

Comp683: Computational Biology

Lecture 9

Feb 12, 2025

Today + Announcements

- Branch preserving visualization for single-cell data
- Single-cell data augmentation

Announcement : Homework 1 is available and due Feb 28.

<https://github.com/natalies-teaching/comp683-2025/tree/main/Homework1>

Review Questions

- ① What is a common data artifact that we have with single-cell RNA-seq?
- ② What is the intuition for how MAGIC corrects these artifacts?

Visualizing what is happening....

As t increases, dense areas of the data result in more possible paths and weights are then concentrated in this area (e.g. closest neighbors remain with high probability). Also, helps to filter out noise.

- v Exponentiate markov matrix

$$[\quad]^t$$



Figure: from van Dijk et al. Cell Systems. 2018

Visualizing Single-Cell Data with PHATE

A Dimensionality Reduction Method for Single Cells

PHATE aims to preserve between-cell population similarities according to the way that cells differentiate.

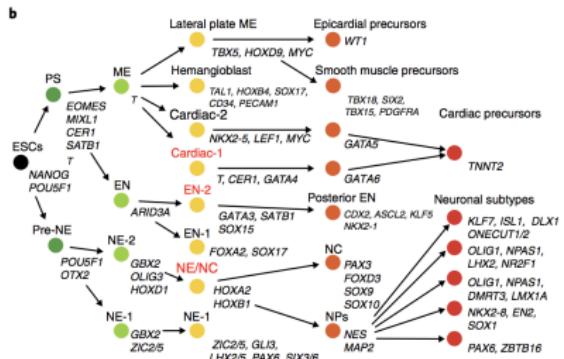


Figure: from Moon et al. Nature Biotechnology. 2019

Latent Structure

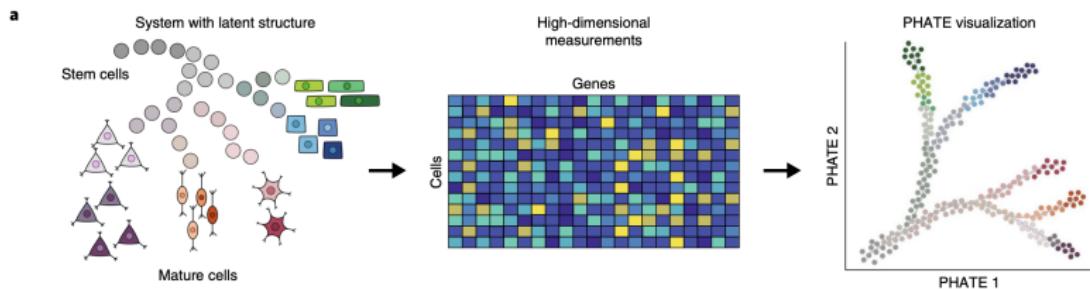


Figure: from Moon *et al.* Nature Biotechnology. 2019

PHATE Overview

Table 1 | General steps in the PHATE algorithm

Input: Data matrix, algorithm parameters (Methods)

Output: The PHATE visualization

- (1) Compute the pairwise distances from the data matrix.
- (2) Transform the distances to affinities to encode local information.
- (3) Learn global relationships via the diffusion process.
- (4) Encode the learned relationships using the potential distance.
- (5) Embed the potential distance information into low dimensions for visualization.

Figure: from Moon *et al.* Nature Biotechnology. 2019

The First Few Steps are Identical to MAGIC

- **Step 1:** Cell \times Cell distance matrix (based on Euclidean distance between cells)
- **Step 2:** Convert to Cell \times Cell row-stochastic affinity matrix \rightarrow, \mathbf{P} , such that $\sum_j \mathbf{P}_{ij} = 1$.
- **Step 3:** Power $\mathbf{P} \rightarrow \mathbf{P}^t$.

A Complementary Method to Determine Optimal t

- The choice of t will determine how much noise you want to filter. This entropy measure will be dominated by large eigenvalues of \mathbf{P} , and will give an idea of signals in the data.
- Let $[\eta(t)]_i = \lambda_i^t / \sum_{j=0}^{N-1} \lambda_j^t$ be the probability distribution defined by the non-negative eigenvalues of \mathbf{P}^t .

