

Deep Learning and Applications

DSA 5204 • Lecture 9
Dr Low Yi Rui (Aaron)
Department of Mathematics



Test (25 Mar 2023, 10:10-10:50am)

- **Format:** 20 multiple choice/selection questions
- **Duration:** 40 minutes
- **Platform:** Canvas Quiz
- **Syllabus:** First 8 slides (everything so far except generative models)
- **Instructions**
 - *Lecture on 25 March will be **fully online***
 - *Log in to zoom session at **10am** to set up*
 - *Join the usual zoom session on your laptop/desktop. Have a second device (handphone) join the zoom session and show via its camera your hands, face and your laptop/desktop screen. **Join Zoom with exactly 2 devices!***
 - ***Attendance will be taken, you need to attend the zoom session for your score to count***

Last Time

Previously, we introduced some unsupervised and supervised learning algorithms

- Autoencoders
- Semi-supervised, multi-task and transfer learning

Things in common: still focused on overcoming missing labels

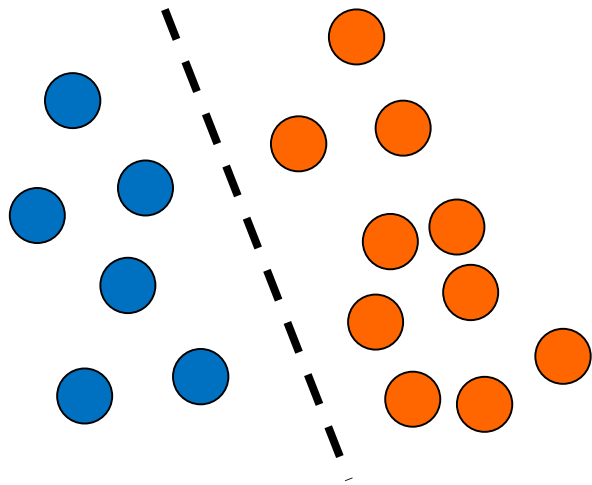
Today and in the next lectures, we will look at very different task – generating new samples by learning distributions.



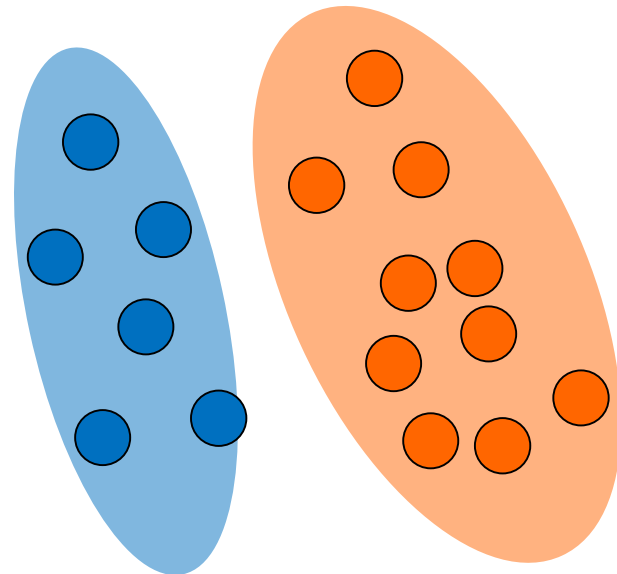
Generative Models Overview

Discriminative vs Generative Models

Discriminative
Modelling $f^*(x)$



Generative
Modelling $p^*(x)$



Modelling Distributions



There are two main ways to model distributions, with important differences.

Density Estimation

Given:

$$\mathcal{D} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}\}$$

$$\mathbf{x}^{(i)} \sim p^* \text{ i.i.d.}$$

Goal:

$$\text{find } \hat{p} \approx p^*$$

Generative Models

Given:

$$\mathcal{D} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}\}$$

$$\mathbf{x}^{(i)} \sim p^* \text{ i.i.d.}$$

Goal:

$$\text{sample new } \tilde{\mathbf{x}} \sim p^* (\text{approximately})$$

Differences

- With \hat{p} , we may not be able to sample from it easily
- Even if we can sample approximately from p^* , we may not have a good representation of it

