

Advanced Machine Learning Summer 2019

Part 19 – Variational Autoencoders II 10.07.2019

Prof. Dr. Bastian Leibe

RWTH Aachen University, Computer Vision Group

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

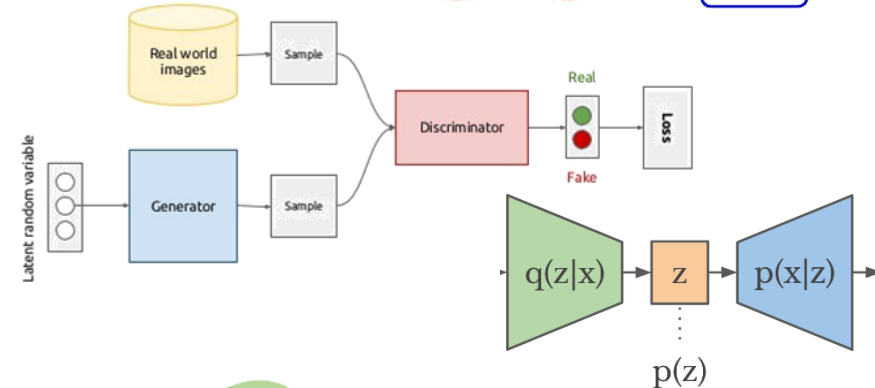
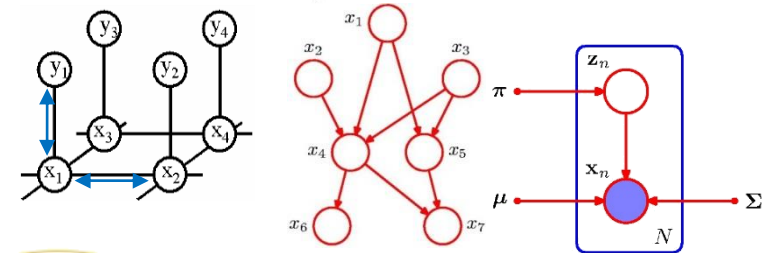
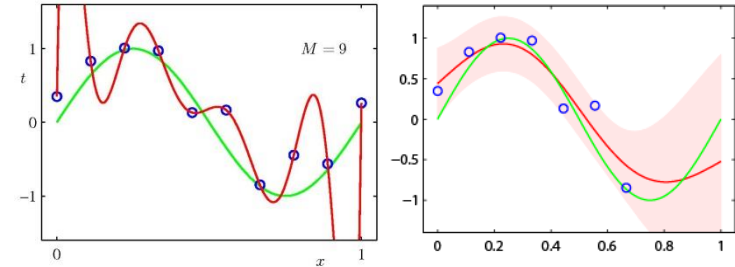


RWTHAACHEN
UNIVERSITY

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**

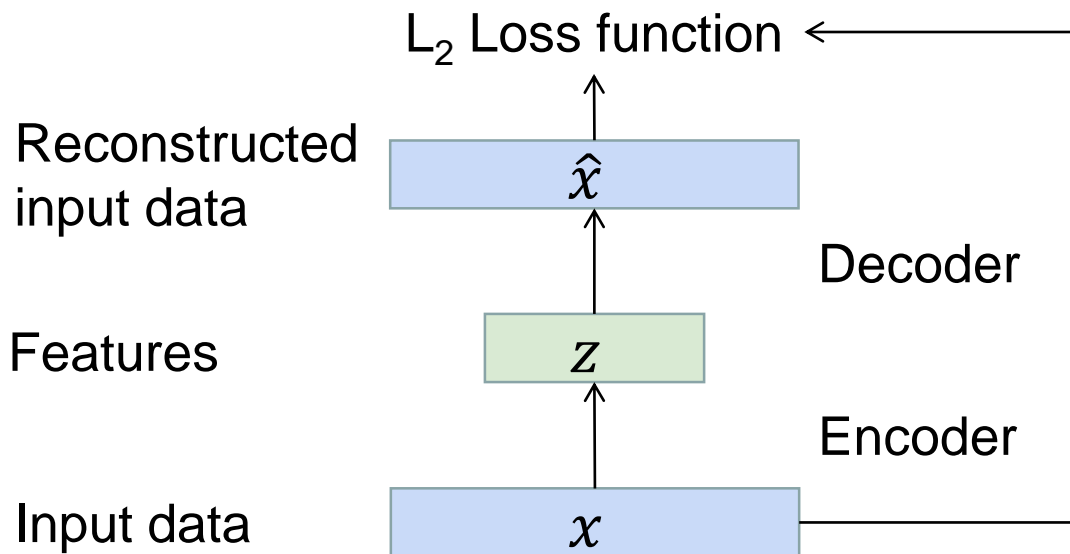
$$f : \mathcal{X} \rightarrow \mathbb{R}$$



