# M2177.003100
# Deep Learning

## [14: Variational Autoencoder (VAE)]

Electrical and Computer Engineering
Seoul National University

(last compiled at 09:23:00 on 2020/11/21)

# Outline

Autoencoders

Variational Autoencoder
  Architecture
  Training
  Remarks

Summary

# References

- *Deep Learning* by Goodfellow, Bengio and Courville ▸ Link
  - ▶ Chapter 14: Autoencoders
  - ▶ Chapter 20: Deep Generative Models

- *Pattern Recognition and Machine Learning* by Bishop
  - ▶ Chapter 10: Approximate Inference

- online resources:
  - ▶ *Stanford CS231n: CNN for Visual Recognition* ▸ Link
  - ▶ *CVPR 2018 GAN Tutorial* ▸ Link
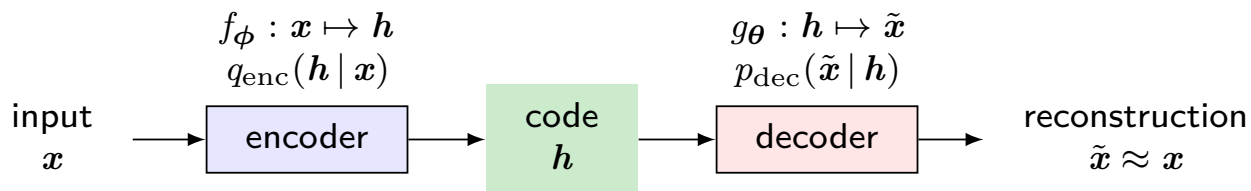  - ▶ *NIPS 2016 Variational Inference Tutorial* ▸ Link

# Outline

**Autoencoders**

Variational Autoencoder

Summary

# Autoencoder (AE)

- neural net trained to reconstruct input $x$ as output $\tilde{x}$ (*i.e.* $x \approx \tilde{x}$)
  - learns <u>code</u> (= efficient representation of input data)
  - unsupervised (no label)

- architecture

| | deterministic | stochastic | parameterized by |
|---|---|---|---|
| encoder | $f_\phi(x) = h$ | $q_{\text{enc}}(h \mid x)$ | $\phi$ |
| decoder | $g_\theta(h) = \tilde{x}$ | $p_{\text{dec}}(\tilde{x} \mid h)$ | $\theta$ |

$$f_\phi : x \mapsto h \qquad\qquad g_\theta : h \mapsto \tilde{x}$$
$$q_{\text{enc}}(h \mid x) \qquad\qquad p_{\text{dec}}(\tilde{x} \mid h)$$

input $x$ $\longrightarrow$ [ encoder ] $\longrightarrow$ [ code $h$ ] $\longrightarrow$ [ decoder ] $\longrightarrow$ reconstruction $\tilde{x} \approx x$

- $\dim(x) > \dim(h)$: undercomplete $\Rightarrow$ capture important features of input
- $\dim(x) < \dim(h)$: overcomplete

- training

$$\phi^*, \theta^* = \underset{\phi, \theta}{\arg\min} \, L(\boldsymbol{x}, \tilde{\boldsymbol{x}})$$

  - $L$: loss (*e.g.* mean squared error)

- applications
  - dimensionality reduction
  - feature learning
  - forefront of *generative* modeling

- linear $f$, $g$ and MSE loss
  - $\boldsymbol{h}$ spans the same subspace as principal component analysis (PCA)
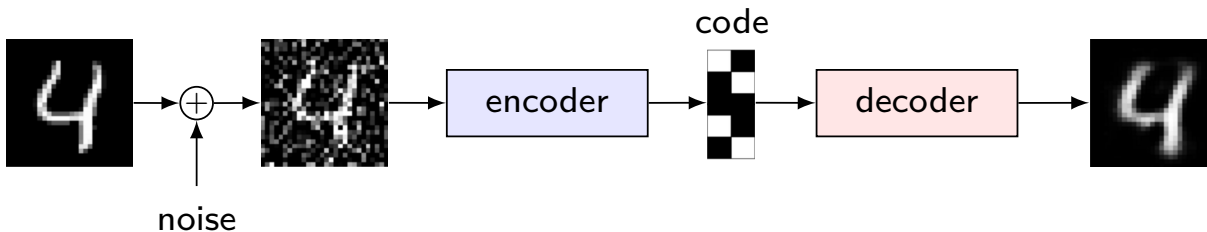
- stacked AE ($=$ deep AE)
  - AE with multiple hidden layers

# Need for regularization

- encoder/decoder with too much <u>**capacity**</u>
    - ▶ simply learn to copy input to output
    - ▶ learn nothing useful about data distribution

- regularizing AE
    - ▶ limit model capacity
    - ▶ use a loss encouraging desirable properties

- these properties
    - ▶ robustness to noise/missing inputs
    - ▶ sparse representation
    - ▶ small derivative of representation

# Denoising autoencoder (DAE)

- one way to force AE to learn useful features
    - ▶ add noise to <u>inputs</u>
    - ▶ then train AE to recover original (noise-free) input



- types of noise
    - ▶ Gaussian noise (shown above)
    - ▶ dropout: randomly switched off inputs

# Sparse autoencoder

- training criterion involves
    1. reconstruction error
    2. sparsity penalty $\Omega(\boldsymbol{h})$ on code $\boldsymbol{h}$

$$J(\boldsymbol{\phi}, \boldsymbol{\theta}) = \underbrace{L(\boldsymbol{x}, \tilde{\boldsymbol{x}})}_{\substack{\text{reconstruction} \\ \text{loss}}} + \underbrace{\Omega(\boldsymbol{h})}_{\substack{\text{sparsity} \\ \text{loss}}}$$

    - ▶ sparsity loss example ▸ Link
        - ▷ $\Omega(\boldsymbol{h}) = \lambda D_{\mathrm{KL}}[p(\boldsymbol{h}; \boldsymbol{\rho}) \,||\, p(\boldsymbol{h}; \hat{\boldsymbol{\rho}})]$    $\leftarrow \lambda$: regularization coefficient
        - ▷ $p(\boldsymbol{h}; \boldsymbol{\rho})$: desired (sparse) distribution of $\boldsymbol{h}$    $\leftarrow$ parameterized by $\boldsymbol{\rho}$
        - ▷ $p(\boldsymbol{h}; \hat{\boldsymbol{\rho}})$: observed distribution of $\boldsymbol{h}$    $\leftarrow$ parameterized by $\hat{\boldsymbol{\rho}}$

- typically used to learn <u>features</u> for another task

  *e.g.* classification

# Contractive autoencoder (CAE)

- introduces explicit regularizer on code $h = f(x)$
  - encourage derivative of $f$ to be as small as possible

$$\Omega(h) = \lambda \left\| \frac{\partial f(x)}{\partial x} \right\|_F^2$$

  - $|| \cdot ||_F^2$ : squared Frobenius norm (= sum of squared elements)

- effect[1]
  - better captures $\underbrace{\text{local directions of variation}}$ dictated by data
    $$\uparrow$$
    = a lower-dimensional non-linear manifold
  - $f$ "contracts" input neighborhood to a smaller output neighborhood
  - *i.e.* similar inputs $\mapsto$ similar codings

---

[1]Rifai, S., Vincent, P., Muller, X., Glorot, X., and Bengio, Y. (2011). Contractive auto-encoders: Explicit invariance during feature extraction. In *Proceedings of the 28th International Conference on International Conference on Machine Learning*, pages 833–840. Omnipress

