

Interpretable Factor Models of Single-cell RNA-Seq via Variational Autoencoders

2019/10/09

Ping-Han Hsieh

Outline

- Background
 - Single-cell RNA-Seq
 - Neural network
 - Autoencoder
 - Variational autoencoder (VAE)
- Methods
 - Linearly decoded variational autoencoder (LDVAE)
- Results
- Discussion

Background

Single-cell RNA-Seq

- Problem
 - scRNA-Seq is useful to analyze relationship between genes depend on cell types.
 - In order to investigate such interaction, we need to learn the manifold of gene expression in different cell types.
- Existing Method
 - Principle Component Analysis: Gaussian likelihood not suitable for RNA-Seq data (Negative Binomial Distribution)
 - ZINB-WaVE: not scalable to big data
 - Variational Autoencoder (scVI): hard to interpret the result

Neural Network – Representation (1)

$$\mathbf{X} = \begin{bmatrix} X_{1,1} & X_{1,2} & \cdots & X_{1,d} \\ X_{2,1} & X_{2,2} & \cdots & X_{2,d} \\ \vdots & \vdots & \ddots & \vdots \\ X_{n,1} & X_{n,2} & \cdots & X_{n,d} \end{bmatrix}$$

Features

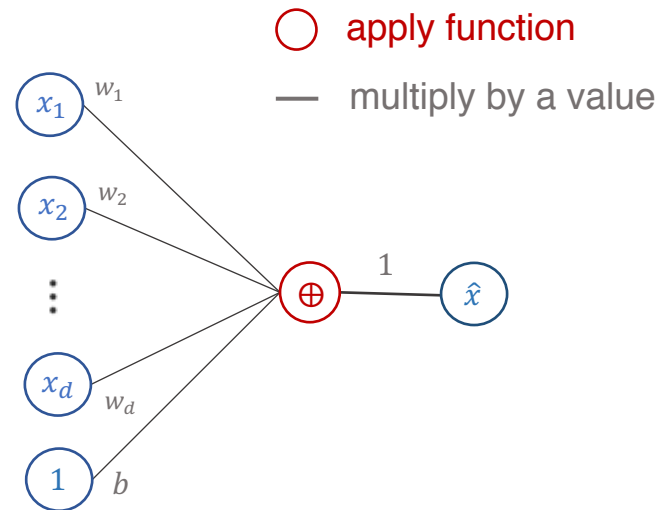
Samples

$$\begin{aligned} x_k &\rightarrow [X_{1,k} \quad X_{2,k} \quad \cdots \quad X_{n,k}] \\ x^{(i)} &\rightarrow [X_{i,1} \quad X_{i,2} \quad \cdots \quad X_{i,d}] \\ \mathbf{w} &= [w_1 \quad w_2 \quad \cdots \quad w_d] \end{aligned}$$

Linear Combination

$$\hat{\mathbf{x}} = \mathbf{X}\mathbf{w}^T + b = w_1\mathbf{x}_1 + w_2\mathbf{x}_2 + \dots + w_d\mathbf{x}_d + b$$

$$w_1 \begin{bmatrix} x_1 \\ X_{1,1} \\ X_{2,1} \\ \vdots \\ X_{n,1} \end{bmatrix} + w_2 \begin{bmatrix} x_2 \\ X_{1,2} \\ X_{2,2} \\ \vdots \\ X_{n,2} \end{bmatrix} + \cdots + w_d \begin{bmatrix} x_d \\ X_{1,d} \\ X_{2,d} \\ \vdots \\ X_{n,d} \end{bmatrix} + b$$



Neural Network – Representation (2)

$$\mathbf{X} = \begin{bmatrix} X_{1,1} & X_{1,2} & \cdots & X_{1,d} \\ X_{2,1} & X_{2,2} & \cdots & X_{2,d} \\ \vdots & \vdots & \ddots & \vdots \\ X_{n,1} & X_{n,2} & \cdots & X_{n,d} \end{bmatrix}$$

Features

Observations

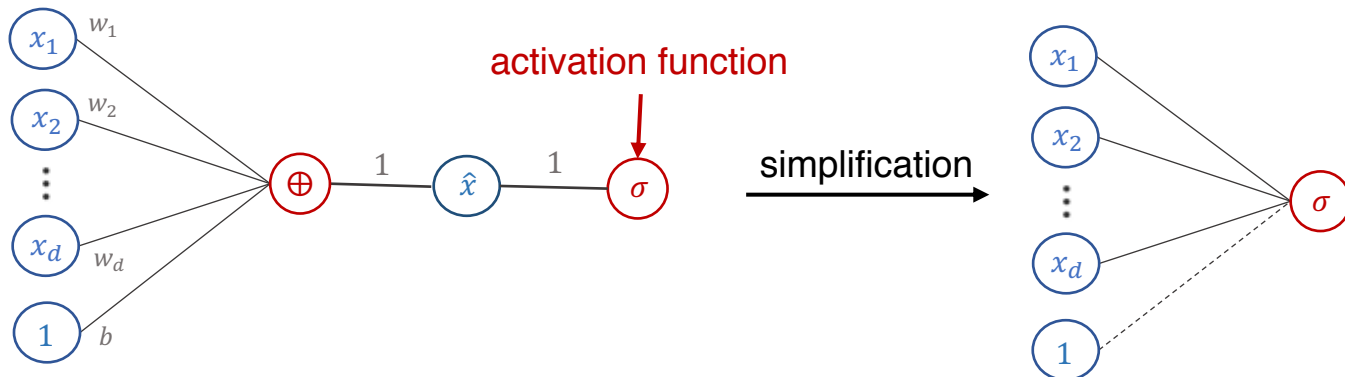
$$\begin{aligned} x_k &\rightarrow [X_{1,k} \quad X_{2,k} \quad \cdots \quad X_{n,k}] \\ x^{(i)} &\rightarrow [X_{i,1} \quad X_{i,2} \quad \cdots \quad X_{i,d}] \\ \mathbf{w} &= [w_1 \quad w_2 \quad \cdots \quad w_d] \end{aligned}$$

Linear combination + activation function

$$\sigma(w_1 \mathbf{x}_1 + w_2 \mathbf{x}_2 + \dots + w_d \mathbf{x}_d + b)$$

○ apply function

— multiply by a value



Neural Network – Gradient Descent

1. Define the loss function

$$L = \sum_{i=1}^n (\hat{y}_i - y_i)^2$$
$$= \sum_{i=1}^n (w_1 X_{i,1} + w_2 X_{i,2} + \dots w_d X_{i,d} - y_i)^2$$

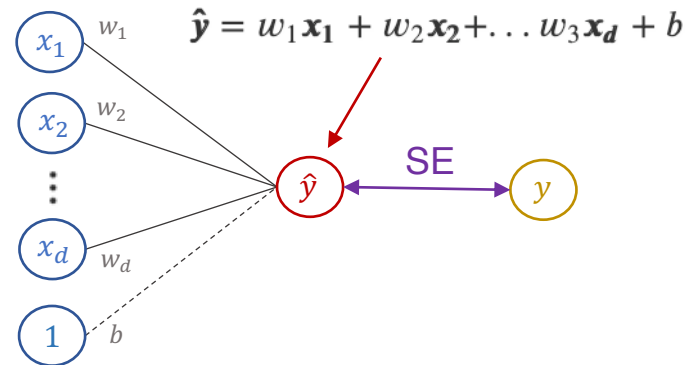
2. The loss function is actually a function of w

