

Autoencoders; Variational Autoencoders

Neural Networks So Far

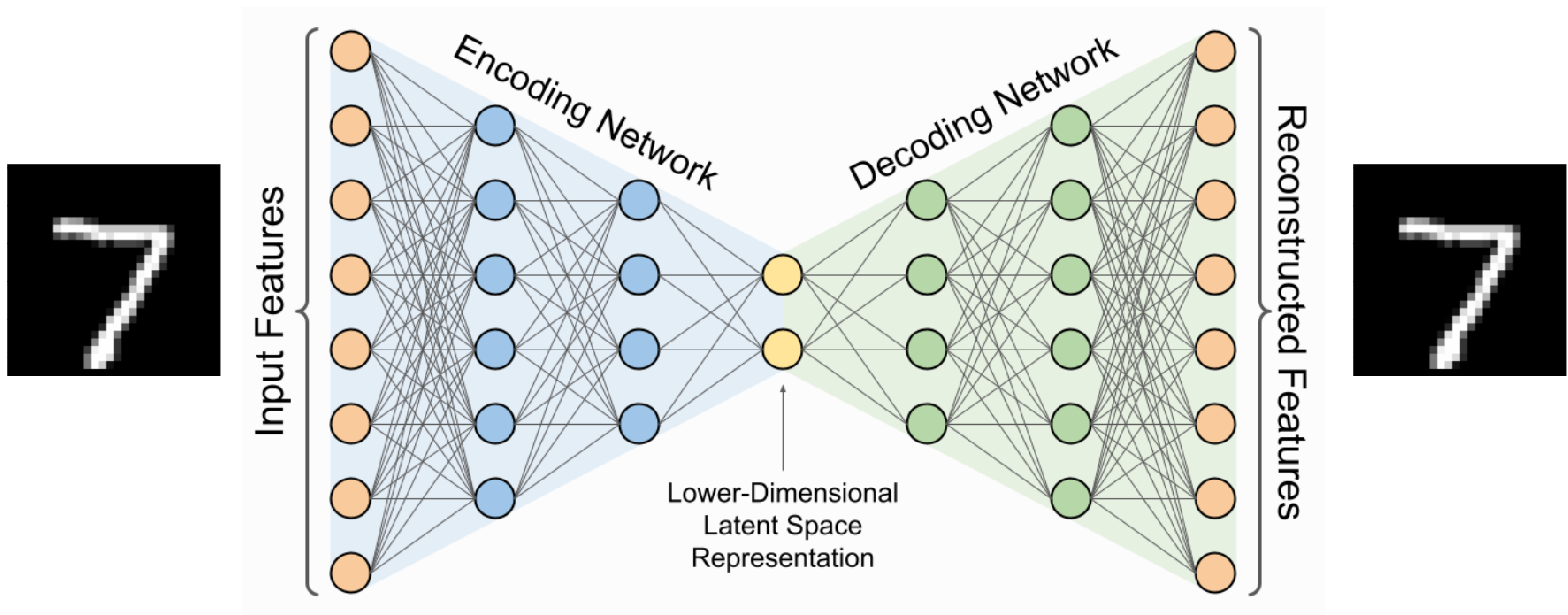
- Multilayer perceptrons
- Convolutional neural networks
- Recurrent neural networks
- Self-attention and Transformers

Today

- Autoencoders
 - dimensionality reduction
 - denoising
 - superresolution
- Variational Autoencoders
 - Generative model: synthesis of new data

Autoencoders

- These are neural networks for dim. reduction / feature extraction
- Idea: Try to reconstruct original data through a bottleneck
- Extracted features = bottleneck layer
- Layers need not be fully connected