# More autoencoder examples

- stacked convolutional autoencoders[2]
  - ▶ learn to extract visual features through conv layers

- generative stochastic network (GSN)[3]
  - ▶ generalization of DAE (added capability to generate data)

- winner-take-all (WTA) autoencoder[4]
  - ▶ only top $k\%$ activation for each neuron are preserved (the rest: set to 0)
  - ⇒ leads to sparse coding

- adversarial autoencoders[5]                     ⇒ AE2 pushes AE1 to learn robust code
  - ▶ AE1: trained to reproduce its input
  - ▶ AE2: trained to find inputs AE1 cannot properly reconstruct

[2] Masci, J., Meier, U., Cireşan, D., and Schmidhuber, J. (2011). Stacked convolutional auto-encoders for hierarchical feature extraction. In *International Conference on Artificial Neural Networks*, pages 52–59. Springer

[3] Bengio, Y., Laufer, E., Alain, G., and Yosinski, J. (2014). Deep generative stochastic networks trainable by backprop. In *International Conference on Machine Learning*, pages 226–234

[4] Makhzani, A. and Frey, B. J. (2015). Winner-take-all autoencoders. In *Advances in Neural Information Processing Systems*, pages 2791–2799

[5] Makhzani, A., Shlens, J., Jaitly, N., Goodfellow, I., and Frey, B. (2015). Adversarial autoencoders. *arXiv preprint arXiv:1511.05644*

# Outline

# Variational autoencoder (VAE)

- a popular approach to unsupervised learning of complex distributions
  - ▶ uses learned approximate inference
  - ▶ can be trained purely with gradient-based methods[6,7]

- appealing because
  - ▶ built on top of standard function approximators (neural nets)
  - ▶ can be trained with *SGD*

---

[6] Kingma, D. P. and Welling, M. (2013). Auto-encoding variational Bayes. *arXiv preprint arXiv:1312.6114*

[7] Rezende, D. J., Mohamed, S., and Wierstra, D. (2014). Stochastic backpropagation and approximate inference in deep generative models. *arXiv preprint arXiv:1401.4082*

# Comparison

- augoregressive models
  - ▶ define a tractable density function

$$p_{\boldsymbol{\theta}}(\boldsymbol{x}) = \prod_{i=1}^{n} p_{\boldsymbol{\theta}}(x_i \mid x_1, \ldots, x_{i-1})$$

- VAE
  - ▶ define an intractable density function with $\underline{latent}$ variable $\mathbf{z}$:

$$p_{\boldsymbol{\theta}}(\boldsymbol{x}) = \int p_{\boldsymbol{\theta}}(\boldsymbol{z}) p_{\boldsymbol{\theta}}(\boldsymbol{x} \mid \boldsymbol{z}) \, d\boldsymbol{z}$$

- GAN
  - ▶ give up explicit $p(\boldsymbol{x})$ modeling
  - ▶ just want to ability to sample from it

# Related models

- autoencoders: reconstruct given input $x$

$$x \xrightarrow[\text{encoder}]{p(z \mid x)} z \xrightarrow[\text{decoder}]{p(x \mid z)} x \qquad \text{(training)}$$

$$x \xrightarrow[\text{encoder}]{p(z \mid x)} z \qquad \text{(main use: feature learner for downstream learning)}$$

- generative models: generate new sample $x$

$$\xrightarrow{p(x)} x$$

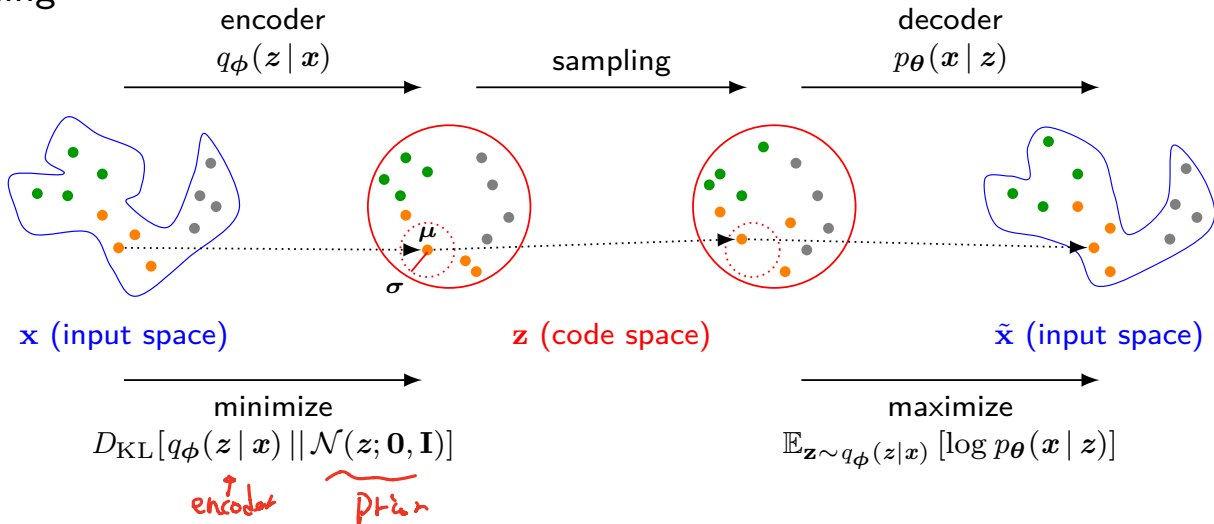  - ▶ one challenge: modeling <u>dependencies</u> between dimensions

- latent generative models: generate new sample $x$ from latent $z$

$$\xrightarrow[\text{prior}]{p(z)} z \xrightarrow{p(x \mid z)} x$$

  - ▶ introducing $z$ improves dependency modeling
  - ▶ how to train? maximize data likelihood $p(x) = \int p(z)p(x \mid z)\,dz$
  - ▶ challenge: integration $\int dz$ is often intractable

# Idea behind VAE

- training



encoder
$q_{\phi}(z \,|\, x)$

sampling

decoder
$p_{\theta}(x \,|\, z)$

x (input space)

z (code space)

$\tilde{x}$ (input space)

minimize
$D_{\mathrm{KL}}[q_{\phi}(z \,|\, x) \,||\, \mathcal{N}(z; \mathbf{0}, \mathbf{I})]$

encoder    prior

maximize
$\mathbb{E}_{z \sim q_{\phi}(z|x)}[\log p_{\theta}(x \,|\, z)]$

- sample generation

  ▶ $z \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$: easy to <u>sample</u> from

decoder
$p_{\theta}(x \,|\, z)$

z (code space)

$\tilde{x}$ (input space)

# Quick look

- goal: find $\boldsymbol{\theta}$ that maximizes likelihood $p_{\boldsymbol{\theta}}(\boldsymbol{x})$ for training data $\boldsymbol{x}$

$$p_{\boldsymbol{\theta}}(\boldsymbol{x}) = \int \underbrace{p_{\boldsymbol{\theta}}(\boldsymbol{z})}_{\substack{\text{simple} \\ \text{Gaussian} \\ \text{prior}}} \underbrace{p_{\boldsymbol{\theta}}(\boldsymbol{x} \mid \boldsymbol{z})}_{\substack{\text{learn with} \\ \text{neural net} \\ \text{``\textbf{decoder}''}}} d\boldsymbol{z} \tag{1}$$

- challenge: cannot optimize (1) directly
  - ▶ $p_{\boldsymbol{\theta}}(\boldsymbol{z})$ and $p_{\boldsymbol{\theta}}(\boldsymbol{x} \mid \boldsymbol{z})$ are okay but the integration is intractable