3. Rewrite the loss function to $L(w)$

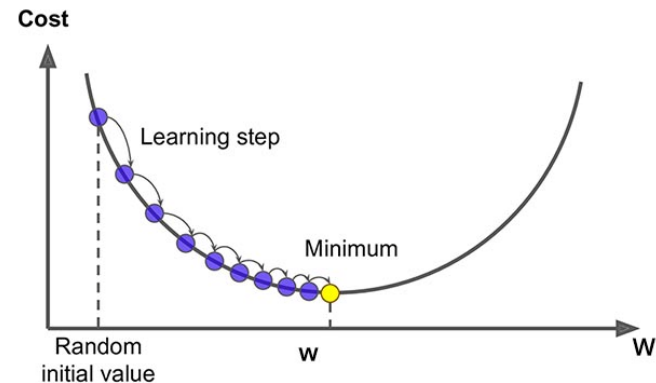
4. We can use the gradient to approximate the minimum loss

$$w^1 = w^0 - \eta \nabla_w L|_{w^0}$$
$$w^2 = w^1 - \eta \nabla_w L|_{w^1}$$
$$w^3 = w^2 - \eta \nabla_w L|_{w^2}$$
$$\vdots$$

Linear regression



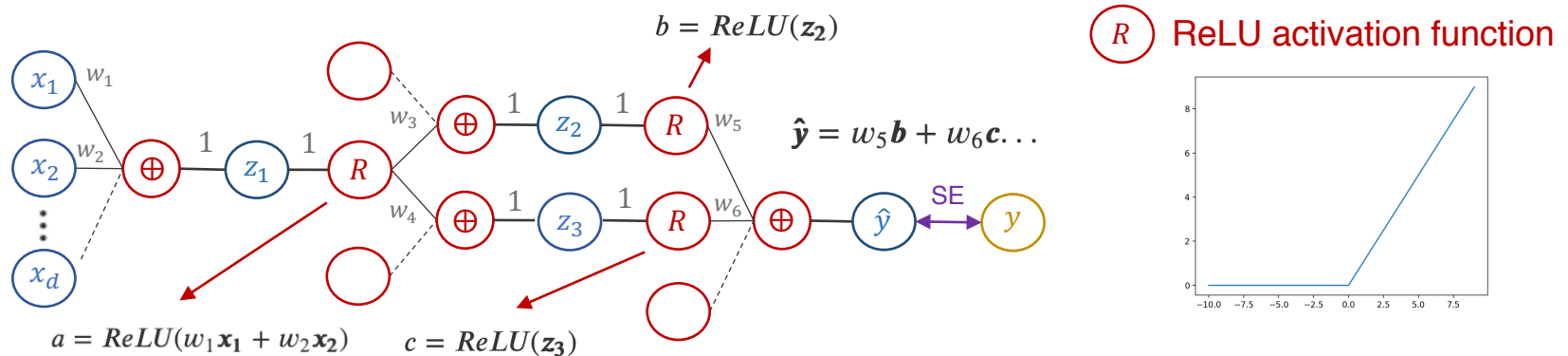
Gradient descent



from Saugat Bhattarai

Neural Network – Backpropagation (1)

Consider more than one hidden layer



1. Feed forward, compute the derivative for \mathbf{z} w.r.t \mathbf{w}

$$\frac{\partial \mathbf{z}_1}{\partial w_1} = \mathbf{x}_1 \quad \frac{\partial \mathbf{z}_2}{\partial w_3} = \mathbf{a}$$

$$\frac{\partial \mathbf{z}_1}{\partial w_2} = \mathbf{x}_2 \quad \frac{\partial \mathbf{z}_3}{\partial w_4} = \mathbf{a}$$

The derivative is the corresponding input of the node

2. Backward pass, compute the derivative for L w.r.t \mathbf{z}

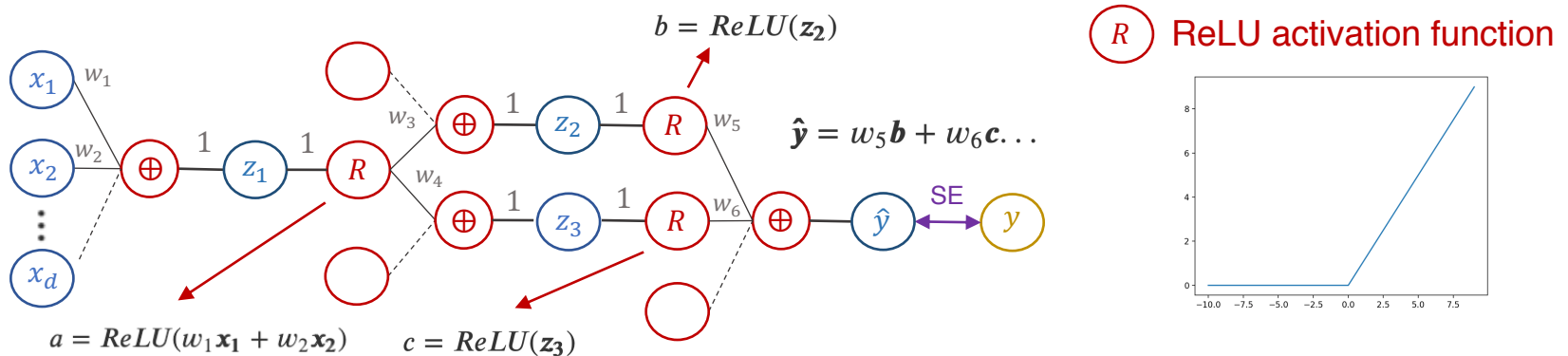
$$\frac{\partial L}{\partial \mathbf{z}_3} = \frac{\partial L}{\partial \mathbf{c}} \frac{\partial \mathbf{c}}{\partial \mathbf{z}_3} = \frac{\partial L}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{c}} \frac{\partial \mathbf{c}}{\partial \mathbf{z}_3}$$

$$\frac{\partial L}{\partial \mathbf{z}_2} = \frac{\partial L}{\partial \mathbf{b}} \frac{\partial \mathbf{b}}{\partial \mathbf{z}_2} = \frac{\partial L}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{b}} \frac{\partial \mathbf{b}}{\partial \mathbf{z}_2}$$

$2\hat{\mathbf{y}} \quad w_6 \quad I\{z_3 > 0\}$
 $2\hat{\mathbf{y}} \quad w_5 \quad I\{z_2 > 0\}$

Neural Network – Backpropagation (2)

Consider more than one hidden layer



$$\frac{\partial z_1}{\partial w_1} = x_1 \quad \frac{\partial z_2}{\partial w_3} = a \quad \frac{\partial L}{\partial z_3} = \frac{\partial L}{\partial c} \frac{\partial c}{\partial z_3} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial c} \frac{\partial c}{\partial z_3} \quad \frac{\partial L}{\partial z_2} = \frac{\partial L}{\partial b} \frac{\partial b}{\partial z_2} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial b} \frac{\partial b}{\partial z_2}$$

$$\frac{\partial z_1}{\partial w_2} = x_2 \quad \frac{\partial z_3}{\partial w_4} = a$$