Topics of This Lecture

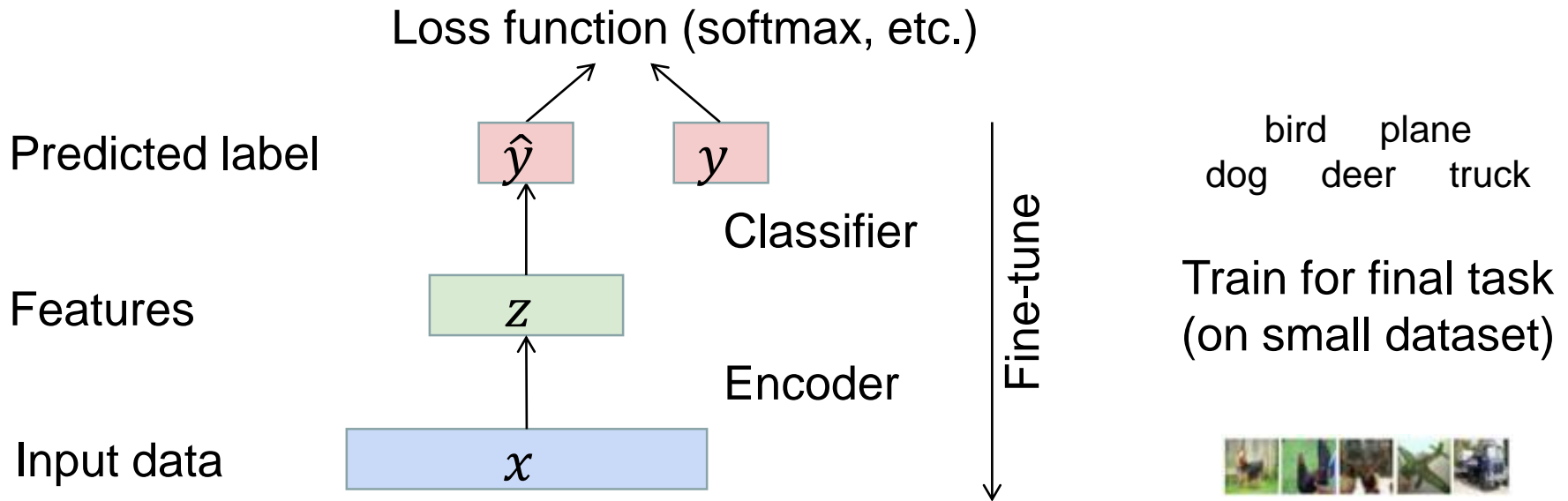
- **Recap: Variational Autoencoders**
 - Autoencoders as Generative Models
 - Intractability
 - Variational Approximation
 - Evidence Lower Bound (ELBO)
- **Applying VAEs**
 - VAE Training
 - VAE Data Generation

Recap: Autoencoders



- How to learn such a feature representation?
 - Unsupervised learning approach for learning a lower-dimensional feature representation z from unlabeled input data x .
 - z usually smaller than x (dimensionality reduction)
 - Want to capture meaningful factors of variation in the data Train such that features can be used to reconstruct original data.

