

# Advanced Machine Learning Summer 2019

## Part 18 – Variational Autoencoders 03.07.2019

Prof. Dr. Bastian Leibe

RWTH Aachen University, Computer Vision Group

<http://www.vision.rwth-aachen.de>

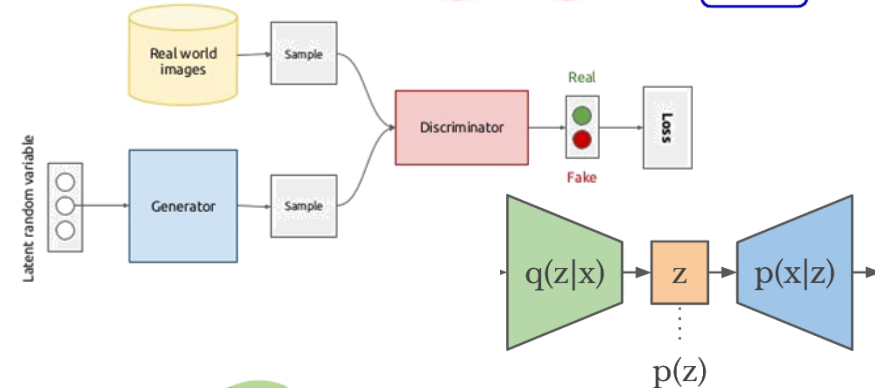
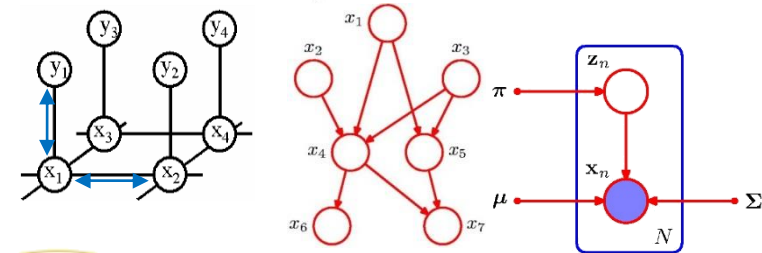
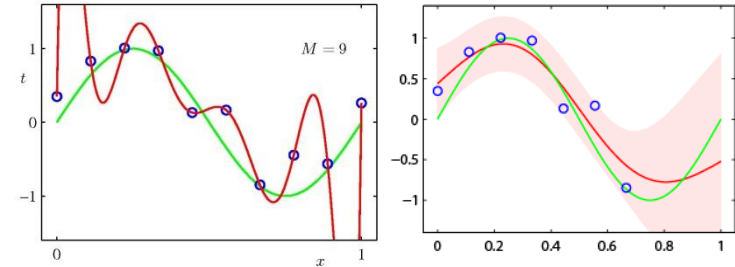


**RWTHAACHEN**  
UNIVERSITY

# Course Outline

- Regression Techniques
  - Linear Regression
  - Regularization (Ridge, Lasso)
  - Kernels (Kernel Ridge Regression)
- Deep Reinforcement Learning
- Probabilistic Graphical Models
  - Bayesian Networks
  - Markov Random Fields
  - Inference (exact & approximate)
  - Latent Variable Models
- **Deep Generative Models**
  - Generative Adversarial Networks
  - **Variational Autoencoders**

$$f : \mathcal{X} \rightarrow \mathbb{R}$$

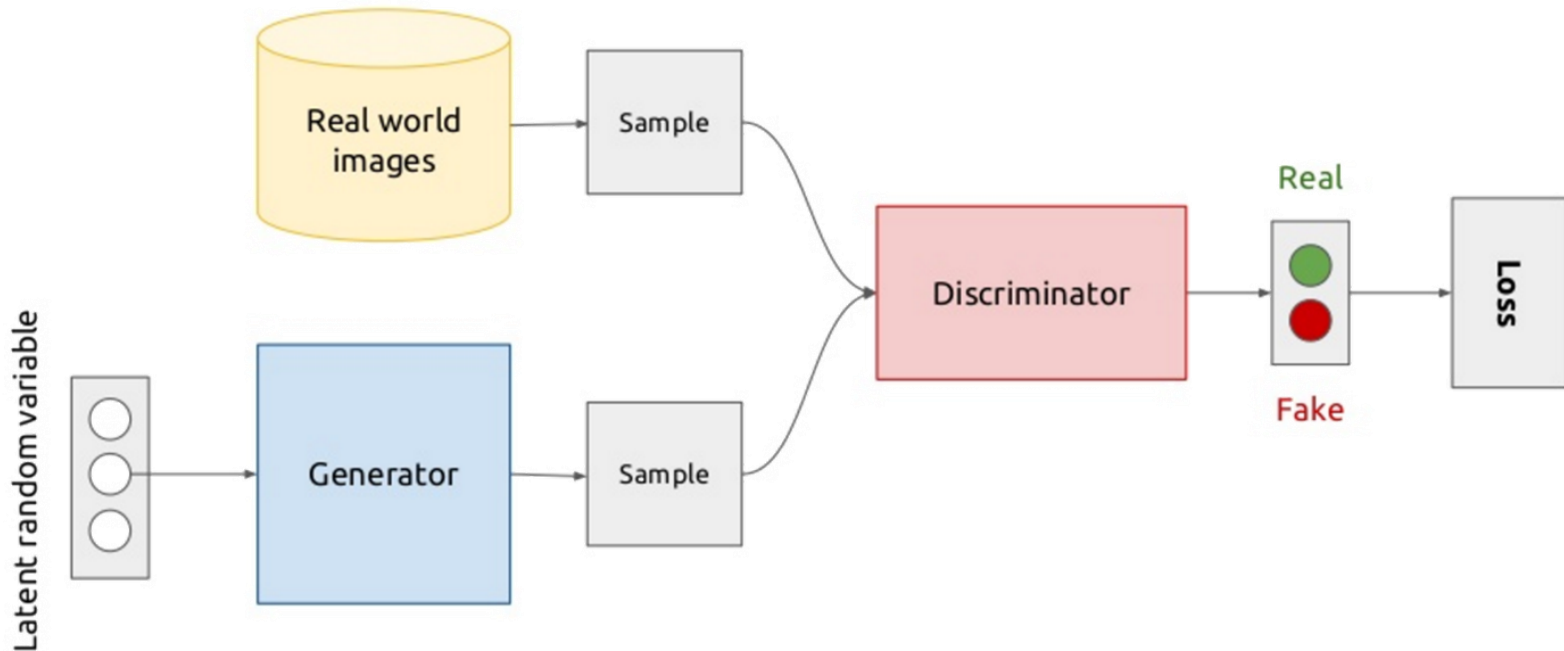


# Topics of This Lecture

- **Recap: GANs**
- **Autoencoders**
  - Motivation
  - Regularized Autoencoder
  - Denoising Autoencoder
- **Variational Autoencoders (VAE)**
  - Autoencoders as Generative Models
  - Intractability
  - Variational Approximation
  - Evidence Lower Bound (ELBO)
- **Application Examples**

# Recap: Generative Adversarial Networks (GANs)

- Conceptual view



- Main idea

- Simultaneously train an image **generator  $G$**  and a **discriminator  $D$** .
- Interpreted as a two-player game

# Recap: GAN Loss Function

- This corresponds to a two-player minimax game:

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log (1 - D(G(\mathbf{z})))]$$

- Explanation

- Train  $D$  to maximize the probability of assigning the correct label to both **training examples** and **samples from  $G$** .
- Simultaneously train  $G$  to minimize  $\log(1 - D(G(\mathbf{z})))$ .

