## Slide 1

**Advanced Machine Learning**
**Summer 2019**

**Part 18 – Variational Autoencoders**
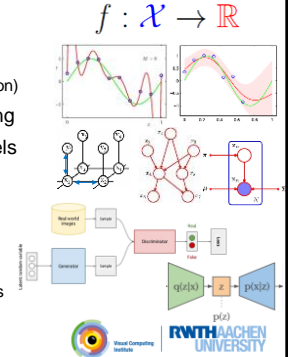**03.07.2019**

Prof. Dr. Bastian Leibe

RWTH Aachen University, Computer Vision Group
http://www.vision.rwth-aachen.de

## Slide 2

### 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

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

Visual Computing Institute | Prof. Dr. Bastian Leibe
Advanced Machine Learning
Part 18 – Variational Autoencoders
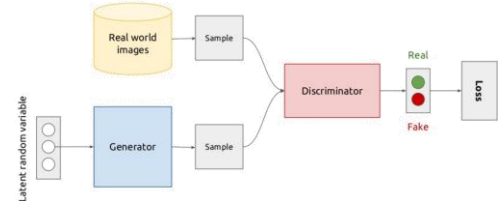
## Slide 3

### 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

Visual Computing Institute | Prof. Dr. Bastian Leibe
Advanced Machine Learning
Part 18 – Variational Autoencoders

## Slide 4

### 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

Visual Computing Institute | Prof. Dr. Bastian Leibe
Advanced Machine Learning
Part 18 – Variational Autoencoders

Image credit: Kevin McGuiness

## Slide 5

### Recap: GAN Loss Function

- This corresponds to a two-player minimax game:

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

- 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 $\log\left(1 - D(G(\boldsymbol{z}))\right)$.

- The Nash equilibrium of this game is achieved at
  - $p_g(\boldsymbol{x}) = p_{data}(\boldsymbol{x}) \quad \forall \boldsymbol{x}$
  - $D(\boldsymbol{x}) = \frac{1}{2} \quad \forall \boldsymbol{x}$

Visual Computing Institute | Prof. Dr. Bastian Leibe
Advanced Machine Learning
Part 18 – Variational Autoencoders

## Slide 6

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

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

  **end for** *(Discriminator updates)*

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

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

*(Generator updates)*

**end for**

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

Visual Computing Institute | Prof. Dr. Bastian Leibe
Advanced Machine Learning
Part 18 – Variational Autoencoders

## Recap: Intuition



discrimi-nator · data · model distrib.

- **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

Features  $z$

Encoder

Input data  $x$

- **Autoencoders**
  - 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

---

## Autoencoders

Features  $z$

Encoder

Input data  $x$

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

---

## Autoencoders

Reconstructed input data  $\hat{x}$

Decoder

Features  $z$

Encoder

Input data  $x$

**Decoder**: 4-layer upconv
**Encoder**: 4-layer conv

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

---

## Autoencoders

$L_2$ Loss function

Reconstructed input data  $\hat{x}$

Decoder

Features  $z$

Encoder

Input data  $x$

**Decoder**: 4-layer upconv
**Encoder**: 4-layer conv

- **How to learn such a feature representation?**
  - Train such that features can be used to reconstruct original data.
  - "Autoencoding" – encoding itself
  - $L_2$ loss function  $\|x - \hat{x}\|^2$
  - Note: this doesn't use any labels!

## Slide 14

### Autoencoders

Reconstructed input data $\hat{x}$

Decoder

Features $z$

Encoder

Input data $x$

**Decoder**: 4-layer upconv
**Encoder**: 4-layer conv

- After training
  - Throw away the decoder part

---

## Slide 15

### Autoencoders

Loss function (softmax, etc.)

Predicted label $\hat{y}$ $y$

Classifier

Features $z$

Encoder

Input data $x$

Fine-tune

bird    plane
dog    deer    truck

Train for final task
(on small dataset)

- 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*

---

## Slide 16

### Variants of Autoencoders

$L_2$ Loss function

Reconstructed input data $\hat{x}$

Decoder

Features $z$

Encoder

Input data $x$

- Analyzing the learning process
  - Learning process minimizes a loss function $L\left(\mathbf{x}, g(f(\mathbf{x}))\right)$
  - 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.

---

## Slide 17

### Variants of Autoencoders

$L_2$ Loss function

Reconstructed input data $\hat{x}$

Decoder

Features $z$

Encoder

Input data $x$

- Analyzing the learning process
  - Learning process minimizes a loss function $L\left(\mathbf{x}, g(f(\mathbf{x}))\right)$
  - *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.