- solution: derive and optimize a tractable <u>lower bound</u> on the likelihood

$$\log p_{\boldsymbol{\theta}}(\boldsymbol{x}) \geq \underbrace{\mathcal{L}(\boldsymbol{x}, \boldsymbol{\theta}, \boldsymbol{\phi})}_{\text{ELBO}}$$

  - ▶ ELBO: depends on $\underbrace{q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x})}_{\substack{\text{learn with} \\ \text{neural net} \\ \text{``\textbf{encoder}''}}} \underbrace{\approx}_{\substack{\text{approximate} \\ \text{using} \\ \text{variational inference}}} \underbrace{p_{\boldsymbol{\theta}}(\boldsymbol{z}|\boldsymbol{x})}_{\substack{\frac{p_{\boldsymbol{\theta}}(\boldsymbol{x}|\boldsymbol{z})\,p_{\boldsymbol{\theta}}(\boldsymbol{z})}{p_{\boldsymbol{\theta}}(\boldsymbol{x})} \\ \nearrow \text{ intractable}}}$

  - ▶ what to learn: **model** parameters $\boldsymbol{\theta}$ and **variational** parameters $\boldsymbol{\phi}$
    decoder                          encoder

# GNU and VAE

- GNU
  - ▶ GNU's Not Unix!                                      (recursive acronym)
  - ▶ an operating system and an extensive collection of computer software



- VAE
  - ▶ VAE's not AutoEncoder!
  - ▶ a generative model to learn complex distributions

image source: https://en.wikipedia.org/wiki/GNU#/media/File:Heckert_GNU_white.svg

# Why are VAEs called "autoencoders"?

- mathematical basis of VAE[8]
  - ▶ has relatively little to do with classical AEs (*e.g.* sparse AE, denoising AE)

- they are called "AEs" only because
  - ▶ training objective has encoder/decoder ($\Rightarrow$ *resembles* traditional AE)

- differences
  - ▶ unlike sparse/denoising AEs
    - ▷ VAE allows us to sample directly from $p(\boldsymbol{x})$ (without doing MCMC)
  - ▶ typical AE: decoders are sometimes removed after training
    - ▷ $p(\boldsymbol{z} \,|\, \boldsymbol{x})$ learned by encoder: used for *e.g.* supervised learning
  - ▶ VAE: encoders are removed after training
    - ▷ $p(\boldsymbol{x} \,|\, \boldsymbol{z})$ learned by decoder: used for data generation

---

[8] Doersch, C. (2016). Tutorial on variational autoencoders. *arXiv preprint arXiv:1606.05908*

# Comparison

- **autoencoders**

$$x \xrightarrow[\text{encoder}]{p(z \mid x)} z \xrightarrow[\text{decoder}]{p(x \mid z)} x \qquad \text{(training)}$$

$$x \xrightarrow[\text{encoder}]{p(z \mid x)} z \qquad \text{(main use: feature learner for downstream learning)}$$

- **generative models**

$$\xrightarrow{p(x)} x$$

- **latent generative models**

$$\xrightarrow[\text{prior}]{p(z)} z \xrightarrow[\text{decoder}]{p(x \mid z)} x$$

- **VAEs**

$$x \xrightarrow[\text{encoder}]{q(z \mid x)} z \xrightarrow[\text{decoder}]{p(x \mid z)} x \qquad \text{(trained as an AE)}$$

$$\xrightarrow[\text{prior}]{p(z)} z \xrightarrow[\text{decoder}]{p(x \mid z)} x \qquad \text{(main use: data \underline{generation})}$$

# Outline

# Problem formulation

- objective: for each $x$ in training set

  ▶ maximize the probability of $x$ under the entire generative process:

$$p(x) = \int p(x \mid z)p(z)\,dz \qquad (2)$$

- to solve (2), VAEs must deal with three problems

  P1 how to define latent variables $\mathbf{z}$
  (*i.e.* decide what information they represent)

  P2 how to define output distribution $p(x \mid z)$

  P3 how to deal with <u>integral</u> over $\mathbf{z}$

# Solution 1: prior $p(z)$

- choose prior $p(z)$ to be <u>simple</u>

  ▶ reasonable for latent attributes (*e.g.* pose, how much smile)

- common choice: unit Gaussian $p(z) = \mathcal{N}(z; \mathbf{0}, \mathbf{I})$

  ▶ diagonal prior $\Rightarrow$ independent latent variables

  ▶ different dimensions of $z$ encode interpretable factors of variation



- varying $z_1$

  ▶ varying *head* pose

- varying $z_2$

  ▶ varying degree of *smile*

image source: Kingma, D. P. and Welling, M. (2013). Auto-encoding variational Bayes. *arXiv preprint arXiv:1312.6114*

# Solution 2: output distribution $p(\boldsymbol{x} \mid \boldsymbol{z})$

- use a differentiable generator ("decoder") network $g(\boldsymbol{z}; \boldsymbol{\theta})$
  - ▶ $\boldsymbol{\theta}$ denotes parameters of decoder neural net
  - ▶ $g(\boldsymbol{z}; \boldsymbol{\theta})$ produces parameters of $p(\boldsymbol{x} \mid \boldsymbol{z})$

$$p(\boldsymbol{x} \mid \boldsymbol{z}) = p(\boldsymbol{x}; g(\boldsymbol{z}; \boldsymbol{\theta})) \triangleq p_{\boldsymbol{\theta}}(\boldsymbol{x} \mid \boldsymbol{z})$$

- common parametric distributions for $p(\boldsymbol{x} \mid \boldsymbol{z})$
  - ▶ continuous $\mathbf{x}$ : a Gaussian parameterized by $g(\boldsymbol{z}; \boldsymbol{\theta})$
  - ▶ binary $\mathbf{x}$ : a Bernoulli parameterized by $g(\boldsymbol{z}; \boldsymbol{\theta})$

- at training time
  - ▶ likelihood of each input $\boldsymbol{x}^{(i)}$ can be computed using $p_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)} \mid \boldsymbol{z})$
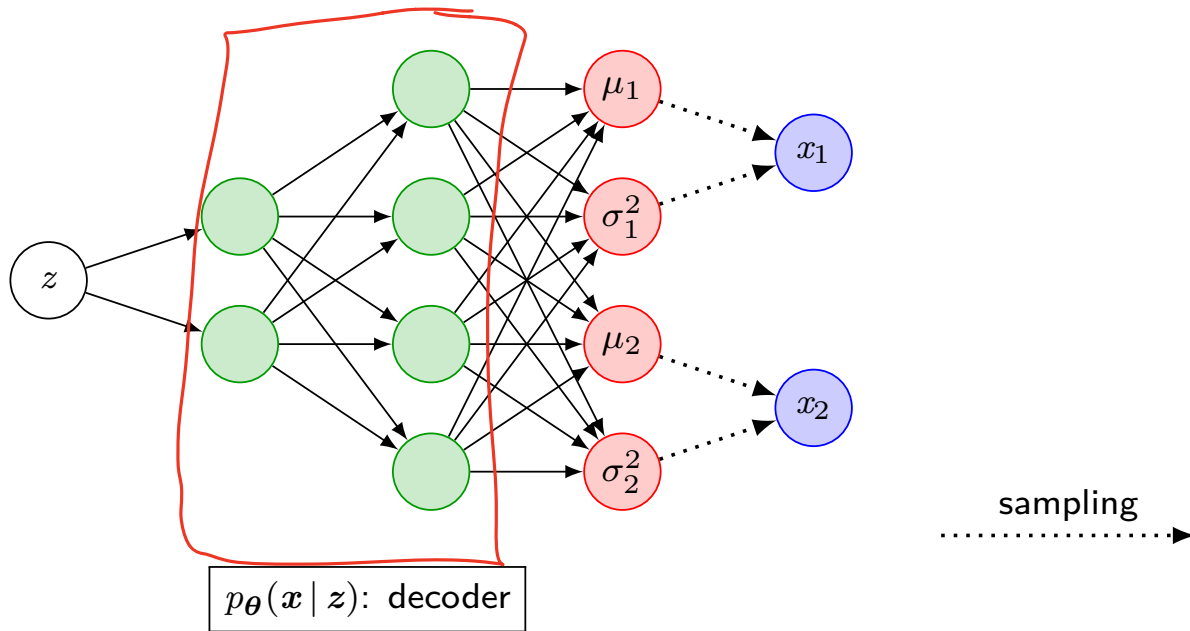  - ⇒ this generates error signal for doing SGD by backprop