2. Backward pass, compute the derivative for L w.r.t z

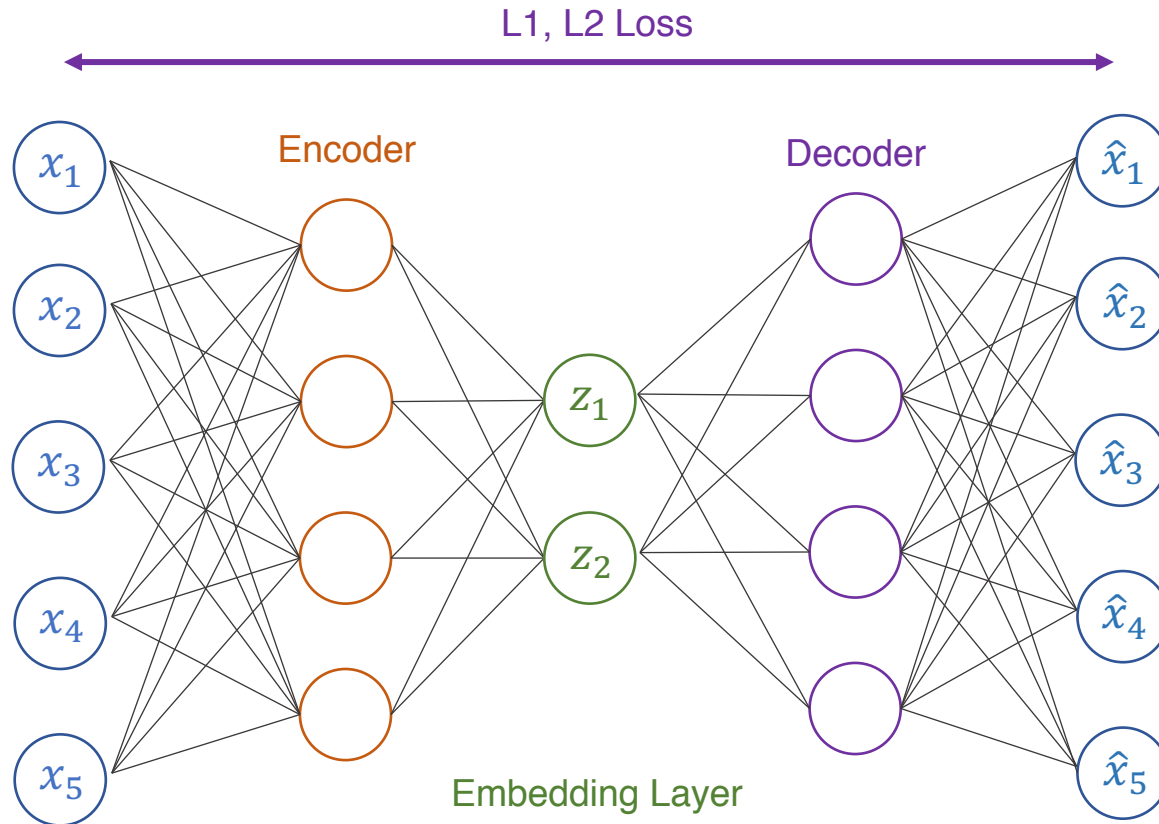
$$\frac{\partial L}{\partial z_1} = \frac{\partial L}{\partial a} \frac{\partial a}{\partial z_1} = \left(\underbrace{\frac{\partial L}{\partial z_2} \frac{\partial z_2}{\partial a}}_{w_3} + \underbrace{\frac{\partial L}{\partial z_3} \frac{\partial z_3}{\partial a}}_{w_4} \right) \frac{\partial a}{\partial z_1}$$

Because L is a function of z_2 and z_3 , and z_2 and z_3 are functions of a

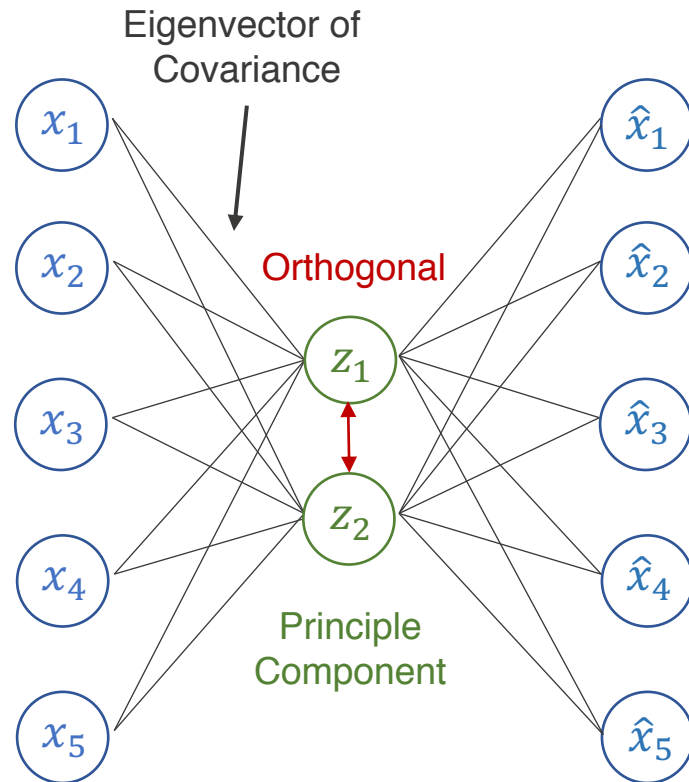
3. Now we can compute the derivative for L w.r.t w