---

## Slide 18

### Variants of Autoencoders

$L_2$ Loss function

Reconstructed input data $\hat{x}$

Decoder

Features $z$

Encoder

Input data $x$

- Regularized Autoencoders
  - Include a regularization term to the loss function: $L\left(\mathbf{x}, g(f(\mathbf{x}))\right) + \Omega(\mathbf{z})$
  - E.g., enforce sparsity by an $L_1$ regularizer   $\Omega(\mathbf{z}) = \lambda \sum_i |z_i|$

---

## Slide 19

### Variants of Autoencoders

$L_2$ Loss function

Reconstructed input data $\hat{x}$

Decoder

Features $z$

Encoder

Input data $x$

- Regularized Autoencoders
  - We can think of the sparse encoder framework as approximating ML training of a generative model with latent variables $\mathbf{z}$.
$$\log p_{model}(\mathbf{x}) = \log \sum_{\mathbf{z}} p_{model}(\mathbf{x}, \mathbf{z})$$
  - The autoencoder approximates this sum with a point estimate for just one highly likely value for $\mathbf{z}$.

## Variants of Autoencoders

Loss function

Reconstructed input data — $\hat{x}$

Decoder

Features — $z$

Encoder

Input data — $x$

- **Denoising Autoencoder** (DAE)
  - Rather than the reconstruction loss, minimize $L\left(\mathbf{x}, g\bigl(f(\tilde{\mathbf{x}})\bigr)\right)$
    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.

Image source: [Goodfellow 2016]

---

## 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?

$\hat{x}$

Decoder

$z$

Encoder

$x$

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

Slide inspired by Feifei Li

---

## Probabilistic Spin on Autoencoders

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

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

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

Slide credit: Feifei Li

---

## Probabilistic Spin on Autoencoders

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

Decoder network

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

- **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 $p(\mathbf{x} \mid \mathbf{z})$ is complex (generates image)
    $\Rightarrow$ Represent with neural network

Slide credit: Feifei Li

4

## Slide 26

### Probabilistic Spin on Autoencoders

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

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

$x$

Decoder network

$z$
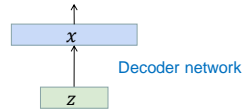
- Idea: Sample the model to generate data
  - We want to estimate the true parameters $\theta^*$ of this generative model.

- How to train the model?
  - Learn model parameters to maximize likelihood of training data
  $$p_\theta(\mathbf{x}) = \int p_\theta(\mathbf{z}) p_\theta(\mathbf{x} \mid \mathbf{z}) d\mathbf{z}$$
  - *What is the problem here?*    Intractable!

Slide credit: Feifei Li

---

## Slide 27

### Variational Autoencoders: Intractibility

- Computing the data likelihood
  $$p_\theta(\mathbf{x}) = \int p_\theta(\mathbf{z}) p_\theta(\mathbf{x} \mid \mathbf{z}) d\mathbf{z}$$

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

  - Posterior density is also intractable
  $$p_\theta(\mathbf{z} \mid \mathbf{x}) = \frac{p_\theta(\mathbf{z}) p_\theta(\mathbf{x} \mid \mathbf{z})}{p_\theta(\mathbf{x})}$$

Slide credit: Feifei Li

---

## Slide 28

### Variational Autoencoders: Intractibility

- Solution
  - In addition to the decoder network modeling $p_\theta(\mathbf{x} \mid \mathbf{z})$, define additional encoder network modeling $q_\phi(\mathbf{z} \mid \mathbf{x})$ that approximates $p_\theta(\mathbf{z} \mid \mathbf{x})$.

  - *We will see that this allows us to derive a lower bound on the data likelihood that is tractable and that we can optimize.*

Slide credit: Feifei Li
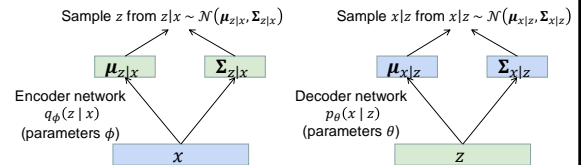
---

## Slide 29

### Variational Autoencoders

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

Sample $z$ from $z|x \sim \mathcal{N}(\boldsymbol{\mu}_{z|x}, \boldsymbol{\Sigma}_{z|x})$