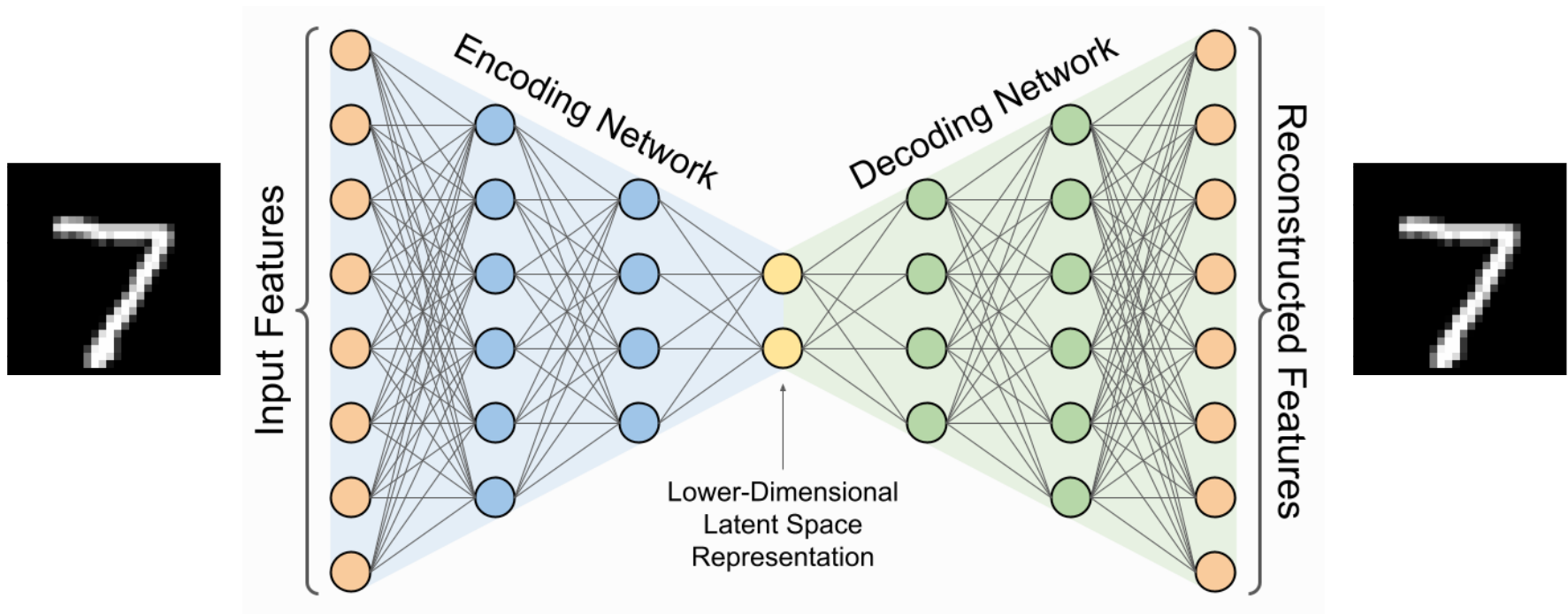


Training Autoencoders

- Consider unlabeled training data $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$, and let $f_{\theta}(\mathbf{x}) \in \mathbb{R}^d$ denote the output of the autoencoder, with weights θ .

- Objective function for learning θ ?

$$\min_{\theta} \sum_{i=1}^n \|\mathbf{x}_i - f_{\theta}(\mathbf{x}_i)\|_2^2$$



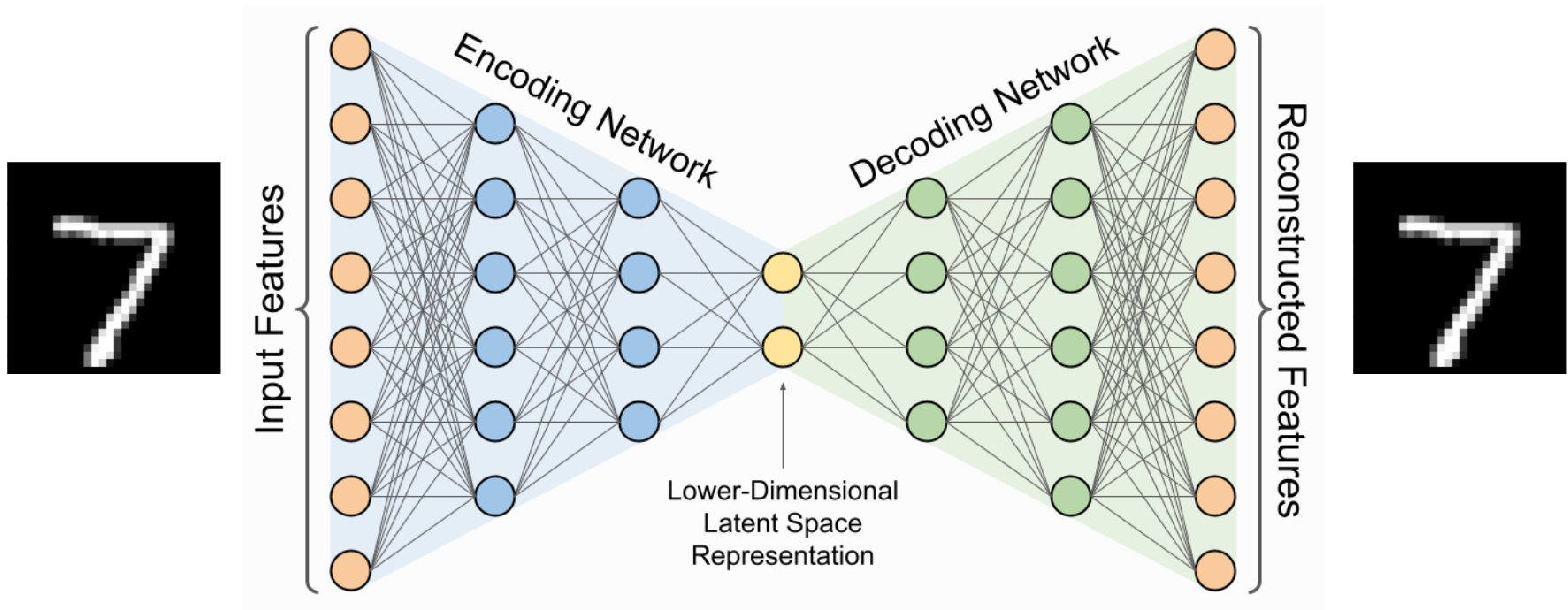
Poll

To optimize the reconstruction error, it suffices to use backpropagation as discussed previously for supervised learning with neural networks