The von Neumann entropy $H(t)$ is then computed based on $[\eta(t)]_i$ and decreases as $t \rightarrow \infty$.

$$H(t) = - \sum_{i=1}^N [\eta(t)]_i \log([\eta(t)]_i) \quad (1)$$

Relationship Between t and $H(t)$

Higher t means less noise, but the implications of t are super powerful!

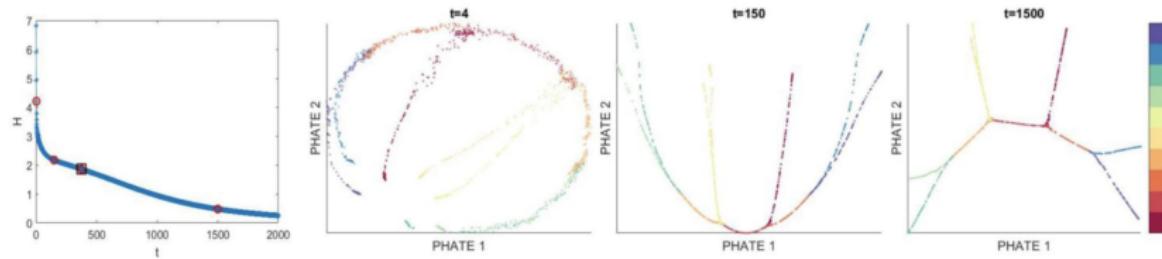


Figure: from Moon et al. Nature Biotechnology. 2019

Computing Potential Distances

- Potential distance is computed as a divergence between the associated diffusion probability distributions of the two cells to all other cells

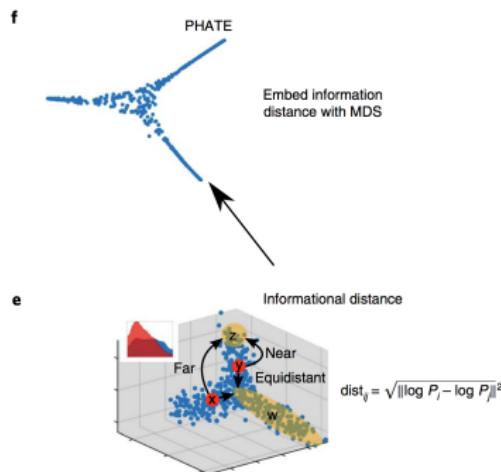


Figure: from Moon et al. Nature Biotechnology. 2019

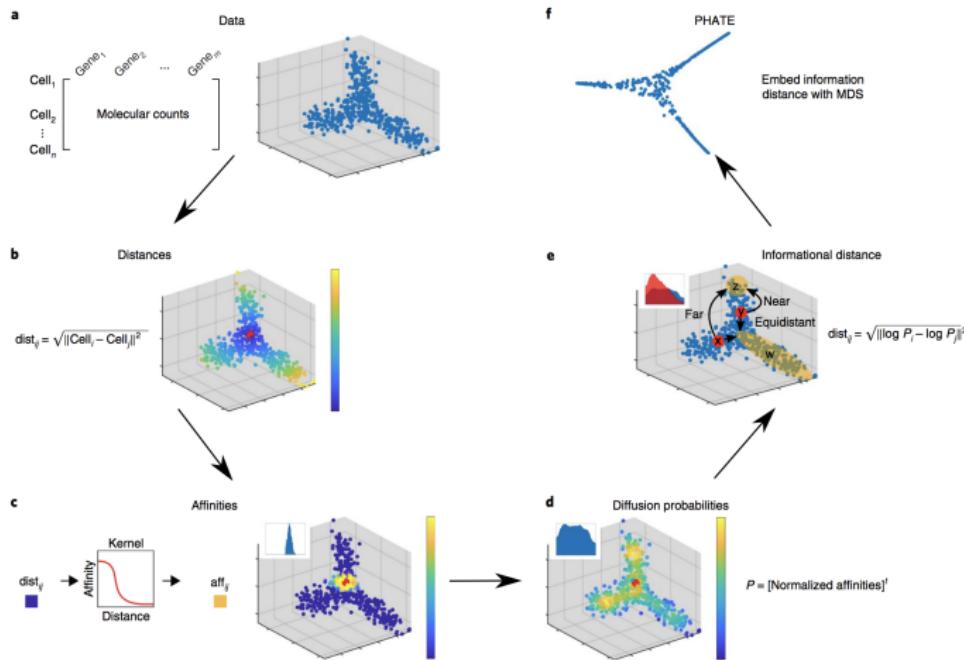
Potential Distances, Written Formally

- $U_x^t = -\log(P_x^t)$
- Then $\text{PD}(\mathbf{x}, \mathbf{y}) = \|U_x^t - U_y^t\|_2$