Sample $x|z$ from $x|z \sim \mathcal{N}(\boldsymbol{\mu}_{x|z}, \boldsymbol{\Sigma}_{x|z})$

$\boldsymbol{\mu}_{z|x}$   $\boldsymbol{\Sigma}_{z|x}$

$\boldsymbol{\mu}_{x|z}$   $\boldsymbol{\Sigma}_{x|z}$

Encoder network
$q_\phi(z \mid x)$
(parameters $\phi$)

Decoder network
$p_\theta(x \mid z)$
(parameters $\theta$)

$x$

$z$

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

Slide credit: Feifei Li

---

## Slide 30

### Variational Autoencoders

- We can now work out the log-likelihood
$$\log p_\theta(x^{(i)}) = \mathbb{E}_{z \sim q_\phi(z|x^{(i)})}\left[\log p_\theta(x^{(i)})\right]$$
($p_\theta(x^{(i)})$ does not depend on $z$)

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

Slide credit: Feifei Li

---

## Slide 31

### Variational Autoencoders

- We can now work out the log-likelihood
$$\log p_\theta(x^{(i)}) = \mathbb{E}_{z \sim q_\phi(z|x^{(i)})}\left[\log p_\theta(x^{(i)})\right] \quad (p_\theta(x^{(i)}) \text{ does not depend on } z)$$
$$= \mathbb{E}_z\left[\log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})}\right] \quad \text{(Bayes' Rule)}$$
$$= \mathbb{E}_z\left[\log \frac{p_\theta(x^{(i)} \mid z) p_\theta(z)}{p_\theta(z \mid x^{(i)})} \frac{q_\phi(z \mid x^{(i)})}{q_\phi(z \mid x^{(i)})}\right] \quad \text{(Multiply by constant)}$$
$$= \mathbb{E}_z\left[\log p_\theta(x^{(i)} \mid z)\right] - \mathbb{E}_z\left[\log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z)}\right] + \mathbb{E}_z\left[\log \frac{q_\phi(z \mid x^{(i)})}{p_\theta(z \mid x^{(i)})}\right]$$
$$= \mathbb{E}_z\left[\log p_\theta(x^{(i)} \mid z)\right] - D_{KL}(q_\phi(z \mid x^{(i)}) \| p_\theta(z)) + D_{KL}(q_\phi(z \mid x^{(i)}) \| p_\theta(z \mid x^{(i)}))$$

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

Slide credit: Feifei Li

5

## Variational Autoencoders

- We can now work out the log-likelihood

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

$$= \mathbb{E}_z\left[\log \frac{p_\theta(x^{(i)}\,|\,z)p_\theta(z)}{p_\theta(z\,|\,x^{(i)})}\right] \qquad \text{(Bayes' Rule)}$$

$$= \mathbb{E}_z\left[\log \frac{p_\theta(x^{(i)}\,|\,z)p_\theta(z)}{p_\theta(z\,|\,x^{(i)})} \frac{q_\phi(z\,|\,x^{(i)})}{q_\phi(z\,|\,x^{(i)})}\right] \qquad \text{(Multiply by constant)}$$

$$= \mathbb{E}_z[\log p_\theta(x^{(i)}\,|\,z)] - \mathbb{E}_z\left[\log \frac{q_\phi(z\,|\,x^{(i)})}{p_\theta(z)}\right] + \mathbb{E}_z\left[\log \frac{q_\phi(z\,|\,x^{(i)})}{p_\theta(z\,|\,x^{(i)})}\right]$$

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

↑
Decoder network gives $p_\theta(\mathbf{x}\,|\,\mathbf{z})$, can compute estimate
of this term through sampling.
(Sampling differentiable through reparametrization trick, see paper)

Visual Computing Institute | Prof. Dr. Bastian Leibe
Advanced Machine Learning
Part 18 – Variational Autoencoders
Slide credit: Feifei Li

---

## Variational Autoencoders

- We can now work out the log-likelihood

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

$$= \mathbb{E}_z\left[\log \frac{p_\theta(x^{(i)}\,|\,z)p_\theta(z)}{p_\theta(z\,|\,x^{(i)})}\right] \qquad \text{(Bayes' Rule)}$$

$$= \mathbb{E}_z\left[\log \frac{p_\theta(x^{(i)}\,|\,z)p_\theta(z)}{p_\theta(z\,|\,x^{(i)})} \frac{q_\phi(z\,|\,x^{(i)})}{q_\phi(z\,|\,x^{(i)})}\right] \qquad \text{(Multiply by constant)}$$

$$= \mathbb{E}_z[\log p_\theta(x^{(i)}\,|\,z)] - \mathbb{E}_z\left[\log \frac{q_\phi(z\,|\,x^{(i)})}{p_\theta(z)}\right] + \mathbb{E}_z\left[\log \frac{q_\phi(z\,|\,x^{(i)})}{p_\theta(z\,|\,x^{(i)})}\right]$$