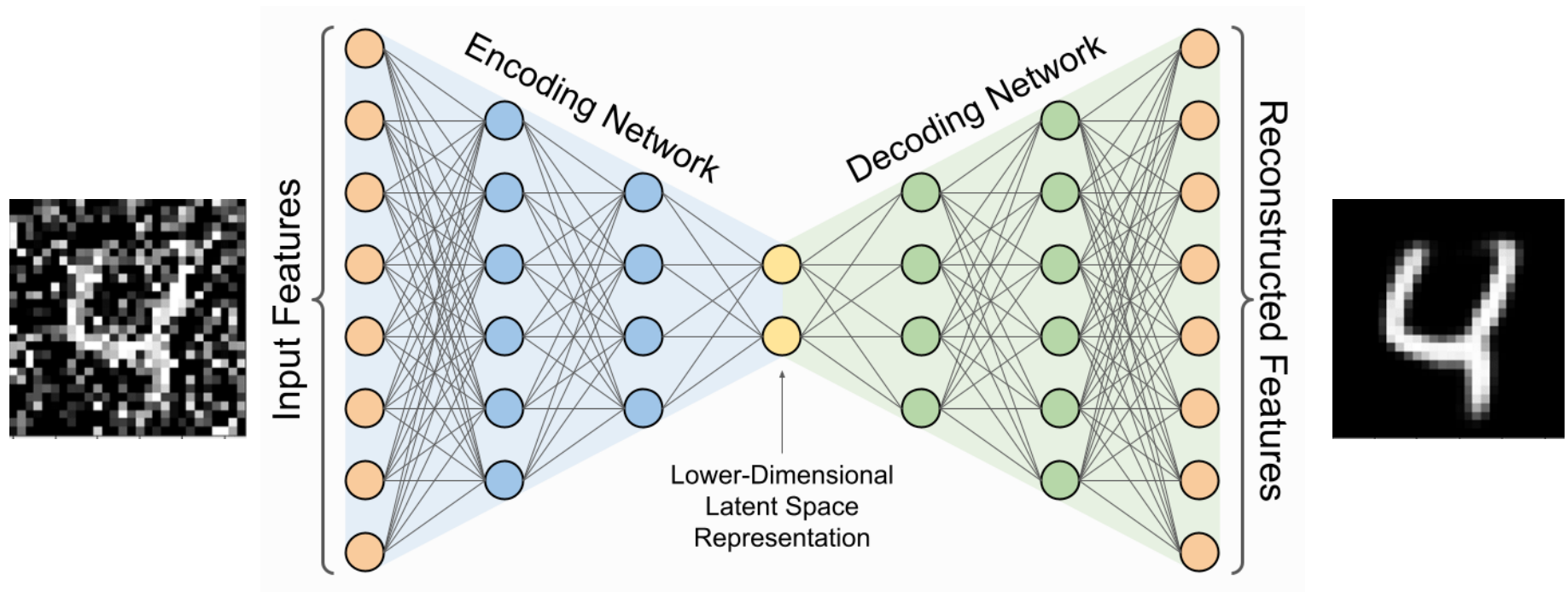
(A) True

(B) False

$$\sum_{i=1}^n \|x_i - f_{\theta}(x_i)\|_2^2$$

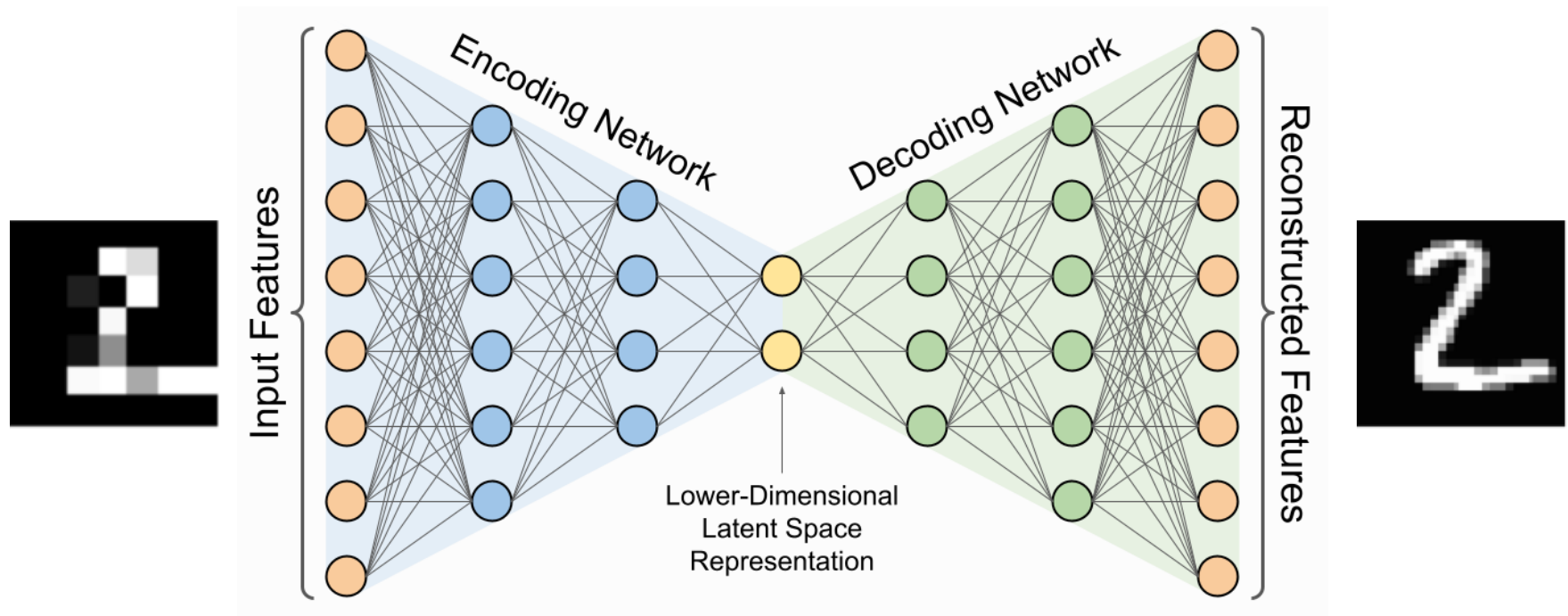


Autoencoders for Denoising