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial w}$$

Autoencoder

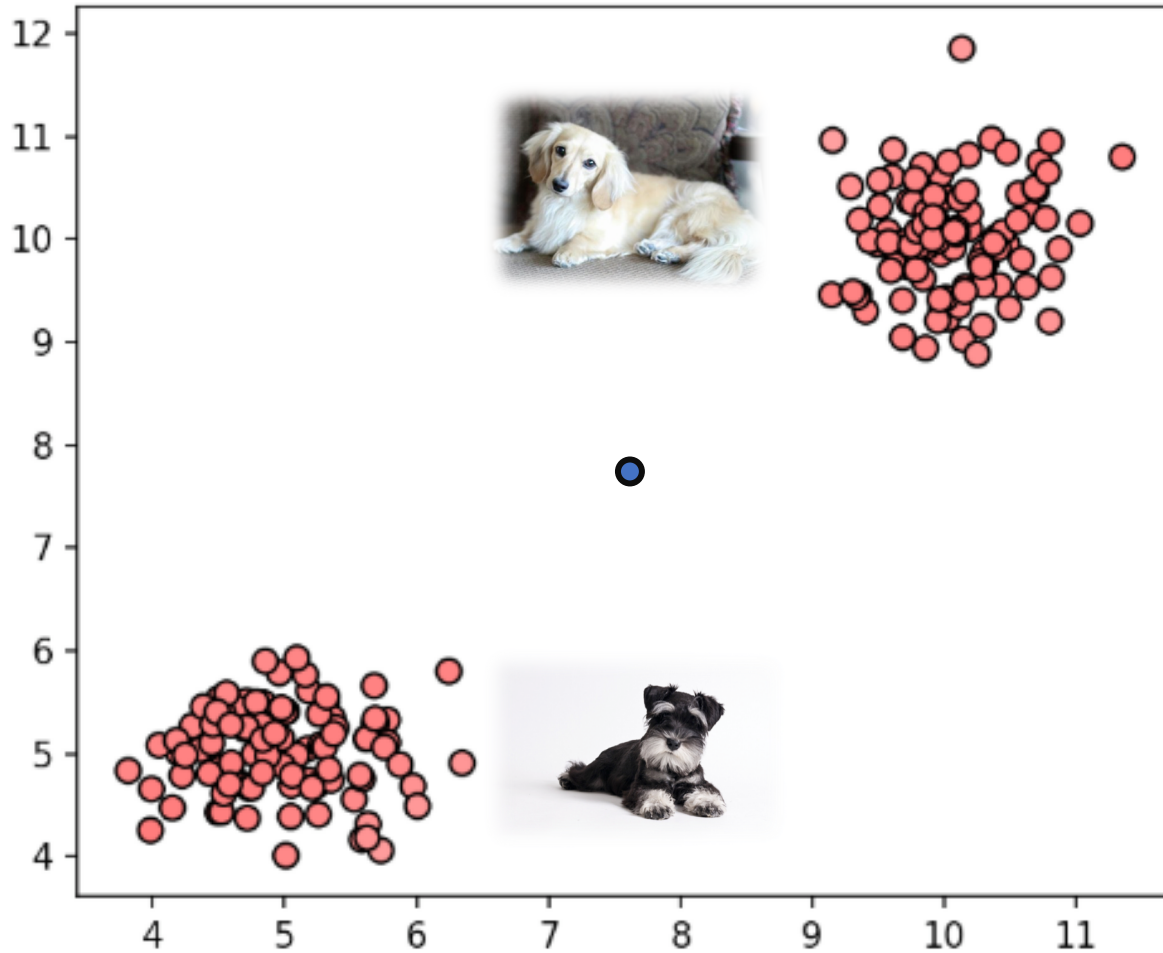


Principle Component Analysis

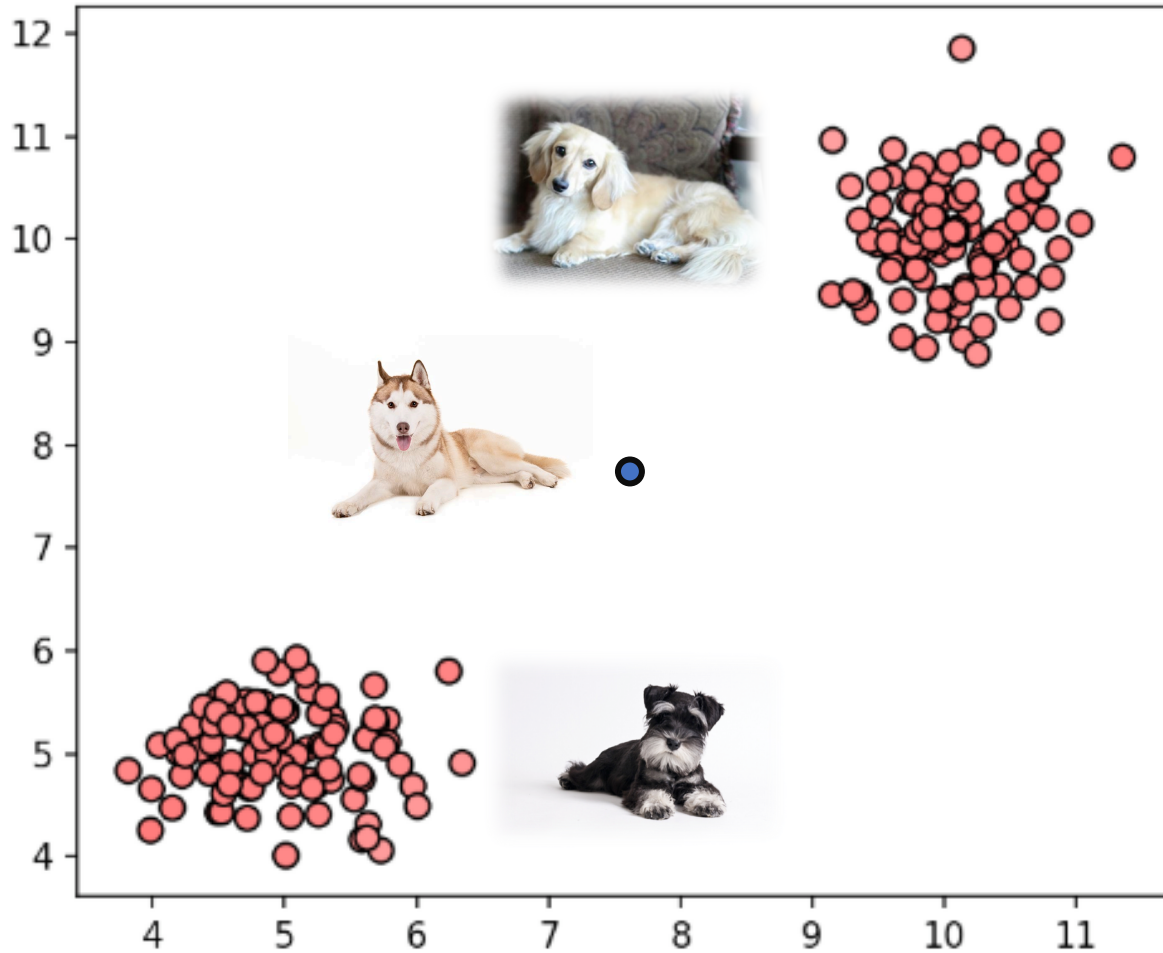


Maximize variance for the projection on PCs

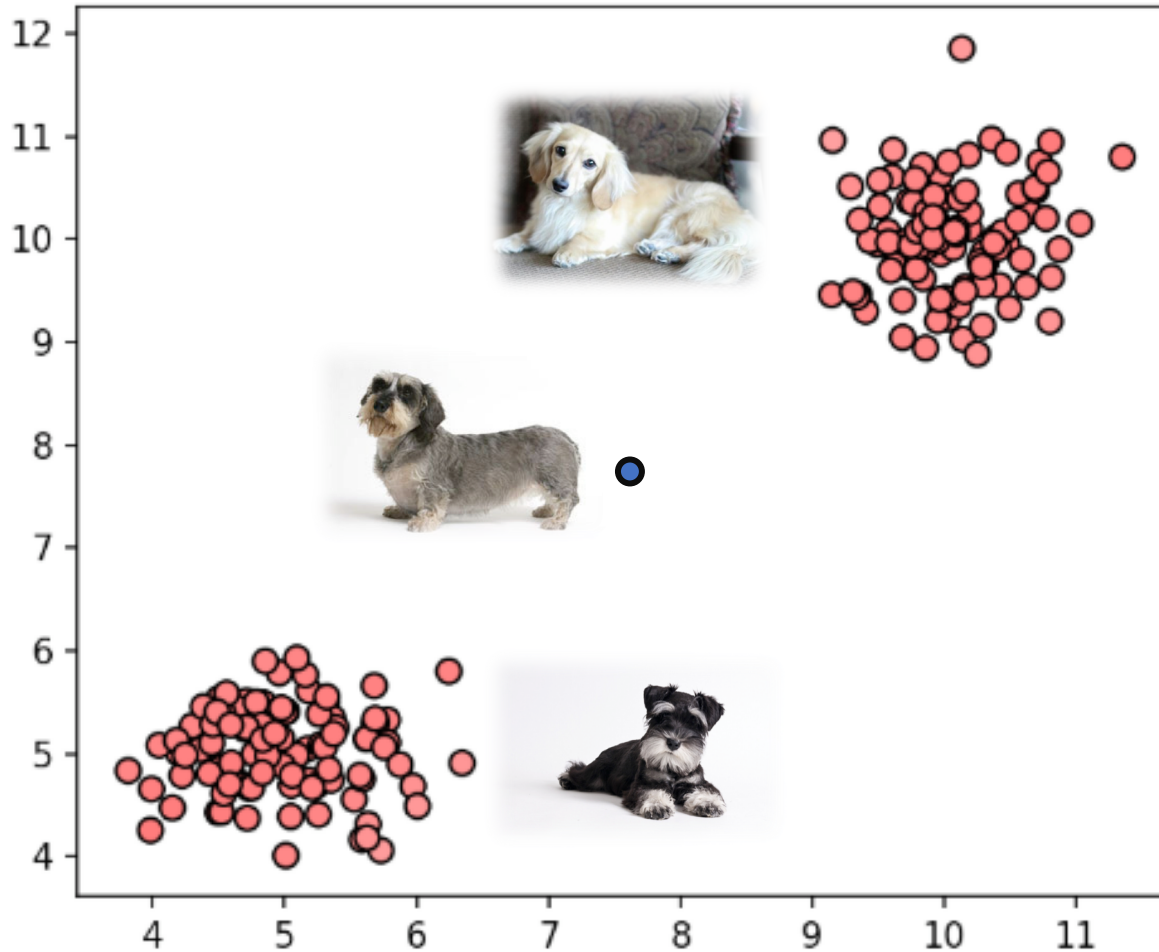
Limitation of Autoencoder



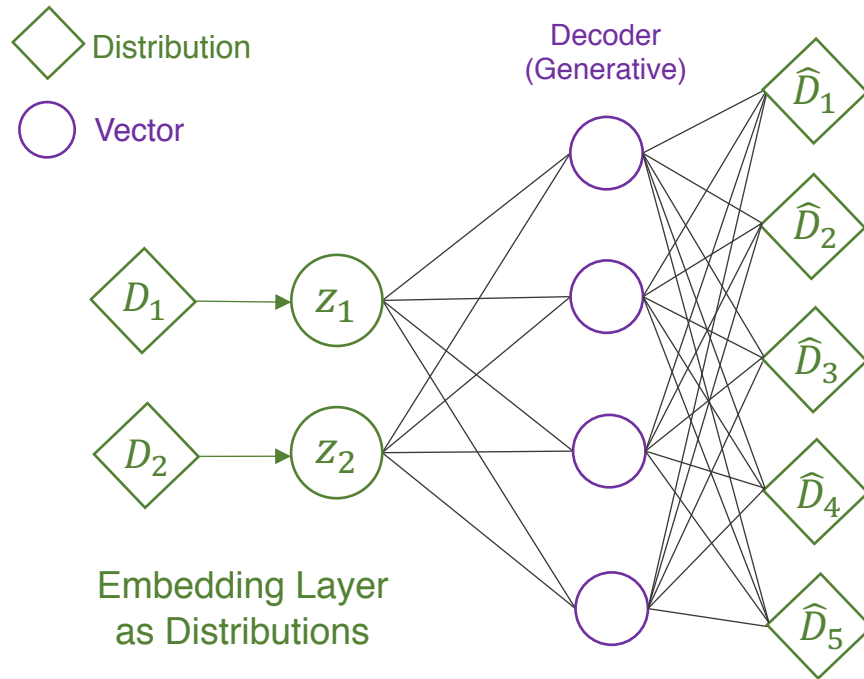
Limitation of Autoencoder



Limitation of Autoencoder



Reconstruct Data from Code



Auto-encoding Variational Bayes

Problem Formation

- Assumptions

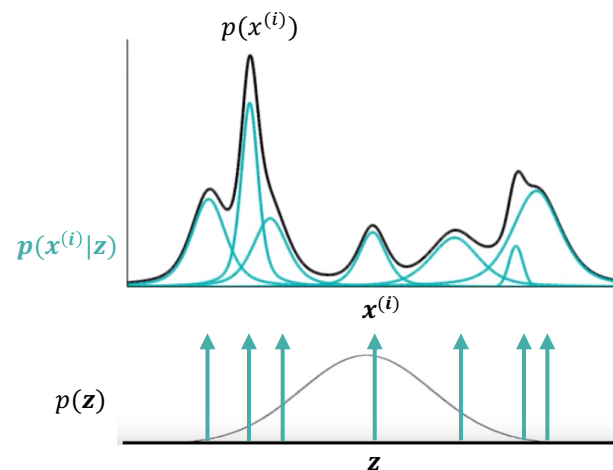
- Suppose $x^{(i)}$ is generated from a random process.
- This random process involve an unobserved continuous random variable z
- z is generated through a random (Gaussian) process $p_{\theta^*}(z)$