- The Nash equilibrium of this game is achieved at

- $p_g(\mathbf{x}) = p_{data}(\mathbf{x}) \quad \forall \mathbf{x}$
- $D(\mathbf{x}) = \frac{1}{2} \quad \forall \mathbf{x}$

# GAN Algorithm

**Algorithm 1** Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator,  $k$ , is a hyperparameter. We used  $k = 1$ , the least expensive option, in our experiments.

**for** number of training iterations **do**

**for**  $k$  steps **do**

- Sample minibatch of  $m$  noise samples  $\{z^{(1)}, \dots, z^{(m)}\}$  from noise prior  $p_g(z)$ .
- Sample minibatch of  $m$  examples  $\{x^{(1)}, \dots, x^{(m)}\}$  from data generating distribution  $p_{\text{data}}(x)$ .
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[ \log D(x^{(i)}) + \log \left( 1 - D(G(z^{(i)})) \right) \right].$$

**Discriminator updates**

**end for**

- Sample minibatch of  $m$  noise samples  $\{z^{(1)}, \dots, z^{(m)}\}$  from noise prior  $p_g(z)$ .
- Update the generator by descending its stochastic gradient:

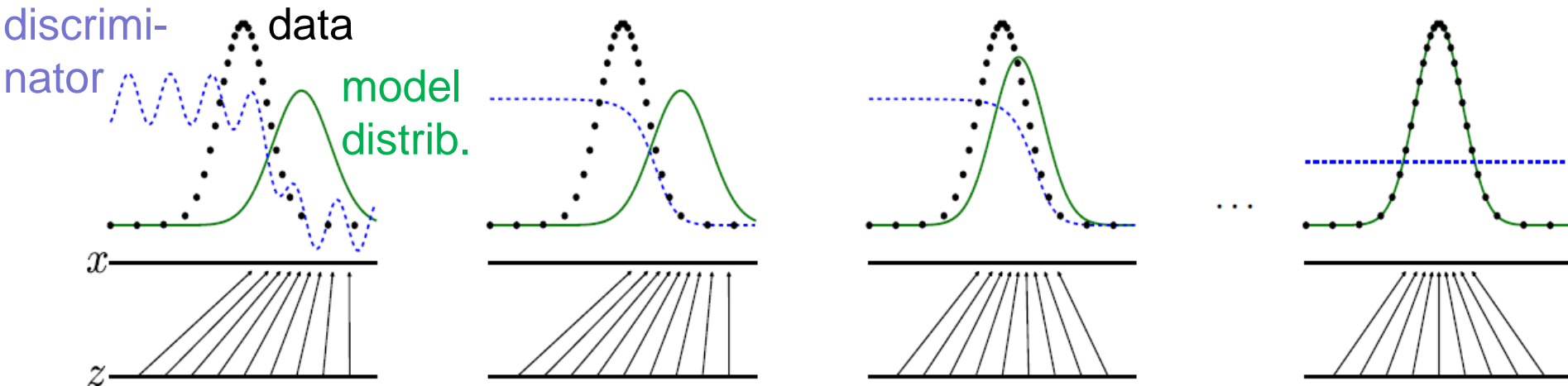
$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log \left( 1 - D(G(z^{(i)})) \right).$$

**Generator updates**

**end for**

The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

# Recap: Intuition



- Behavior near convergence

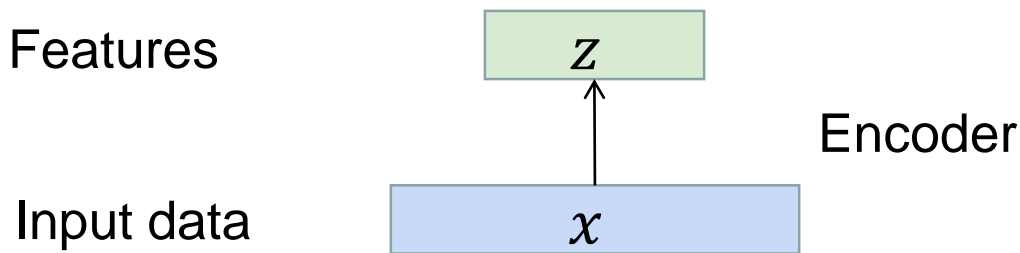
- In the inner loop,  $D$  is trained to discriminate samples from data.
- Gradient of  $D$  guides  $G$  to flow to regions that are more likely to be classified as data.
- After several steps of training,  $G$  and  $D$  will reach a point at which they cannot further improve, because  $p_g = p_{data}$ .
- Now, the discriminator is unable to differentiate between the two distributions, i.e.,  $D(x) = 0.5$ .

# Topics of This Lecture

- Recap: GANs
- **Autoencoders**
  - Motivation
  - Regularized Autoencoder
  - Denoising Autoencoder
- Variational Autoencoders (VAE)
  - Autoencoders as Generative Models
  - Intractability
  - Variational Approximation
  - Evidence Lower Bound (ELBO)
- Application Examples



# Autoencoders

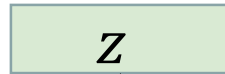


- **Autoencoders**

- Unsupervised learning approach for learning a lower-dimensional feature representation  $\mathbf{z}$  from unlabeled input data  $\mathbf{x}$ .
- $\mathbf{z}$  usually smaller than  $\mathbf{x}$  (dimensionality reduction)
- Want to capture meaningful factors of variation in the data