Then use a classic algorithm called multi-dimensional scaling to project all points in two dimensions. Apply to,

$$\mathfrak{V}^t(\mathbf{x}, \mathbf{y}) = \|U_x^t - U_y^t\|_2, \mathbf{x}, \mathbf{y} \in \mathcal{X}$$

Recap



PHATE preserves branch structure

As expected tSNE is preserving local similarities, but not necessarily distances between groups of points.

b

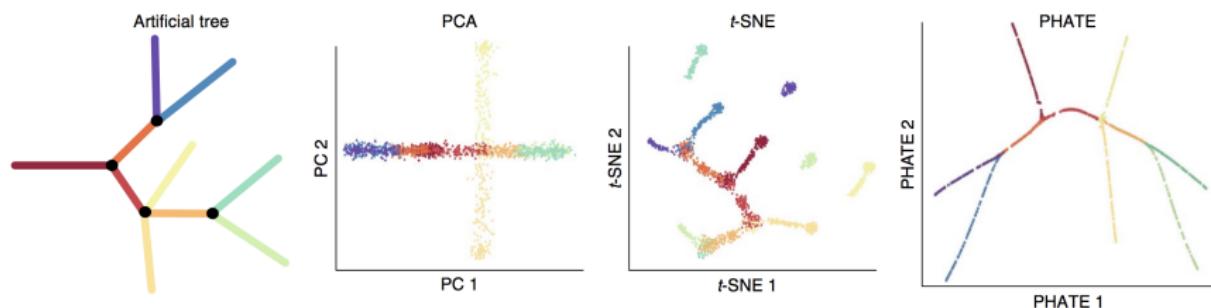


Figure: from Moon et al. Nature Biotechnology. 2019

And Preserves Cellular Differentiation Patterns!