- $x^{(i)}$ is generated from $p_{\theta^*}(x|z)$ (The distribution of x , given the latent variable z)
- z and θ^* are unknown
- $p_{\theta}(x) = \int p_{\theta}(z) p_{\theta}(x|z) dz$ is **intractable** since we have infinite possibility of z
- Therefore the posterior $p_{\theta}(z|x) = \int p_{\theta}(x|z) p_{\theta}(z) / p_{\theta}(x) dz$ is also **intractable**

- Problems

- Estimate the true parameter of θ^* efficiently
- But we don't know how to optimize the problem

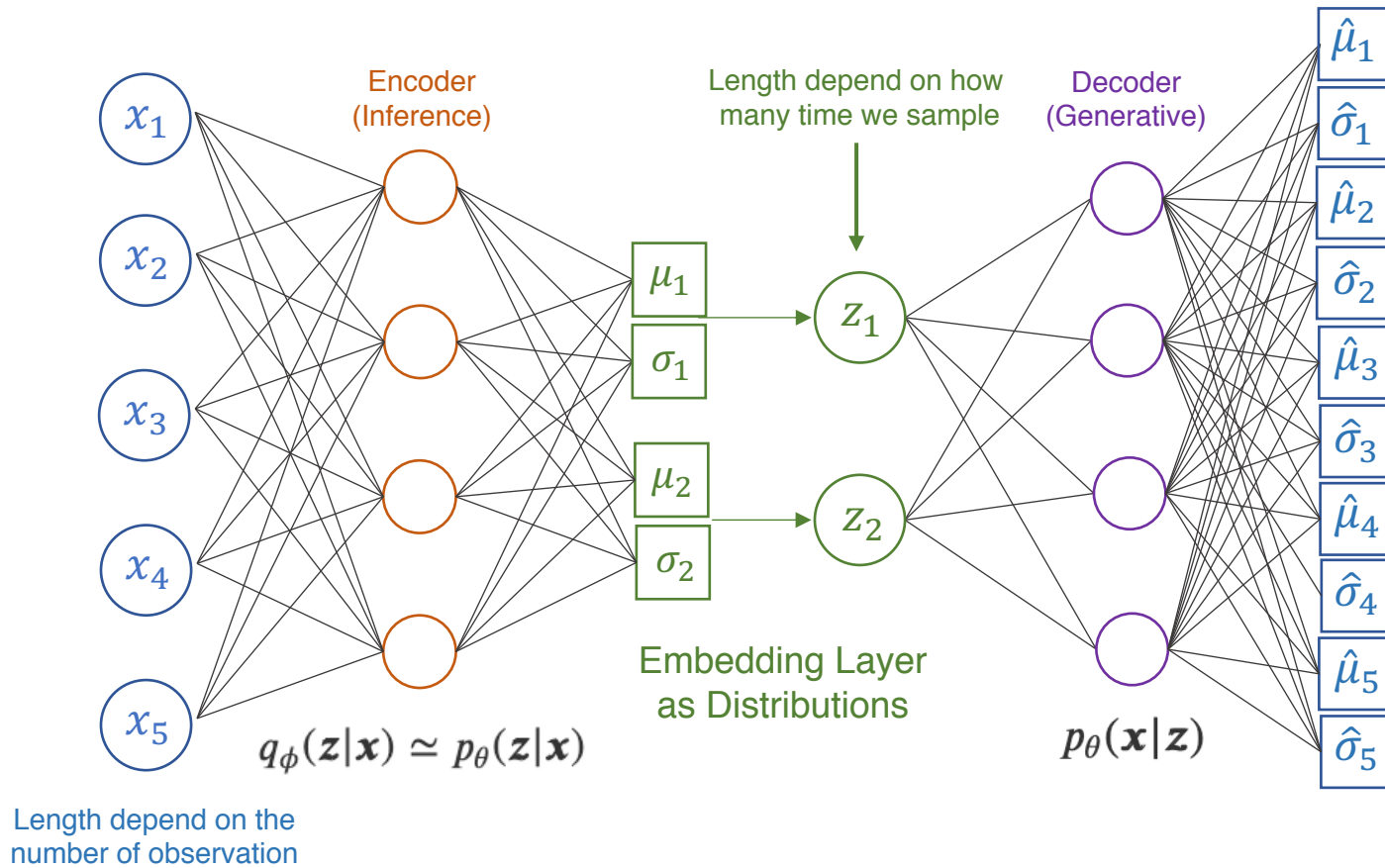
- Consider using autoencoder

- Let encoder function to be $q_{\phi} \cong p_{\theta}(z|x)$
- Let decoder function to be $p_{\theta}(x|z)$



from Hung-Yi Lee

Embedding Layer as Distributions



Auto-encoding Variational Bayes

Methods (1)