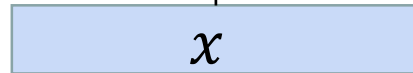
# Autoencoders

Features



Encoder

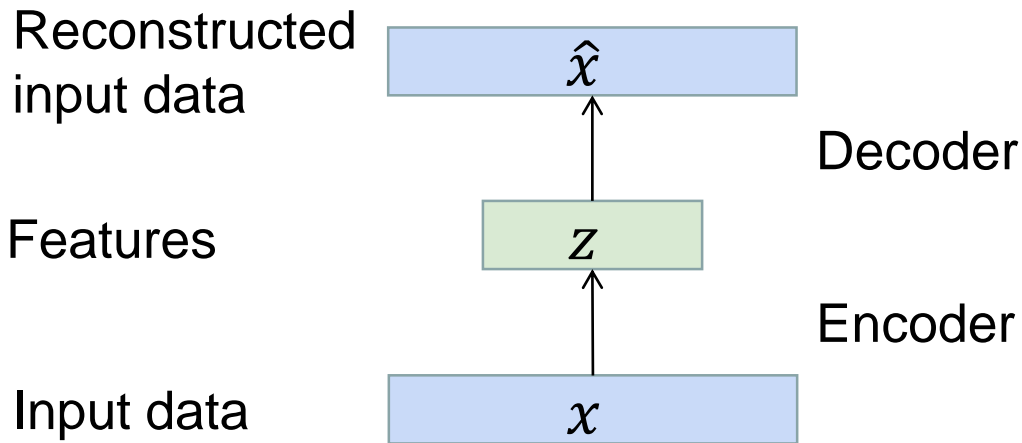
Input data



- **Encoder**

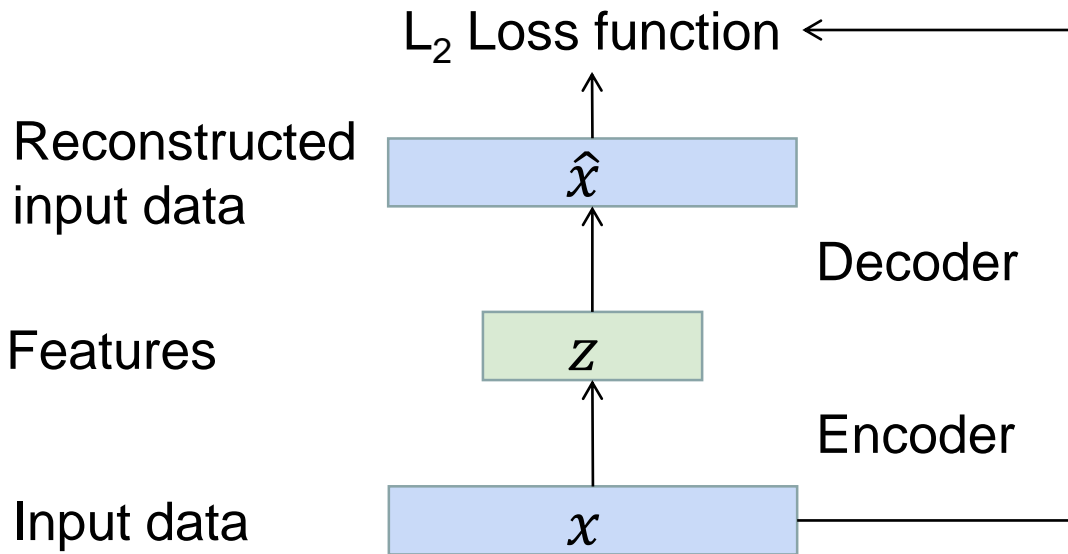
- Originally: shallow function (linear + sigmoid)
- Later: Deep, fully-connected
- Later: ReLU CNN

# Autoencoders



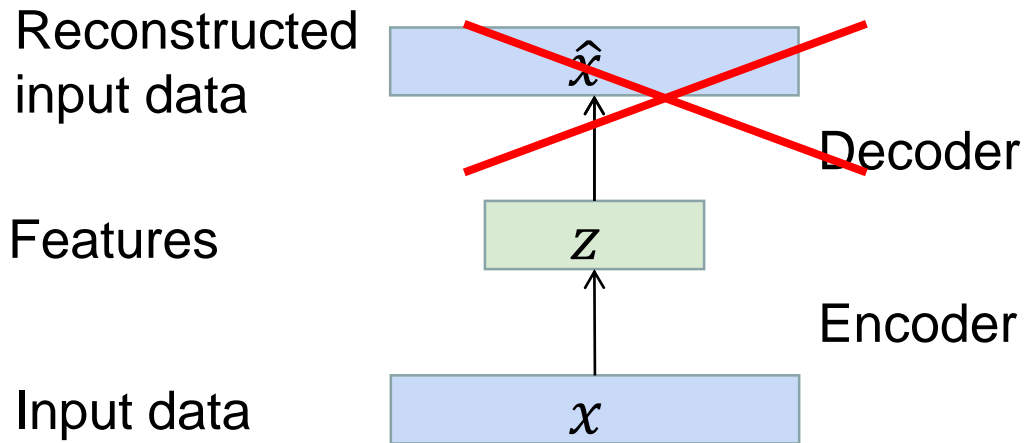
- How to learn such a feature representation?
  - Train such that features can be used to reconstruct original data.
  - “Autoencoding” – encoding itself

# Autoencoders



- How to learn such a feature representation?
  - Train such that features can be used to reconstruct original data.
  - “Autoencoding” – encoding itself
  - L<sub>2</sub> loss function  $\|x - \hat{x}\|^2$
  - Note: this doesn't use any labels!