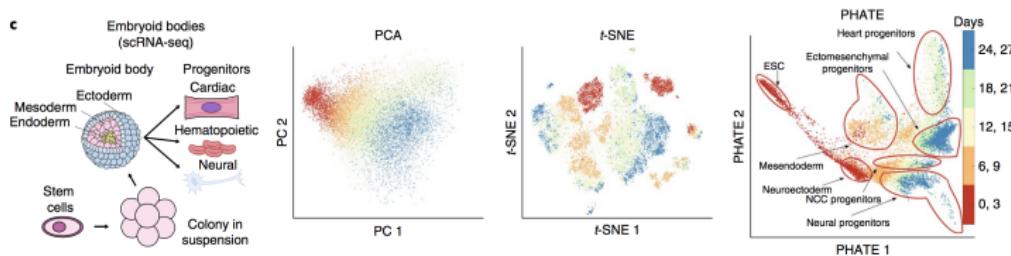


Figure: from Moon et al. Nature Biotechnology. 2019

Compare Phate to Other Dimension Reduction Techniques

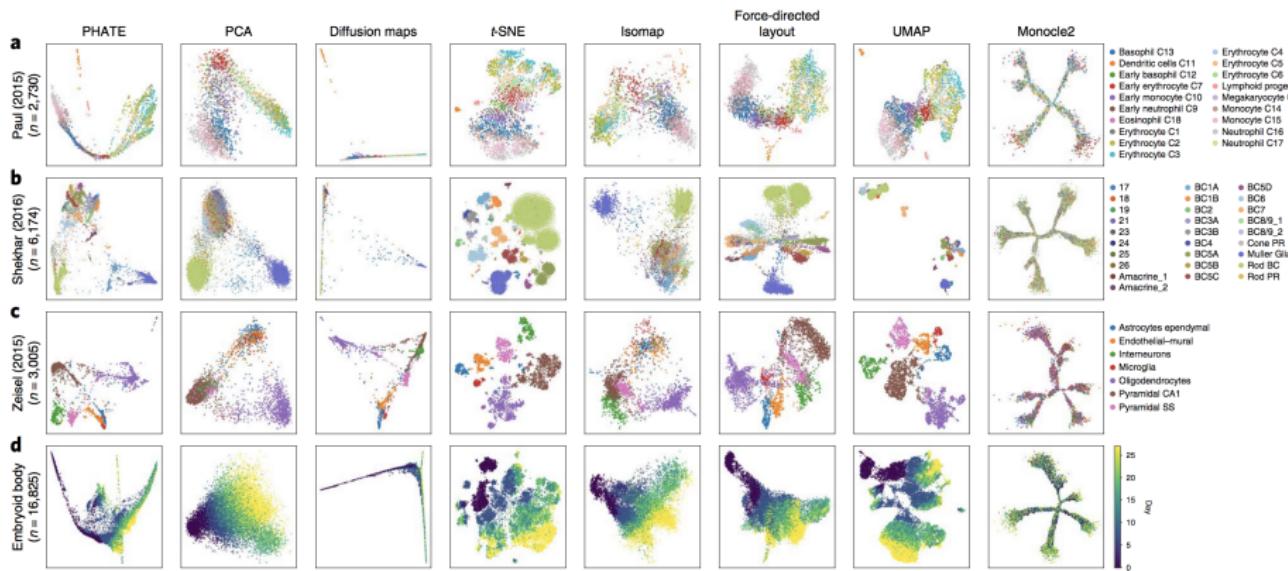


Figure: from Moon *et al.* Nature Biotechnology. 2019

Rare + Sparse Cell Populations

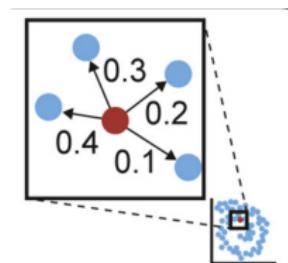


Figure: from MAGIC paper

- ‘Dense’ cell-populations vs ‘Sparse’ cell-populations.
- Low-frequency cell-populations and data artifacts can fail to be under-represented and not accurately reflect the underlying biology.
- Are members of a particular cell-type not there or are they just very infrequent?

General Problems with Downstream Tasks from Sparse Data

- Imbalanced classes can affect classification accuracy (**classifying cells**)
- In clustering, imbalanced ‘ground-truth’ clusters can cause distortion or mis-representation of the clusters in the data (**frequency features**)
- Too few samples/observations can cause mis-representation of the dependencies or correlations between features.

Welcome SUGAR

Generate points uniformly from intrinsic data geometry / manifold:

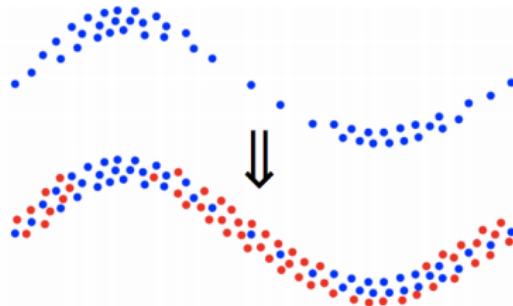


Figure: from Lindenbaum *et al.* NeurIPS. 2018

- Most generative modeling approaches seek to learn and replicate the data density
- SUGAR is a data-generation approach that uses the underlying geometry of the data (e.g. a random walk based approach)

Common Data Generation Approaches

- Older Techniques
 - **Parametric Models :** Specify a model and optimize the parameters through maximum likelihood optimization.
 - Use the learned model to generate new data
 - Use a histogram or kernel to estimate generating distributions
- Newer techniques to generate additional points from complicated data distributions
 - GAN (generative adversarial network) →
<https://papers.nips.cc/paper/2014/file/5ca3e9b122f61f8f06494c97b1afccf3-Paper.pdf>
 - VAE (variational autoencoders) →
<https://arxiv.org/abs/1606.05908>

Remembering our Good Friend Gaussian Kernel

For a pair of nodes, i and j , the strength of their connection can be computed as $\mathcal{K}(\mathbf{x}_i, \mathbf{x}_j)$, or

$$K_{ij} = \exp\left(-\frac{||\mathbf{x}_i - \mathbf{x}_j||^2}{2\sigma^2}\right) \quad (2)$$

We also remember our row-stochastic markov matrix, \mathbf{P} computed from \mathbf{K} .