- after training
  - ▶ a new $\boldsymbol{x}$ is sampled from $p_{\boldsymbol{\theta}}(\boldsymbol{x} \mid \boldsymbol{z})$

- *e.g.* Gaussian with diagonal covariance: $p(\boldsymbol{x} \,|\, z) = \mathcal{N}(\boldsymbol{x}; \boldsymbol{\mu}(z), \boldsymbol{\Sigma}(z))$

  ▸ $\boldsymbol{x} \in \mathbb{R}^2$, $z \in \mathbb{R}$

$$\boldsymbol{\mu}(z) = \begin{bmatrix} \mu_1(z) \\ \mu_2(z) \end{bmatrix}, \quad \boldsymbol{\Sigma}(z) = \begin{bmatrix} \sigma_1^2(z) & 0 \\ 0 & \sigma_2^2(z) \end{bmatrix}$$



$p_{\boldsymbol{\theta}}(\boldsymbol{x} \,|\, \boldsymbol{z})$: decoder

sampling

# Solution 3: integration over $z$

- from Solutions 1 & 2:

$$p(\boldsymbol{x}) = \int \underbrace{p(\boldsymbol{z})}_{\substack{\text{simple} \\ \text{Gaussian} \\ \text{prior}}} \underbrace{p(\boldsymbol{x} \mid \boldsymbol{z})}_{\substack{\text{learn with} \\ \text{neural net} \\ \text{"decoder"}}} d\boldsymbol{z}$$

  - ▶ $p_{\boldsymbol{\theta}}(\boldsymbol{z})$ and $p_{\boldsymbol{\theta}}(\boldsymbol{x} \mid \boldsymbol{z})$ are okay but integration remains <u>intractable</u>

- conceptually: easy to compute $p(\boldsymbol{x})$ approximately
  - ▶ sample a large number of $\boldsymbol{z}$ values $\{\boldsymbol{z}_1, \ldots, \boldsymbol{z}_n\}$ and compute

$$p(\boldsymbol{x}) \approx \frac{1}{n} \sum_{i=1}^{n} p(\boldsymbol{x} \mid \boldsymbol{z}_i)$$

- challenge: in high-dimensional space
  - ▶ $n$ : extremely large before having an accurate estimate of $p(\boldsymbol{x})$

- in practice: $p(\boldsymbol{x}\,|\,\boldsymbol{z}) \approx 0$ for most $\boldsymbol{z} \rightarrow$ contribute nothing to estimate of $p(\boldsymbol{x})$

- key idea behind VAE #1:
  - ▸ try to sample values of $\boldsymbol{z}$ that are likely to have produced $\boldsymbol{x}$, and
  - ▸ compute $p(\boldsymbol{x})$ just from those

- this means: we need a new function $q(\boldsymbol{z}\,|\,\boldsymbol{x})$, which can
  - ▸ take a value of $\boldsymbol{x}$, and
  - ▸ give us $\boxed{\text{a distribution over } \boldsymbol{z} \text{ that are likely to produce } \boldsymbol{x}}$

- why not exact posterior $p(\boldsymbol{z}\,|\,\boldsymbol{x})$?
  - ▸ also $\underline{\text{intractable}}$: $p(\boldsymbol{z}\,|\,\boldsymbol{x}) = p(\boldsymbol{x}\,|\,\boldsymbol{z})p(\boldsymbol{z})/p(\boldsymbol{x})$

- key idea #2: resort to variational inference, which allows us to
  1. find an approximate posterior $q(\boldsymbol{z}\,|\,\boldsymbol{x}) \approx p(\boldsymbol{z}\,|\,\boldsymbol{x})$
  2. $\underbrace{\text{maximize } \mathrm{ELBO}(q)}_{\text{tractable}}$ rather than $\underbrace{\text{maximizing } p(\boldsymbol{x})}_{\text{intractable}}$ directly
  - ▸ issue: $\boxed{\text{traditional VI using optimization is slow/requires closed form}}$

- key idea #3: use a $\boxed{\text{differentiable inference ("encoder") network } f(\boldsymbol{x}; \boldsymbol{\phi})}$

  - ▶ $\phi$ denotes parameters of encoder neural net
  - ▶ $f(\boldsymbol{x}; \boldsymbol{\phi})$ produces parameters of $q(\boldsymbol{z} \mid \boldsymbol{x})$

  $$q(\boldsymbol{z} \mid \boldsymbol{x}) = q(\boldsymbol{z}; f(\boldsymbol{x}; \boldsymbol{\phi})) \triangleq q_{\boldsymbol{\phi}}(\boldsymbol{z} \mid \boldsymbol{x})$$

  - ▶ training the encoder: often easier than optimization in traditional VI
    - ▷ backed by various efficient techniques (*e.g.* SGD, backprop)

- common parametric distributions for $q(\boldsymbol{z} \mid \boldsymbol{x})$
  - ▶ a Gaussian parameterized by $f(\boldsymbol{x}; \boldsymbol{\phi})$
  - ▶ its covariance $\boldsymbol{\Sigma}$ : constrained to be diagonal (for computational issues)

- z is sampled from
  - ▶ $q_{\boldsymbol{\phi}}(\boldsymbol{z} \mid \boldsymbol{x})$ in training
  - ▶ $p(\boldsymbol{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$ after training
  - $\Rightarrow$ and then the sampled $\mathbf{z}$ is fed into decoder to generate $\mathbf{x}$

- *e.g.* Gaussian $q(z \mid \boldsymbol{x}) = \mathcal{N}(z; \mu(\boldsymbol{x}), \sigma^2(\boldsymbol{x}))$