- For any PDF, we can maximize the likelihood of this function with available data

$$\begin{aligned}\log p(x) &= \log p(x) \int_z q(z|x) dz \quad (\text{Integral of any PDF equals to 1}) \\ &= \int_z q(z|x) \log p(x) dz \\ &= \int_z q(z|x) \log \left(\frac{p(z,x)}{p(z|x)} \right) dz \quad (\text{Bayes theorem}) \\ &= \int_z q(z|x) \log \left(\frac{p(z,x)}{q(z|x)} \frac{q(z|x)}{p(z|x)} \right) dz \\ &= \int_z q(z|x) \log \left(\frac{p(z,x)}{q(z|x)} \right) dz + \int_z q(z|x) \log \left(\frac{q(z|x)}{p(z|x)} \right) dz\end{aligned}$$

- Let $L_b = \int_z q(z|x) \log \left(\frac{p(z,x)}{q(z|x)} \right) dz$

$$D_{KL}(q(z|x) || p(z|x)) = \int_z q(z|x) \log \left(\frac{q(z|x)}{p(z|x)} \right) dz \quad \leftarrow \text{Intractable, but } D_{KL} \text{ always } \geq 0$$

- Then $\log p(x) = L_b + D_{KL}(q(z|x) || p(z|x)) \geq L_b$ Evidence lower bound (ELBO)

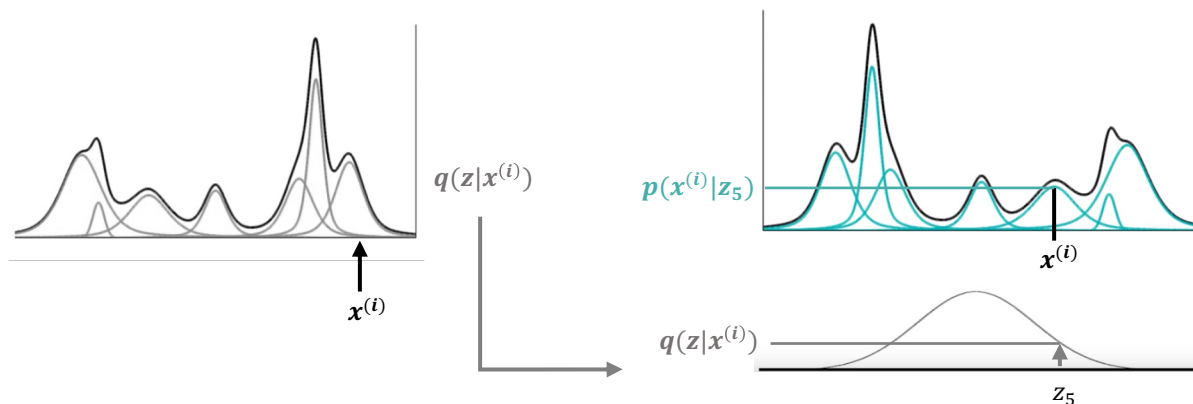
- Rewrite $L_b = \int_z q(z|x) \log \left(\frac{p(z,x)}{q(z|x)} \right) dz$
$$\begin{aligned}&= \int_z q(z|x) \log \left(\frac{p(x|z)p(z)}{q(z|x)} \right) dz \\ &= \int_z q(z|x) \log p(x|z) dz + \int_z q(z|x) \log \left(\frac{p(z)}{q(z|x)} \right) dz \\ &= \int_z q(z|x) \log p(x|z) dz - D_{KL}(q(z|x) || p(z))\end{aligned}$$

Auto-encoding Variational Bayes

Methods (2)

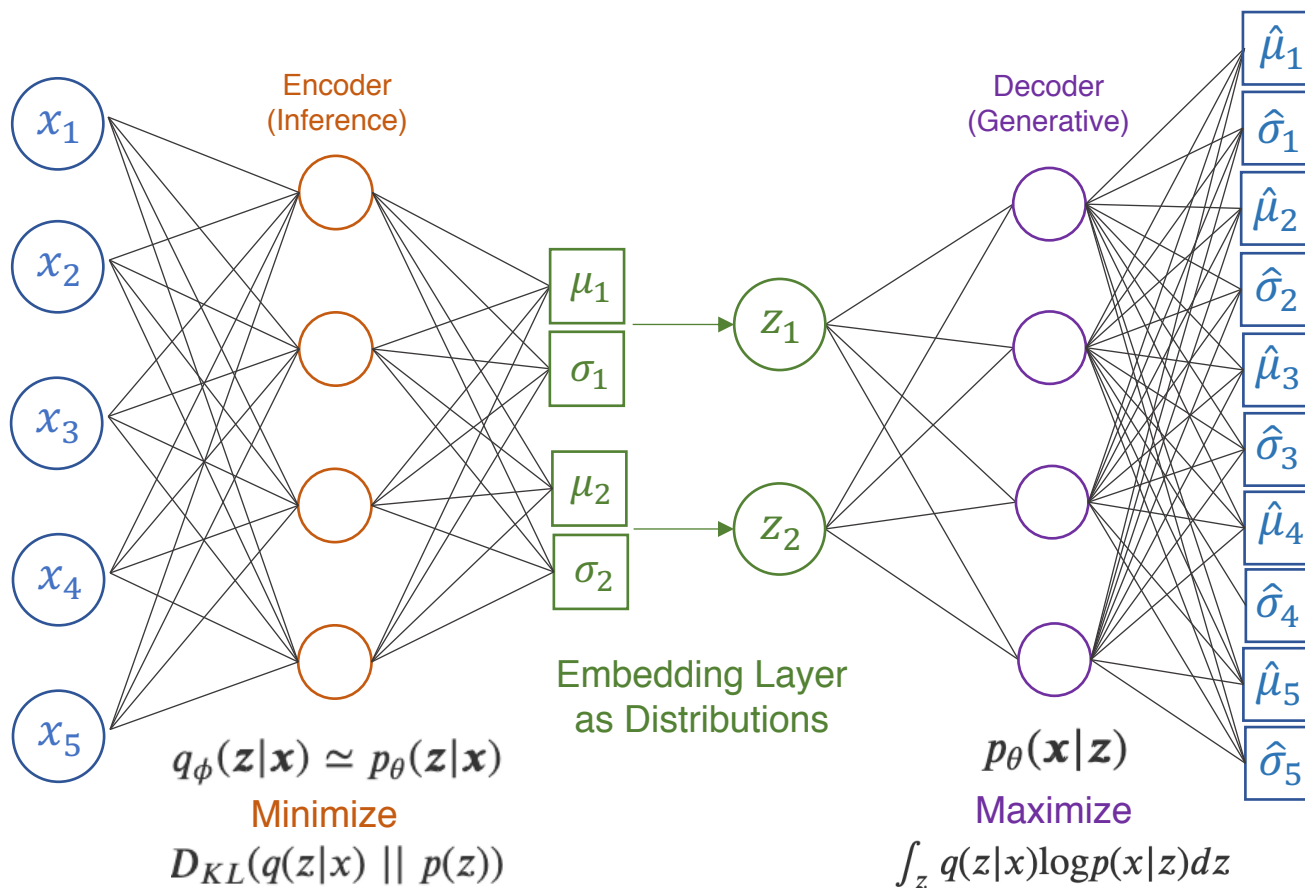
Evidence lower bound

- Rewrite $\log p(x) = \int_z q(z|x) \log p(x|z) dz - D_{KL}(q(z|x) || p(z)) + D_{KL}(q(z|x) || p(z|x))$
 - ↑ The contribution of this paper
 - ↑ Has solution (*)
 - ↑ Intractable
- We can maximize the evidence lower bound and “attempt to” maximize the likelihood of $p(x)$
- $-D_{KL}(q(z|x) || p(z)) = \sum_{j=1}^J (\sigma_j)^2 - (1 + \log(\sigma_j)) + (\mu_j)^2$ (Appendix B)
- Try Monte Carlo estimator to maximize $\int_z q(z|x) \log p(x|z) dz \rightarrow$ Exhibit high variance