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

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

33 Visual Computing Institute | Prof. Dr. Bastian Leibe
Advanced Machine Learning
Part 18 – Variational Autoencoders
Slide credit: Feifei Li

---

## Variational Autoencoders

- We can now work out the log-likelihood

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

$$= \mathbb{E}_z\left[\log \frac{p_\theta(x^{(i)}\,|\,z)p_\theta(z)}{p_\theta(z\,|\,x^{(i)})}\right] \qquad \text{(Bayes' Rule)}$$

$$= \mathbb{E}_z\left[\log \frac{p_\theta(x^{(i)}\,|\,z)p_\theta(z)}{p_\theta(z\,|\,x^{(i)})} \frac{q_\phi(z\,|\,x^{(i)})}{q_\phi(z\,|\,x^{(i)})}\right] \qquad \text{(Multiply by constant)}$$

$$= \mathbb{E}_z[\log p_\theta(x^{(i)}\,|\,z)] - \mathbb{E}_z\left[\log \frac{q_\phi(z\,|\,x^{(i)})}{p_\theta(z)}\right] + \mathbb{E}_z\left[\log \frac{q_\phi(z\,|\,x^{(i)})}{p_\theta(z\,|\,x^{(i)})}\right]$$

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

↑
$p_\theta(\mathbf{z}\,|\,\mathbf{x})$ intractable (as seen earlier),
can't compute this KL term ☹
But we know KL divergence always $\geq 0$.

34 Visual Computing Institute | Prof. Dr. Bastian Leibe
Advanced Machine Learning
Part 18 – Variational Autoencoders
Slide credit: Feifei Li

---

## Variational Autoencoders

- We can now work out the log-likelihood

Want to
maximize
data
likelihood

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

$$= \mathbb{E}_z\left[\log \frac{p_\theta(x^{(i)}\,|\,z)p_\theta(z)}{p_\theta(z\,|\,x^{(i)})}\right] \qquad \text{(Bayes' Rule)}$$

$$= \mathbb{E}_z\left[\log \frac{p_\theta(x^{(i)}\,|\,z)p_\theta(z)}{p_\theta(z\,|\,x^{(i)})} \frac{q_\phi(z\,|\,x^{(i)})}{q_\phi(z\,|\,x^{(i)})}\right] \qquad \text{(Multiply by constant)}$$

$$= \mathbb{E}_z[\log p_\theta(x^{(i)}\,|\,z)] - \mathbb{E}_z\left[\log \frac{q_\phi(z\,|\,x^{(i)})}{p_\theta(z)}\right] + \mathbb{E}_z\left[\log \frac{q_\phi(z\,|\,x^{(i)})}{p_\theta(z\,|\,x^{(i)})}\right]$$

$$= \underbrace{\mathbb{E}_z[\log p_\theta(x^{(i)}\,|\,z)] - D_{KL}(q_\phi(z\,|\,x^{(i)})\|p_\theta(z))}_{\mathcal{L}(x^{(i)},\theta,\phi)} + \underbrace{D_{KL}(q_\phi(z\,|\,x^{(i)})\|p_\theta(z\,|\,x^{(i)}))}_{\geq 0}$$

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

35 Visual Computing Institute | Prof. Dr. Bastian Leibe
Advanced Machine Learning
Part 18 – Variational Autoencoders
Slide credit: Feifei Li

---

## Variational Autoencoders

- Variational Lower Bound ("ELBO")

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

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

- Training: Maximize lower bound

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

36 Visual Computing Institute | Prof. Dr. Bastian Leibe
Advanced Machine Learning
Part 18 – Variational Autoencoders
Slide adapted from Feifei Li

---

## 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

37 Visual Computing Institute | Prof. Dr. Bastian Leibe
Advanced Machine Learning
Part 18 – Variational Autoencoders

## Application Examples



32x32 CIFAR-10



Labeled Faces in the Wild

## References

• Variational Auto-Encoders
  – D. Kingma, M. Welling, Auto-Encoding Variational Bayes, ICLR 2014.