

Similarity-assisted Variational Autoencoder for Nonlinear Dimension Reduction with Application to Single-cell RNA Sequencing Data

Gwangwoo Kim and Hyonho Chun

Department of Mathematical Sciences
Korea Advanced Institute of Science and Technology



Outline

1 Introduction

2 Preliminary

3 Method

4 Experiment

5 Conclusion

Introduction

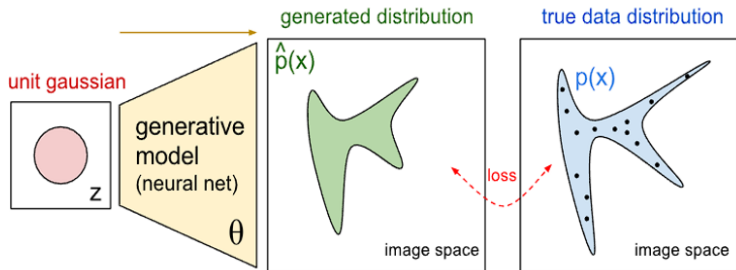
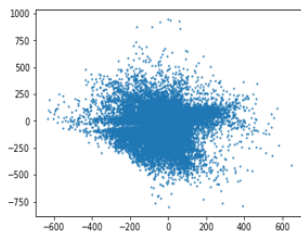


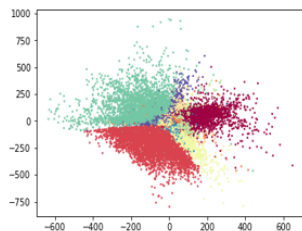
Figure: The generative modeling process

- A generative model describes how a dataset is generated.

Introduction



(a) Latent space without true labels

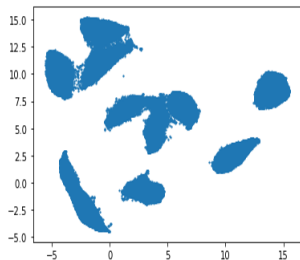


(b) Latent space without true labels

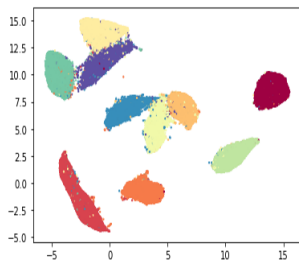
Figure: The latent encoding of the MNIST data from variational autoencoder (VAE)

- The inferred space is hard to interpret without further information.

Introduction



(a) Latent space without true labels



(b) Latent space without true labels

Figure: The visualization of the same data from a dimension reduction method

- Incorporate the similarity information to better illustrate the latent structure.

Introduction

Deep generative models

- **Variational autoencoder (VAE)**
- Generative adversarial networks (GAN)
- Flow-based generative model

Dimension Reduction Algorithms

- **Uniform manifold approximation and projection (UMAP)**
- t-distributed stochastic neighbor embedding (t-SNE)
- Large-scale Dimensionality Reduction Using Triplets (TriMAP)

Outline

- 1 Introduction
- 2 Preliminary
- 3 Method
- 4 Experiment
- 5 Conclusion

Preliminary

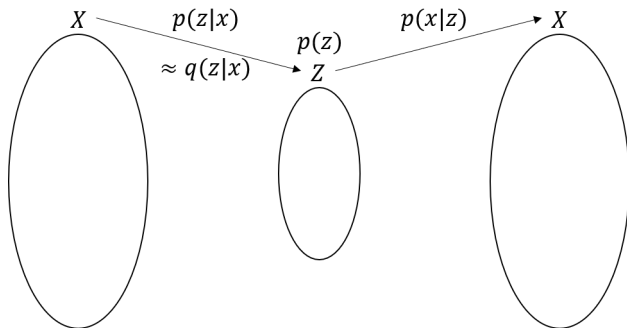


Figure: Variational autoencoder (VAE)

- VAE (Kingma and Welling, 2014) are generative models that attempt to describe data generation through a probabilistic distribution.

Preliminary

Evidence Lower Bound (ELBO)

The evidence lower bound (ELBO) is an useful lower bound on the log-likelihood of some observed data.