  ▶ $\boldsymbol{x} \in \mathbb{R}^2$, $z \in \mathbb{R}$



$q_{\boldsymbol{\phi}}(\boldsymbol{z} \mid \boldsymbol{x})$: encoder

sampling

# VAE architecture

- encoder and decoder are stacked



- we can learn parameters $\phi$ and $\boldsymbol{\theta}$ via _____

# Outline

# Training objective

- as in VI, decomposing $\log p(\boldsymbol{x})$ reveals it:

$$
\begin{aligned}
\log p(\boldsymbol{x}) &= \mathbb{E}_{\mathbf{z} \sim q(z|x)}[\log p(\boldsymbol{x})] \\
&= \mathbb{E}_{\mathbf{z}}\left[\log \frac{p(\boldsymbol{x} \mid \boldsymbol{z})p(\boldsymbol{z})}{p(\boldsymbol{z} \mid \boldsymbol{x})}\right] \\
&= \mathbb{E}_{\mathbf{z}}\left[\log \frac{p(\boldsymbol{x} \mid \boldsymbol{z})p(\boldsymbol{z})}{p(\boldsymbol{z} \mid \boldsymbol{x})}\frac{q(\boldsymbol{z} \mid \boldsymbol{x})}{q(\boldsymbol{z} \mid \boldsymbol{x})}\right] \\
&= \mathbb{E}_{\mathbf{z}}\left[\log p(\boldsymbol{x} \mid \boldsymbol{z})\right] - \mathbb{E}_{\mathbf{z}}\left[\log \frac{q(\boldsymbol{z} \mid \boldsymbol{x})}{p(\boldsymbol{z})}\right] + \mathbb{E}_{\mathbf{z}}\left[\log \frac{q(\boldsymbol{z} \mid \boldsymbol{x})}{p(\boldsymbol{z} \mid \boldsymbol{x})}\right] \\
&= \underbrace{\mathbb{E}_{\mathbf{z}}\left[\log p_{\boldsymbol{\theta}}(\boldsymbol{x} \mid \boldsymbol{z})\right] - D_{\mathrm{KL}}\left[q_{\boldsymbol{\phi}}(\boldsymbol{z} \mid \boldsymbol{x}) \,\|\, p(\boldsymbol{z})\right]}_{\text{ELBO: } \mathcal{L}(\boldsymbol{x}, \boldsymbol{\theta}, \boldsymbol{\phi})} + \underbrace{D_{\mathrm{KL}}\left[q_{\boldsymbol{\phi}}(\boldsymbol{z} \mid \boldsymbol{x}) \,\|\, p(\boldsymbol{z} \mid \boldsymbol{x})\right]}_{\geq 0 \text{ (intractable)}}
\end{aligned}
$$

- note: *any* distribution over $\mathbf{z}$ can be proxy $q$ for posterior $p(\boldsymbol{z} \mid \boldsymbol{x})$
  - ▶ EM/VI: using $q(\boldsymbol{z})$ is more common
  - ▶ VAE: using $q(z|x)$ is more common[9]

---

[9] since we want to infer $p(\boldsymbol{x})$, it makes sense to construct $q$ that does depend on $\boldsymbol{x}$ [see Doersch, C. (2016). Tutorial on variational autoencoders. *arXiv preprint arXiv:1606.05908*]

# Interpretation[10]

- VAE: trained by maximizing variational lower bound $\mathcal{L}(\boldsymbol{x}, \boldsymbol{\theta}, \boldsymbol{\phi})$

$$\log p(\boldsymbol{x}) \geq \mathcal{L}(\boldsymbol{x}, \boldsymbol{\theta}, \boldsymbol{\phi})$$

$$= \underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x})} \left[ \log p(\boldsymbol{z}, \boldsymbol{x}) \right]}_{\textcircled{1}} + \underbrace{H[q_{\boldsymbol{\phi}}(\boldsymbol{z}\,|\,\boldsymbol{x})\,]}_{\textcircled{2}\ \text{entropy}}$$

$$= \underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x})} \left[ \log p_{\boldsymbol{\theta}}(\boldsymbol{x}\,|\,\boldsymbol{z}) \right]}_{\textcircled{3}\ \text{reconstruction}} - \underbrace{D_{\mathrm{KL}}[q_{\boldsymbol{\phi}}(\boldsymbol{z}\,|\,\boldsymbol{x})\,||\,p(\boldsymbol{z})\,]}_{\textcircled{4}\ \text{regularizer}}$$

① joint log-likelihood of visible/hidden variables

   ▶ under approximate posterior over latent variables

② entropy of approximate posterior: maximizing this encourages to

   ▶ place high prob mass on <u>many</u> $z$ values that could have generated $\boldsymbol{x}$

   ▶ rather than collapsing to a single point estimate of the most likely value

③ reconstruction log-likelihood found in other autoencoders

④ tries to make approximate posterior $q_{\boldsymbol{\phi}}(\boldsymbol{z}\,|\,\boldsymbol{x})$ and model prior $p(\boldsymbol{z})$ close

   ▶ tries to reflect prior knowledge $\Rightarrow$ $\boxed{\text{regularizer}}$

---

[10] Goodfellow, I., Bengio, Y., and Courville, A. (2016). *Deep learning.* MIT Press

# Interpretation[11]

- rearranging decomposition of $\log p(\boldsymbol{x})$ gives

$$\underbrace{\log p(\boldsymbol{x})}_{①} - \underbrace{D_{\mathrm{KL}}[q_\phi(\boldsymbol{z}\,|\,\boldsymbol{x})\,||\,p(\boldsymbol{z}\,|\,\boldsymbol{x})]}_{②}$$

$$= \underbrace{\mathbb{E}_{\mathbf{z}\sim q_\phi(z|x)}[\log p_\theta(\boldsymbol{x}\,|\,\boldsymbol{z})]}_{③ \text{ decoding objective}} - \underbrace{D_{\mathrm{KL}}[q_\phi(\boldsymbol{z}\,|\,\boldsymbol{x})\,||\,p(\boldsymbol{z})]}_{④ \text{ encoding objective}}$$

- **left hand side**

  ① quantity we want to maximize

  ② error term (this will become small if $q$ is high-capacity)

     ▷ makes $q$ produce $\boldsymbol{z}$'s that can reproduce a given $\boldsymbol{x}$

- **right hand side**: something we can optimize via SGD (given right choice of $q$)

  ▶ suddenly takes a form which looks like an <u>autoencoder</u>

     ▷ $q$ is "encoding" $\boldsymbol{x}$ into $\boldsymbol{z}$ (quality measured by ④)

     ▷ $p$ is "decoding" it to reconstruct $\boldsymbol{x}$ (quality measured by ③)

---

[11]Doersch, C. (2016). Tutorial on variational autoencoders. *arXiv preprint arXiv:1606.05908*

# Training

- perform SGD to maximize the objective:

$$J = \mathcal{L}(\boldsymbol{x}, \boldsymbol{\theta}, \boldsymbol{\phi})$$
$$= \underbrace{-D_{\mathrm{KL}}[q_{\boldsymbol{\phi}}(\boldsymbol{z}\,|\,\boldsymbol{x})\,||\,p(\boldsymbol{z})]}_{\text{\textcircled{1} encoding objective}} + \underbrace{\mathbb{E}_{\mathbf{z}\sim q_{\boldsymbol{\phi}}(z|x)}\left[\log p_{\boldsymbol{\theta}}(\boldsymbol{x}\,|\,\boldsymbol{z})\right]}_{\text{\textcircled{2} decoding objective}}$$

- usual choice: $q(\boldsymbol{z}\,|\,\boldsymbol{x}) = \mathcal{N}(\boldsymbol{z}; \boldsymbol{\mu}(\boldsymbol{x}), \boldsymbol{\Sigma}(\boldsymbol{x}))$

  ▶ $\boldsymbol{\mu}$, $\boldsymbol{\Sigma}$ : learned/represented via neural nets

  ▶ $\boldsymbol{\Sigma}$ : a diagonal matrix

① $D_{\mathrm{KL}}$ between two multivariate Gaussians

  ▶ can be computed in closed form

② a bit more tricky

  ▶ involves random <u>sampling</u> $\Rightarrow$ **how to backprop** gradient for SGD?

