types, there are $d = 16655$ genes and $n = 4316$ cells in the dataset. Also, we merge CD8+ T-cells and CD4+ T-cells in one type named T-cells resulting in $k = 4$ cell types. The ground truth cell types are provided by SingleR annotation after marker gene verification in github.com/SingleR.

Details about the pre-processing of data sets can be found in S2 Text. For the following UMAP and t -SNE results, Linnorm normalization [83] was applied without denoising, as this

Table 4. The results of clustering accuracy (ARI) for scRNAseq data.

Method	RNA1	RNA2	TMLung	Beta	TMPanc	BaronPanc	PBMC4K	CellMix
SC3	0.637	0.827	0.798	0.969	0.894	0.767	0.889	1
scanpy	0.620	0.825	0.796	0.910	0.615	0.966	0.977	1
RaceID3	0.730	0.520	0.900	0.714	0.751	0.651	0.763	1
SIMLR	0.878	0.792	0.727	0.975	0.599	0.698	0.705	1
Seurat	0.792	0.667	0.843	0.891	0.547	0.941	0.889	0.993
Seurat_def	0.714	0.785	0.764	0.919	0.798	0.971	0.975	1
<i>k</i> -means	0.921	0.786	0.848	0.969	0.840	0.662	0.747	1
DBSCAN	0.952	0.826	0.587	0.568	0.734	0.724	0.889	1
UMAP+DBSCAN	0.926	0.892	0.619	0.565	0.893	0.848	0.974	1
<i>t</i> -SNE+ <i>k</i> -means	0.943	0.915	0.753	0.969	0.620	0.641	0.596	0.878
PM _{1.5}	0.939	0.924	0.888	0.969	0.626	0.804	0.754	1
PM ₂	0.939	0.973	0.808	0.969	0.918	0.969	0.757	1
PM ₄	0.939	0.939	0.731	0.921	0.775	0.853	0.978	1

<https://doi.org/10.1371/journal.pcbi.1012014.t004>

normalization gave the best results. Note Seurat_def refers to the results of the entire Seurat pipeline, whereas Seurat refers to the result of using Seurat clustering on data with the same processing and normalization as for PM. The embedding dimension r selected by scPMP ranged from 3 to 7 for PM_{1.5} and PM₂, and from 3 to 11 for PM₄.

Table 4 reports the clustering accuracy regarding ARI achieved by scPMP and other methods; see Table C in S2 Text and Table D in S2 Text for ECP and ECA. The path metric methods perform equally well or better than the rest of the methods. Once again PM₂ exhibits the best overall performance, with a high ARI ($\geq 90\%$) on all data sets except TM lung and PBMC4K; the next best method is PM₄, which achieves a high ARI on all but 3 data sets. Seurat_def and PM_{1.5} had a low ARI for 4 of 8 data sets; scanpy, *k*-means, UMAP+DBSCAN and *t*-SNE+*k*-means had a low ARI on 5 of the 8 data sets; SC3, RaceId3, SIMLR and Seurat had a low ARI ($< 90\%$) on 6 of the 8 data sets. These results indicate that incorporating both density-based and geometric information when determining similarity generally leads to more robust results for scRNA-seq data. Moreover, PM₂ achieves the best median ECP and median ECA values across all RNA data sets. Although the optimal balance depends on the data set (for example PBMC4K does best with $p = 4$, while TMLung does best with $p = 1.5$), path metrics with a moderate p exhibit the best performance across a wide range of data sets.

For BaronPanc, we observe that Seurat_def achieves a slightly higher ARI than all the reported path metric methods ($p = 1.5, 2, 4$). However, a significant advantage of scPMP over Seurat is the high clustering performance on a wide range of sample sizes. To demonstrate our claim we compare the ARI results in different down-sampled versions of BaronsPanc. We selected a stratified sample of 50%, 25% and 10% of the cells of the BaronPanc data set. The results can be found in Table E in S2 Text. We observed no ARI deterioration for scPMP for the 50% and 25% down-sampled data set and only a moderate decrease for the 10% down-sampled dataset (ARI of 0.67 at 10% downsampling for $p = 1.5$). On the contrary, there is significant ARI deterioration both for Seurat and Seurat_def; in particular, at 10% downsampling the ARI deteriorates to 0.405 for Seurat and to 0.185 for Seurat_def. Notice that in the 10% down-sampled data set, we use regular *k*-means for PM₂ to allow for the prediction of smaller sized clusters.

We also investigated whether we could learn the ground truth number of clusters by optimizing the silhouette criterion in the scPMP embedding, and compared this with the number

Table 5. Geometric perturbation for RNA data.