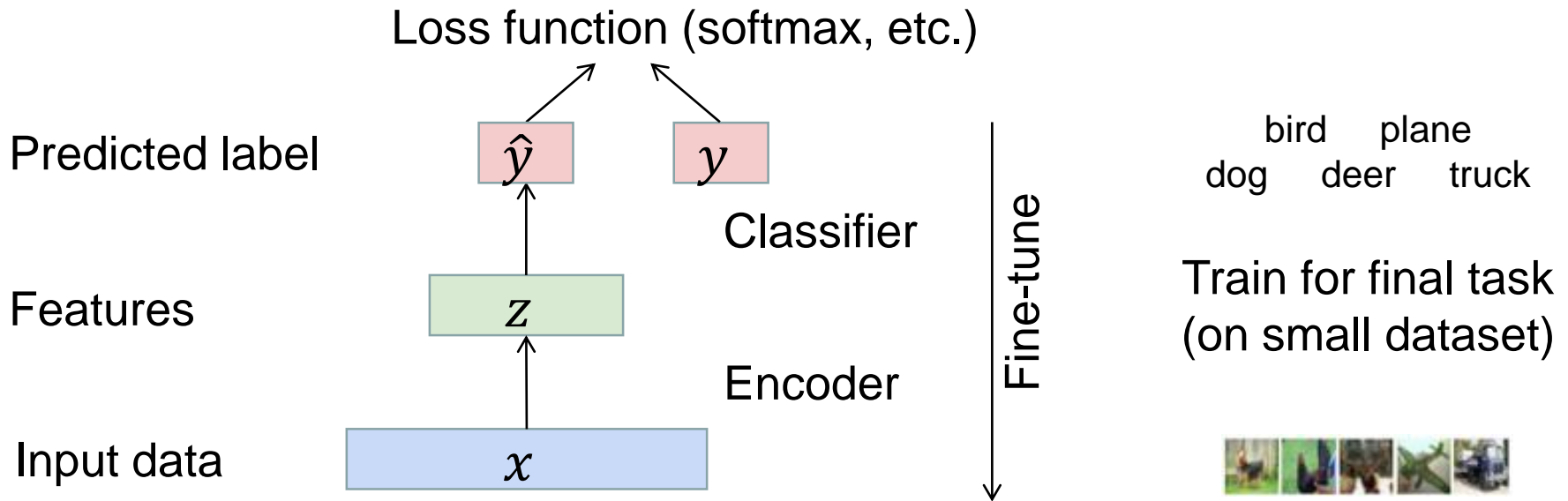
# Autoencoders



- After training
  - Throw away the decoder part



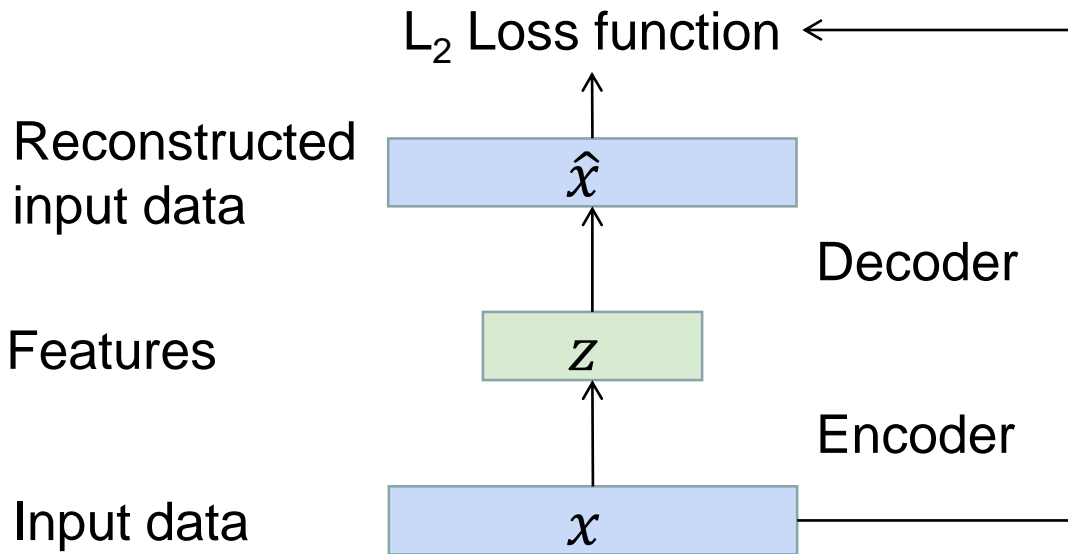
# Autoencoders



- After training

- Throw away the decoder part
- Encoder can be used to initialize a supervised model
- Fine-tune encoder jointly with supervised model
- *Idea used in the 90s and early 2000s to pre-train deeper models*

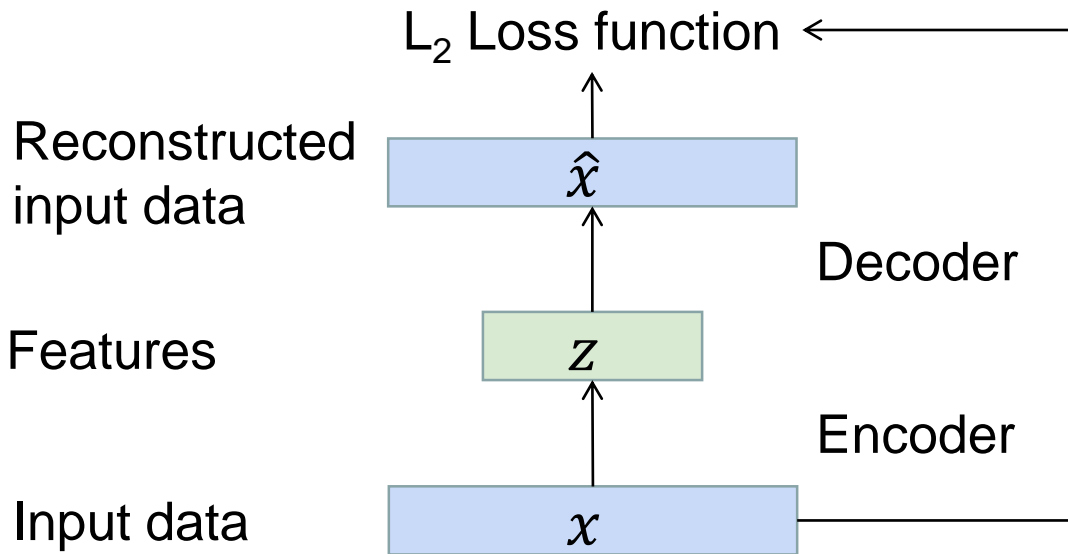
# Variants of Autoencoders



- Analyzing the learning process

- Learning process minimizes a loss function  $L(\mathbf{x}, g(f(\mathbf{x})))$
- Linear decoder +  $L_2$  loss: Autoencoder learns **PCA** subspace
- Autoencoders with nonlinear encoder and decoder functions thus learn a more powerful **nonlinear generalization of PCA**.

# Variants of Autoencoders

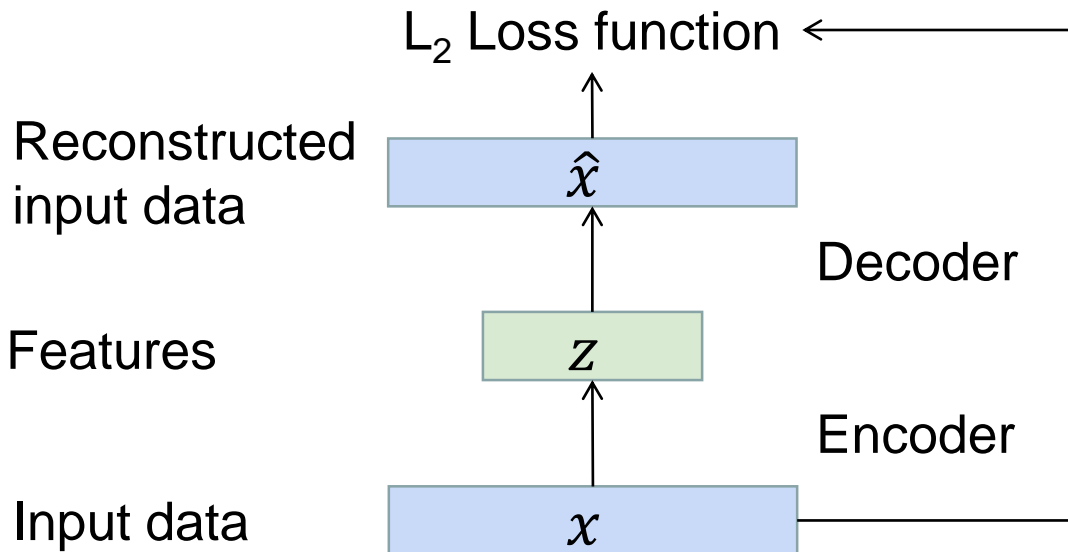


- Analyzing the learning process

- Learning process minimizes a loss function  $L(\mathbf{x}, g(f(\mathbf{x})))$
- *Unfortunately, if the encoder and decoder are too powerful, they can learn to perform the copying task without learning useful information.*
- E.g., learn a 1D code to memorize each training example.