Train with pairs (noisy image, clean image)

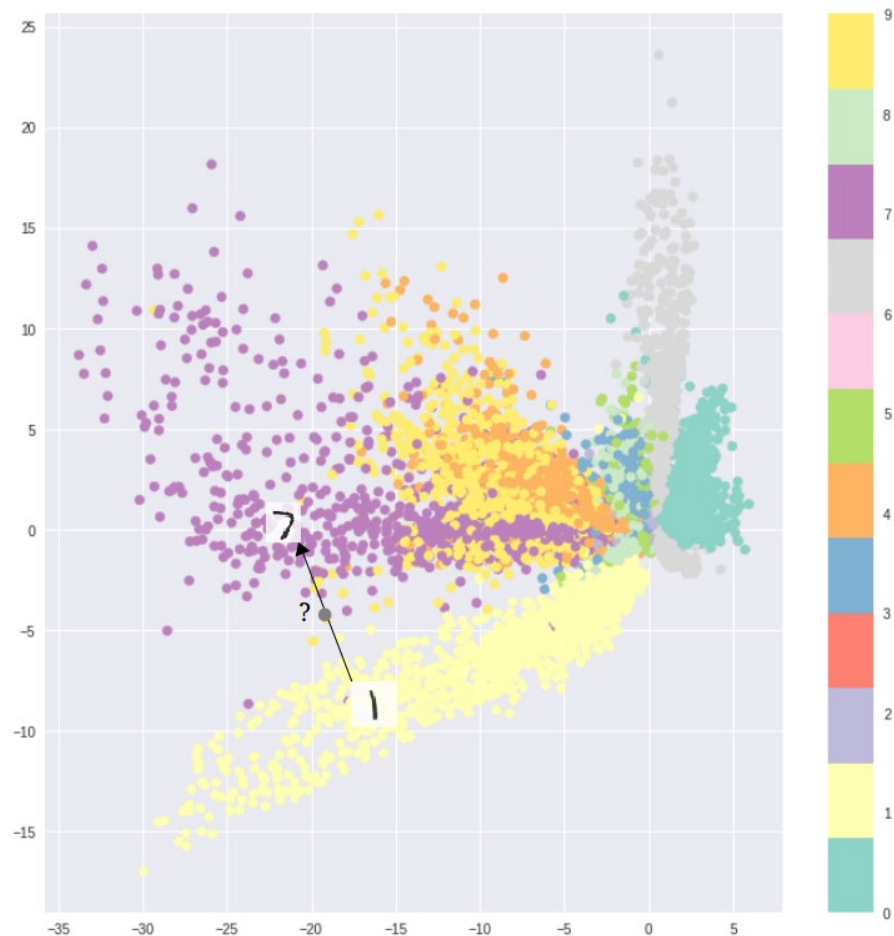
Autoencoders for Superresolution



Train with pairs (blurry image, sharp image)

Autoencoders for Simulation?