Method	RNA1	RNA2	TMLung	Beta	TMPanc	BaronPanc	PBMC4k	CellMix
2d UMAP	0.122	0.142	0.057	0.025	0.064	0.115	0.015	0.090
rd UMAP	0.160	0.131	0.092	0.026	0.036	0.129	0.027	0.050
2d <i>t</i> -SNE	0.059	0.054	0.042	0.024	0.048	0.206	0.038	0.061
rd <i>t</i> -SNE	0.035	0.054	0.027	0.016	0.040	0.229	0.050	0.033
2d PM _{1.5}	0.010	0.013	0.046	0.002	0.076	0.067	0.028	0.098
rd PM _{1.5}	0.017	0.009	0.006	0	0.019	0.006	0.007	0.007
2d PM ₂	0.040	0.040	0.085	0.003	0.150	0.103	0.050	0.101
rd PM ₂	0.048	0.036	0.029	0.003	0.051	0.010	0.013	0.008
2d PM ₄	0.108	0.135	0.246	0.016	0.265	0.193	0.069	0.107
rd PM ₄	0.100	0.082	0.083	0.015	0.099	0.027	0.029	0.008

For rd UMAP $r = 7, 6, 5, 3, 5, 9, 3, 4$ for the various data sets, which was the maximum of the PM_{1.5} dimension and the PM₂ dimension. For rd *t*-SNE $r = 3$.

<https://doi.org/10.1371/journal.pcbi.1012014.t005>

of clusters obtained from Seurat using the default resolution; see Table F in S2 Text. For 4 out of the 8 RNA data sets evaluated in this article (RNAmix1, RNAmix2, BaronPanc, and CellMix), this procedure on PM₂ yielded an estimate for k which matched the number of distinct annotated labels. On the other hand, Seurat correctly estimates the number of clusters for only 2 out of the 8 RNA data sets (RNAmix1 and TMLung).

Table 5 reports the geometric perturbation. We see that increasing p increases the geometric perturbation, with PM_{1.5} yielding the smallest geometric perturbation on all data sets. Although PM_{1.5} is the clear winner in terms of this metric, PM₂ still performed favorably with respect to UMAP and *t*-SNE. Indeed, rd PM₂ had lower geometric perturbation than UMAP on all but one data set (TMPanc), and lower geometric perturbation than *t*-SNE on the majority of data sets. Fig 4 shows the PCA, PM₂, UMAP, and *t*-SNE embeddings of the Cell Mix data set, as well as a tree structure on the clusters. The tree structure was obtained by first computing the cluster means in the embedding and then applying hierarchical clustering with average linkage to the means. The PCA tree (Fig 4(E)) was computed using 40 PCs so that it accurately reflects the global geometry of the clusters. Interestingly path metrics recover the same hierarchical structure on the clusters as PCA: the cell types HCC827 and H1975 are the most similar, and H838 is the most distinct. This is what one would expect given more extensive biological information about the cell types, since H838 is the only cell line here derived from metastatic site Lymph node on a male patient, while both HCC827 and H1975 originated from the primary site of female lung cancer patients. However, neither UMAP or *t*-SNE give the correct hierarchical representation of the clusters, because both methods struggle to preserve global geometric structure as observed in numerous studies [84, 85]. We note that in Fig 4(B) the clusters appear elongated in the PM₂ embedding; such elongated cluster shapes occur when clusters living in nearly orthogonal subspaces (due for example to different genetic signatures) are projected into a lower-dimensional space; see S3 Text for an example illustrating how this phenomenon occurs. While this is also the case for PCA, the PM embedding exaggerates the elongation by shrinking noisy directions. Although 2 dimensions is generally not sufficient to visualize the true cluster shapes, the PM embedding is able to simultaneously denoise the clusters while preserving their global layout.

Fig 5 records the runtime for processing and clustering (in minutes) of the Baron's Pancreatic ($n = 1844$) and PBMC4K ($n = 4316$) data sets. For PBMC4k (our largest data set), we use the landmark-based approximation of path distances for scalability. All the PM methods run in less than a minute on BaronPanc and less than 6 minutes on PBMC4k; RaceID3, scanpy,