Measure Based Gaussian Correlation (MGC)

- Given your data, \mathbf{X} , define a set of reference points, \mathbf{r} with $r \in \mathbf{X}$.
- Define $\mu(\mathbf{r})$ as some measure over the reference points

Then define a new kernel, $\hat{\mathcal{K}}(\mathbf{x}_i, \mathbf{x}_j)$ as,

$$\hat{\mathcal{K}}(\mathbf{x}_i, \mathbf{x}_j) = \sum_{\mathbf{r} \in \mathbf{X}} \mathcal{K}(\mathbf{x}_i, \mathbf{r}) \mathcal{K}(\mathbf{x}_j, \mathbf{r}) \mu(\mathbf{r}) \quad (3)$$

For our purposes, $\mu(\mathbf{r})$ will be some value that relates to sparsity and is the inverse of node r 's degree.

Problem Formulation for Data Generation Problem

- Let \mathcal{M} be a d -dimensional manifold such that $\mathcal{M} \in \mathbb{R}^d$
- Let $\mathbf{X} \subset \mathcal{M}$ be a subset of points samples from \mathcal{M} with $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\} \in \mathbf{X}$
- Assuming that instances, \mathbf{X} are unevenly sampled from \mathcal{M} , we seek a set of new points, $\mathbf{Y} = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_M\}$ that also lie of \mathcal{M} and that $\mathbf{X} \cup \mathbf{Y}$ is uniform.

Synthesis Using Geometrically Aligned Random Walks

Algorithm 1 SUGAR: Synthesis Using Geometrically Aligned Random-walks

Input: Dataset $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}, \mathbf{x}_i \in \mathbb{R}^D$.

Output: Generated set of points $\mathbf{Y} = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_M\}, \mathbf{y}_i \in \mathbb{R}^D$.

- 1: Compute the diffusion geometry operators \mathbf{K} , \mathbf{P} , and degrees $\hat{d}(i)$, $i = 1, \dots, N$ (see Sec. 3)
- 2: Define a sparsity measure $\hat{s}(i), i = 1, \dots, N$ (Eq. 2).
- 3: Estimate a local covariance Σ_i , $i = 1, \dots, N$, using k nearest neighbors around each \mathbf{x}_i .
- 4: For each point $i = 1, \dots, N$ draw $\hat{\ell}(i)$ vectors (see Sec. 4.3) from a Gaussian distribution $\mathcal{N}(\mathbf{x}_i, \Sigma_i)$. Let $\hat{\mathbf{Y}}_0$ be a matrix with these $M = \sum_{i=1}^N \hat{\ell}(i)$ generated vectors as its rows.

Figure: from Lindenbaum *et al.* NeurIPS. 2018. Let's walk through steps 1-4 for now.

Synthesis Using Geometrically Aligned Random Walks

- **Specify Kernel:** Initialized by forming a Gaussian Kernel over the input data, \mathbf{X} , \mathbf{G}_x .
 - Use \mathbf{G}_x to estimate the degree of each node i , $d(i)$.
 - The sparsity of each point, $s(i)$ is defined as the inverse degree of node i or $s(i) = 1/d(i)$.
- **Sample According to Each Point** Next, sample $\ell(i)$ points around each $\mathbf{x}_i \in \mathbf{X}$ from a localized gaussian distribution (e.g. k -nearest points around i).
 - $G_i = \mathcal{N}(\mathbf{x}_i, \Sigma_i)$

Practical : Sampling from MV Gaussians

`numpy.random.multivariate_normal`

```
random.multivariate_normal(mean, cov, size=None,  
check_valid='warn', tol=1e-8)
```

Draw random samples from a multivariate normal distribution.

The multivariate normal, multinormal or Gaussian distribution is a generalization of the one-dimensional normal distribution to higher dimensions. Such a distribution is specified by its mean and covariance matrix. These parameters are analogous to the mean (average or “center”) and variance (standard deviation, or “width,” squared) of the one-dimensional normal distribution.

Back to SUGAR

- Let \mathbf{Y}_0 be the set of all new $M = \sum_i \ell(i)$ points generated around each i , with $\mathbf{Y}_0 = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_M\}$
- **MGC Kernel:** Now use affinities between points in \mathbf{X} and \mathbf{Y}_0 . Here, points in \mathbf{X} are used as reference.

$$\hat{\mathcal{K}}(\mathbf{y}_i, \mathbf{y}_j) = \sum_r \mathcal{K}(\mathbf{y}_i, \mathbf{x}_r) \mathcal{K}(\mathbf{x}_r, \mathbf{y}_j) s(r) \quad (4)$$