Gaussian Mixture Models

For Gaussian mixture models (GMM), we consider a parametric family

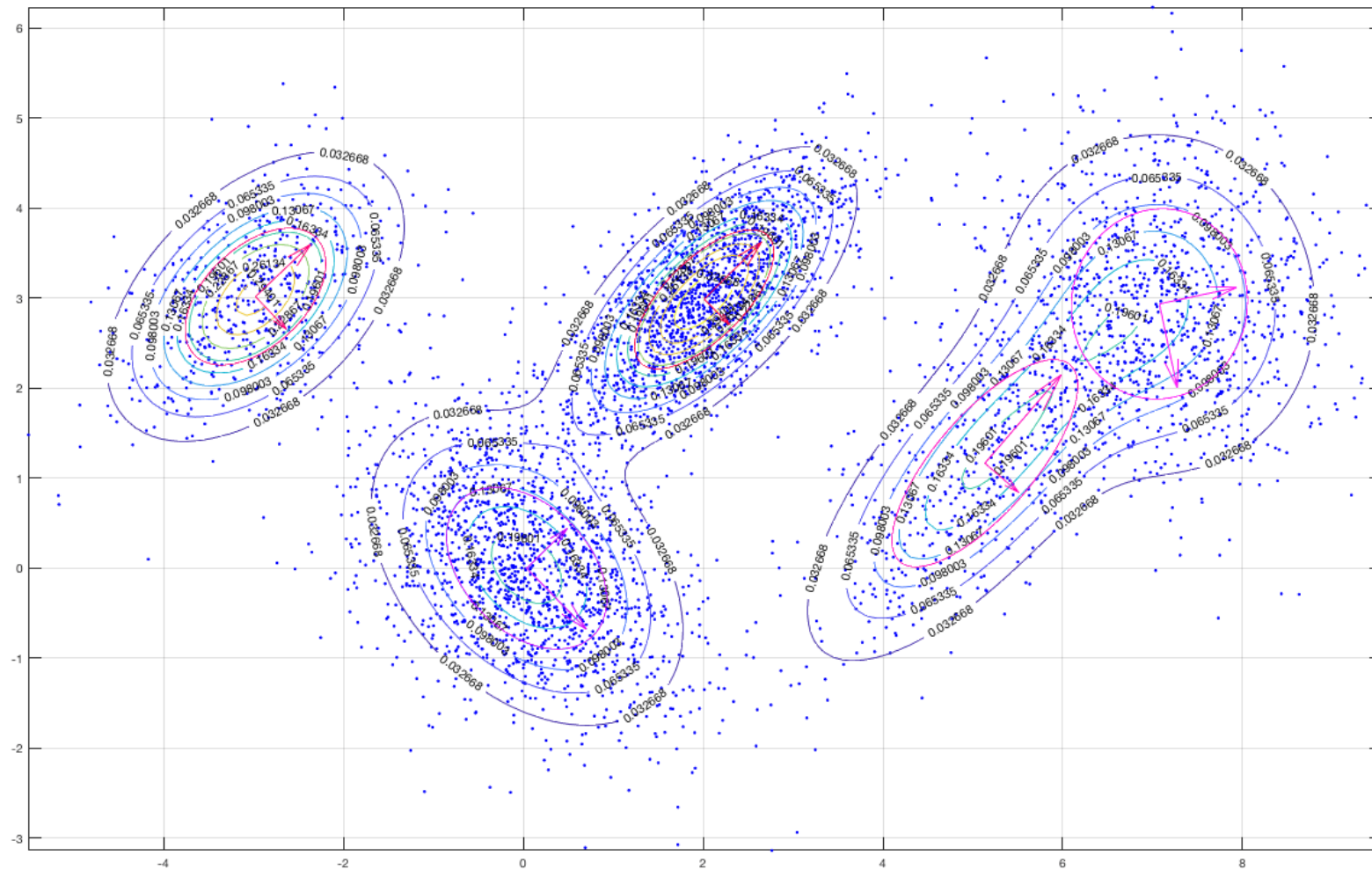
$$p_{\theta}(\mathbf{x}) = \sum_j \alpha_j p_{\mu_j, \Sigma_j}(\mathbf{x}) \quad \alpha_j \geq 0 \quad \sum_j \alpha_j = 1$$

where $p_{\mu, \Sigma}(\mathbf{x})$ is the PDF of Gaussian RVs with mean μ and covariance matrix Σ

Once fitted, we can generate samples from p_{θ} , by

- Generate index i from distribution $\{\alpha_j\}$
- Generate from p_{μ_i, Σ_i}

Gaussian Mixture Models



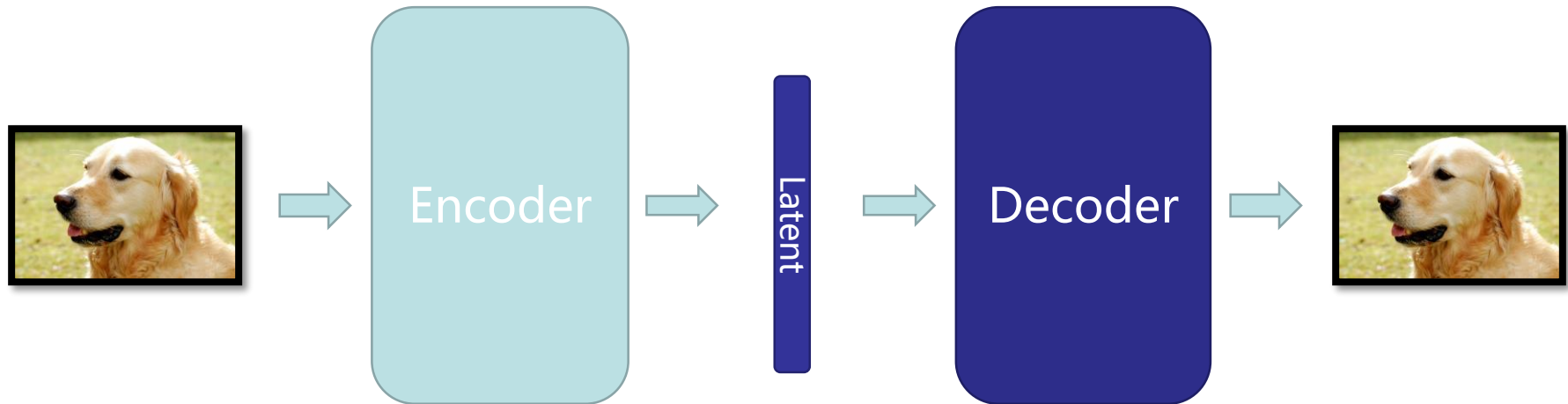
Limitations of GMM

- **A simple mixture of Gaussian is rarely enough to model high dimensional data**
- **Number of Gaussians is a hyper-parameter that is hard to tune in practice**
- **Cannot use domain knowledge (e.g. translation invariance, CNNs)**

Deep Learning Based Generative Model?



Recall that we introduced the autoencoder architecture



Idea: generating something in latent space, and then decode?