# Variants of Autoencoders

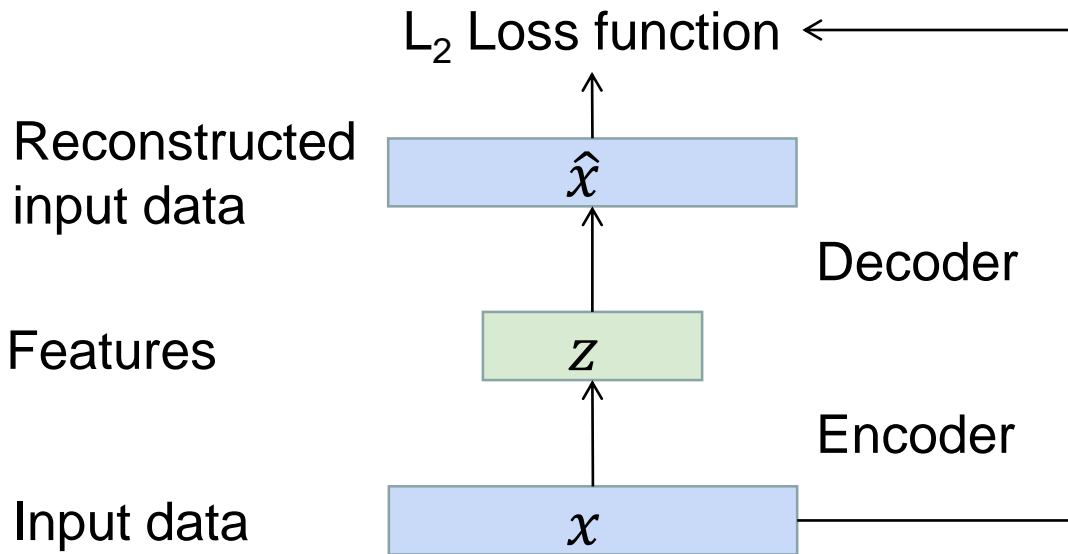


- **Regularized Autoencoders**

- Include a regularization term to the loss function:  $L(\mathbf{x}, g(f(\mathbf{x}))) + \Omega(\mathbf{z})$

- E.g., enforce sparsity by an L<sub>1</sub> regularizer  $\Omega(\mathbf{z}) = \lambda \sum_i |z_i|$

# Variants of Autoencoders



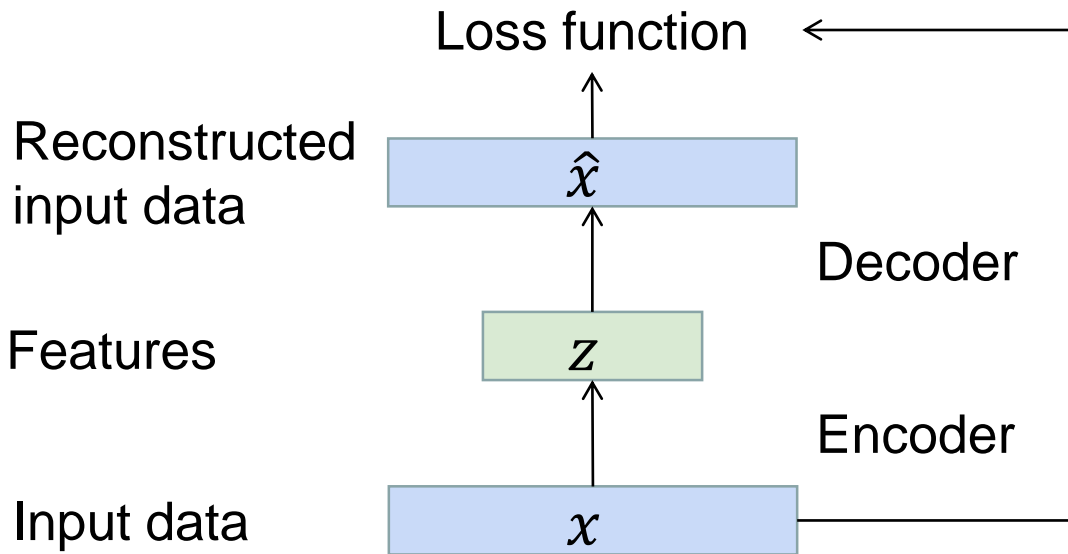
- Regularized Autoencoders

- We can think of the sparse encoder framework as approximating ML training of a generative model with latent variables  $z$ .

$$\log p_{model}(\mathbf{x}) = \log \sum_{\mathbf{z}} p_{model}(\mathbf{x}, \mathbf{z})$$

- The autoencoder approximates this sum with a point estimate for just one highly likely value for  $z$ .

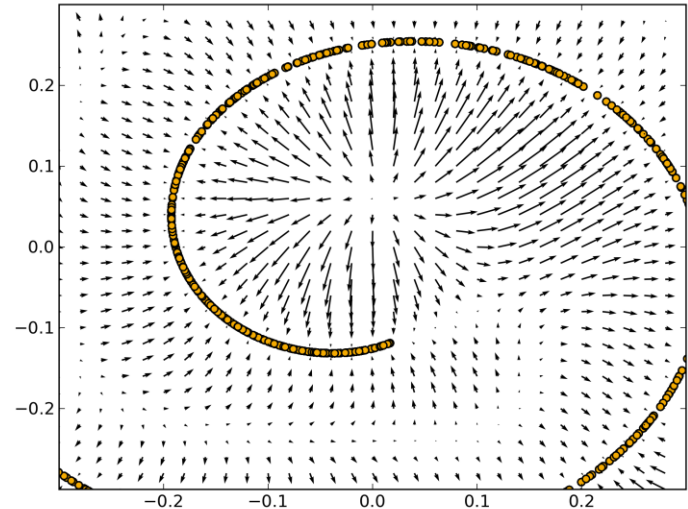
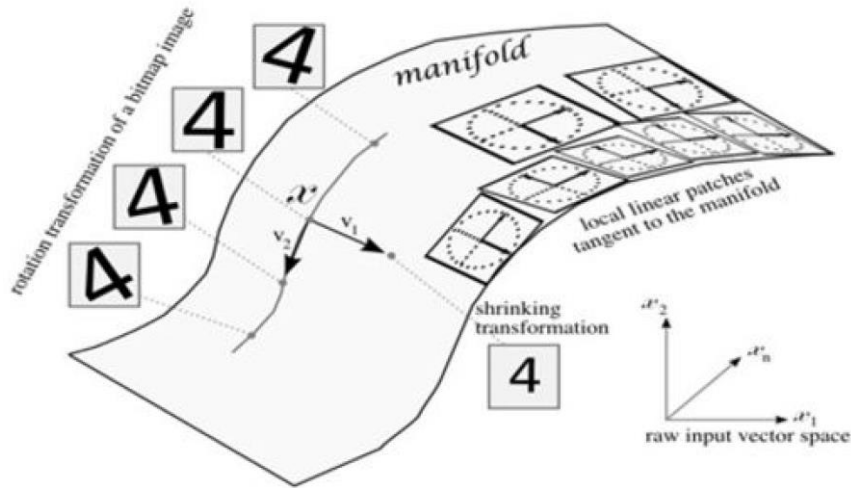
# Variants of Autoencoders



- **Denoising Autoencoder (DAE)**

- Rather than the reconstruction loss, minimize  $L(\mathbf{x}, g(f(\tilde{\mathbf{x}})))$  where  $\tilde{\mathbf{x}}$  is a copy of  $\mathbf{x}$  that has been corrupted by some noise.
- Denoising forces  $f$  and  $g$  to implicitly learn the structure of  $p_{data}(\mathbf{x})$ .