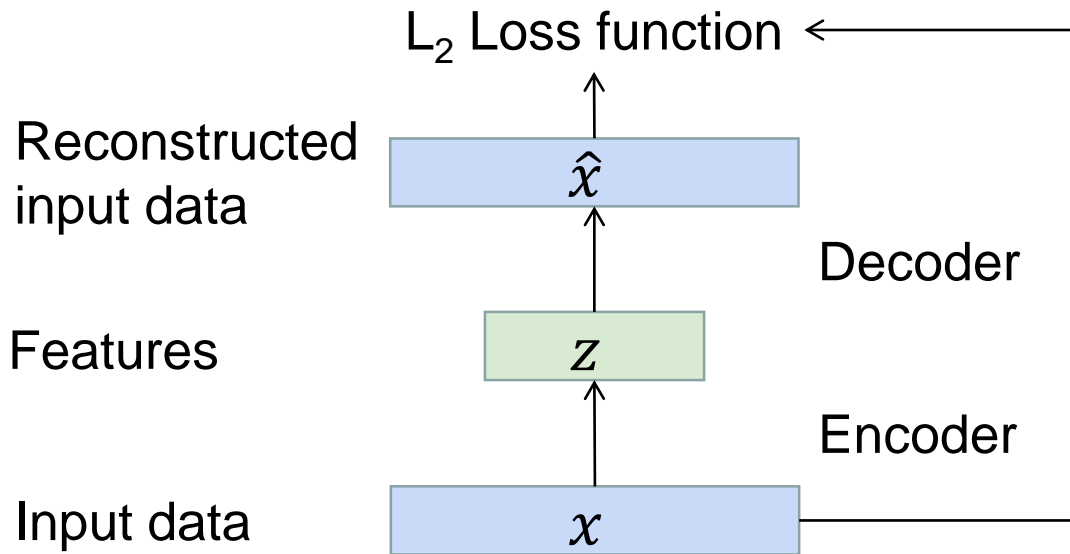
Recap: Autoencoders



- 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*

Recap: Variants of Autoencoders

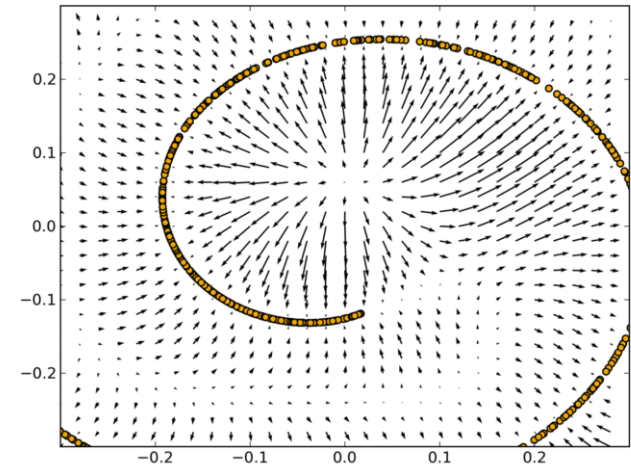
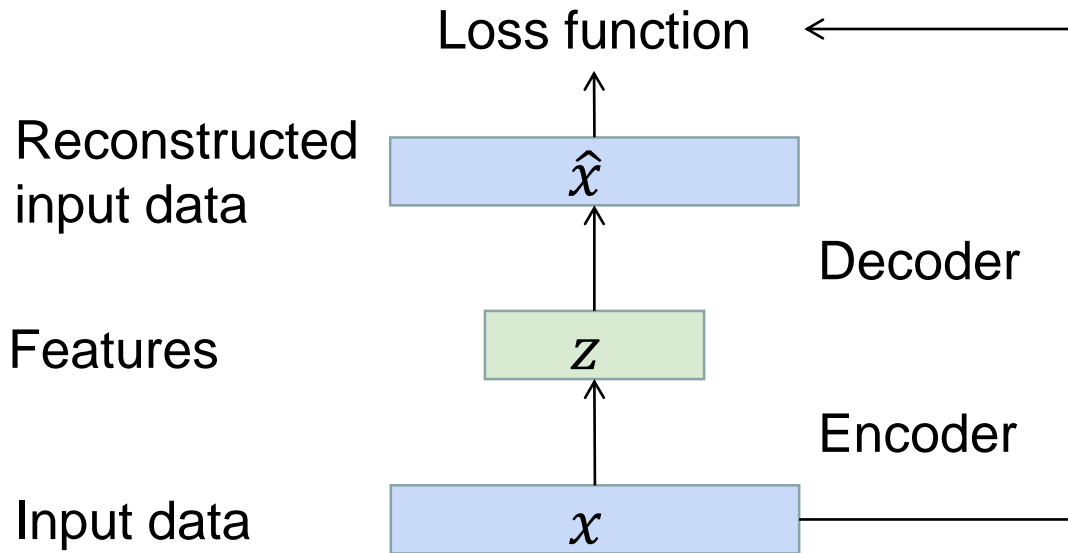


- **Regularized Autoencoders**

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

- E.g., enforce sparsity by an L₁ regularizer $\Omega(\mathbf{z}) = \lambda \sum_i |z_i|$

Recap: 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 $p_{data}(\mathbf{x})$.

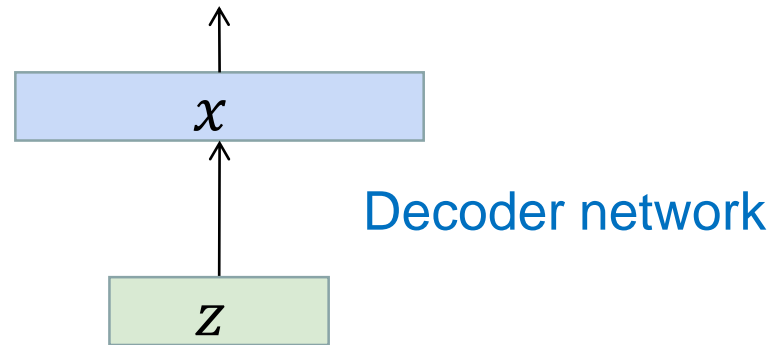
Recap: Probabilistic Spin on Autoencoders