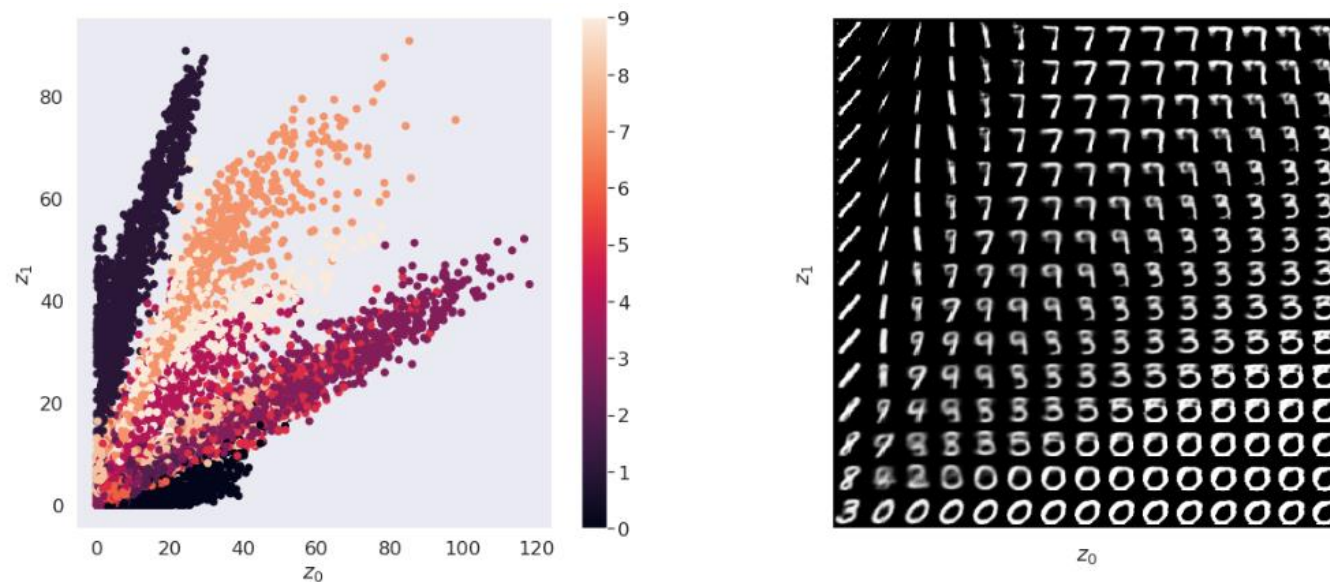


Demo: Generating Samples using Autoencoder

The issue with AEs



The latent space representation is not “continuous”



Therefore, to improve on this, we should try to make the latent space correspond smoothly with the output!



Variational Autoencoder

The Maximum Likelihood Approach

Given $p^*(x)$, consider

- A parametric model $p_{\theta}(x)$
- A dataset $\mathcal{D} = \{x^{(1)}, x^{(2)}, \dots, x^{(N)}\}$ with $x^{(i)} \sim p^*(x)$.

One way to estimate p^* using p_{θ} is the *maximum likelihood estimator (MLE)*:

$$\hat{\theta} = \operatorname{argmax}_{\theta} \prod_i p_{\theta}(x^{(i)})$$

This is equivalent to maximizing the log-likelihood

$$\hat{\theta} = \operatorname{argmax}_{\theta} \sum_i \log(p_{\theta}(x^{(i)}))$$

Example: MLE for Gaussian Family

<Lecture Notebook>

Consider $\theta = (\mu, \sigma)$ and the parametric Gaussian family in 1D

$$p_{\theta}(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} (x - \mu)^2\right)$$

Applying MLE to a dataset $\{x^{(1)}, x^{(2)}, \dots, x^{(N)}\}$

$$\hat{\theta} = (\hat{\mu}, \hat{\sigma}) = \arg \max_{\theta} \sum_{i=1}^N \log p_{\theta}(x^{(i)})$$

We find

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x^{(i)} \quad \hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N (x^{(i)} - \hat{\mu})^2$$

Sampling using MLEs

Once we found $\hat{\theta}$, we can generate samples
 $x' \sim p_{\hat{\theta}}(x)$

Limitations of this approach

- If we use a simple parametric form p_{θ} (e.g. Gaussian), then it may not be a good model for the data
- If we use a complex parametric form (e.g. NNs), then there is often no easy way to perform optimization and sampling...

This motivates another way to build generative models using *latent random variables*

Latent Generative Models

An alternative to direct modelling is to consider *latent* generative models.

Here, we write

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z})p_{\theta}(\mathbf{z})d\mathbf{z}$$

where

- \mathbf{z} : *latent* variable
- $p_{\theta}(\mathbf{z})$: latent variable *prior* distribution
- $p_{\theta}(\mathbf{x}|\mathbf{z})$: *generative/conditional* distribution

Example: Latent Generative Model

Let

- z be a uniform random variable on the interval $[a, b]$
- x be conditional random normal with mean μ and variance z^2

Then

$$p_{\theta}(z) = \frac{1}{b-a} \mathbb{I}_{a \leq z \leq b}$$
$$p_{\theta}(x|z) = \frac{1}{\sqrt{2\pi z^2}} e^{-\frac{(x-\mu)^2}{2z^2}}$$

with $\theta = (a, b, \mu)$

Direct vs Latent Approach



Given a data point x

Direct Approach:

- Learn $p_{\theta}(x)$
- Sample $x' \sim p_{\theta}(x)$

Latent Variable Approach

- Learn posterior latent distribution $p_{\theta}(z|x)$
- Learn generative distribution $p_{\theta}(x|z)$
- Sample from $z \sim p_{\theta}(z|x)$, and then $x' \sim p_{\theta}(x|z)$