Pulling \mathbf{Y}_0 towards sparser regions of \mathcal{M}

Pasted psuedo code for the second part of SUGAR

- 5: Compute the sparsity based diffusion operator $\hat{\mathbf{P}}$ (see Sec 4.2).
- 6: Apply the operator $\hat{\mathbf{P}}$ at time instant t to the new generated points in $\hat{\mathbf{Y}}_0$ to get diffused points as rows of $\mathbf{Y}_t = \hat{\mathbf{P}}^t \cdot \mathbf{Y}_0$.
- 7: Rescale \mathbf{Y}_t to get the output $\mathbf{Y}[\cdot, j] = \mathbf{Y}_t[\cdot, j] \cdot \frac{\text{percentile}(\mathbf{X}_{[\cdot, j], 99})}{\max_{t[\cdot, j]} \mathbf{Y}_t[\cdot, j]}$, $j = 1, \dots, D$, in order to fit the original range of feature values in the data.

Figure: from Lindenbaum *et al.* NeurIPS. 2018.

Diffusion Operator Again

- Take affinities from $\hat{\mathcal{K}}$ and convert them to \mathbf{P} , the row-normalized Markov matrix.
- This will allow us to correct points in \mathbf{Y}_0 according to neighborhood regions

As you will all recognize, powering \mathbf{P} as \mathbf{P}^t estimates the probability of successfully traveling between nodes with t steps. The transformed matrix, \mathbf{Y}_t is computed as

$$\mathbf{Y}_t = \mathbf{P}^t \times \mathbf{Y}_o \quad (5)$$

Same Story Regarding t

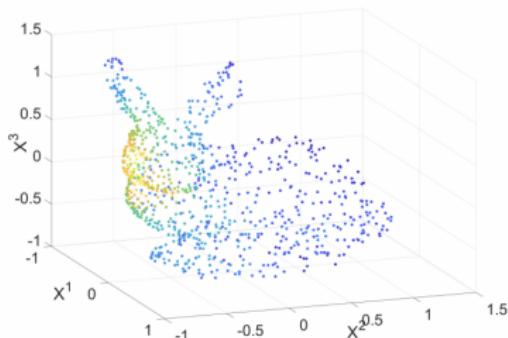
Remember in PHATE, t was chosen according to the knee point of the Von-Neumann entropy of the normalized eigenvalues of \mathbf{P}^t .

- Let $[\eta(t)]_i = \lambda_i^t / \sum_{j=0}^{N-1} \lambda_j^t$ be the probability distribution defined by the non-negative eigenvalues of \mathbf{P}^t .

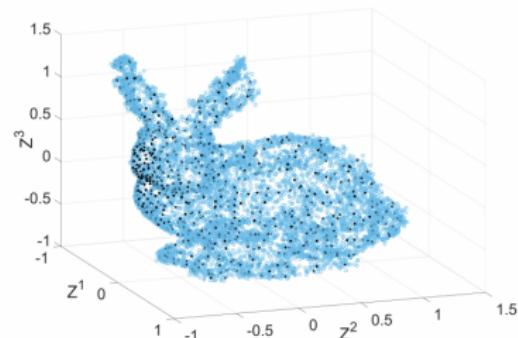
$$H(t) = - \sum_{i=1}^N [\eta(t)]_i \log([\eta(t)]_i) \quad (6)$$

Experiments : Synthetic and Biological

Stanford Bunny



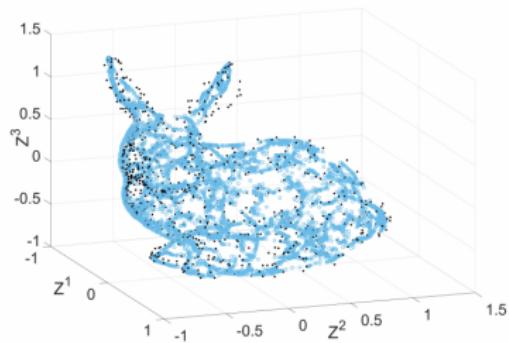
(a)



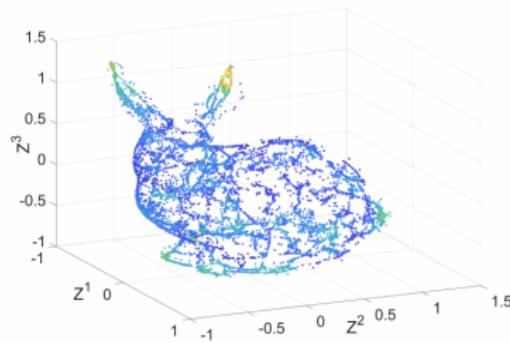
(b)

Figure: from Lindenbaum *et al.* NeurIPS. 2018. (a). The original set, \mathbf{X} of points, colored by node degree. (b) \mathbf{Y}_0 are generated points (blue) and original points, \mathbf{X} (black).