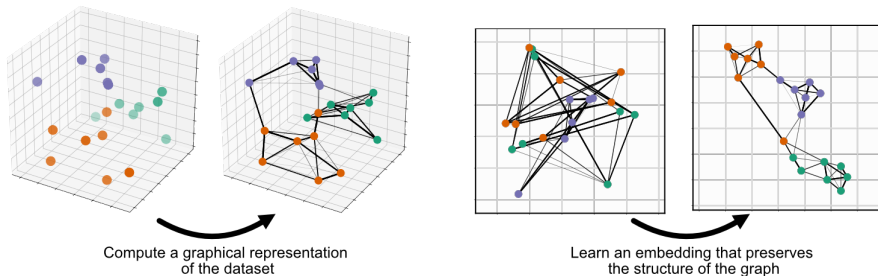
$$\mathbb{E}_{q_{\phi}(z|x)} \log \frac{p_{\theta}(z, x)}{q_{\phi}(z|x)} \leq \log p_{\theta}(x)$$

Remark

The ELBO can be decomposed into two components:

$$\begin{aligned} \mathbb{E}_{q_{\phi}(z|x)} \log \frac{p_{\theta}(z, x)}{q_{\phi}(z|x)} &= \mathbb{E}_{q_{\phi}(z|x)} \log p_{\theta}(x|z) - \mathbb{E}_{q_{\phi}(z|x)} \log \frac{q_{\phi}(z|x)}{p_{\theta}(z)} \\ &= \mathbb{E}_{q_{\phi}(z|x)} \log p_{\theta}(x|z) - D_{\text{KL}}(q_{\phi}(z|x) \parallel p_{\theta}(z)). \end{aligned}$$

Preliminary



- Uniform manifold approximation and projection (McInnes et al., 2018) is an algorithm for dimension reduction based on manifold learning techniques.

Graph construction

- 1 For an observation x_i , denote the k -nearest neighborhoods by $\{x_{i_j}\}_{j=1}^k$. Define the weight function by

$$w(x_i, x_j) = \begin{cases} \exp\left(\frac{-\max(0, d(x_i, x_j) - \rho_i)}{\sigma_i}\right) & \text{for } j \in \{i_1, \dots, i_k\} \\ 0 & \text{otherwise,} \end{cases}$$

where ρ_i is the minimum positive distance from x_i and σ_i is a normalizing constant.

- 2 The symmetrized similarity can be computed as

$$\mu_{ij} = w(x_i, x_j) + w(x_j, x_i) - w(x_i, x_j)w(x_j, x_i).$$

Layout optimization

- 1 Let y_i be an initial embedding of x_i .
- 2 The similarities among the embeddings are approximations to μ_{ij} , which can be defined as

$$\nu_{ij} := \frac{1}{1 + a\|y_i - y_j\|^{2b}},$$

where a and b are pre-defined numbers.

- 3 Optimize the cross entropy loss of μ_{ij} and ν_{ij} , i.e.,

$$-\sum_i \sum_{j \neq i} \mu_{ij} \ln(\nu_{ij}) + (1 - \mu_{ij}) \ln(1 - \nu_{ij}).$$

Outline

- 1 Introduction
- 2 Preliminary
- 3 Method**
- 4 Experiment
- 5 Conclusion

Contributions

There are three big deals to integrate the mentioned algorithms.

- Their optimization targets are different.
- Their batch schemes conflict each other.
- It is difficult to strike a balance between them.

Method

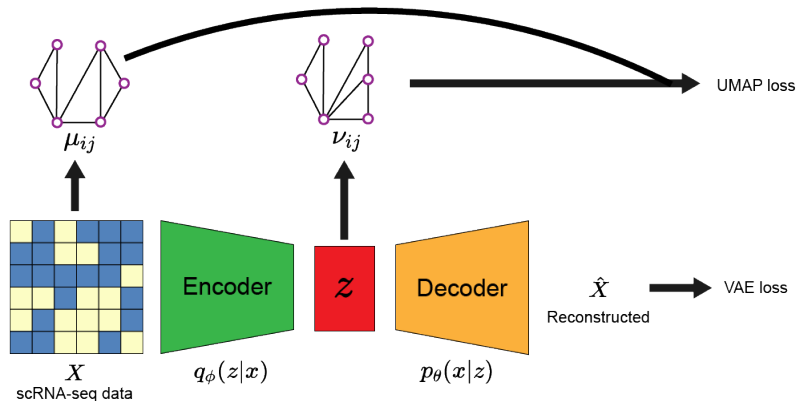


Figure: Overview of saVAE

- Our saVAE connects the VAE and UMAP objective functions by utilizing the expected UMAP loss function,