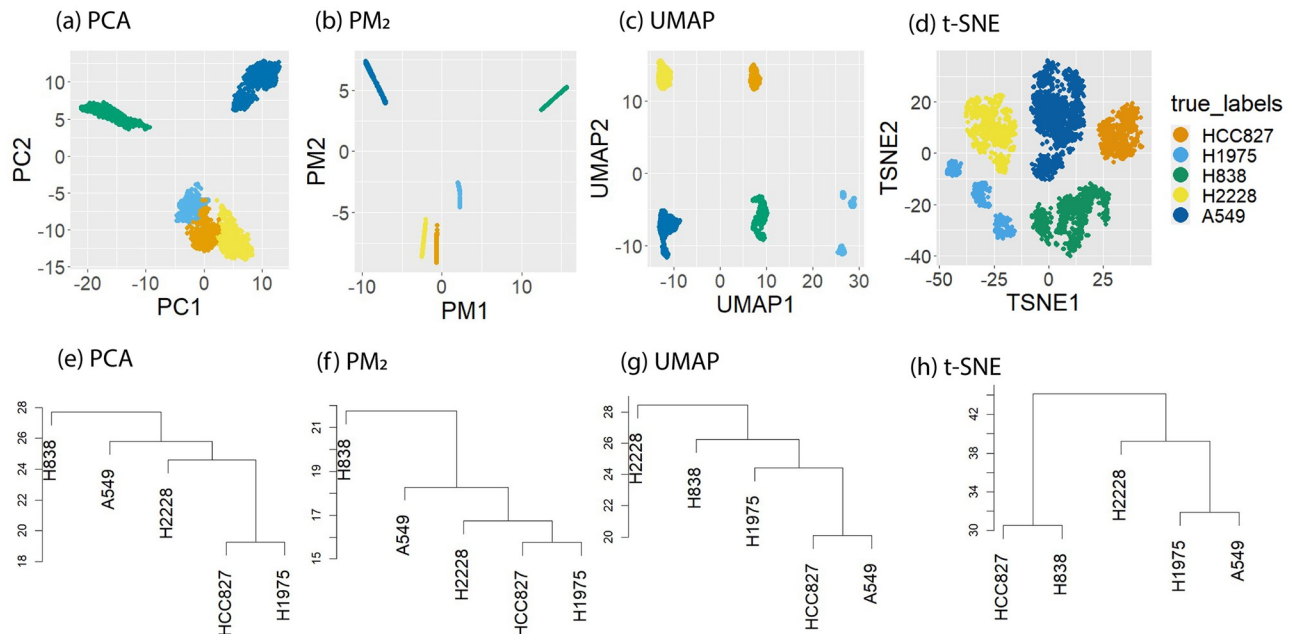


Fig 4. Comparison of cluster structure preservation on PCA, UMAP and *t*-SNE embeddings. Top row: 2d PCA, PM₂, UMAP, and *t*-SNE embeddings of Cell Mix data set, colored by true cell type. Bottom row: average linkage dendrograms of cluster means for the *rd* embeddings, where *r* = 40 for PCA, *r* = 4 for PM₂ and UMAP, and *r* = 3 for *t*-SNE.

<https://doi.org/10.1371/journal.pcbi.1012014.g004>

and Seurat were also fast. SC3 and SIMLR had long runtimes, requiring 37.9 and 91.1 minutes respectively for PBMC4k.

3.3 Determining the parameter *p*

In this section, we explore the clustering performance of scPMP for different values of the parameter *p*. We record the ARI achieved by scPMP for each real data set for *p* ranging from 1 to 10 in increments of 0.5. Fig 6(A) plots the corresponding distributions of ARI; *p* = 2 is the clear winner across various *p*, achieving the highest median ARI with the smallest spread of values.

Furthermore, for each RNA data set we determined the *p* value maximizing the data set's ARI and investigated whether there was a correlation between the best *p* and the degree of data elongation. We define an elongation score for each data set by computing the skewness coefficient of *k*th nearest neighbor distances for $k = 10 \log(n)$. More specifically, letting d_i^k denote the Euclidean distance of x_i from its *k*th nearest neighbor, we define the data elongation score as the following measure of skewness:

$$G_1 = \frac{n}{(n-1)(n-2)} \sum_{i=1}^n \left(\frac{d_i^k - \bar{d}}{s} \right)^3,$$

where $\bar{d} = \frac{1}{n} \sum_{i=1}^n d_i^k$ and *s* is the standard deviation of the $\{d_i^k\}_{i=1}^n$.

We observe a moderately strong linear relationship ($r = 0.866$) between the elongation score of a data set and the value of *p* achieving the best ARI as in Fig 6(B). Overall these results support using *p* = 2 as a default, but increasing *p* if the data set exhibits strong elongation; the elongation score is a completely unsupervised statistic, and can thus be computed without access to data labels.