- Given: Examples $\mathbf{x}_1, \dots, \mathbf{x}_n$ of a probability distribution
- Goal: Simulate new instances from the distribution
- First try: Use the decoder half of an autoencoder with random inputs
- Issue: unrealistic synthetic images



Variational Autoencoders

- An autoencoder that functions as a generative model
- A type of latent variable model
- Approximate maximum likelihood estimation by *variational inference*

Deep Latent Variable Models

- Given: Examples $\mathbf{x}_1, \dots, \mathbf{x}_n$ of a probability distribution
- Goal: Simulate new instances from the distribution
- Idea: Latent variable model

$$z \sim p(z)$$
$$\mathbf{x} | z \sim p(\mathbf{x} | z; \theta)$$

where $p(\mathbf{x} | z; \theta)$ is parametrized by a neural network.

DLVM Example: Binary Outputs

- Suppose

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} \in \{0, 1\}^d.$$

- DLVM:

$$\begin{aligned} \mathbf{z} &\sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \\ x_j | \mathbf{z} &\stackrel{\text{ind}}{\sim} \text{Bernoulli}(p_j), \end{aligned} \quad \begin{bmatrix} p_1 \\ \vdots \\ p_d \end{bmatrix} = \sigma(\mathbf{d}_\theta(\mathbf{z}))$$

where

- $\mathbf{d}_\theta : \mathbb{R}^k \rightarrow \mathbb{R}^d$ is a neural network with weights θ
- σ is a sigmoid function, such as

$$\sigma(t) = \frac{1}{1 + \exp(-t)}$$

DLVM Example: Continuous Outputs

- Suppose $\mathbf{x} \in \mathbb{R}^d$.
- DLVM:

$$\begin{aligned} \mathbf{z} &\sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \\ \mathbf{x} | \mathbf{z} &\sim \mathcal{N}(\mathbf{d}_{\boldsymbol{\theta}}(\mathbf{z}), c\mathbf{I}) \end{aligned}$$

where

- $\mathbf{d}_{\boldsymbol{\theta}} : \mathbb{R}^k \rightarrow \mathbb{R}^d$ is a neural network with weights $\boldsymbol{\theta}$
- c is a hyperparameter.

Parameter Estimation

- It is common to refer to $p(\mathbf{z})$ as the *prior* and $p(\mathbf{z}|\mathbf{x})$ as the *posterior*
- Maximization of the log-likelihood

$$\begin{aligned}\log p(\mathbf{x}; \boldsymbol{\theta}) &= \log \left(\int p(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta}) p(\mathbf{z}) d\mathbf{z} \right) \\ &= \log \left(\mathbb{E}_{p(\mathbf{z})} [p(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta})] \right)\end{aligned}$$

is intractable:

- No analytic formula (because of the neural network)
- No viable stochastic estimate of gradient (because of the log)
- For this reason, the posterior

$$p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}) = \frac{p(\mathbf{z})p(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta})}{p(\mathbf{x}; \boldsymbol{\theta})}$$

is also intractable.

Variational Approximation

- To facilitate an approximation of the posterior, introduce a parametric family of distributions

$$q(\mathbf{z}|\mathbf{x}; \phi)$$

- The idea is to approximate $p(\mathbf{z}|\mathbf{x}; \theta)$ with a member of this family:

$$q(\mathbf{z}|\mathbf{x}; \phi) \approx p(\mathbf{z}|\mathbf{x}; \theta).$$

- We will also model $q(\mathbf{z}|\mathbf{x}; \phi)$ with a neural network. For example:

$$\begin{aligned} (\mu, \nu) &= \mathbf{e}_\phi(\mathbf{x}) \\ \sigma &= \exp(\nu) \quad (\text{elementwise}) \\ q(\mathbf{z}|\mathbf{x}; \phi) &= \mathcal{N}(\mathbf{z}; \mu, \text{diag}(\sigma)) \end{aligned}$$

where $\mathbf{e}_\phi : \mathbb{R}^d \rightarrow \mathbb{R}^k$ is a neural network with weights ϕ