saVAE

Similarity-assisted VAE (saVAE) unifies VAE and UMAP, and its objective function is given by

$$-(\text{ELBO}) - \lambda \sum_i \sum_{j \neq i} \mu_{ij} \ln(\tilde{\nu}_{ij}) + (1 - \mu_{ij}) \ln(1 - \tilde{\nu}_{ij}),$$

where $\tilde{\nu}_{ij} = \left[\mathbb{E}_{q_\phi(z|x_i)} \mathbb{E}_{q_\phi(z'|x_j)} \frac{1}{1+a\|z-z'\|^{2b}} \right]$ and λ is a weight.

Method

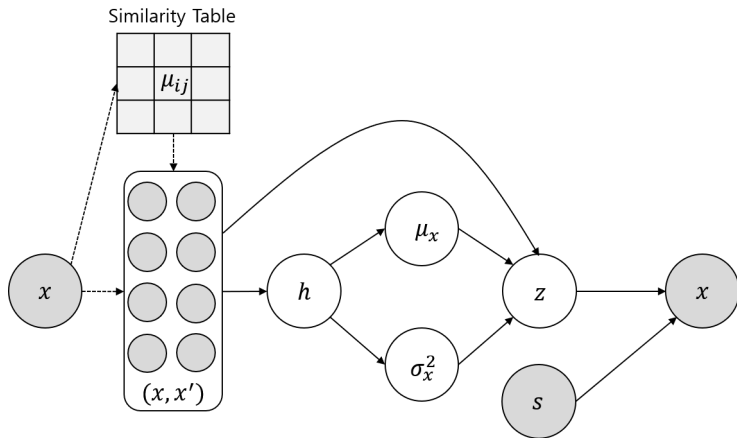


Figure: Graphical model of our saVAE

Method

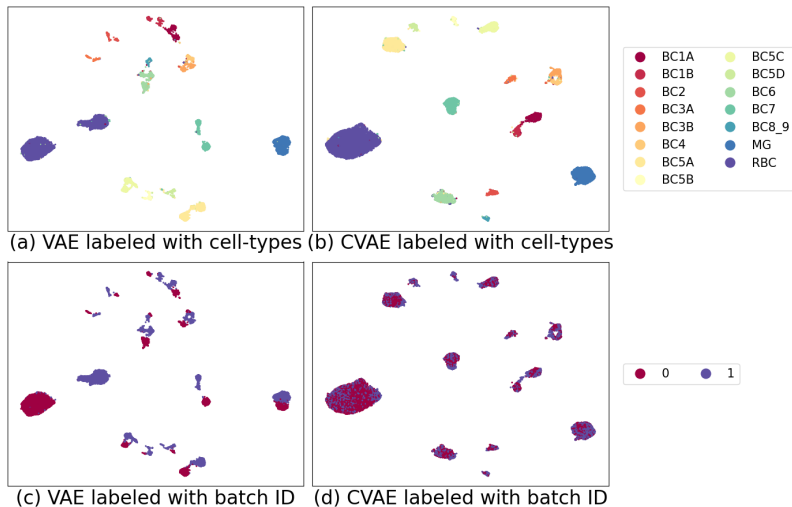


Figure: Embeddings from retina dataset

saCVAE

For additional information s , our objective function of the saCVAE is

$$\begin{aligned} & -\mathbb{E}_{q_\phi(z|x)} \log p_\theta(x|z, \textcolor{red}{s}) + \mathbb{E}_{q_\phi(z|x)} \log \frac{q_\phi(z|x)}{p_\theta(z)} \\ & -\lambda \sum_i \sum_{j \neq i} \textcolor{red}{\tilde{\mu}_{ij}} \ln(\tilde{\nu}_{ij}) + (1 - \textcolor{red}{\tilde{\mu}_{ij}}) \ln(1 - \tilde{\nu}_{ij}), \end{aligned}$$

where $\tilde{\mu}_{ij}$ is the weight computed from the covariate-adjusted embeddings and $\tilde{\nu}_{ij}$ and λ are defined in the previous section.