Sample from true conditional

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

Sample from true prior

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



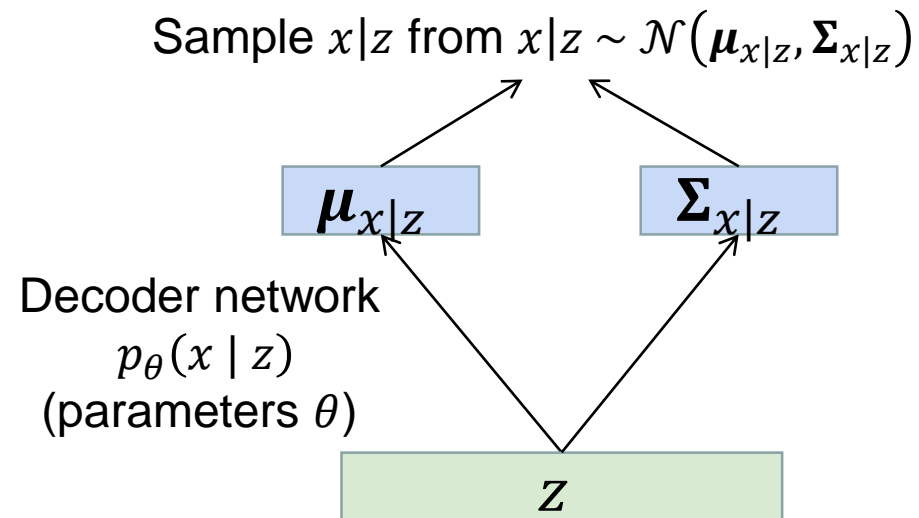
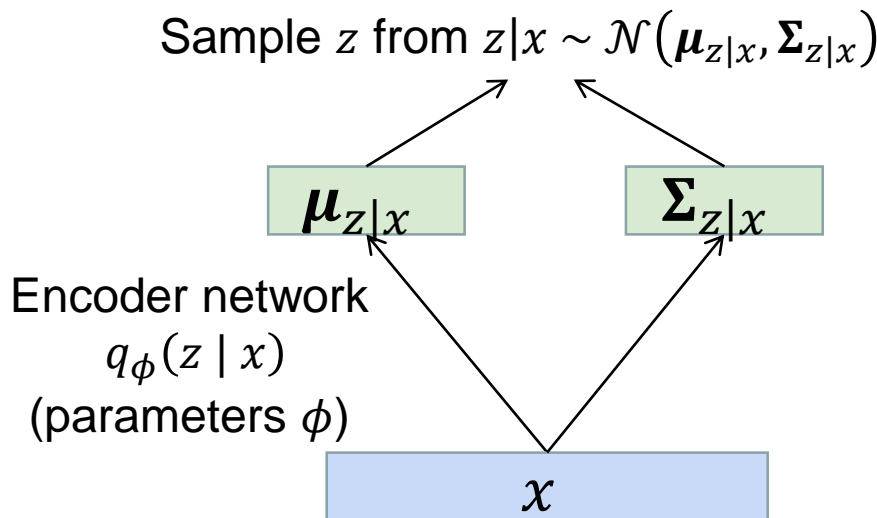
- Idea: Sample the model to generate data
 - We want to estimate the true parameters θ^* of this generative model.
- How should we represent the model?
 - Choose prior $p(\mathbf{z})$ to be simple, e.g., Gaussian
 - Conditional $p(\mathbf{x} | \mathbf{z})$ is complex (generates image)
 - ⇒ Represent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

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

Intractable!

Recap: Variational Autoencoders

- Define additional encoder network $q_{\phi}(z | x)$
 - 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

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

Recap: 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

Recap: 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: Variational Autoencoders
 - Autoencoders as Generative Models
 - Intractability
 - Variational Approximation
 - Evidence Lower Bound (ELBO)
- **Applying VAEs**
 - VAE Training
 - VAE Data Generation

Applying Variational Autoencoders

- Putting it all together...
 - Maximizing the likelihood lower bound

$$\underbrace{\mathbb{E}_z[\log p_\theta(x^{(i)} | z)] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

- Let's look at computing the bound for a given minibatch of input data (forward pass)...

Input data

x

Applying Variational Autoencoders

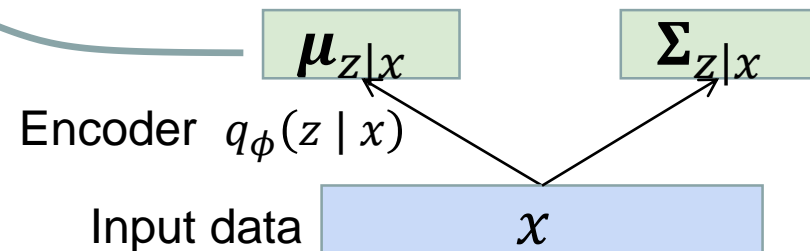
- Putting it all together...