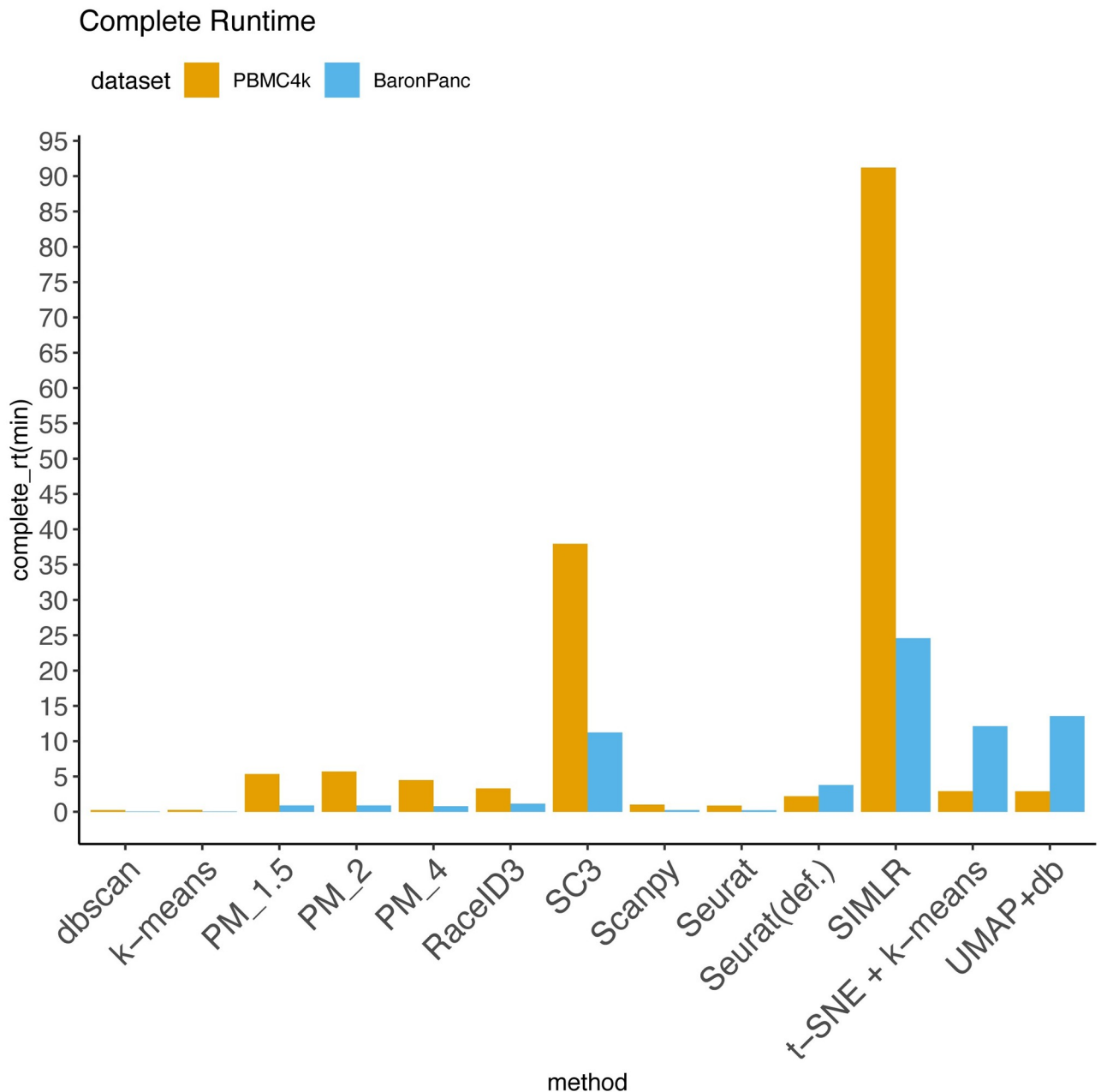


Fig 5. Processing and clustering time for PBMC4K and Baron's Pancreatic data sets.

<https://doi.org/10.1371/journal.pcbi.1012014.g005>

4 Discussion

This article introduces a new theoretical framework to analyze single-cell RNA-seq data based on the computation of optimal paths. Specifically, path metrics encode both geometric and density-based information, and the resulting low-dimensional embeddings simultaneously preserve density-based cluster structure as well as global cluster orientation. Thus, our method with theoretical guarantees addresses the inherent challenge of balancing the preservation of local cluster structures and the global data geometry, a common limitation in existing scRNAseq

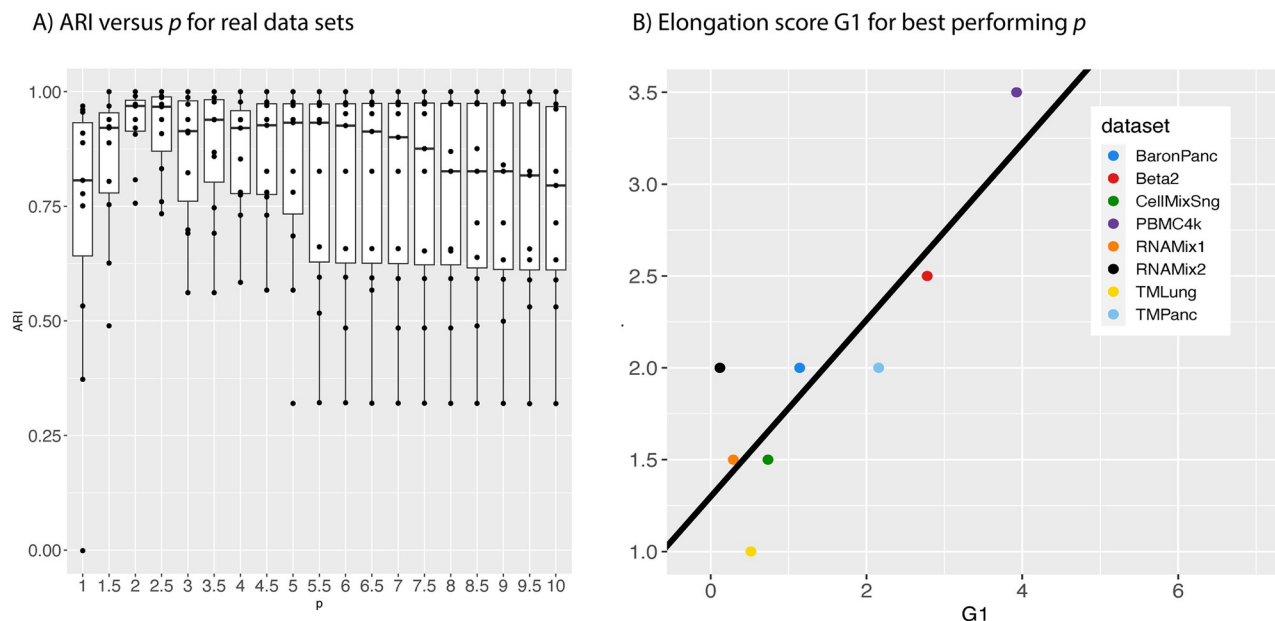


Fig 6. Clustering performance for different values of p .

<https://doi.org/10.1371/journal.pcbi.1012014.g006>

clustering and visualization methods such as DBSCAN, SC3, scanpy, and Seurat. The flexibility in choosing the parameter p allows researchers to adjust the balance between density sensitivity and geometry preservation, tailoring the analysis to their dataset's specific characteristics, such as noise level and elongation. Compared to deep learning-based methods, such as CellPLM, scMoFormer, and scMoGNN, scPMP based on path metrics offers greater interpretability making it easier to derive biological insights. More importantly, scPMP is more robust to smaller datasets than deep learning-based methods since it has fewer parameters to be trained.