Posterior Model for $p_{\theta}(z|x)$

Computationally
Intractable

By Bayes theorem, we have

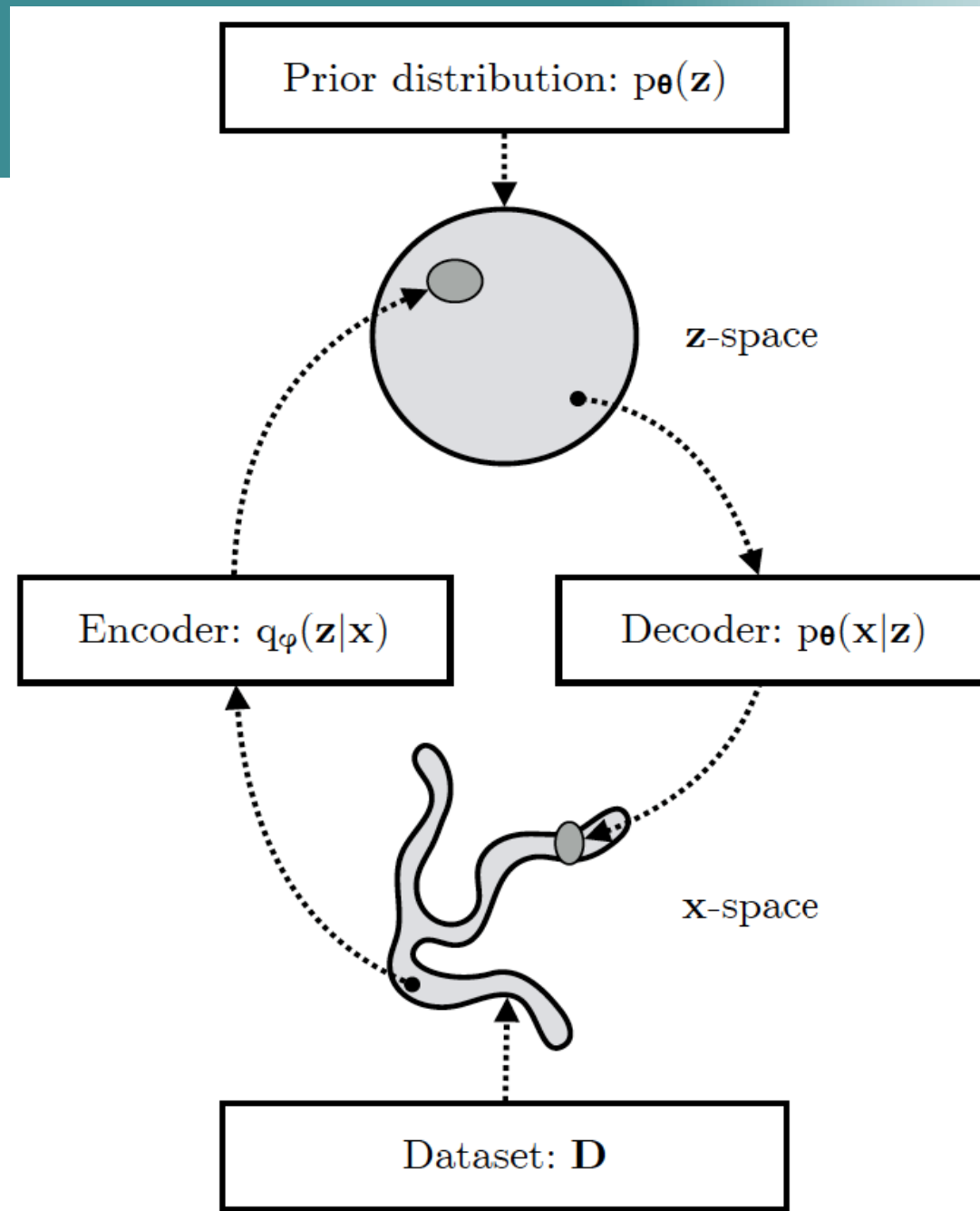
$$p_{\theta}(z|x) = \frac{p_{\theta}(x|z)p_{\theta}(z)}{\int p_{\theta}(x, z') dz'}$$

$\rightarrow p_{\theta}(x)$

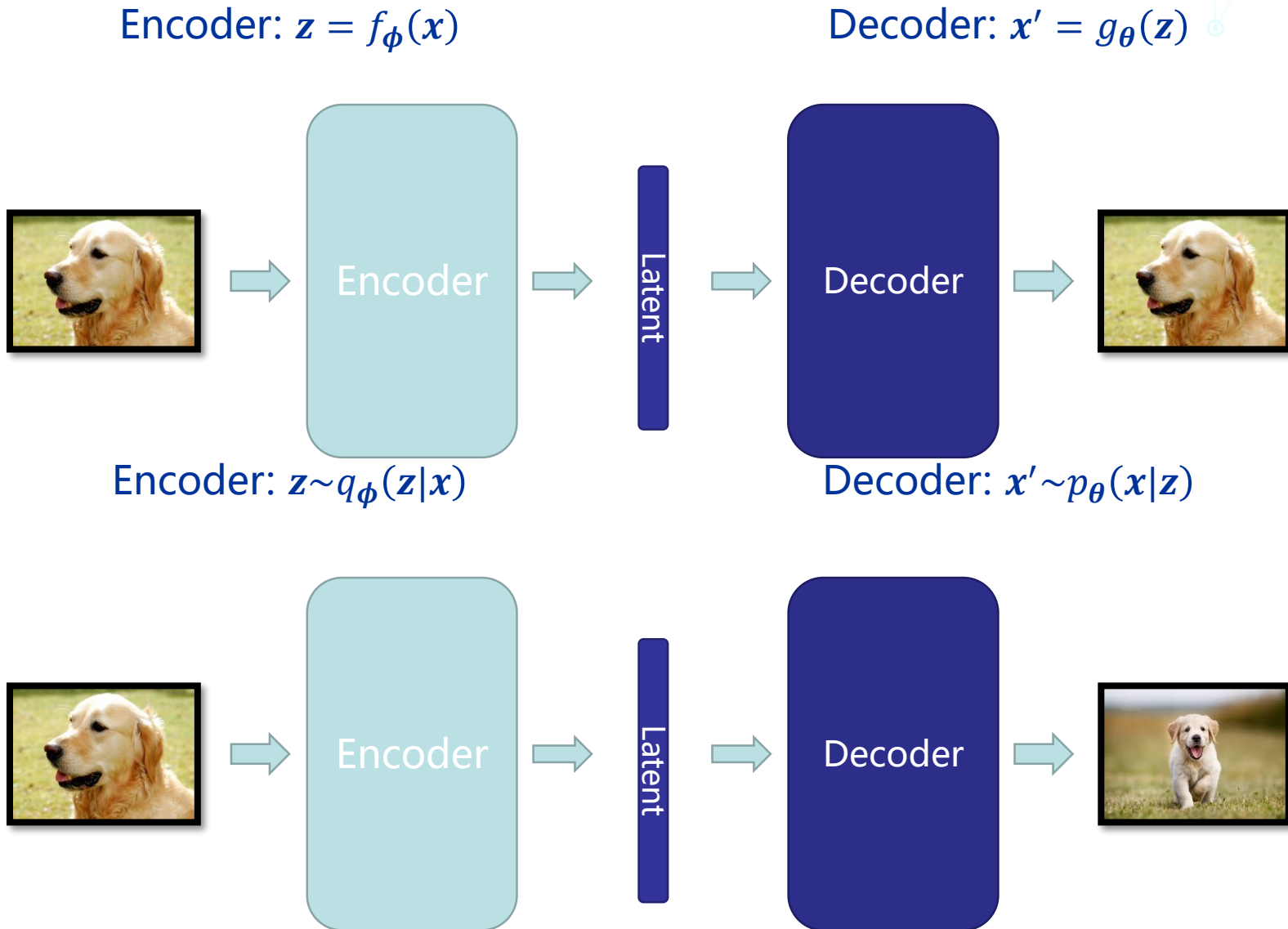
Instead of computing this, we will resort to a *posterior model*