- Maximizing the likelihood lower bound

$$\mathbb{E}_z [\log p_\theta(x^{(i)} | z)] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))$$

$$\mathcal{L}(x^{(i)}, \theta, \phi)$$

Make approximate
posterior distribution
close to prior



Applying Variational Autoencoders

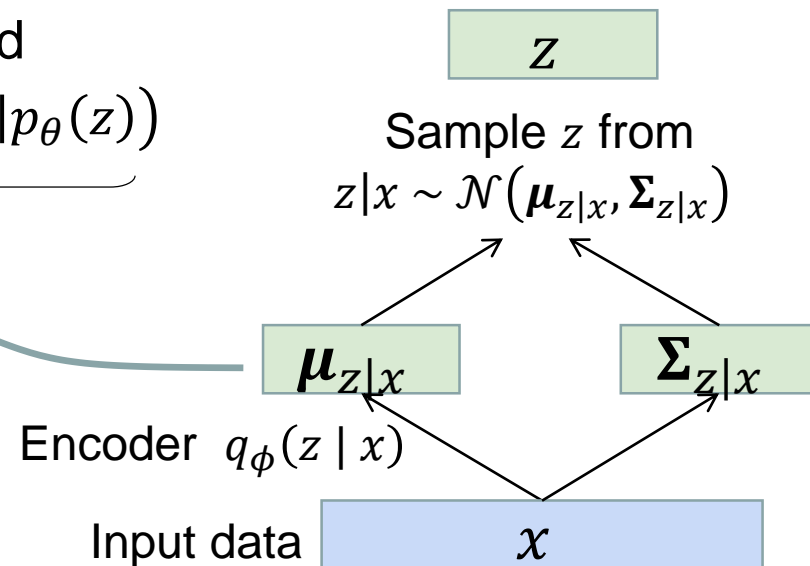
- Putting it all together...

- Maximizing the likelihood lower bound

$$\mathbb{E}_z [\log p_\theta(x^{(i)} | z)] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))$$

$$\mathcal{L}(x^{(i)}, \theta, \phi)$$

Make approximate
posterior distribution
close to prior



Applying Variational Autoencoders

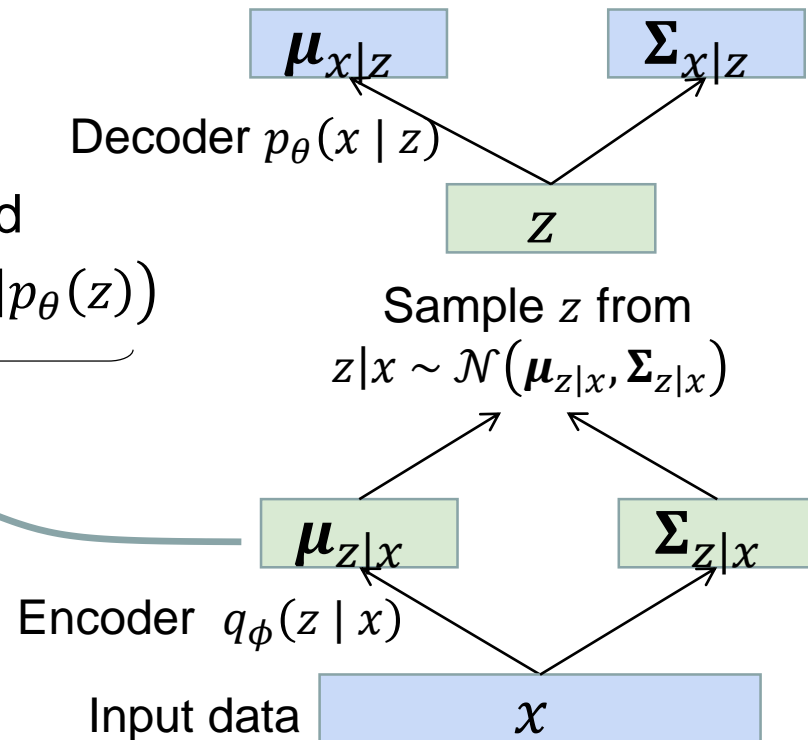
- Putting it all together...

- Maximizing the likelihood lower bound

$$\mathbb{E}_z [\log p_\theta(x^{(i)} | z)] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))$$

$$\mathcal{L}(x^{(i)}, \theta, \phi)$$

Make approximate
posterior distribution
close to prior



Applying Variational Autoencoders

- Putting it all together...