# Training: encoding objective

- optimizing ①: encoding objective $-D_{\mathrm{KL}}[q_\phi(z\,|\,x)\,||\,p(z)]$

  ▶ in general ($K$: dimensionality of variable)

$$D_{\mathrm{KL}}[\mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)\,||\,\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)]$$
$$= \frac{1}{2}\left(\mathrm{tr}\left(\boldsymbol{\Sigma}_1^{-1}\boldsymbol{\Sigma}_0\right) + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)^\top \boldsymbol{\Sigma}_1^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) - K + \ln\left(\frac{\det \boldsymbol{\Sigma}_1}{\det \boldsymbol{\Sigma}_0}\right)\right)$$

  ▶ in VAE this simplifies[12] to

$$D_{\mathrm{KL}}[\mathcal{N}(\boldsymbol{\mu}(x), \boldsymbol{\Sigma}(x))\,||\,\mathcal{N}(\mathbf{0}, \mathbf{I})] = \frac{1}{2}\left(\mathrm{tr}\left(\boldsymbol{\Sigma}(x)\right) + \boldsymbol{\mu}(x)^\top \boldsymbol{\mu}(x) - K - \ln\det \boldsymbol{\Sigma}(x)\right)$$
$$= \frac{1}{2}\left(\sum_{k=1}^{K} \sigma_k^2(x) + \sum_{k=1}^{K} \mu_k^2(x) - \sum_{k=1}^{K} 1 - \ln\prod_{k=1}^{K} \sigma_k^2(x)\right)$$

$$\boxed{= \frac{1}{2}\sum_{k=1}^{K}\left(\sigma_k^2(x; \phi) + \mu_k^2(x; \phi) - 1 - \ln\sigma_k^2(x; \phi)\right)}$$

  ▶ differentiable $\Rightarrow$ can be minimized wrt $\phi$ via SGD [13]

---

[12] for diagonal $\boldsymbol{\Sigma} = \mathrm{diag}[\sigma_1^2, \sigma_2^2, \ldots]$, $\mathrm{tr}(\boldsymbol{\Sigma}) = \sum_k \sigma_k^2$ and $\det \boldsymbol{\Sigma} = \prod_k \sigma_k^2$

[13] to avoid clutter, revealed dependence on $\phi$ only in the last equation

# Training: decoding objective

- forward prop to compute ②: decoding objective $\mathbb{E}_{\mathbf{z} \sim q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x})} \left[ \log p_{\boldsymbol{\theta}}(\boldsymbol{x} \,|\, \boldsymbol{z}) \right]$

  - ▶ use Monte Carlo sampling: for each input $\boldsymbol{x}^{(i)}$

    1. $\boldsymbol{z}^{(i,l)}$ is sampled from encoder $q_{\boldsymbol{\phi}}(\boldsymbol{z} \,|\, \boldsymbol{x}^{(i)})$

    2. $\boldsymbol{z}^{(i,l)}$ is fed to decoder $\Rightarrow$ gives $\log p_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)} \,|\, \boldsymbol{z}^{(i,l)})$

    3. repeat above steps $L$ times to approximate

$$\mathbb{E}_{\mathbf{z} \sim q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x})} \left[ \log p_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)} \,|\, \boldsymbol{z}^{(i,l)}) \right] \approx \frac{1}{L} \sum_{l=1}^{L} \log p_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)} \,|\, \boldsymbol{z}^{(i,l)})$$

  - ▶ in practice, often $L = 1$ (as in SGD; the same motivation)

- however, to maximize ② via SGD with backprop[14]

  - ▶ we need a trick called <u>*reparameterization*</u>

  - ▶ without it, we cannot backprop gradient through the sampling process

    - ▷ sampling: not continuous $\Rightarrow$ not differntiable

---

[14]SGD via backprop can handle stochastic inputs but not stochastic units within the network
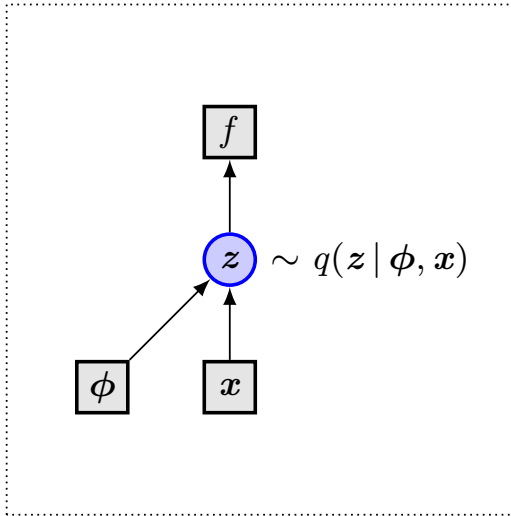
# Reparameterization trick

- idea: move sampling to an <u>input layer</u>

- given $\boldsymbol{\mu}(\boldsymbol{x})$ and $\boldsymbol{\Sigma}(\boldsymbol{x})$ (*i.e.* mean and covariance of $q(\boldsymbol{z} \,|\, \boldsymbol{x})$):
  - ▶ instead of sampling $\boldsymbol{z}^{(i,l)} \sim \mathcal{N}(\boldsymbol{\mu}(\boldsymbol{x}^{(i)}), \boldsymbol{\Sigma}(\boldsymbol{x}^{(i)}))$ directly
  - ▶ we can compute $\boldsymbol{z}^{(i,l)}$ by
    1. first sampling $\boldsymbol{\epsilon}^{(l)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
    2. then computing

    $$\boldsymbol{z}^{(i,l)} = \boldsymbol{\mu}(\boldsymbol{x}^{(i)}) + \boldsymbol{\sigma}(\boldsymbol{x}^{(i)}) \odot \boldsymbol{\epsilon}^{(l)}$$
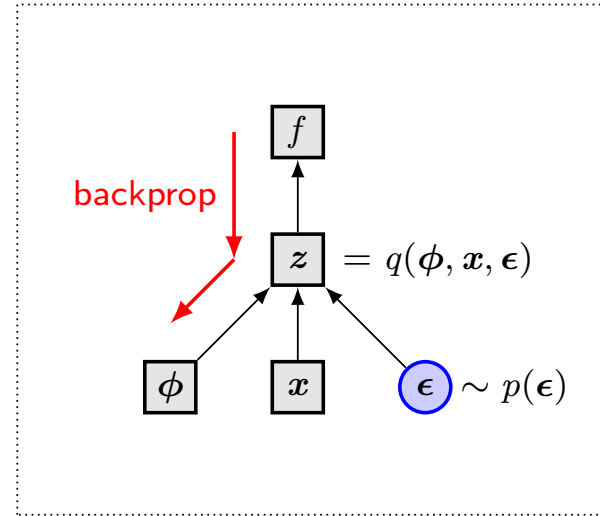
    ▷ $\odot$: element-wise multiplication
  - ▶ $\boldsymbol{z}^{(i,l)}$ has the same distribution as before, but now we can do backprop

original form

reparameterized form

$z \sim q(z \mid \phi, x)$

backprop

$z = q(\phi, x, \epsilon)$

$\epsilon \sim p(\epsilon)$

deterministic node

random node

# Training: combined objectives

- finally, training objective for datapoint $\boldsymbol{x}^{(i)}$:

$$\mathcal{L}(\boldsymbol{x}^{(i)}, \boldsymbol{\theta}, \boldsymbol{\phi}) = \underbrace{-D_{\mathrm{KL}}[q_{\boldsymbol{\phi}}(\boldsymbol{z} \mid \boldsymbol{x}) \,\|\, p(\boldsymbol{z})]}_{\textbf{encoding} \text{ objective}} + \underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\boldsymbol{\phi}}(\boldsymbol{z} \mid \boldsymbol{x})}\left[\log p_{\boldsymbol{\theta}}(\boldsymbol{x} \mid \boldsymbol{z})\right]}_{\textbf{decoding} \text{ objective}} \tag{3}$$

$$\simeq \frac{1}{2} \sum_{k=1}^{K} \left(1 + \ln \sigma_k^2(\boldsymbol{x}^{(i)}; \boldsymbol{\phi}) - \mu_k^2(\boldsymbol{x}^{(i)}; \boldsymbol{\phi}) - \sigma_k^2(\boldsymbol{x}^{(i)}; \boldsymbol{\phi})\right) + \frac{1}{L} \sum_{l=1}^{L} \log p_{\boldsymbol{\theta}}\left(\boldsymbol{x}^{(i)} \mid \boldsymbol{z}^{(i,l)}\right)$$