$$q_{\phi}(z|x) \approx p_{\theta}(z|x)$$

whose parameters ϕ can be trained to make this approximation good.



Comparing with Autoencoders



Variational Autoencoders

In this sense, the latent variable generative process is very much like an autoencoder. Thus, we call it a variational autoencoder (VAE).

The term “variational” refers to our approximation of the posterior distribution $p_{\theta}(z|x)$ model by $q_{\phi}(z|x)$.

It remains to ask:

- How do we parameterize $p_{\theta}(x|z)$ (generative distribution) and $q_{\phi}(z|x)$ (posterior/variational distribution)?
- How do we train θ and ϕ ?



The Loss for VAEs

Training Autoencoders vs VAEs

Recall that in AEs, we can simply train the network using inputs as labels

$$\min_{\theta, \phi} L(x, x') = L\left(x, g_{\theta}\left(f_{\phi}(x)\right)\right)$$

In the probabilistic case of VAEs, x' and x are *not required* (rather, *required not*) to be the same!

Instead, we train their *distributions* q_{ϕ}, p_{θ} . So, how can we do training?

Preliminaries I: KL Divergence



Given two densities $p(x)$ and $q(x)$. The Kullback-Leibler (KL) Divergence is a way to measure their “distance” from one another.

For continuous RVs, we have

$$D_{KL}(p||q) = \int p(x) \log \frac{p(x)}{q(x)} dx = \mathbb{E}_{p(x)} \left[\log \left(\frac{p(x)}{q(x)} \right) \right]$$

For discrete RVs, replace integral by sum

Important Properties:

- $D_{KL}(p||q) \geq 0$ with equality if and only if $p = q$
- $D_{KL}(p||q)$ is convex in both p and q
- $D_{KL}(p||q)$ is not symmetric (in particular, not a metric)

Example: KL Divergence of Gaussians

<Lecture Notebook>

Consider two Gaussian PDFs in 1D

$$p_{\theta_1}(x) = \frac{1}{\sqrt{2\pi\sigma_1^2}} \exp\left(-\frac{1}{2\sigma_1^2}(x - \mu_1)^2\right)$$
$$p_{\theta_2}(x) = \frac{1}{\sqrt{2\pi\sigma_2^2}} \exp\left(-\frac{1}{2\sigma_2^2}(x - \mu_2)^2\right)$$

Then,

$$D_{KL}(p_{\theta_1} \| p_{\theta_2}) = \log \frac{\sigma_2}{\sigma_1} + \frac{\sigma_1^2}{2\sigma_2^2} - \frac{1}{2} + \frac{(\mu_1 - \mu_2)^2}{2\sigma_2^2}$$

Note the properties of D_{KL}

- $D_{KL}(p_{\theta_1} \| p_{\theta_2}) \geq 0$. Note $-\log u + \frac{1}{2}u^2 - \frac{1}{2} \geq 0$ for $u > 0$
- $D_{KL}(p_{\theta_1} \| p_{\theta_2}) = 0$ iff $\mu_1 = \mu_2$ and $\sigma_1 = \sigma_2$
- $D_{KL}(p_{\theta_1} \| p_{\theta_2}) \neq D_{KL}(p_{\theta_2} \| p_{\theta_1})$

Preliminaries II: Monte-Carlo Estimation of Expectations/Integrals



Given a random variable $x \in \mathbb{R}^d$ with density p and a function $f: \mathbb{R}^d \rightarrow \mathbb{R}$, we wish to compute the expectation

$$\mathbb{E}[f(x)] = \int f(x)p(x)dx$$

An efficient way (especially for large d) is *Monte-Carlo approximation*, where

$$\mathbb{E}[f(x)] \approx \frac{1}{N} \sum_{i=1}^N f(x^{(i)}) \quad x^{(i)} \sim p \text{ (i. i. d.)}$$

which becomes exact in the limit of $N \rightarrow \infty$.

Evidence Lower Bound (ELBO)

<Lecture Notebook>



Let us now apply MLE estimation, but now on the latent model

$$\log p_{\theta}(x) = \log \int p_{\theta}(x|z)p_{\theta}(z)dz$$

The problem is that we can't easily compute the integral, so we will use some tricks.