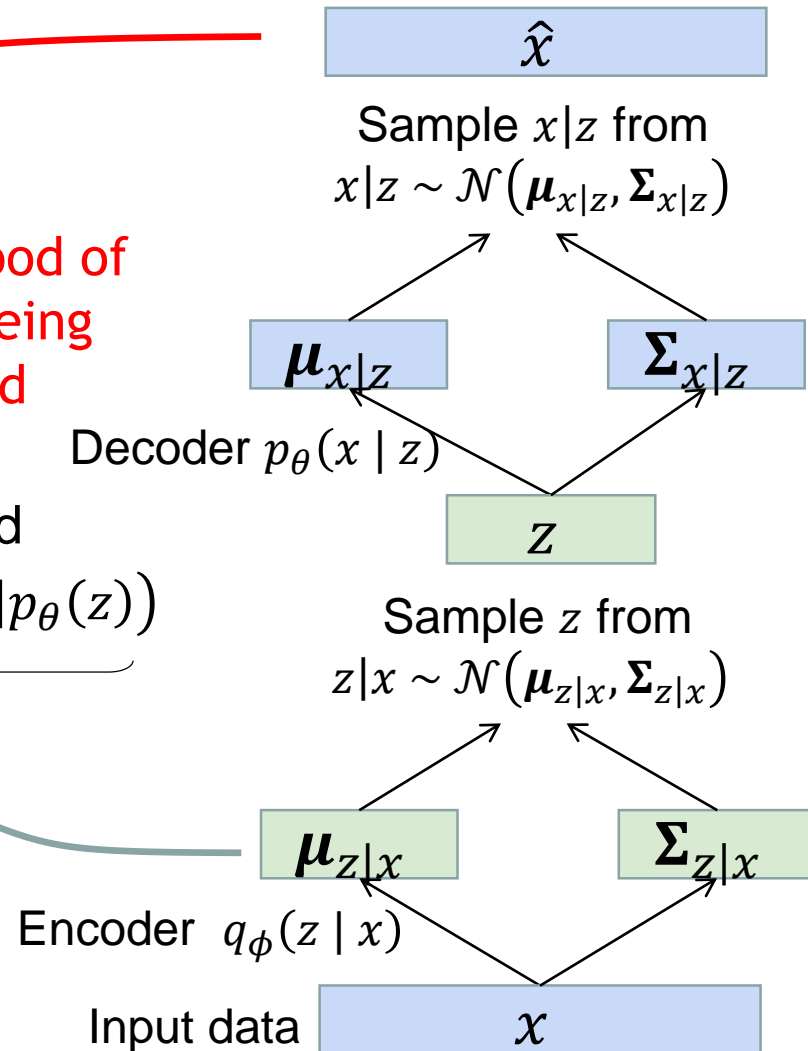
Maximize likelihood of original input being reconstructed

- Maximizing the likelihood lower bound

$$\mathbb{E}_z [\log p_\theta(x^{(i)} | z)] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))$$

$$\mathcal{L}(x^{(i)}, \theta, \phi)$$

Make approximate posterior distribution close to prior



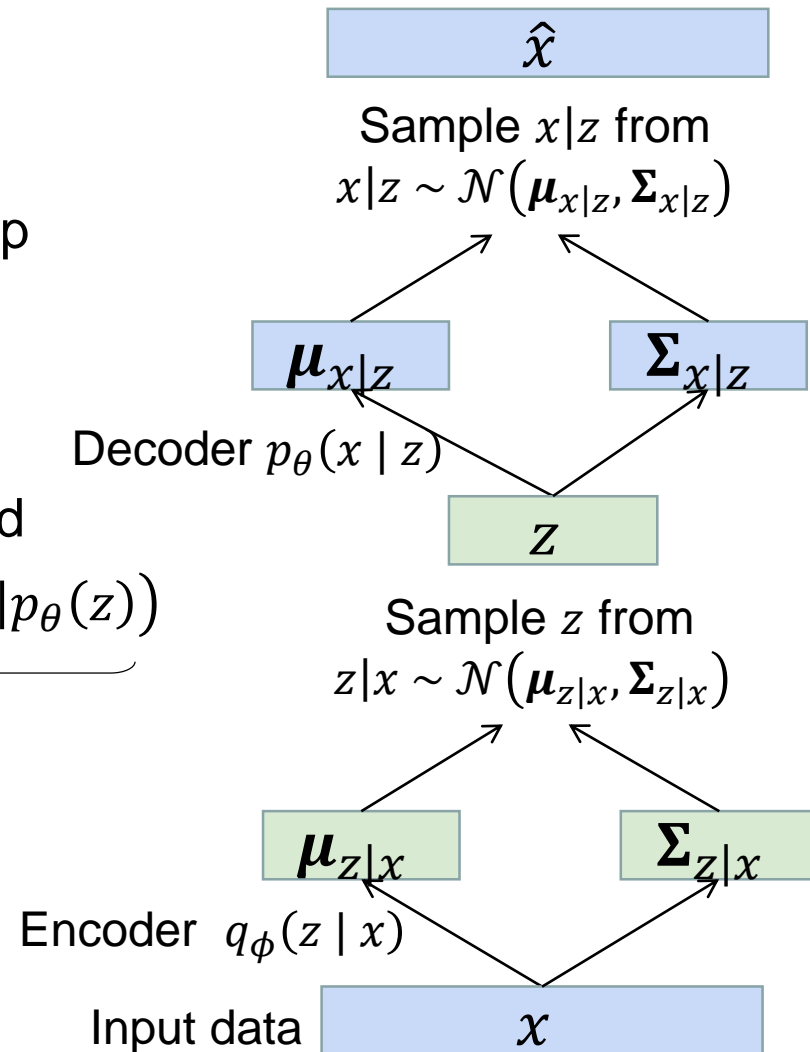
Applying Variational Autoencoders

- Putting it all together...
 - Compute this forward pass for every minibatch of input data, then backprop