# Variants of Autoencoders



- **Denoising Autoencoder (DAE)**

- Assumption: Natural data actually lies in a (low-dimensional) manifold of the high-dimensional space of input data  $\mathbf{x}$ .
- By corrupting the input data with noise, we force the DAE to learn a vector field that pushes towards this low-dimensional manifold.

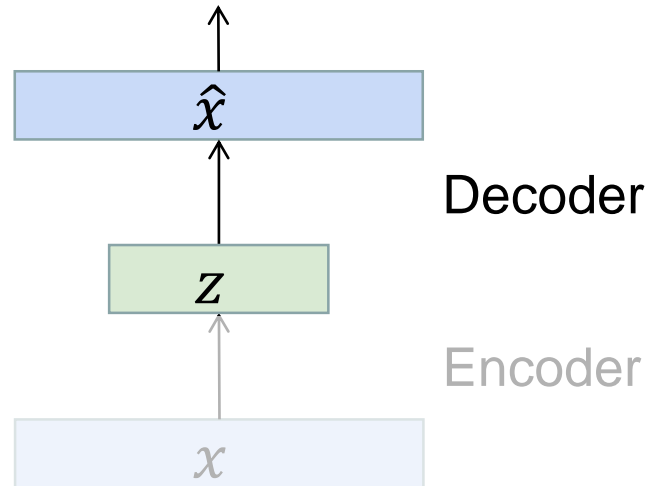
# Topics of This Lecture

- Recap: GANs
- Autoencoders
  - Motivation
  - Regularized Autoencoder
  - Denoising Autoencoder
- **Variational Autoencoders (VAE)**
  - Autoencoders as Generative Models
  - Intractability
  - Variational Approximation
  - Evidence Lower Bound (ELBO)
- Application Examples

# Autoencoders as Data Generators

- Autoencoders

- Can reconstruct data and can learn features to initialize a supervised model
- Features capture factors of variation in training data
- Can we generate new images from an autoencoder?



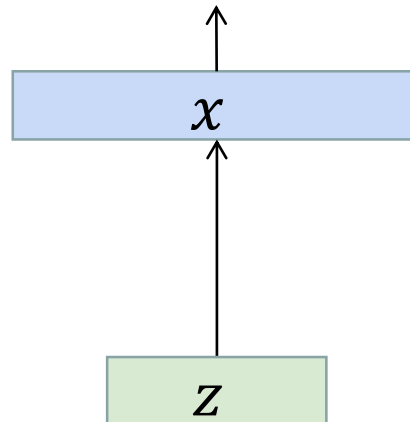
- For this we need to generate samples from the data manifold. How?

# Probabilistic Spin on Autoencoders

- Idea: Sample the model to generate data
  - Assume training data  $\{\mathbf{x}^{(i)}\}_{i=1}^N$  is generated from underlying latent representation  $\mathbf{z}$ .

Sample from  
true conditional  
 $p_{\theta^*}(\mathbf{x}|\mathbf{z}^{(i)})$

Sample from  
true prior  
 $p_{\theta^*}(\mathbf{z})$



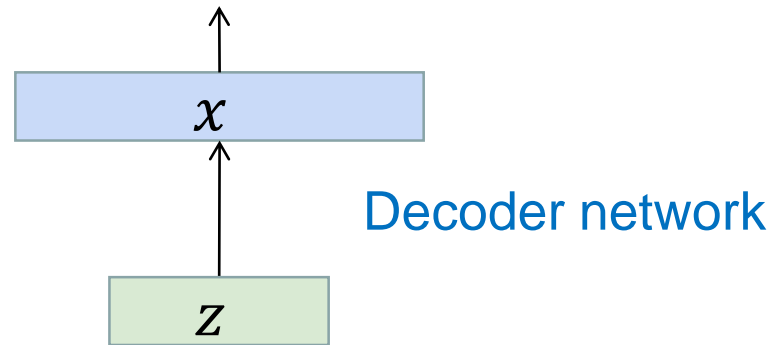
# Probabilistic Spin on Autoencoders

Sample from true conditional

$$p_{\theta^*}(\mathbf{x}|\mathbf{z}^{(i)})$$

Sample from true prior

$$p_{\theta^*}(\mathbf{z})$$



- Idea: Sample the model to generate data
  - We want to estimate the true parameters  $\theta^*$  of this generative model.
- How should we represent the model?
  - Choose prior  $p(\mathbf{z})$  to be simple, e.g., Gaussian
  - Conditional  $p(\mathbf{x} | \mathbf{z})$  is complex (generates image)
    - ⇒ Represent with neural network



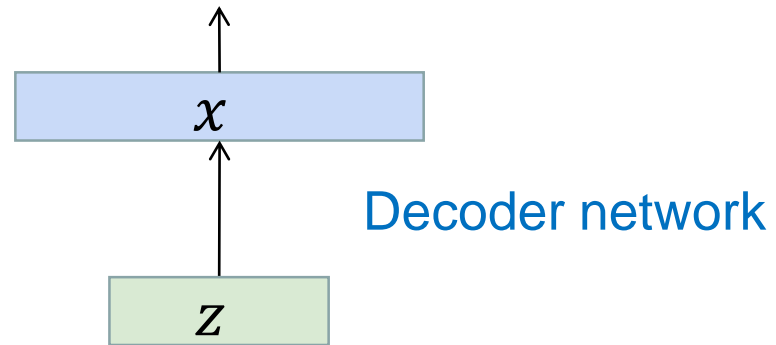
# Probabilistic Spin on Autoencoders

Sample from true conditional

$$p_{\theta^*}(\mathbf{x}|\mathbf{z}^{(i)})$$

Sample from true prior

$$p_{\theta^*}(\mathbf{z})$$



- Idea: Sample the model to generate data
  - We want to estimate the true parameters  $\theta^*$  of this generative model.
- How to train the model?
  - Learn model parameters to maximize likelihood of training data