The method exhibits competitive performance when applied to numerous benchmarks, and the implementation is scalable to large data sets. Although we investigated other choices of p , we found that $p = 2$ performed well on a wide range of RNA data sets, indicating that $p = 2$ is an appropriate balance between density and geometry for this application. Future research will explore ways to make the method more robust to noise, tools for better visualization of the PM embeddings, and adapting the method to the semi-supervised context.

Supporting information

S1 Text. Data preprocessing.

(PDF)

S2 Text. Additional clustering results.

(PDF)

S3 Text. Clustering visualizations on PCA and scPMP embedding.

(PDF)

Author Contributions

Conceptualization: Anna Little, Yuying Xie.

Data curation: Andriana Manousidaki, Anna Little.

Formal analysis: Andriana Manousidaki.

Funding acquisition: Anna Little, Yuying Xie.

Investigation: Andriana Manousidaki, Anna Little.

Methodology: Andriana Manousidaki, Anna Little, Yuying Xie.

Project administration: Andriana Manousidaki, Anna Little, Yuying Xie.

Software: Andriana Manousidaki, Anna Little.

Supervision: Anna Little, Yuying Xie.

Validation: Andriana Manousidaki, Anna Little.

Visualization: Andriana Manousidaki, Anna Little.

Writing – original draft: Anna Little.

Writing – review & editing: Andriana Manousidaki, Anna Little, Yuying Xie.

References

1. Saliba AE, Westermann AJ, Gorski SA, Vogel J. Single-cell RNA-seq: advances and future challenges. *Nucleic Acids Research*. 2014; 42(14):8845–8860. <https://doi.org/10.1093/nar/gku555> PMID: 25053837
2. Eberwine J, Yeh H, Miyashiro K, Cao Y, Nair S, Finnell R, et al. Analysis of gene expression in single live neurons. *Proceedings of the National Academy of Sciences*. 1992; 89(7):3010–3014. <https://doi.org/10.1073/pnas.89.7.3010> PMID: 1557406
3. Tang F, Barbacioru C, Wang Y, Nordman E, Lee C, Xu N, et al. mRNA-Seq whole-transcriptome analysis of a single cell. *Nature methods*. 2009; 6(5):377–382. <https://doi.org/10.1038/nmeth.1315> PMID: 19349980
4. Kiselev VY, Kirschner K, Schaub MT, Andrews T, Yiu A, Chandra T, et al. SC3: consensus clustering of single-cell RNA-seq data. *Nature Methods*. 2017; 14:483–486. <https://doi.org/10.1038/nmeth.4236> PMID: 28346451
5. Wang B, Zhu J, Pierson E, Ramazzotti D, Batzoglou S. Visualization and analysis of single-cell RNA-seq data by kernel-based similarity learning. *Nature Methods*. 2017; 14:414–416. <https://doi.org/10.1038/nmeth.4207> PMID: 28263960
6. Herman JS, Grün D, et al. FateID infers cell fate bias in multipotent progenitors from single-cell RNA-seq data. *Nature methods*. 2018; 15(5):379. <https://doi.org/10.1038/nmeth.4662> PMID: 29630061
7. Lin P, Troup M, Ho JWK. CIDR: Ultrafast and accurate clustering through imputation for single-cell RNA-seq data. *Genome Biology*. 2017; 18. <https://doi.org/10.1186/s13059-017-1188-0> PMID: 28351406
8. Z A, Muñoz-Manchado AB, Codeluppi S, L P, LM G, J A, et al. Brain structure. Cell types in the mouse cortex and hippocampus revealed by single-cell RNA-seq. *Science (New York, NY)*. 2015; 347:1138–1142. <https://doi.org/10.1126/science.aaa1934>
9. Žurauskienė J, Yau C. pcaReduce: hierarchical clustering of single cell transcriptional profiles. *BMC Bioinformatics*. 2016; 17. <https://doi.org/10.1186/s12859-016-0984-y> PMID: 27005807
10. CLevine J, Simonds E, Bendall S, Davis K, Amir Ea, Tadmor M, et al. Data-Driven Phenotypic Dissection of AML Reveals Progenitor-like Cells that Correlate with Prognosis. *Cell*. 2015.
11. Xu C, Su Z. Identification of cell types from single-cell transcriptomes using a novel clustering method. *Bioinformatics*. 2015; 31(12):1974–1980. <https://doi.org/10.1093/bioinformatics/btv088> PMID: 25805722
12. Stuart T, Butler A, Hoffman P, Hafemeister C, Papalexi E, III WMM, et al. Comprehensive Integration of Single-Cell Data. *Cell*. 2019; 177:1888–1902. <https://doi.org/10.1016/j.cell.2019.05.031> PMID: 31178118
13. Zhu X, Zhang J, Xu Y, Wang J, Peng X, Li HD. Single-Cell Clustering Based on Shared Nearest Neighbor and Graph Partitioning. *Interdisciplinary Sciences: Computational Life Sciences*. 2020; 12:117–130. PMID: 32086753