Outline

- 1 Introduction
- 2 Preliminary
- 3 Method
- 4 Experiment**
- 5 Conclusion

Experiment

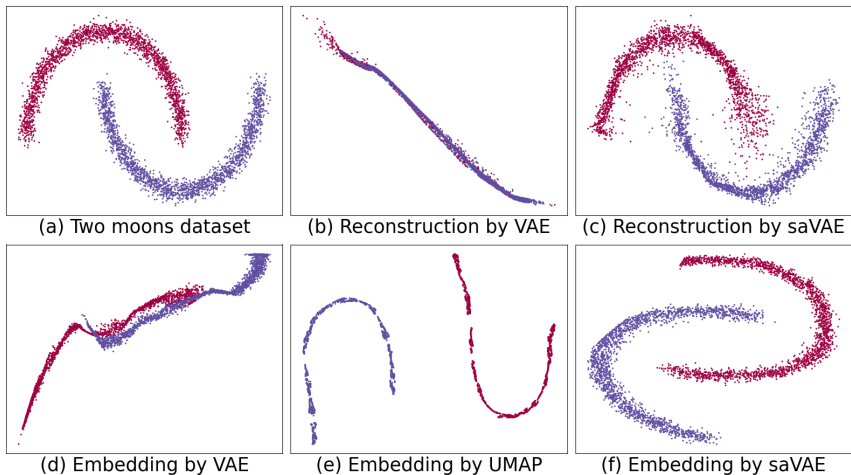


Figure: Method comparison using two moons dataset

Experiment

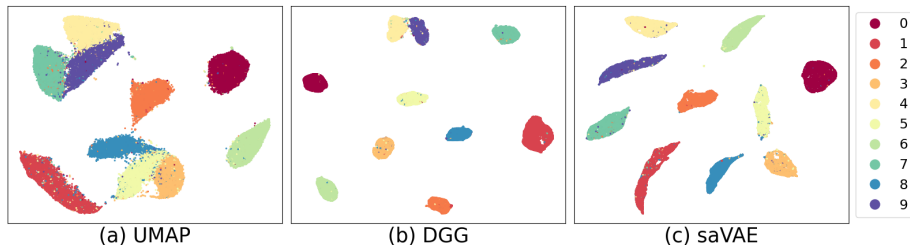


Figure: Embeddings from MNIST dataset, which is a set of handwritten digits images from 0 to 9.