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

– *What is the problem here?* **Intractable!**

# Variational Autoencoders: Intractability

- Computing the data likelihood

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

- $p(\mathbf{z})$  is a simple Gaussian prior.  $\Rightarrow$  ok.
- $p(\mathbf{x} | \mathbf{z})$  is a decoder Neural network.  $\Rightarrow$  ok.
- **But is is intractable to compute  $p(\mathbf{x} | \mathbf{z})$  for every  $\mathbf{z}$ !**
- Posterior density is also intractable

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

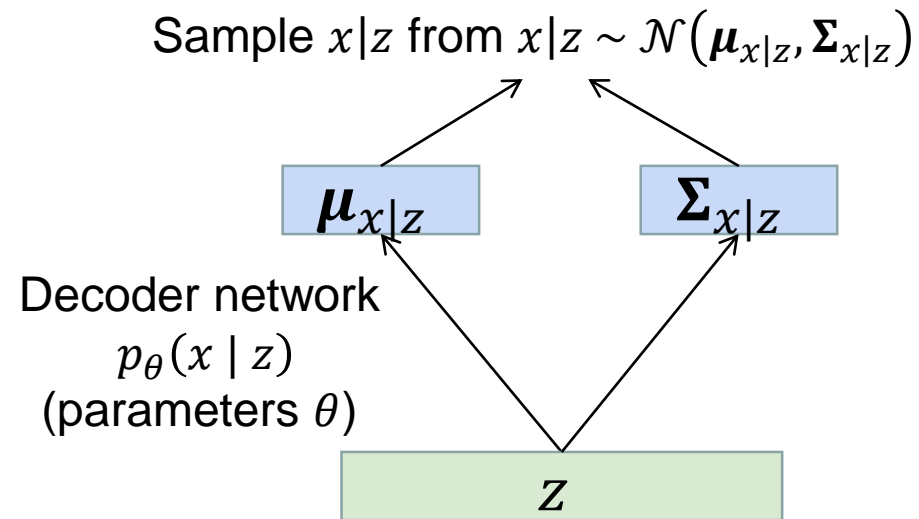
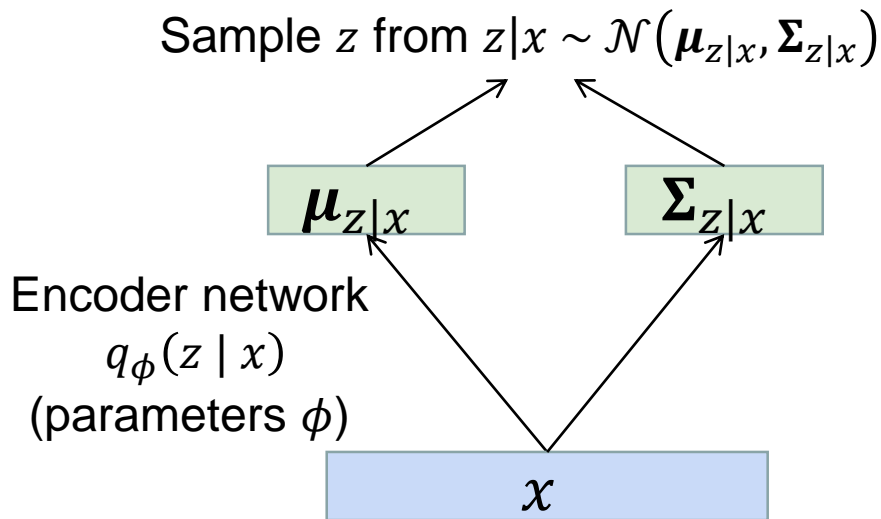
# Variational Autoencoders: Intractability

- Solution

- In addition to the decoder network modeling  $p_{\theta}(\mathbf{x} | \mathbf{z})$ , define additional encoder network modeling  $q_{\phi}(\mathbf{z} | \mathbf{x})$  that approximates  $p_{\theta}(\mathbf{z} | \mathbf{x})$ .
- *We will see that this allows us to derive a lower bound on the data likelihood that is tractable and that we can optimize.*

# Variational Autoencoders

- Since we are modelling probabilistic generation of data, encoder and decoder networks are probabilistic



- Encoder and decoder networks are also called **recognition/inference** and **generation** networks

# Variational Autoencoders

- We can now work out the log-likelihood

$$\log p_{\theta}(x^{(i)}) = \mathbb{E}_{z \sim q_{\phi}(z|x^{(i)})} [\log p_{\theta}(x^{(i)})] \quad (p_{\theta}(x^{(i)}) \text{ does not depend on } z)$$

Taking expectation w.r.t.  $z$   
(using encoder network)  
will come in handy later

# Variational Autoencoders

- We can now work out the log-likelihood

$$\begin{aligned}\log p_{\theta}(x^{(i)}) &= \mathbb{E}_{z \sim q_{\phi}(z|x^{(i)})} [\log p_{\theta}(x^{(i)})] && (p_{\theta}(x^{(i)}) \text{ does not depend on } z) \\ &= \mathbb{E}_z \left[ \log \frac{p_{\theta}(x^{(i)} | z) p_{\theta}(z)}{p_{\theta}(z | x^{(i)})} \right] && \text{(Bayes' Rule)} \\ &= \mathbb{E}_z \left[ \log \frac{p_{\theta}(x^{(i)} | z) p_{\theta}(z)}{p_{\theta}(z | x^{(i)})} \frac{q_{\phi}(z | x^{(i)})}{q_{\phi}(z | x^{(i)})} \right] && \text{(Multiply by constant)} \\ &= \mathbb{E}_z [\log p_{\theta}(x^{(i)} | z)] - \mathbb{E}_z \left[ \log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z)} \right] + \mathbb{E}_z \left[ \log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z | x^{(i)})} \right] \\ &= \mathbb{E}_z [\log p_{\theta}(x^{(i)} | z)] - D_{KL}(q_{\phi}(z | x^{(i)}) \| p_{\theta}(z)) + D_{KL}(q_{\phi}(z | x^{(i)}) \| p_{\theta}(z | x^{(i)}))\end{aligned}$$

The expectation w.r.t  $z$   
(using encoder network) lets  
us write nice KL terms

# Variational Autoencoders

- We can now work out the log-likelihood

$$\begin{aligned}\log p_{\theta}(x^{(i)}) &= \mathbb{E}_{z \sim q_{\phi}(z|x^{(i)})} [\log p_{\theta}(x^{(i)})] && (p_{\theta}(x^{(i)}) \text{ does not depend on } z) \\ &= \mathbb{E}_z \left[ \log \frac{p_{\theta}(x^{(i)} | z) p_{\theta}(z)}{p_{\theta}(z | x^{(i)})} \right] && (\text{Bayes' Rule}) \\ &= \mathbb{E}_z \left[ \log \frac{p_{\theta}(x^{(i)} | z) p_{\theta}(z)}{p_{\theta}(z | x^{(i)})} \frac{q_{\phi}(z | x^{(i)})}{q_{\phi}(z | x^{(i)})} \right] && (\text{Multiply by constant}) \\ &= \mathbb{E}_z [\log p_{\theta}(x^{(i)} | z)] - \mathbb{E}_z \left[ \log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z)} \right] + \mathbb{E}_z \left[ \log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z | x^{(i)})} \right] \\ &= \mathbb{E}_z [\log p_{\theta}(x^{(i)} | z)] - D_{KL}(q_{\phi}(z | x^{(i)}) \| p_{\theta}(z)) + D_{KL}(q_{\phi}(z | x^{(i)}) \| p_{\theta}(z | x^{(i)}))\end{aligned}$$



Decoder network gives  $p_{\theta}(x | z)$ , can compute estimate of this term through sampling.