We can show that

$$\log p_{\theta}(x) = \underbrace{\mathbb{E}_{q_{\phi}(z|x)} \left[\log \left(\frac{p_{\theta}(x, z)}{q_{\phi}(z|x)} \right) \right]}_{L(x; \theta, \phi)} + \underbrace{D_{KL} \left(q_{\phi}(z|x) || p_{\theta}(z|x) \right)}_{\substack{K(x; \theta, \phi) \\ \text{KL Divergence}}}$$

Evidence Lower Bound (ELBO)

Optimizing the ELBO

Note that the KL divergence is non-negative, hence

$$\log p_{\theta}(x) = L(x; \theta, \phi) + K(x; \theta, \phi) \geq L(x; \theta, \phi)$$

In other words, L is a lower bound for the log-likelihood.

Therefore, we can replace

$$\max_{\theta} \log p_{\theta}(x) \rightarrow \max_{\theta, \phi} L(x; \theta, \phi)$$

Key point: the latter can be computed without computing the exact posterior $p_{\theta}(z|x)$!

Multiple Datapoints and GD

The case for multiple i.i.d. datapoints is completely analogous

$$\begin{aligned} L(\mathcal{D}; \theta, \phi) &= \frac{1}{N} \sum_{i=1}^N L(\mathbf{x}^{(i)}; \theta, \phi) \\ &= \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})} \left[\log \left(\frac{p_{\theta}(\mathbf{x}^{(i)}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})} \right) \right] \end{aligned}$$

As always, we optimize L via (stochastic) gradient descent on θ, ϕ , which requires the computation of the gradients

$$\nabla_{\theta} L \quad \text{and} \quad \nabla_{\phi} L$$

Gradient with respect to θ

The gradient with respect to θ is easy

$$\begin{aligned}\nabla_{\theta} L &= \nabla_{\theta} \mathbb{E}_{q_{\phi}(z|x)} \left[\log \left(\frac{p_{\theta}(x, z)}{q_{\phi}(z|x)} \right) \right] \\ &= \mathbb{E}_{q_{\phi}(z|x)} \left[\nabla_{\theta} \log \left(\frac{p_{\theta}(x, z)}{q_{\phi}(z|x)} \right) \right] \\ &= \mathbb{E}_{q_{\phi}(z|x)} [\nabla_{\theta} \log(p_{\theta}(x, z))]\end{aligned}$$

This can be computed by Monte-Carlo approximation provided we can sample from $q_{\phi}(\cdot | x)$.

Gradient with respect to ϕ

The gradient with respect to ϕ is more involved, since

$$\begin{aligned}\nabla_{\phi} L &= \nabla_{\phi} \mathbb{E}_{q_{\phi}(z|x)} \left[\log \left(\frac{p_{\theta}(x, z)}{q_{\phi}(z|x)} \right) \right] \\ &\neq \mathbb{E}_{q_{\phi}(z|x)} \left[\nabla_{\phi} \log \left(\frac{p_{\theta}(x, z)}{q_{\phi}(z|x)} \right) \right]\end{aligned}$$

Hence, we need to rewrite this in order to

- Apply Monte-Carlo approximation
- Apply back-propagation

The Reparameterization Trick

Suppose we have a random variable z with a parameterized conditional distribution

$$z \sim q_{\phi}(z|x)$$

We can reparametrize z as a *deterministic transformation of a fixed random variable*

$$z = g(u, \phi, x)$$

where

- g is a deterministic function depending on ϕ
- $u \sim p_0(u)$ is fixed and independent of ϕ

Example: Reparametrizing Gaussians

A d -dimensional Gaussian random variable with mean $\mu + x$ and covariance matrix C can be parameterized by

$$q_{\phi}(z|x) = (2\pi)^{-\frac{d}{2}} \det(C)^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(z - \mu - x)^T C^{-1}(z - \mu - x)\right)$$

where $\phi = (\mu, C)$

Using the reparameterization trick, we can write

$$z = g(u, \phi, x)$$

where $u \sim \mathcal{N}(0, I)$ and

$$g(u, \phi, x) = \mu + x + C^{\frac{1}{2}}u \sim \mathcal{N}(\mu + x, C)$$

Why Reparametrize?

Consider any function f . We have

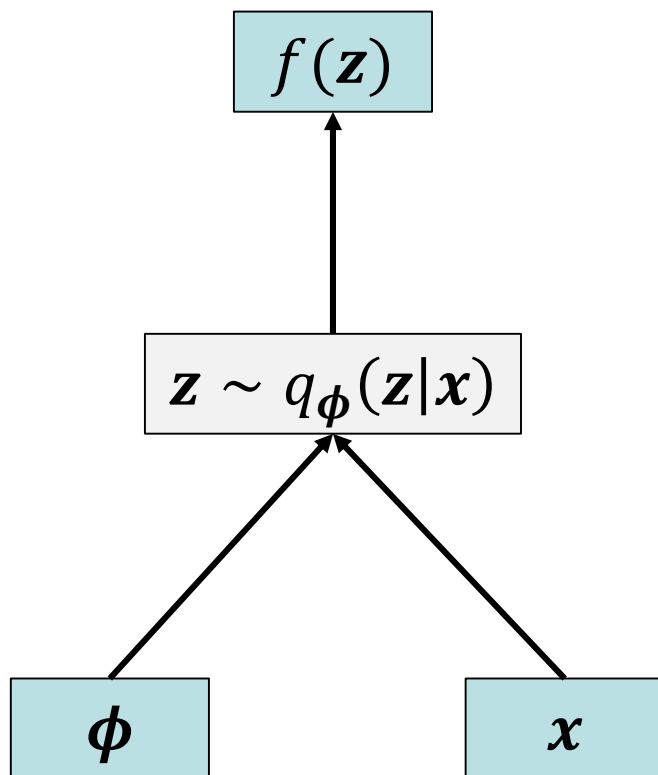
$$\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|x)}[f(\mathbf{z})] = \mathbb{E}_{\mathbf{u} \sim p_0(\mathbf{u})}[f(g(\mathbf{u}, \phi, x))]$$