  - ▶ where

$$\boldsymbol{z}^{(i,l)} = \boldsymbol{\mu}(\boldsymbol{x}^{(i)}) + \boldsymbol{\sigma}(\boldsymbol{x}^{(i)}) \odot \boldsymbol{\epsilon}^{(l)} \text{ and } \boldsymbol{\epsilon}^{(l)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$L$ : Monte Carlo sample size (often $L = 1$)

$K$ : dimensionality of latent variable $\mathbf{z}$

- given a fixed $\boldsymbol{x}$ and $\boldsymbol{\epsilon}$
  - ▶ (3) is deterministic and continuous in $\boldsymbol{\theta}$ and $\boldsymbol{\phi}$
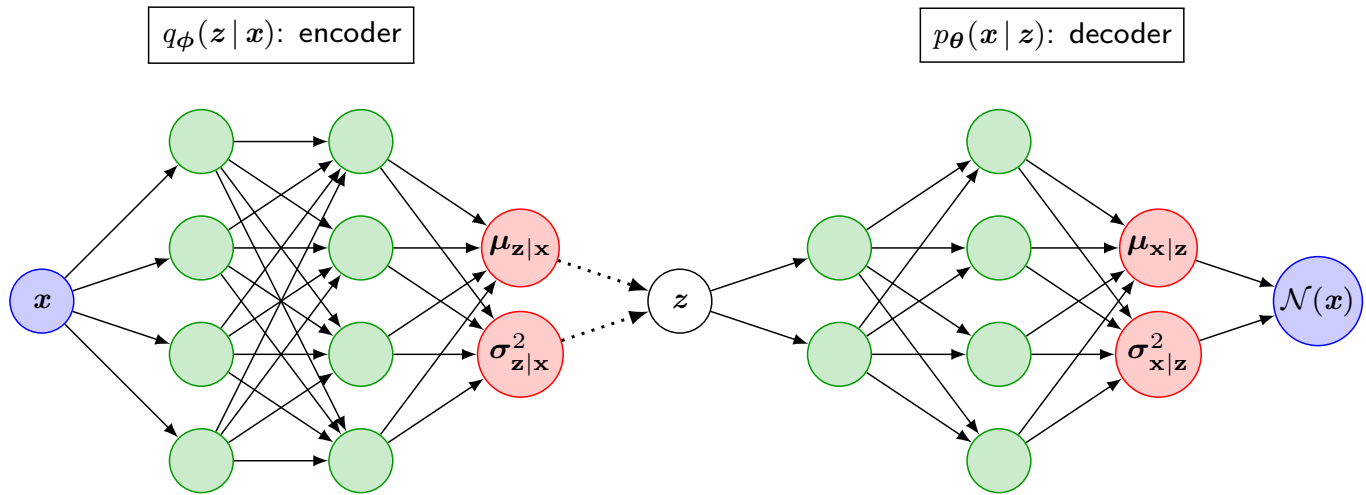  - ⇒ backprop can compute a gradient that will work for SGD

# Training summary

**training in variational autoencoders (VAEs):**

$$\boldsymbol{\theta}^*, \boldsymbol{\phi}^* = \underset{\boldsymbol{\theta}, \boldsymbol{\phi}}{\operatorname{argmax}} \sum_{i=1}^{N} \mathcal{L}(\boldsymbol{x}^{(i)}, \boldsymbol{\theta}, \boldsymbol{\phi})$$

- $\boldsymbol{\theta}$ : parameters of generator network ("decoder")
- $\boldsymbol{\phi}$ : parameters of inference network ("encoder")
- $N$ : number of training samples ($\{\boldsymbol{x}^{(1)}, \dots, \boldsymbol{x}^{(N)}\}$)
- $\mathcal{L}(\boldsymbol{x}^{(i)}, \boldsymbol{\theta}, \boldsymbol{\phi})$ : variational lower bound (see below)

$$\mathcal{L}(\boldsymbol{x}^{(i)}, \boldsymbol{\theta}, \boldsymbol{\phi}) = \underbrace{-D_{\mathrm{KL}}[q_{\boldsymbol{\phi}}(\boldsymbol{z} \mid \boldsymbol{x}) \,||\, p(\boldsymbol{z})]}_{\textbf{encoding } \text{objective}} + \underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x})}[\log p_{\boldsymbol{\theta}}(\boldsymbol{x} \mid \boldsymbol{z})]}_{\textbf{decoding } \text{objective}}$$

$$\simeq \frac{1}{2} \sum_{k=1}^{K} \left( 1 + \ln \sigma_k^2(\boldsymbol{x}^{(i)}; \boldsymbol{\phi}) - \mu_k^2(\boldsymbol{x}^{(i)}; \boldsymbol{\phi}) - \sigma_k^2(\boldsymbol{x}^{(i)}; \boldsymbol{\phi}) \right) + \frac{1}{L} \sum_{l=1}^{L} \log p_{\boldsymbol{\theta}}\left( \boldsymbol{x}^{(i)} | \boldsymbol{z}^{(i,l)} \right)$$

$$\boldsymbol{z}^{(i,l)} = \boldsymbol{\mu}(\boldsymbol{x}^{(i)}) + \boldsymbol{\sigma}(\boldsymbol{x}^{(i)}) \odot \boldsymbol{\epsilon}^{(l)} \quad \text{where} \quad \boldsymbol{\epsilon}^{(l)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

- VAE with ~~Gaussian~~ encoder/decoder:

$q_{\boldsymbol{\phi}}(\boldsymbol{z}\,|\,\boldsymbol{x})$: encoder

$p_{\boldsymbol{\theta}}(\boldsymbol{x}\,|\,\boldsymbol{z})$: decoder



$$\mathcal{L}(\boldsymbol{x}^{(i)}, \boldsymbol{\theta}, \boldsymbol{\phi}) = \underbrace{-D_{\mathrm{KL}}[q_{\boldsymbol{\phi}}(\boldsymbol{z}\,|\,\boldsymbol{x})\,||\,p(\boldsymbol{z})]}_{\text{encoding objective}} + \underbrace{\mathbb{E}_{\mathbf{z}\sim q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x})}[\log p_{\boldsymbol{\theta}}(\boldsymbol{x}\,|\,\boldsymbol{z})]}_{\text{decoding objective}}$$

$$\simeq \frac{1}{2}\sum_{k=1}^{K}\left(1 + \ln\sigma^2_{\mathbf{z}|\mathbf{x},k}(\boldsymbol{x}^{(i)}; \boldsymbol{\phi}) - \mu^2_{\mathbf{z}|\mathbf{x},k}(\boldsymbol{x}^{(i)}; \boldsymbol{\phi}) - \sigma^2_{\mathbf{z}|\mathbf{x},k}(\boldsymbol{x}^{(i)}; \boldsymbol{\phi})\right) + \frac{1}{L}\sum_{l=1}^{L}\log p_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)}|\boldsymbol{z}^{(i,l)})$$

$$\log p_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)}|\boldsymbol{z}^{(i,l)}) = \log\mathcal{N}\left(\boldsymbol{x}^{(i)}; \boldsymbol{\mu}_{\mathbf{x}|\mathbf{z}}(\boldsymbol{z}^{(i,l)}), \boldsymbol{\sigma}^2_{\mathbf{x}|\mathbf{z}}(\boldsymbol{z}^{(i,l)})\mathbf{I}\right)$$

$$\boldsymbol{z}^{(i,l)} = \boldsymbol{\mu}_{\mathbf{z}|\mathbf{x}}(\boldsymbol{x}^{(i)}) + \boldsymbol{\sigma}_{\mathbf{z}|\mathbf{x}}(\boldsymbol{x}^{(i)}) \odot \boldsymbol{\epsilon}^{(l)} \quad \text{where} \quad \boldsymbol{\epsilon}^{(l)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