(Sampling differentiable through reparametrization trick, see paper)

# Variational Autoencoders

- We can now work out the log-likelihood

$$\begin{aligned}\log p_{\theta}(x^{(i)}) &= \mathbb{E}_{z \sim q_{\phi}(z|x^{(i)})} [\log p_{\theta}(x^{(i)})] && (p_{\theta}(x^{(i)}) \text{ does not depend on } z) \\ &= \mathbb{E}_z \left[ \log \frac{p_{\theta}(x^{(i)} | z) p_{\theta}(z)}{p_{\theta}(z | x^{(i)})} \right] && \text{(Bayes' Rule)} \\ &= \mathbb{E}_z \left[ \log \frac{p_{\theta}(x^{(i)} | z) p_{\theta}(z)}{p_{\theta}(z | x^{(i)})} \frac{q_{\phi}(z | x^{(i)})}{q_{\phi}(z | x^{(i)})} \right] && \text{(Multiply by constant)} \\ &= \mathbb{E}_z [\log p_{\theta}(x^{(i)} | z)] - \mathbb{E}_z \left[ \log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z)} \right] + \mathbb{E}_z \left[ \log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z | x^{(i)})} \right] \\ &= \mathbb{E}_z [\log p_{\theta}(x^{(i)} | z)] - D_{KL}(q_{\phi}(z | x^{(i)}) \| p_{\theta}(z)) + D_{KL}(q_{\phi}(z | x^{(i)}) \| p_{\theta}(z | x^{(i)}))\end{aligned}$$



This KL term (between  
Gaussians for encoder/prior)  
has a nice closed-form solution



# Variational Autoencoders

- We can now work out the log-likelihood

$$\begin{aligned}\log p_{\theta}(x^{(i)}) &= \mathbb{E}_{z \sim q_{\phi}(z|x^{(i)})} [\log p_{\theta}(x^{(i)})] && (p_{\theta}(x^{(i)}) \text{ does not depend on } z) \\ &= \mathbb{E}_z \left[ \log \frac{p_{\theta}(x^{(i)} | z) p_{\theta}(z)}{p_{\theta}(z | x^{(i)})} \right] && (\text{Bayes' Rule}) \\ &= \mathbb{E}_z \left[ \log \frac{p_{\theta}(x^{(i)} | z) p_{\theta}(z)}{p_{\theta}(z | x^{(i)})} \frac{q_{\phi}(z | x^{(i)})}{q_{\phi}(z | x^{(i)})} \right] && (\text{Multiply by constant}) \\ &= \mathbb{E}_z [\log p_{\theta}(x^{(i)} | z)] - \mathbb{E}_z \left[ \log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z)} \right] + \mathbb{E}_z \left[ \log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z | x^{(i)})} \right] \\ &= \mathbb{E}_z [\log p_{\theta}(x^{(i)} | z)] - D_{KL}(q_{\phi}(z | x^{(i)}) \| p_{\theta}(z)) + D_{KL}(q_{\phi}(z | x^{(i)}) \| p_{\theta}(z | x^{(i)}))\end{aligned}$$



$p_{\theta}(z | x)$  intractable (as seen earlier),  
can't compute this KL term ☹

But we know KL divergence always  $\geq 0$ .

# Variational Autoencoders

- We can now work out the log-likelihood

Want to maximize data likelihood

$$\begin{aligned}\log p_{\theta}(x^{(i)}) &= \mathbb{E}_{z \sim q_{\phi}(z|x^{(i)})} [\log p_{\theta}(x^{(i)})] && (p_{\theta}(x^{(i)}) \text{ does not depend on } z) \\ &= \mathbb{E}_z \left[ \log \frac{p_{\theta}(x^{(i)} | z) p_{\theta}(z)}{p_{\theta}(z | x^{(i)})} \right] && (\text{Bayes' Rule}) \\ &= \mathbb{E}_z \left[ \log \frac{p_{\theta}(x^{(i)} | z) p_{\theta}(z)}{p_{\theta}(z | x^{(i)})} \frac{q_{\phi}(z | x^{(i)})}{q_{\phi}(z | x^{(i)})} \right] && (\text{Multiply by constant}) \\ &= \mathbb{E}_z [\log p_{\theta}(x^{(i)} | z)] - \mathbb{E}_z \left[ \log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z)} \right] + \mathbb{E}_z \left[ \log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z | x^{(i)})} \right] \\ &= \underbrace{\mathbb{E}_z [\log p_{\theta}(x^{(i)} | z)] - D_{KL}(q_{\phi}(z | x^{(i)}) \| p_{\theta}(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)} + \underbrace{D_{KL}(q_{\phi}(z | x^{(i)}) \| p_{\theta}(z | x^{(i)}))}_{\geq 0}\end{aligned}$$

Tractable lower bound, which we can take gradient of and optimize

# Variational Autoencoders

- Variational Lower Bound (“ELBO”)

$$\begin{aligned}\log p_{\theta}(x^{(i)}) &\geq \mathcal{L}(x^{(i)}, \theta, \phi) \\ &= \underbrace{\mathbb{E}_z[\log p_{\theta}(x^{(i)} | z)]}_{\text{“Reconstruct the input data”}} - \underbrace{D_{KL}(q_{\phi}(z | x^{(i)}) \| p_{\theta}(z))}_{\text{“Make approximate posterior distribution close to prior”}}\end{aligned}$$

- Training: Maximize lower bound

$$\theta^*, \phi^* = \arg \max_{\theta, \phi} \sum_{i=1}^N \mathcal{L}(x^{(i)}, \theta, \phi)$$

# Topics of This Lecture

- Recap: GANs
- Autoencoders
  - Motivation
  - Regularized Autoencoder
  - Denoising Autoencoder
- Variational Autoencoders (VAE)
  - Autoencoders as Generative Models
  - Intractability
  - Variational Approximation
  - Evidence Lower Bound (ELBO)
- **Application Examples**

# Application Examples



32x32 CIFAR-10



Labeled Faces in the Wild

# References

- Variational Auto-Encoders
  - D. Kingma, M. Welling, [Auto-Encoding Variational Bayes](#), ICLR 2014.