Stanford Bunny Part II



(c)



(d)

Figure: from Lindenbaum *et al.* NeurIPS. 2018. (c). Original and generated points, before MGC diffusion. (d) Generated points (\mathbf{Y}) after MGC diffusion. Points are colored by degree.

Application : Augmented Clustering

SUGAR was used to generate additional data points to improve the quality of clusters identified with k -means.

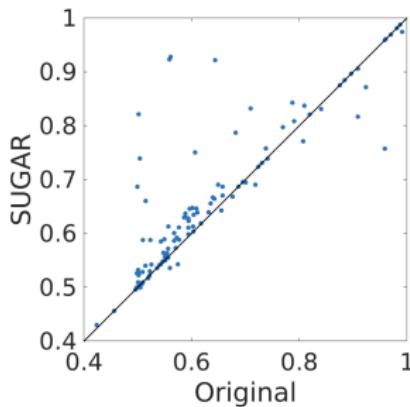


Figure: from Lindenbaum *et al.* NeurIPS. 2018. In 119 datasets, the data were augmented with additional datapoints using sugar. Adjusted Rand Index was computed for the original data vs data + SUGAR.

SUGAR on Single-Cell

SUGAR was used to augment cells in a single-cell RNA-seq dataset.

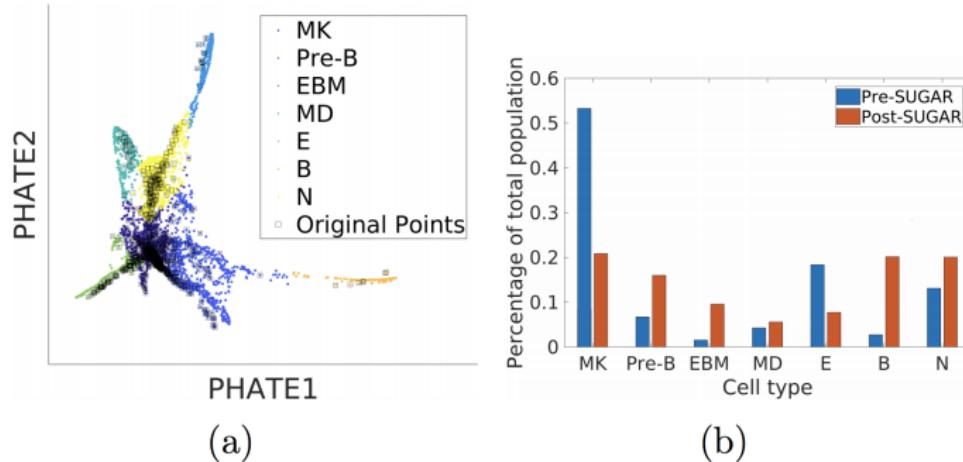


Figure: from Lindenbaum *et al.* NeurIPS. 2018

Maintaining Intra-Module Marker Co-Expression

Among cells assigned to the same module or cluster, after SUGAR led to higher intra-module between-marker correlation.

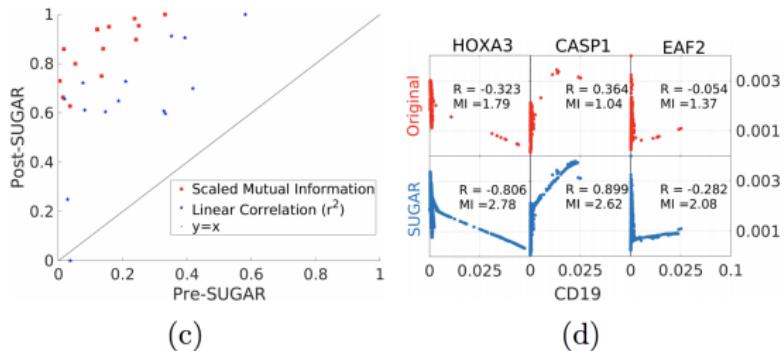


Figure: from Lindenbaum *et al.* NeurIPS. 2018