Variational Approximation

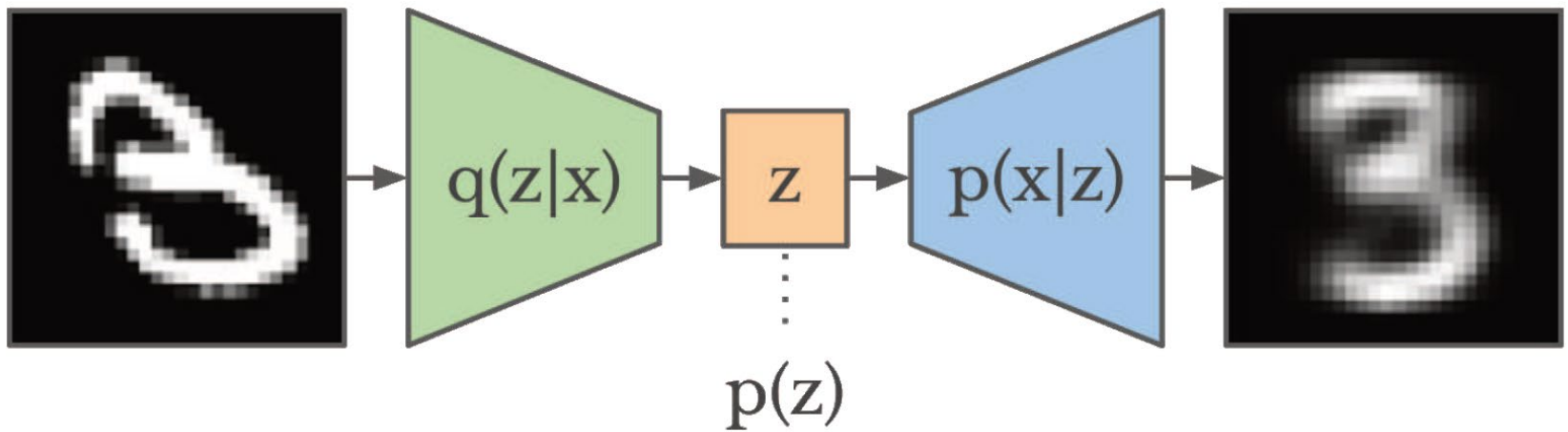


Figure credit: Murphy, *Probabilistic Machine Learning: Advanced Topics*

Variational Approximation

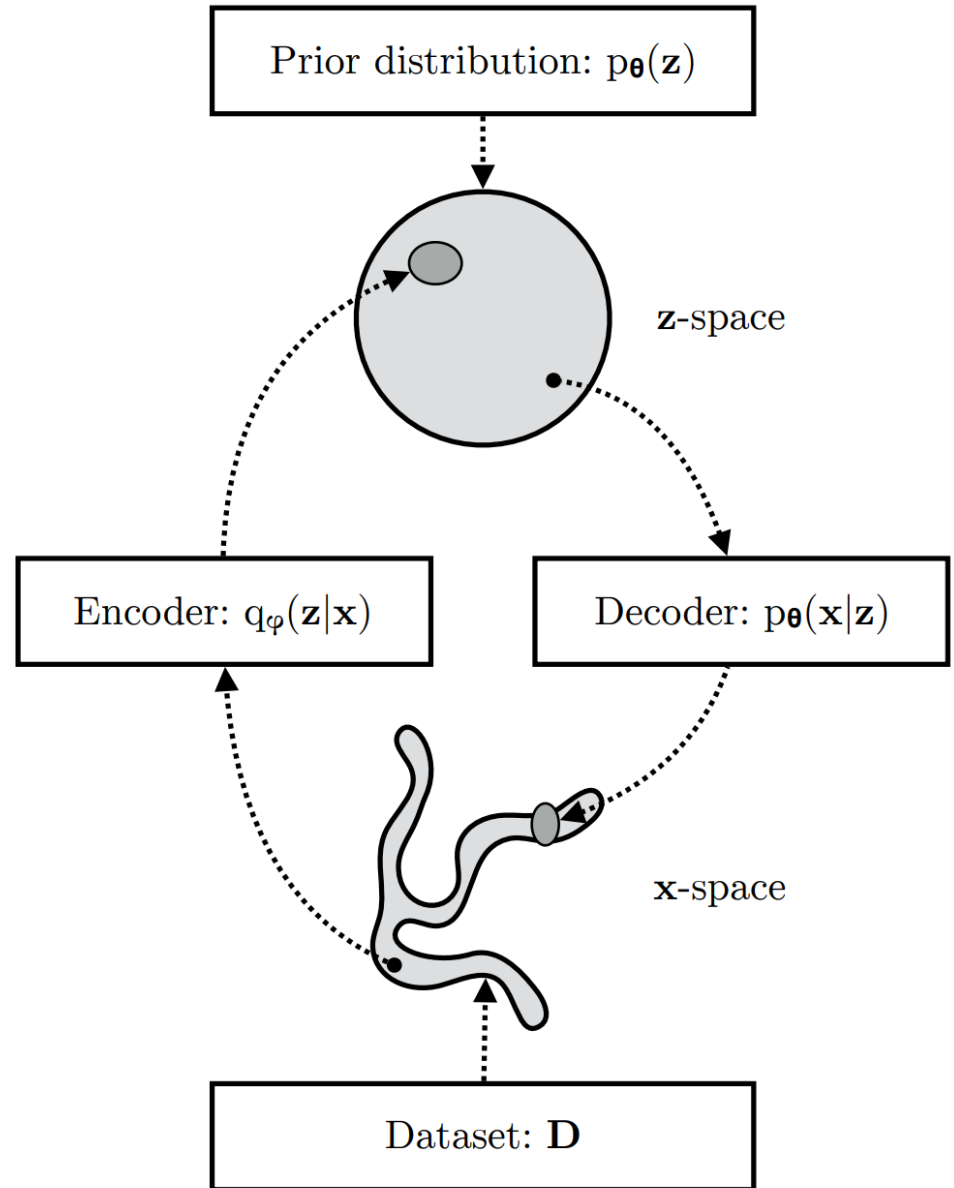


Figure credit: Kingma and Welling,
An Introduction to Variational
Autoencoders

Variational Inference

- So why did we introduce the approximate distribution $q(\mathbf{z}|\mathbf{x}; \phi)$?
 - It is tractable (we know how to calculate it)
 - Thanks to the encoder network, it can achieve a good approximation of the posterior
 - It can be leveraged in an inference (parameter estimation) strategy known as *variational inference*