Then, we have

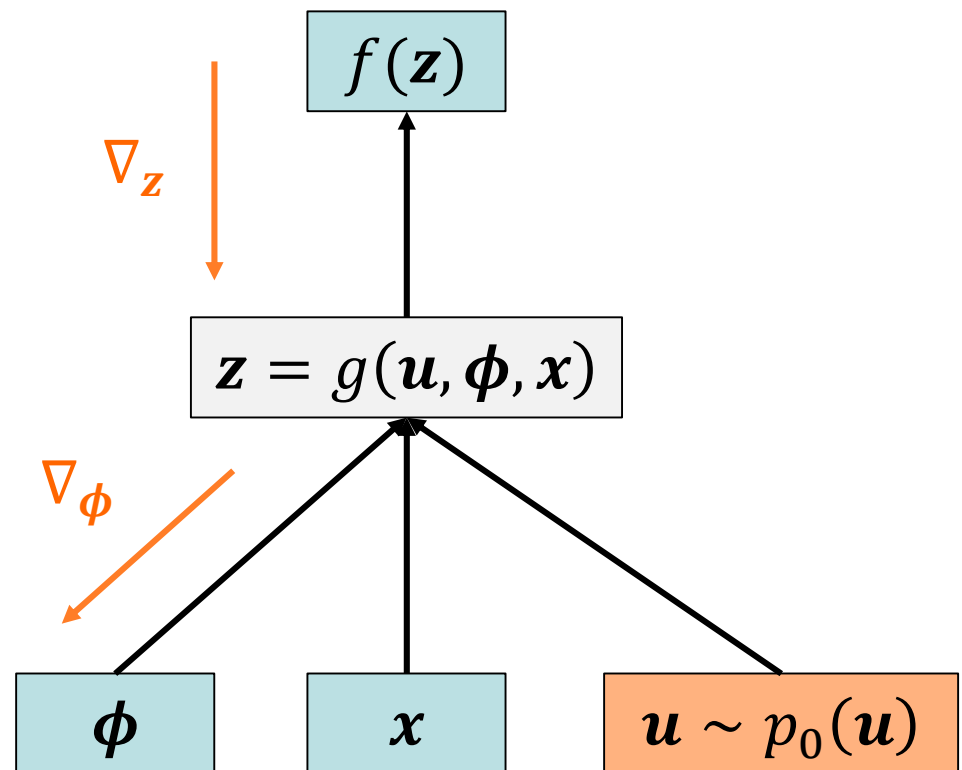
$$\nabla_{\phi} \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|x)}[f(\mathbf{z})] = \mathbb{E}_{\mathbf{u} \sim p_0(\mathbf{u})}[\nabla_{\phi} f(g(\mathbf{u}, \phi, x))]$$

The latter can be approximated by Monte-Carlo averages!

Original Form



Reparametrized Form



Gradient with respect to ϕ revisited

Under reparameterization, we have

$$\begin{aligned} L(x; \theta, \phi) &= \mathbb{E}_{q_{\phi}(z|x)} \left[\log p_{\theta}(x, z) - \log \left(q_{\phi}(z|x) \right) \right] \\ &= \mathbb{E}_{p_0(u)} \left[\log p_{\theta}(x, g(u, \phi, x)) - \log \left(q_{\phi}(g(u, \phi, x)|x) \right) \right] \end{aligned}$$

And hence

$$\nabla_{\phi} L(x; \theta, \phi) = \mathbb{E}_{p_0(u)} \left[\nabla_{\phi} \log p_{\theta}(x, g(u, \phi, x)) - \nabla_{\phi} \log \left(q_{\phi}(g(u, \phi, x)|x) \right) \right]$$

This can be computed by Monte-Carlo approximation.

Neural Network Representations

A decorative network diagram in the top right corner, consisting of a series of interconnected nodes (small circles) and edges (thin lines), forming a complex web-like structure.

Since the prior $p_{\theta}(z)$ is not conditioned on any data, we may pick it to be some simple form, e.g. factorized/iid Gaussians

Next, we need to discuss how to represent the densities
 $p_{\theta}(x|z)$ and $q_{\phi}(z|x)$

In practice, we should do it in such a way that

- They can be easily evaluated
- They can be easily differentiated (e.g. by back-prop)
- We can easily sample from them

Representing $p_{\theta}(x|z)$

Suppose $x \in \{0,1\}^d$ (binary data). Then a common choice is the *factorized Bernoulli model*.

Recall: Let x be a scalar Bernoulli random variable, then

$$x = \begin{cases} 1 & \text{with prob } s \\ 0 & \text{with prob } 1 - s \end{cases}$$

Here $s \in [0,1]$ is the probability of successful outcome

The distribution function then depends on s , and is given by

$$p_s(x) = s^x (1 - s)^{1-x}$$

Representing $p_{\theta}(x|z)$

In general, we may represent $p_{\theta}(x|z)$ as *factorized Bernoulli model* with the success probability vector s being a deterministic function of the latent variable z

$$s = \text{DecodingNN}(z; \theta) \quad s \in [0,1]^d$$

$$p_{\theta}(x|z) = \prod_j s_j^{x_j} (1 - s_j)^{1-x_j}$$

$$\log p_{\theta}(x|z) = \sum_j x_j \log s_j + (1 - x_j) \log(1 - s_j)$$

The last term is just the negative of the binary cross-entropy loss!