- Maximizing the likelihood lower bound

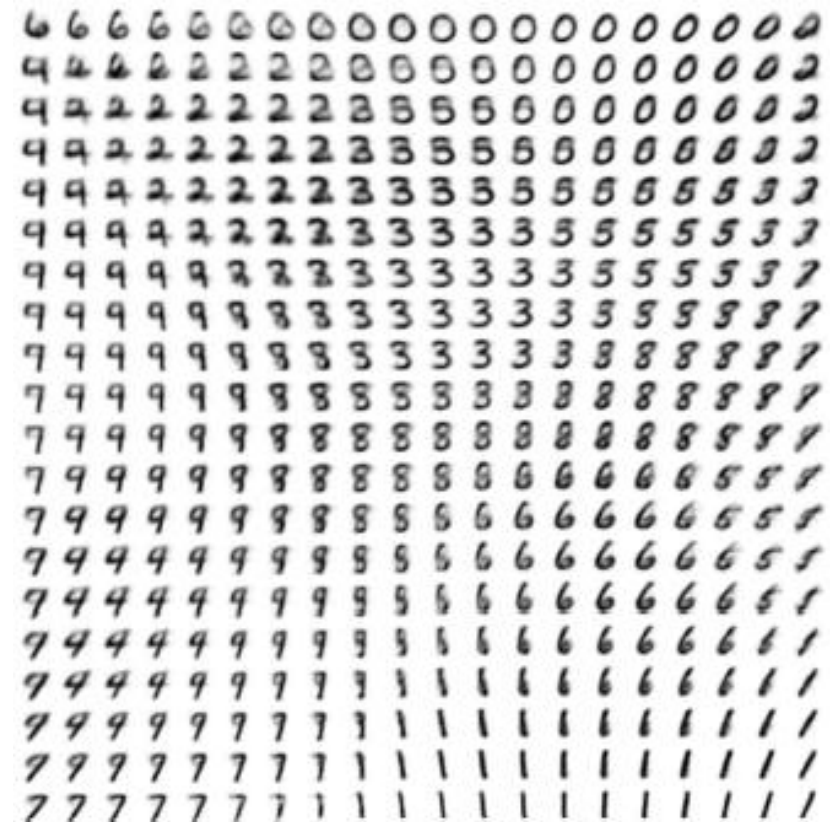
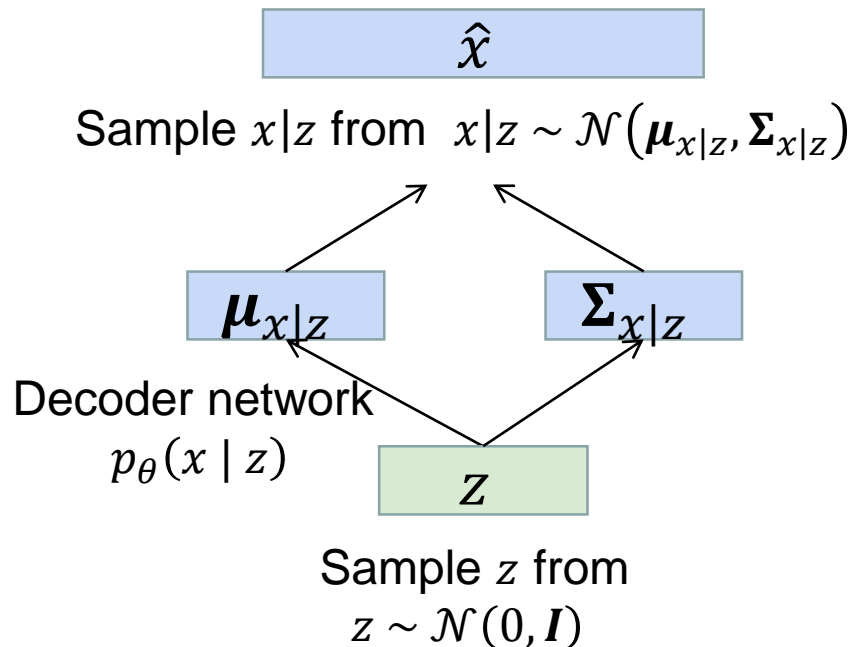
$$\mathbb{E}_z [\log p_\theta(x^{(i)} | z)] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))$$