Auto-encoding Variational Bayes

Methods (3)



Auto-encoding Variational Bayes

Methods (4)

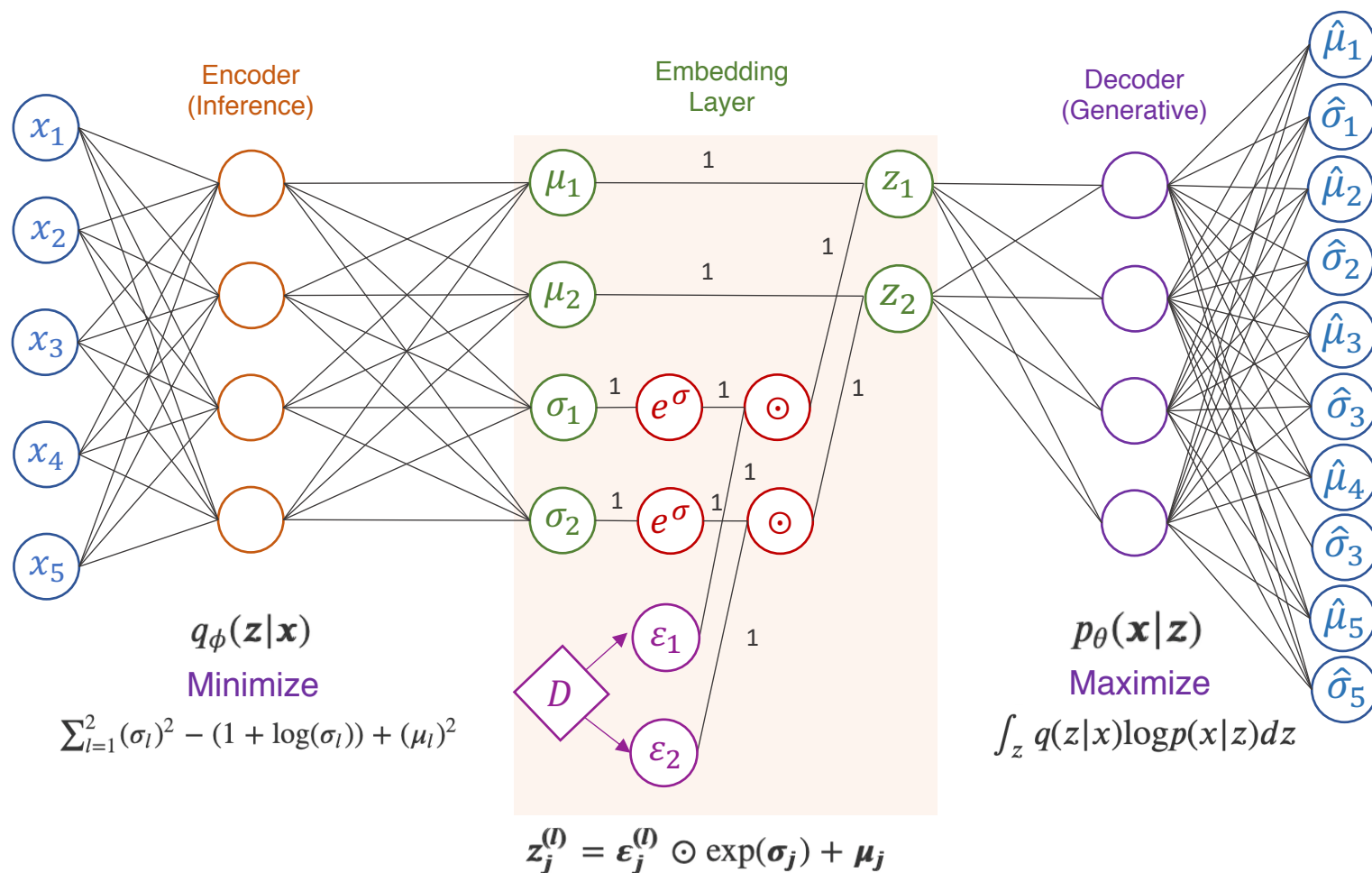
- Instead of output the conditional distribution from the encoder function, we can output a function for \mathbf{z} .

$$\mathbf{z}_j^{(l)} = \boldsymbol{\varepsilon}_j^{(l)} \odot \exp(\boldsymbol{\sigma}_j) + \boldsymbol{\mu}_j, \text{ where } \boldsymbol{\varepsilon} \sim N(0, I)$$

- The variables σ_j and μ_j are the output from the encoder function. In other word, they are a function of x
- ε is a random variable sampled from a standard distribution (now it is independent from the parameter φ)
- l is the index of the sampling process.
- The reparameterization alleviate the problem of high variance using Monte Carlo estimator.
- The authors did not mention why it help.

Auto-encoding Variational Bayes

Methods (5)



In practice we only sample ϵ once for each observation

Methods

Methods

$$\mathbf{X} = \begin{bmatrix} X_{1,1} & X_{1,2} & \cdots & X_{1,d} \\ X_{2,1} & X_{2,2} & \cdots & X_{2,d} \\ \vdots & \vdots & \ddots & \vdots \\ X_{n,1} & X_{n,2} & \cdots & X_{n,d} \end{bmatrix}$$

Genes

Cell types

- Model the observed expression with a conditional negative binomial distribution
- Recall that negative binomial distribution can be written as Gamma-Poisson (mixture) distribution:

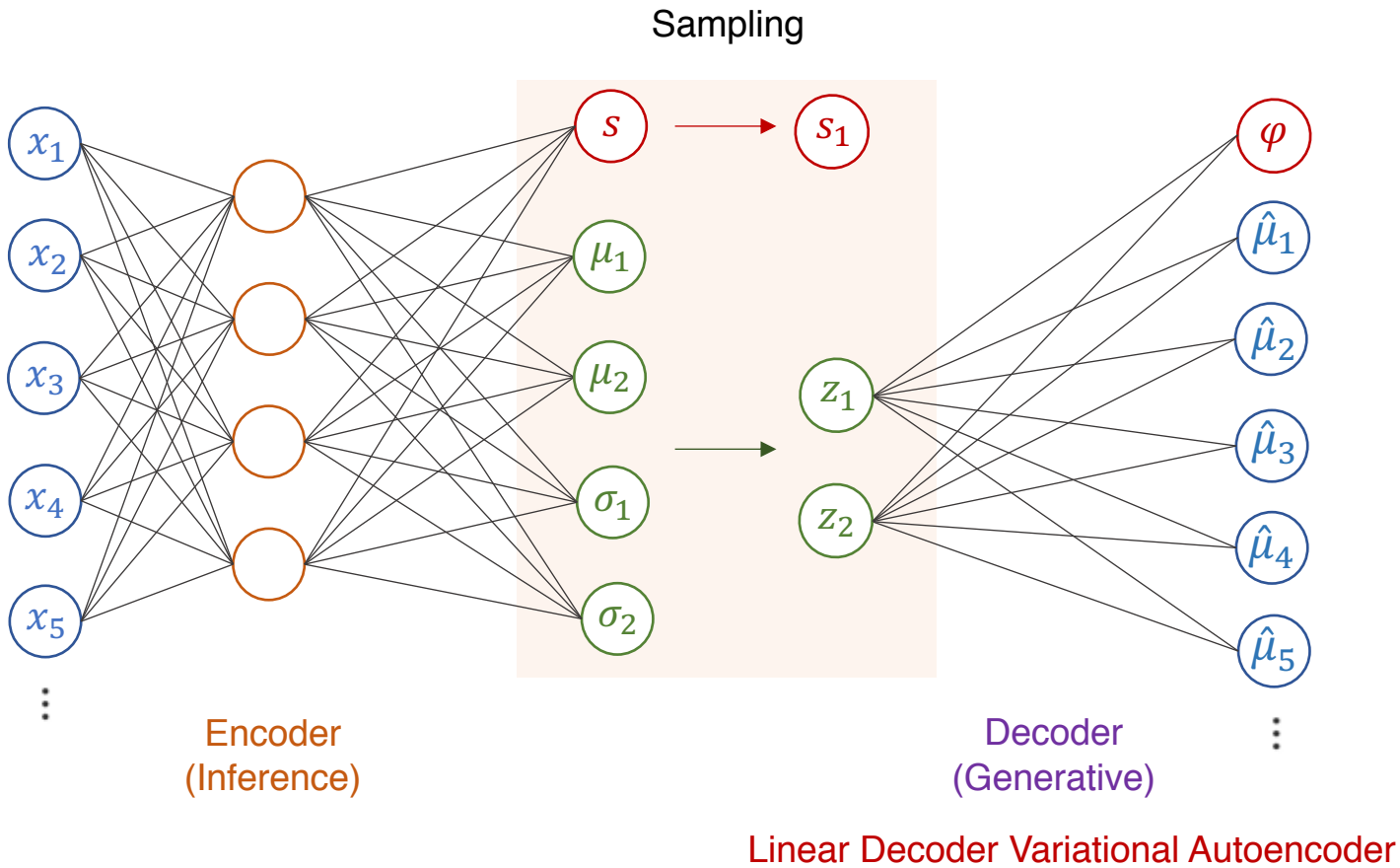
$$NB(k|r, p) = \int \underset{\text{sampling}}{f_{Poisson}(\lambda)(k)} \underset{\text{technical noise}}{f_{Gamma}(r, \frac{1-p}{p})(\lambda)} d\lambda$$