Representing $q_{\phi}(\mathbf{z}|\mathbf{x})$

Similarly, we can pick the simplest case of *factorized Gaussian model* for the encoding distribution

$$(\boldsymbol{\mu}, \log \boldsymbol{\sigma}) = \text{EncodingNN}(\mathbf{x}; \boldsymbol{\phi})$$

$$q_{\phi}(\mathbf{z}|\mathbf{x}) = \prod_j (2\pi\sigma_j^2)^{-\frac{1}{2}} \exp\left(-\frac{1}{2\sigma_j^2} (z_j - \mu_j)^2\right)$$

Under the reparameterization trick, we have

$$\mathbf{z} = g(\mathbf{u}, \boldsymbol{\phi}, \mathbf{x}) = \boldsymbol{\mu} + \boldsymbol{\sigma} \circ \mathbf{u} \quad \mathbf{u} \sim p_0 = \mathcal{N}(\mathbf{0}, I)$$

Giving

$$\log q_{\phi}(\mathbf{z}|\mathbf{x}) = \sum_j -\frac{1}{2} \log(2\pi\sigma_j^2) - \frac{1}{2} u_j^2 \rightarrow u_j = \frac{z_j - \mu_j}{\sigma_j}$$

Simplifying the ELBO Loss

<Lecture Notebook>

Under the current choices, we can dramatically simplify the ELBO loss

$$\begin{aligned}\mathbb{E}_{q_{\phi}(z|x)} \left[\log \left(\frac{p_{\theta}(x, z)}{q_{\phi}(z|x)} \right) \right] &= \mathbb{E}_{q_{\phi}(z|x)} \left[\log \left(\frac{p_{\theta}(x|z)p_{\theta}(z)}{q_{\phi}(z|x)} \right) \right] \\ &= \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z) + \log p_{\theta}(z) - \log q_{\phi}(z|x)] \\ &= \mathbb{E}_{p_0(u)} [\log p_{\theta}(x|g(u, \phi, x)) + \log p_{\theta}(g(u, \phi, x)) - \log q_{\phi}(g(u, \phi, x)|x)]\end{aligned}$$

We have

$$\begin{aligned}\log p_{\theta}(x|z) &= \sum_j x_j \log s_j + (1 - x_j) \log(1 - s_j) \\ \log p_{\theta}(z) &= -\frac{1}{2} \sum_j [z_j^2 + \log(2\pi)] = -\frac{1}{2} \|u\|^2 - \frac{1}{2} \|\sigma\|^2 \\ \log q_{\phi}(z|x) &= -\frac{1}{2} \sum_j [u_j^2 + \log(2\pi) + 2 \log(\sigma_j)]\end{aligned}$$



The only term that requires further attention is

$$\log p_{\theta}(\mathbf{z}) = -\frac{1}{2} \sum_j [z_j^2 + \log(2\pi)] \quad z_j = \mu_j + u_j \sigma_j$$

Since we have to get rid of \mathbf{z} in order to let the gradient flow through the reparameterization trick.

We can show that (dropping constants)

$$\begin{aligned} \mathbb{E}_{p_0(u)}[\log p_{\theta}(\mathbf{z})] &= \mathbb{E}_{\mathbf{u} \sim \mathcal{N}(\mathbf{0}, I)} \left[-\frac{1}{2} \sum_j (\mu_j + \sigma_j u_j)^2 \right] \\ &= -\frac{1}{2} \sum_j (\mu_j^2 + \sigma_j^2) \end{aligned}$$

Summary of VAE (Basic Form)



Input dimension = d . **Latent dimension** = m .

Encoder:

$$\begin{aligned} (\underbrace{\mathbf{y}_1}_{\boldsymbol{\mu}}, \underbrace{\mathbf{y}_2}_{\log \sigma}) &= \text{EncodingNN}(\mathbf{x}; \boldsymbol{\phi}) \\ \mathbf{z} &= \mathbf{y}_1 + e^{\mathbf{y}_2} \circ \mathbf{u} \quad \mathbf{u} \sim \mathcal{N}(\mathbf{0}, I) \quad (\mathbf{u} \in \mathbb{R}^m) \end{aligned}$$

Decoder:

$$\begin{aligned} \mathbf{s} &= \text{DecodingNN}(\mathbf{z}; \boldsymbol{\theta}) \quad (\mathbf{s} \in \mathbb{R}^d) \\ \text{Inference: } \mathbf{x}' &\sim \text{Bernoulli}(\mathbf{s}) \end{aligned}$$

Loss:

$$-L(\mathbf{x}; \boldsymbol{\theta}, \boldsymbol{\phi}) = \text{BCE}(\mathbf{x}, \mathbf{s}) + \frac{1}{2} \|\mathbf{y}_1\|^2 + \frac{1}{2} \|e^{\mathbf{y}_2}\|^2 - \sum_j y_{2,j}$$

Some Remarks

- The loss is usually split into two parts

$$-L(x; \theta, \phi) = \underbrace{\text{BCE}(x, s)}_{\text{Reconstruction Loss}} + \underbrace{\left[\frac{1}{2} \|y_1\|^2 + \frac{1}{2} \|e^{y_2}\|^2 - \sum_j y_{2,j} \right]}_{\text{KL-divergence loss}}$$

Why this name? Observe

$$\begin{aligned} -L &= -\mathbb{E}_{q_\phi(z|x)} \log p_\theta(x|z) + \mathbb{E}_{q_\phi(z|x)} \log \left(\frac{q_\phi(z|x)}{p_\theta(z)} \right) \\ &= \text{BCE}(x, s) + D_{KL} \left(q_\phi(z|x) \| p_\theta(z) \right) \end{aligned}$$

- The multiple-sample loss sums over all single sample losses
- Further reading on VAEs
(<https://arxiv.org/abs/1312.6114>)



Demo: Variational Autoencoder

Summary



In this lecture, we introduced a useful class of generative models known as variational autoencoders

Some key take-aways

- Unlike AEs, VAEs can learn “continuous” representations in latent space, giving rise to useful generative models
- VAEs also belongs to a powerful class of generative models making use of *latent variables*