Experiment

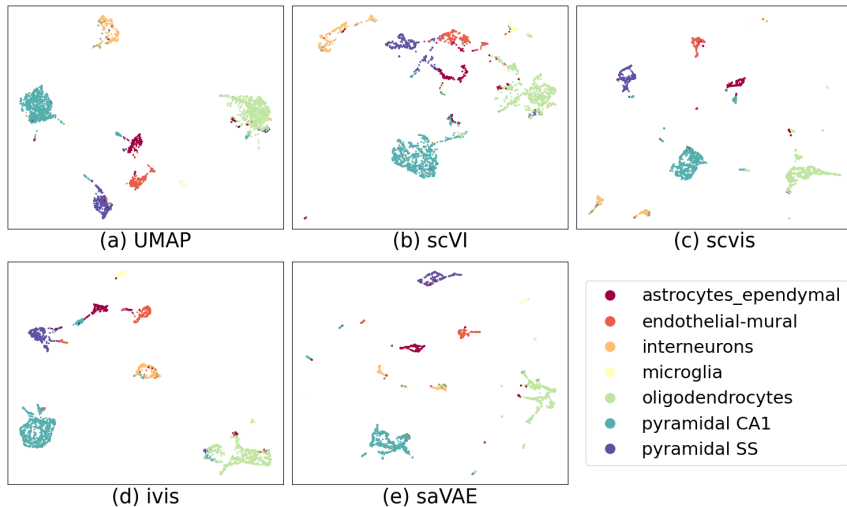


Figure: Embeddings from cortex dataset

Experiment

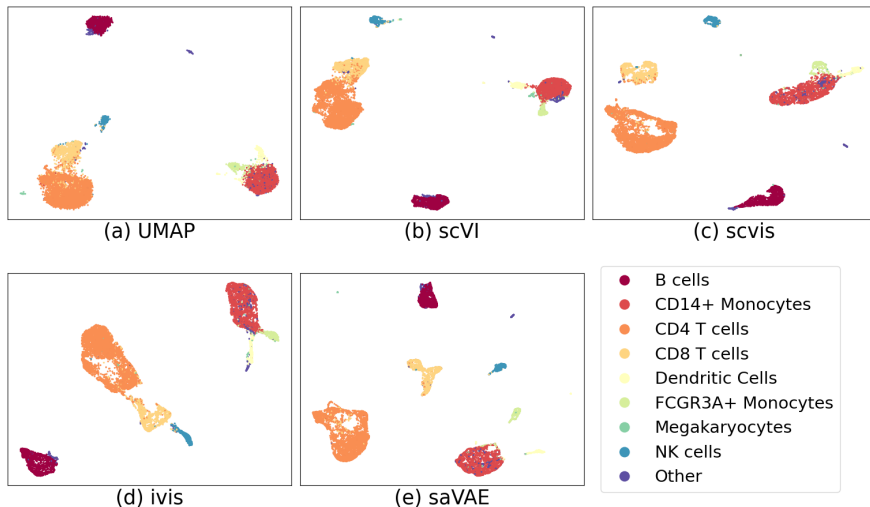


Figure: Embeddings from PBMC dataset

Experiment

Dataset	Method	ARI	NMI	Dataset	Method	ARI	NMI
cortex	UMAP	0.79	0.78	PBMC	UMAP	0.84	0.80
	scvis	0.79	0.77		scvis	0.82	0.79
	Ivis	0.70	0.72		Ivis	0.63	0.69
	scVI	0.75	0.72		scVI	0.53	0.72
	saVAE	0.80	0.77		saVAE	0.84	0.81

Table: Performance comparison. Our saVAE performs better than UMAP, scvis, Ivis, and scVI. All values are reported as the highest out of 10 trials.

Experiment

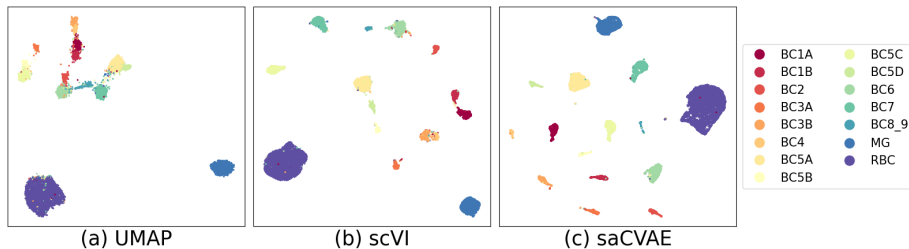


Figure: Embeddings from retina dataset

Experiment

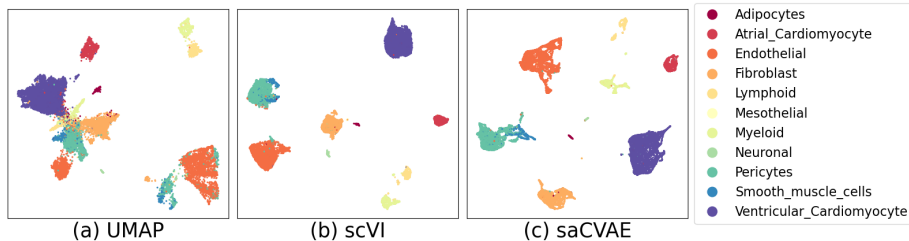


Figure: Embeddings from heart cell atlas dataset

Experiment

Dataset	Method	ARI	NMI	Dataset	Method	ARI	NMI
retina	UMAP	0.90	0.85	heart cell atlas	UMAP	0.64	0.69
	scVI	0.56	0.84		scVI	0.72	0.80
	saVAE	0.98	0.95		saVAE	0.93	0.92

Table: Performance comparison. Our saCVAE formulation effectively eliminates the covariate effects, resulting in a better cell-type clustering performance. All values are reported as the highest out of 10 trials.

Outline

- 1 Introduction
- 2 Preliminary
- 3 Method
- 4 Experiment
- 5 Conclusion**

Conclusions

- We propose a framework to combine VAE based and similarity based approaches to reflect intrinsic structures in the data.
- Computational conflicts are resolved via applying a sampling idea to a normal mini-batch.
- By introducing CVAE, it is possible to adjust the covariate effect of transcriptomes, with one framework.

- Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. In Yoshua Bengio and Yann LeCun, editors, *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, 2014.
- Leland McInnes, John Healy, Nathaniel Saul, and Lukas Grossberger. UMAP: Uniform Manifold Approximation and Projection. *The Journal of Open Source Software*, 3:861, 2018.