Evidence Lower Bound (ELBO)

$$\begin{aligned}\log p(\mathbf{x}; \boldsymbol{\theta}) &= \mathbb{E}_{q(\mathbf{z}|\mathbf{x}; \phi)} [\log p(\mathbf{x}; \boldsymbol{\theta})] \\&= \mathbb{E}_{q(\mathbf{z}|\mathbf{x}; \phi)} \left[\log \left[\frac{p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})}{p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta})} \right] \right] \\&= \mathbb{E}_{q(\mathbf{z}|\mathbf{x}; \phi)} \left[\log \left[\frac{p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})}{q(\mathbf{z}|\mathbf{x}; \phi)} \frac{q(\mathbf{z}|\mathbf{x}; \phi)}{p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta})} \right] \right] \\&= \underbrace{\mathbb{E}_{q(\mathbf{z}|\mathbf{x}; \phi)} \left[\log \left[\frac{p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})}{q(\mathbf{z}|\mathbf{x}; \phi)} \right] \right]}_{\text{ELBO}(\mathbf{x}; \boldsymbol{\theta}, \phi)} + \underbrace{\mathbb{E}_{q(\mathbf{z}|\mathbf{x}; \phi)} \left[\log \left[\frac{q(\mathbf{z}|\mathbf{x}; \phi)}{p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta})} \right] \right]}_{D_{\text{KL}}(q(\mathbf{z}|\mathbf{x}; \phi) \| p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}))} \\&\quad \geq 0\end{aligned}$$

Lower bound on log-likelihood:

$$l(\boldsymbol{\theta}; \underline{\mathbf{x}}) = \sum_{i=1}^n \log p(\mathbf{x}_i; \boldsymbol{\theta}) \geq \sum_{i=1}^n \text{ELBO}(\mathbf{x}_i; \boldsymbol{\theta}, \phi) = \underline{\ell}(\boldsymbol{\theta}, \phi; \underline{\mathbf{x}})$$

$\underline{\mathbf{x}} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$

Variational Inference

- Basic idea is to maximize the lower bound with respect to both θ and ϕ using stochastic gradient ascent
- Need stochastic estimates of gradient wrt both θ and ϕ
- For θ , we have

$$\begin{aligned}\nabla_{\theta} \text{ELBO}(\theta, \phi; \mathbf{x}) &= \nabla_{\theta} \mathbb{E}_{q(\mathbf{z}|\mathbf{x};\phi)} [\log p(\mathbf{x}, \mathbf{z}; \theta) - \log q(\mathbf{z}|\mathbf{x}; \phi)] \\ &= \mathbb{E}_{q(\mathbf{z}|\mathbf{x};\phi)} [\nabla_{\theta} (\log p(\mathbf{x}, \mathbf{z}; \theta) - \log q(\mathbf{z}|\mathbf{x}; \phi))] \\ &= \mathbb{E}_{q(\mathbf{z}|\mathbf{x};\phi)} [\nabla_{\theta} (\log p(\mathbf{x}, \mathbf{z}; \theta))] \\ &\approx \nabla \log p(\mathbf{x}, \mathbf{z}; \theta)\end{aligned}$$

$\approx \frac{1}{n} \nabla_{\theta} (\log p(\mathbf{x}, \mathbf{z}_i; \theta))$

where in the last line, \mathbf{z} is a realization of

$$\mathbf{z} \sim \mathbb{E}_{q(\mathbf{z}|\mathbf{x};\phi)}$$

$\mathbf{z}_1, \dots, \mathbf{z}_n \sim p(\mathbf{z})$

Variational Inference

- For ϕ , let's reparametrize the probability model

$$\begin{aligned} (\boldsymbol{\mu}, \log \boldsymbol{\sigma}) &= \mathbf{e}_\phi(\mathbf{x}) \\ q(\mathbf{z}|\mathbf{x}; \phi) &= \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma})) \end{aligned} \iff \begin{aligned} \boldsymbol{\epsilon} &\sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \\ (\boldsymbol{\mu}, \log \boldsymbol{\sigma}) &= \mathbf{e}_\phi(\mathbf{x}) \\ \mathbf{z} &= \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{\epsilon} \end{aligned}$$