14. Wolf FA, Angerer P, Theis FJ. SCANPY: large-scale single-cell gene expression data analysis. *Genome Biology*. 2018; 19. <https://doi.org/10.1186/s13059-017-1382-0> PMID: 29409532
15. Wang J, Ma A, Chang Y, Gong J, Jiang Y, Qi R, et al. scGNN is a novel graph neural network framework for single-cell RNA-Seq analyses. *Nature communications*. 2021; 12(1):1882. <https://doi.org/10.1038/s41467-021-22197-x> PMID: 33767197
16. Lopez R, Regier J, Cole MB, Jordan MI, Yosef N. Deep generative modeling for single-cell transcriptomics. *Nature methods*. 2018; 15(12):1053–1058. <https://doi.org/10.1038/s41592-018-0229-2> PMID: 30504886
17. Tian T, Wan J, Song Q, Wei Z. Clustering single-cell RNA-seq data with a model-based deep learning approach. *Nature Machine Intelligence*. 2019; 1(4):191–198. <https://doi.org/10.1038/s42256-019-0037-0>
18. Ding J, Wen H, Tang W, Liu R, Li Z, Venegas J, et al. DANCE: A Deep Learning Library and Benchmark for Single-Cell Analysis. *bioRxiv*. 2022; p. 2022–10.
19. Ciortan M, Defrance M. GNN-based embedding for clustering scRNA-seq data. *Bioinformatics*. 2021; 38(4):1037–1044. <https://doi.org/10.1093/bioinformatics/btab787>
20. Zeng Y, Zhou X, Rao J, Lu Y, Yang Y. Accurately Clustering Single-cell RNA-seq data by Capturing Structural Relations between Cells through Graph Convolutional Network. In: 2020 IEEE International Conference on Bioinformatics and Biomedicine (BIBM); 2020. p. 519–522.
21. Tian T, Zhang J, Lin X, Wei Z, Hakonarson H. Model-based deep embedding for constrained clustering analysis of single cell RNA-seq data. *Nature communications*. 2021; 12(1):1873. <https://doi.org/10.1038/s41467-021-22008-3> PMID: 33767149
22. Li X, Wang K, Lyu Y, Pan H, Zhang J, Stambolian D, et al. Deep learning enables accurate clustering with batch effect removal in single-cell RNA-seq analysis. *Nature communications*. 2020; 11(1):2338. <https://doi.org/10.1038/s41467-020-15851-3> PMID: 32393754
23. Tran D, Nguyen H, Tran B, La Vecchia C, Luu HN, Nguyen T. Fast and precise single-cell data analysis using a hierarchical autoencoder. *Nature communications*. 2021; 12(1):1029. <https://doi.org/10.1038/s41467-021-21312-2> PMID: 33589635
24. Chen L, Wang W, Zhai Y, Deng M. Deep soft K-means clustering with self-training for single-cell RNA sequence data. *NAR genomics and bioinformatics*. 2020; 2(2):lqaa039. <https://doi.org/10.1093/nargab/lqaa039> PMID: 33575592
25. Gan Y, Huang X, Zou G, Zhou S, Guan J. Deep structural clustering for single-cell RNA-seq data jointly through autoencoder and graph neural network. *Briefings in Bioinformatics*. 2022; 23(2):bbac018. <https://doi.org/10.1093/bib/bbac018> PMID: 35172334
26. Wen H, Tang W, Dai X, Ding J, Jin W, Xie Y, et al. CellPLM: Pre-training of Cell Language Model Beyond Single Cells. *bioRxiv*. 2023; p. 2023–10.
27. Tang W, Liu R, Wen H, Dai X, Ding J, Li H, et al. A General Single-Cell Analysis Framework via Conditional Diffusion Generative Models. *bioRxiv*. 2023; p. 2023–10.
28. Wen H, Ding J, Jin W, Wang Y, Xie Y, Tang J. Graph neural networks for multimodal single-cell data integration. In: *Proceedings of the 28th ACM SIGKDD conference on knowledge discovery and data mining*; 2022. p. 4153–4163.
29. Tang W, Wen H, Liu R, Ding J, Jin W, Xie Y, et al. Single-Cell Multimodal Prediction via Transformers. *arXiv preprint arXiv:230300233*. 2023;.
30. Yu Z, Lu Y, Wang Y, Tang F, Wong KC, Li X. Zinb-based graph embedding autoencoder for single-cell rna-seq interpretations. In: *Proceedings of the AAAI conference on artificial intelligence*; 2022. p. 4671–4679.
31. Molho D, Ding J, Tang W, Li Z, Wen H, Wang Y, et al. Deep learning in single-cell analysis. *ACM Transactions on Intelligent Systems and Technology*. 2022;.
32. McInnes L, Healy J, Melville J. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:180203426*. 2018;.
33. Van der Maaten L, Hinton G. Visualizing data using t-SNE. *Journal of machine learning research*. 2008; 9(11).
34. Moon KR, van Dijk D, Wang Z, Gigante S, Burkhardt DB, Chen WS, et al. Visualizing structure and transitions in high-dimensional biological data. *Nature Biotechnology*. 2019; 37(12):1482–1492. <https://doi.org/10.1038/s41587-019-0336-3> PMID: 31796933
35. Tang J, Liu J, Zhang M, Mei Q. Visualizing large-scale and high-dimensional data. In: *Proceedings of the 25th international conference on world wide web*; 2016. p. 287–297.
36. Tian L, Dong X, Freytag S, Lê Cao KA, Su S, JalalAbadi A, et al. Benchmarking single cell RNA-sequencing analysis pipelines using mixture control experiments. *Nature methods*. 2019; 16(6):479–487. <https://doi.org/10.1038/s41592-019-0425-8> PMID: 31133762