- Combine $y \sim Poisson(v, s)$ and $v \sim Gamma(\exp(\mu), \frac{1}{\varphi})$
 $y \sim NB(\exp(\mu), s, \frac{1}{\varphi})$

- Here, $r = \exp(\mu_n^g)$ and $p = \varphi$ are latent variables
- Therefore, the generative model condition on latent variable $\mu_n^g, \varphi^g, s_n^g$
- Also, the authors set the decoder to be a linear function $\mu_n^g = z_n W^T$

W is the weight in decoder function

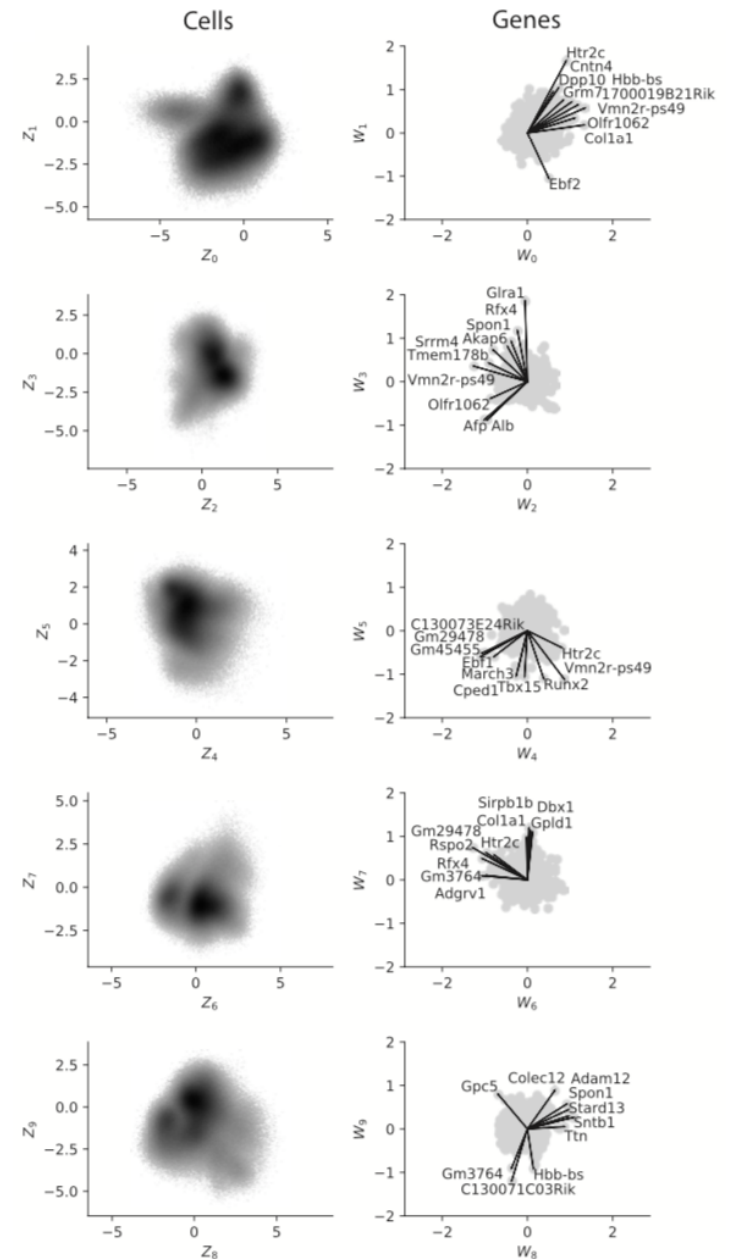
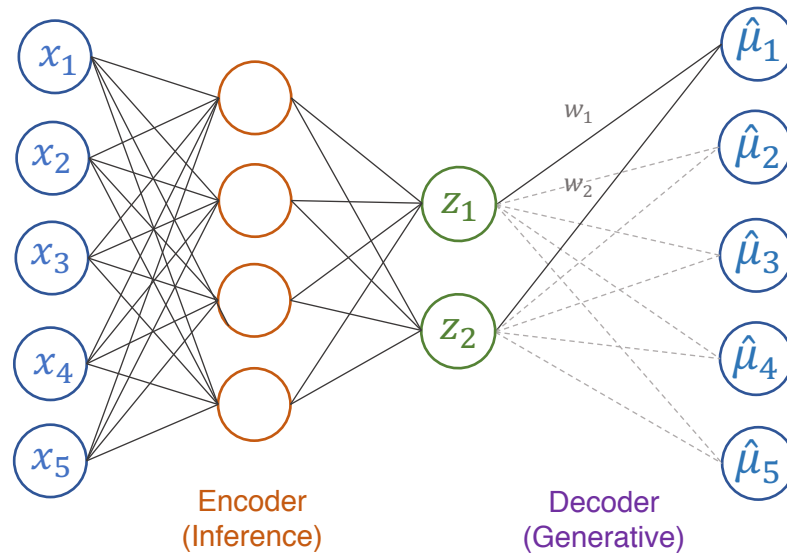
Methods



Results

Results

- Apply on developing mouse embryos in different stages of development (Cao *et al.*, 2019)
- The variation in latent variable z can directly related to variation in gene expression



Discussion

Discussion

- Manuscripts
 - This is an additional feature added to scVI package
 - Promote the concept of extend functions on existing frameworks in bioinformatics
 - Enable interpretable analysis of data at massive scale
- Something more
 - The effect of using linear function as decoder
 - Improve the generative model with InfoGAN?

References and Resources

- Neural Network
 - Numerical Optimization by Ben Frederickson
 - Back Propagation by 3Blue1Brown
- Variational Autoencoder
 - Auto-Encoding Variational Bayes, arXiv 2014.
 - Stanford CS231n: Convolutional Neural Networks for Visual Recognition, Generative Model
 - Stanford CS228: Probabilistic Graphical Models, Variational Autoencoder
 - The Reparameterization Tricks by JP Zhang
- Generative Model on scRNA
 - Deep generative modeling for single-cell transcriptomics, Nature Methods 2018.

Thanks