- Then

$$\begin{aligned} \nabla_\phi \text{ELBO}(\boldsymbol{\theta}, \phi; \mathbf{x}) &= \nabla_\phi \mathbb{E}_{q(\mathbf{z}|\mathbf{x}; \phi)} [\log p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) - \log q(\mathbf{z}|\mathbf{x}; \phi)] \\ &= \nabla_\phi \mathbb{E}_{p(\boldsymbol{\epsilon})} [\log p(\mathbf{x}, \mathbf{z}(\boldsymbol{\epsilon}; \phi); \boldsymbol{\theta}) - \log q(\mathbf{z}|\mathbf{x}; \phi)] \\ &= \mathbb{E}_{p(\boldsymbol{\epsilon})} [\nabla_\phi (\log p(\mathbf{x}, \mathbf{z}(\boldsymbol{\epsilon}; \phi); \boldsymbol{\theta}) - \log q(\mathbf{z}(\boldsymbol{\epsilon}; \phi)|\mathbf{x}; \phi))] \\ &\approx [\nabla_\phi (\log p(\mathbf{x}, \mathbf{z}(\boldsymbol{\epsilon}; \phi); \boldsymbol{\theta}) - \log q(\mathbf{z}(\boldsymbol{\epsilon}; \phi)|\mathbf{x}; \phi))] \end{aligned}$$

where in the last line, $\boldsymbol{\epsilon}$ is a realization $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

- For additional details, see <https://arxiv.org/pdf/1906.02691.pdf>

Poll

True or False: The EM algorithm is another algorithm to maximize the ELBO

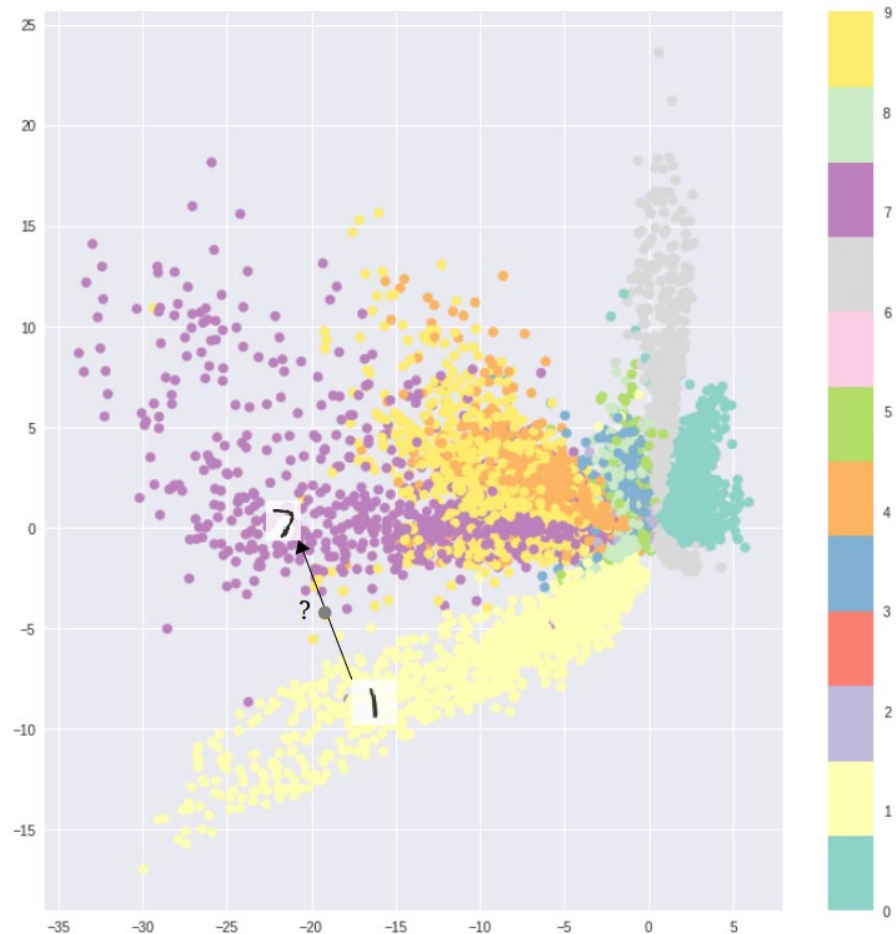
(A) True

(B) False

$$\text{ELBO}(x; \theta, \phi) = \mathbb{E}_{q(z|x; \phi)} \left[\log \left(\frac{p(x, z; \theta)}{q(z|x; \phi)} \right) \right]$$

Variational Autoencoders

- Given: Examples $\mathbf{x}_1, \dots, \mathbf{x}_n$ of a probability distribution
- Goal: Simulate new instances from the distribution
- First try: Use the decoder half of an autoencoder with random inputs
- Issue: unrealistic synthetic images



Variational Autoencoders

- VAEs generate realistic images while the latent variable is varied continuously