37. Ester M, Kriegel HP, Sander J, Xu X, et al. A density-based algorithm for discovering clusters in large spatial databases with noise. In: Kdd. vol. 96; 1996. p. 226–231.
38. Xu X, Ester M, Kriegel HP, Sander J. A distribution-based clustering algorithm for mining in large spatial databases. In: Proceedings 14th International Conference on Data Engineering. IEEE; 1998. p. 324–331.
39. Tabula Muris Consortium Lcea Overall coordination. Single-cell transcriptomics of 20 mouse organs creates a Tabula Muris. *Nature*. 2018; 562:367–372. <https://doi.org/10.1038/s41586-018-0590-4>
40. Vincent P, Bengio Y. Density-sensitive metrics and kernels. In: Snowbird Learning Workshop; 2003.
41. Bousquet O, Chapelle O, Hein M. Measure based regularization. In: NIPS; 2004. p. 1221–1228.
42. Sajama, Orlitsky A. Estimating and computing density based distance metrics. In: ICML; 2005. p. 760–767.
43. Chang H, Yeung DY. Robust path-based spectral clustering. *Pattern Recognition*. 2008; 41(1):191–203. <https://doi.org/10.1016/j.patcog.2007.04.010>
44. Bijral AS, Ratliff N, Srebro N. Semi-supervised Learning with density based distances. In: UAI; 2011. p. 43–50.
45. Moscovich A, Jaffe A, Nadler B. Minimax-optimal semi-supervised regression on unknown manifolds. In: AISTATS; 2017. p. 933–942.
46. McKenzie D, Damelin S. Power weighted shortest paths for clustering Euclidean data. *Foundations of Data Science*. 2019; 1(3):307. <https://doi.org/10.3934/fods.2019014>
47. Little A, Maggioni M, Murphy JM. Path-Based Spectral Clustering: Guarantees, Robustness to Outliers, and Fast Algorithms. *Journal of Machine Learning Research*. 2020; 21(6):1–66.
48. Fernández X, Borghini E, Mindlin G, Groisman P. Intrinsic persistent homology via density-based metric learning. *Journal of Machine Learning Research*. 2023; 24(75):1–42.
49. Fischer B, Zöller T, Buhmann JM. Path based pairwise data clustering with application to texture segmentation. In: International Workshop on Energy Minimization Methods in Computer Vision and Pattern Recognition. Springer; 2001. p. 235–250.
50. Zhang S, Murphy JM. Hyperspectral image clustering with spatially-regularized ultrametrics. *Remote Sensing*. 2021; 13(5):955. <https://doi.org/10.3390/rs13050955>
51. Pollack M. Letter to the Editor: The Maximum Capacity Through a Network. *Operations Research*. 1960; 8(5):733–736. <https://doi.org/10.1287/opre.8.5.733>
52. Hu TC. Letter to the Editor: The Maximum Capacity Route Problem. *Operations Research*. 1961; 9(6):898–900. <https://doi.org/10.1287/opre.9.6.898>
53. Camerini PM. The min-max spanning tree problem and some extensions. *Information Processing Letters*. 1978; 1(10–14). [https://doi.org/10.1016/0020-0190\(78\)90030-3](https://doi.org/10.1016/0020-0190(78)90030-3)
54. Gabow H, Tarjan RE. Algorithms for Two Bottleneck Optimization Problems. *Journal of Algorithms*. 1988; 9:411–417. [https://doi.org/10.1016/0196-6774\(88\)90031-4](https://doi.org/10.1016/0196-6774(88)90031-4)
55. Tenenbaum JB, Silva VD, Langford JC. A global geometric framework for nonlinear dimensionality reduction. *Science*. 2000; 290(5500):2319–2323. <https://doi.org/10.1126/science.290.5500.2319> PMID: 11125149
56. Lee JM. Introduction to Riemannian manifolds. Springer; 2018.
57. Hwang SJ, Damelin SB, Hero A. Shortest path through random points. *The Annals of Applied Probability*. 2016; 26(5):2791–2823. <https://doi.org/10.1214/15-AAP1162>
58. Groisman P, Jonckheere M, Sapienza F. Nonhomogeneous Euclidean first-passage percolation and distance learning. *Bernoulli*. 2022; 28(1):255–276. <https://doi.org/10.3150/21-BEJ1341>
59. Fernández X, Borghini E, Mindlin G, Groisman P. Intrinsic Persistent Homology via Density-based Metric Learning. *Journal of Machine Learning Research*. 2023; 24(75):1–42.
60. Chu T, Miller G, Sheehy D. Exploration of a graph-based density sensitive metric. *arXiv preprint arXiv:170907797*. 2017;.
61. Little A, McKenzie D, Murphy JM. Balancing geometry and density: Path distances on high-dimensional data. *SIAM Journal on Mathematics of Data Science*. 2022; 4(1):72–99. <https://doi.org/10.1137/20M1386657>
62. Groisman P, Jonckheere M, Sapienza F. Nonhomogeneous Euclidean first-passage percolation and distance learning. *Bernoulli*. 2022; 28(1):255–276. <https://doi.org/10.3150/21-BEJ1341>
63. Chu T, Miller GL, Sheehy DR. Exact computation of a manifold metric, via Lipschitz Embeddings and Shortest Paths on a Graph. In: SODA; 2020. p. 411–425.
64. Sniedovich M. Dijkstra's algorithm revisited: the dynamic programming connexion. *Control and cybernetics*. 2006; 35(3):599–620.