$$\mathcal{L}(x^{(i)}, \theta, \phi)$$



Variational Autoencoders: Generating Data

- Use decoder network
 - Now sample x from prior

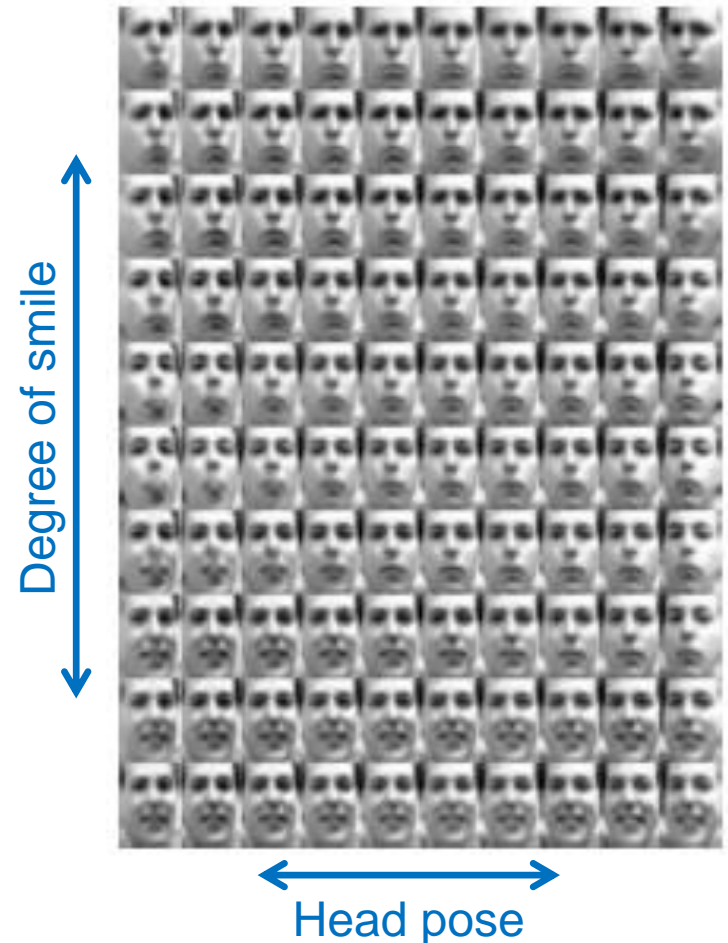


Latent MNIST manifold

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

Variational Autoencoders: Generating Data

- Another example
 - Learning a face manifold
- Comments
 - Diagonal prior on \mathbf{z}
⇒ Independent latent variables
 - Different dimensions of \mathbf{z} encode interpretable factors of variation



Some More Learned Manifolds



32x32 CIFAR-10

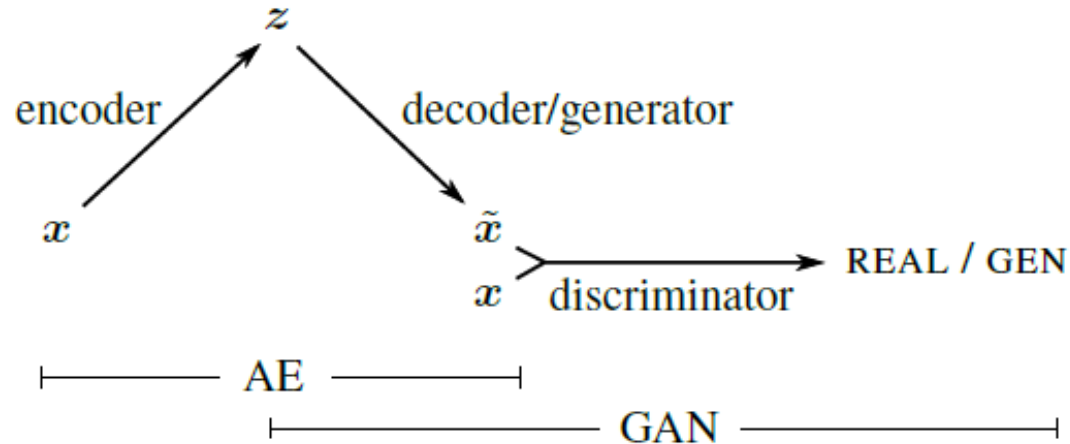


Labeled Faces in the Wild

Summary: Variational Autoencoders

- Idea
 - Probabilistic Spin on traditional autoencoders
 - Intractable density \Rightarrow derive & optimize a variational lower bound
- Pros
 - Principled approach to generative models
 - Allows inference of $q_{\phi}(z | x)$, can be useful feature representation for other tasks
- Cons
 - Only maximizes lower bound of likelihood
 - Samples blurrier and lower quality compared to state-of-the-art (GANs)
- Active area of research
 - More flexible approximations, e.g., GMMs instead of diagonal Gaussian

Combinations



- Attempts at combining the advantages
 - Use learned feature representations in the GAN discriminator as basis for the VAE reconstruction objective
 - Replacing element-wise errors with feature-wise errors to better capture the data distribution

A. Larsen, S. Sonderby, H. Larochelle, O. Winther, [Autoencoding beyond Pixels using a Learned Similarity Metric](#), arXiv 1512.09300

Results



Samples from different generative models



Reconstructions from different autoencoders

- VAE_{Disl}: Train a GAN first, then use the discriminator to train a VAE
- VAE/GAN: VAE and GAN trained together

References

- Variational Auto-Encoders

- D. Kingma, M. Welling, [Auto-Encoding Variational Bayes](#), ICLR 2014.
- A. Larsen, S. Sonderby, H. Larochelle, O. Winther, [Autoencoding beyond Pixels using a Learned Similarity Metric](#), arXiv:1512.09300, 2015.