# Sample generation

- when we want to generate a new sample
    - simply input a sample $z \sim \underline{N(0,I)}$ into decoder $= \; p(z) \text{ is prior}$
        - ▷ do not use encoder at test time
    - decoder will produce parameters of $p(x \,|\, z)$
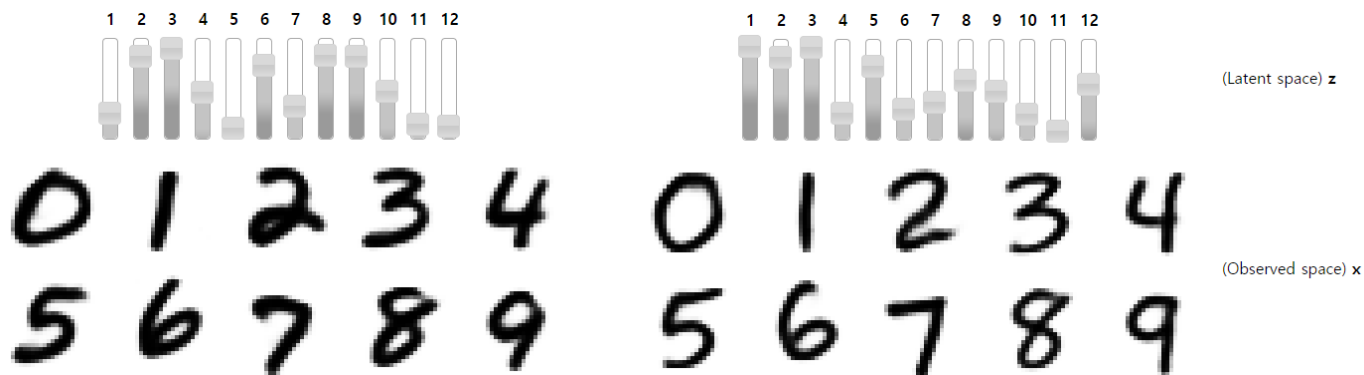    - we then sample a new $x$ from $p(x \,|\, z)$



image source: http://www.dpkingma.com/sgvb_mnist_demo/demo.html

# Outline

# VAE in a nutshell

- a latent generative model using **neural nets** and **variational inference**

  ▶ relies on variational inference to resolve the challenge of inference

  $$\log p(\boldsymbol{x}) \geq \underbrace{\text{ELBO}(\boldsymbol{x}, \boldsymbol{\theta}, \boldsymbol{\phi})}_{\substack{\boldsymbol{\theta}: \text{ parameters for decoder} \\ \boldsymbol{\phi}: \text{ parameters for encoder}}} = \underbrace{\mathbb{E}_{\mathbf{z} \sim q(\boldsymbol{z}|\boldsymbol{x})}[\log p(\boldsymbol{x} \,|\, \boldsymbol{z})]}_{\substack{\text{expected log likelihood} \\ \text{of the data}}} - \underbrace{D_{\text{KL}}[q(\boldsymbol{z} \,|\, \boldsymbol{x}) \,||\, p(\boldsymbol{z})]}_{\substack{D_{\text{KL}} \text{ between} \\ \text{prior } p(\boldsymbol{z}) \text{ and } q(\boldsymbol{z}|\boldsymbol{x})}}$$

  ▶ add "encoder" to learn parametric $q(\boldsymbol{z} \,|\, \boldsymbol{x})$ that approximates $p(\boldsymbol{z} \,|\, \boldsymbol{x})$

  $$\boldsymbol{x} \xrightarrow[\text{encoder}]{q(\boldsymbol{z} \,|\, \boldsymbol{x})} \boldsymbol{z} \xrightarrow[\text{decoder}]{p(\boldsymbol{x} \,|\, \boldsymbol{z})} \boldsymbol{x} \qquad\qquad \text{(trained as an AE)}$$

    ▷ $p(\boldsymbol{x} \,|\, \boldsymbol{z})$: now viewed as "decoder"  (encoder/decoder: neural nets)
    ▷ encoder learns distribution parameters of $q(\boldsymbol{z} \,|\, \boldsymbol{x})$  (e.g. $\boldsymbol{\mu}, \boldsymbol{\Sigma}$)
    ▷ $\boldsymbol{z}$ is sampled from $q(\boldsymbol{z} \,|\, \boldsymbol{x})$  (e.g. $\boldsymbol{z} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$)

  ▶ train VAE = maximize $\text{ELBO}$ (tractable)
     = learn $\boldsymbol{\theta}$ and $\boldsymbol{\phi}$ by backprop (+ <u>reparameterization</u> trick)

  ▶ after training: encoder is removed  (c.f. decoder is removed in AEs)

  $$\xrightarrow[\text{prior}]{p(\boldsymbol{z})} \boldsymbol{z} \xrightarrow[\text{decoder}]{p(\boldsymbol{x} \,|\, \boldsymbol{z})} \boldsymbol{x}$$

    ▷ use a simple prior $p(\boldsymbol{z}) \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  ($\Rightarrow D_{\text{KL}}$ above: "regularizer")

# Manifold learning by VAE

- nice property of VAE

  ▶ simultaneous training of a parametric encoder with decoder

  ⇒ forces to learn a coordinate system the encoder can capture

  ⇒ makes VAE an excellent manifold learning algorithm



- example of low-dim manifolds learned by VAE

  ▶ 2D map of MNIST manifold

# Limitations

- main drawback
  - ▶ image samples from VAE: *blurry*
  - ▶ cause: some explanations possible but not yet completely known

# Outline

# Summary

- autoencoders: stacked/denoising/sparse/contractive and many more
  - ▶ copy input to output learning useful representation of input

- variational autoencoder (VAE)
  - ▶ probabilistic spin to traditional AEs $\Rightarrow$ allows generating data
  - ▶ defines an intractable density $\Rightarrow$ derive and optimize a (variational) lower bound

- VAE advantages: principled approach to generative models
  - ▶ allows inference of $q(z \,|\, x)$
  - $\Rightarrow$ can be useful feature representation for other tasks

- VAE limitations
  - ▶ maximizes lower bound of likelihood (not as good evaluation as PixelCNN)
  - ▶ samples blurrier and lower quality compared to state-of-the-art (GANs)

- active areas of research in VAE
  - ▶ more flexible approximations (*e.g.* richer approximate posterior than diagonal $\mathcal{N}$)
  - ▶ incorporating structure in latent variables