65. García Trillos N, Sanz-Alonso D, Yang R. Local Regularization of Noisy Point Clouds: Improved Global Geometric Estimates and Data Analysis. *Journal of Machine Learning Research*. 2019; 20(136):1–37.
66. Ghoghogh B, Ghodsi A, Karray F, Crowley M. Multidimensional scaling, sammon mapping, and isomap: Tutorial and survey; 2020.
67. Lam C, Yao Q. Factor modeling for high-dimensional time series: inference for the number of factors. *The Annals of Statistics*. 2012; p. 694–726.
68. Williams C, Seeger M. Using the Nyström method to speed up kernel machines. In: *Proceedings of the 14th annual conference on neural information processing systems*. CONF; 2001. p. 682–688.
69. Ghoghogh B, Ghodsi A, Karray F, Crowley M. Multidimensional scaling, Sammon mapping, and Isomap: Tutorial and survey. *arXiv preprint arXiv:200908136*. 2020;.
70. Platt J. Fastmap, metricmap, and landmark mds are all nyström algorithms. In: *International Workshop on Artificial Intelligence and Statistics*. PMLR; 2005. p. 261–268.
71. Yu H, Zhao X, Zhang X, Yang Y. ISOMAP using Nyström method with incremental sampling. *Advances in Information Sciences & Service Sciences*. 2012; 4(12).
72. Civril A, Magdon-Ismael M, Bocek-Rivele E. SSDE: Fast graph drawing using sampled spectral distance embedding. In: *International Symposium on Graph Drawing*. Springer; 2006. p. 30–41.
73. Shamaï G, Zibulevsky M, Kimmel R. Efficient Inter-Geodesic Distance Computation and Fast Classical Scaling. *IEEE Transactions on Pattern Analysis and Machine Intelligence*. 2020; 42(1):74–85. <https://doi.org/10.1109/TPAMI.2018.2877961> PMID: 30369438
74. Kaufman L, Rousseeuw P. *Finding Groups in Data: An Introduction to Cluster Analysis*; 2009.
75. Maechler M, Rousseeuw P, Struyf A, Hubert M, Hornik K. *cluster: Cluster Analysis Basics and Extensions*; 2021. Available from: <https://CRAN.R-project.org/package=cluster>.
76. Von Luxburg U. A tutorial on spectral clustering. *Statistics and computing*. 2007; 17(4):395–416. <https://doi.org/10.1007/s11222-007-9033-z>
77. Grün D, et al. Revealing Dynamics of Gene Expression Variability in Cell State Space. *Nature methods*. 2018; 17:45–49.
78. Huang M, Wang J, Torre E, Dueck H, Shaffer S, Bonasio R, et al. SAVER: gene expression recovery for single-cell RNA sequencing. *Nature methods*. 2018; 15(7):539–542. <https://doi.org/10.1038/s41592-018-0033-z> PMID: 29941873
79. Li WV, Li JJ. An accurate and robust imputation method scImpute for single-cell RNA-seq data. *Nature communications*. 2018; 9(1):1–9. <https://doi.org/10.1038/s41467-018-03405-7> PMID: 29520097
80. Baron M, et al. A Single-Cell Transcriptomic Map of the Human and Mouse Pancreas Reveals Inter- and Intra-cell Population Structure. *Cell Systems*. 2016; 3(4):346–360. <https://doi.org/10.1016/j.cels.2016.08.011> PMID: 27667365
81. Hafemeister C, Satija R. Normalization and variance stabilization of single-cell RNA-seq data using regularized negative binomial regression. *Genome Biology*. 2019; 20(1). <https://doi.org/10.1186/s13059-019-1874-1> PMID: 31870423
82. Choudhary S, Satija R. Comparison and evaluation of statistical error models for scRNA-seq. *Genome Biology*. 2022; 23. <https://doi.org/10.1186/s13059-021-02584-9> PMID: 35042561
83. Yip SH, Wang P, Kocher JPA, Sham PC, Wang J. Linnorm: improved statistical analysis for single cell RNA-seq expression data. *Nucleic Acids Research*. 2017; 45(22):e179–e179. <https://doi.org/10.1093/nar/gkx828> PMID: 28981748
84. Kobak D, Berens P. The art of using t-SNE for single-cell transcriptomics. *Nature Communications*. 2019; 10:2041–1723. <https://doi.org/10.1038/s41467-019-13056-x> PMID: 31780648
85. Cooley SM, Hamilton T, Aragonés SD, Ray JCJ, Deeds EJ. A novel metric reveals previously unrecognized distortion in dimensionality reduction of scRNA-Seq data. *Biorxiv*. 2019; p. 689851.