

# CS-E4890: Deep Learning Deep autoencoders

Alexander Ilin

#### Motivation

• Supervised learning problems: datasets consist of input-output pairs

 $(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), \dots, (\mathbf{x}^{(n)}, \mathbf{y}^{(n)})$ 

- Deep learning: supervised learning solved.
- Unsupervised learning: Make computers learn from unlabeled data

$$\mathbf{x}^{(1)},\ldots,\mathbf{x}^{(n)}$$

- Unsupervised learning seems important for building intelligent systems that can learn quickly. We humans learn a lot from unlabeled data.
- Unsupervised learning can be useful in practical applications:
  - detect samples that look different from training population (novelty/anomaly detection)
  - visualize data, discover patterns (information visualization)
  - generate new samples which look similar to the training data (generative models)

- We can use unlabeled data to do representation learning.
- Representation learning: extract features that may be useful for future (downstream) tasks

 $\mathbf{x} \xrightarrow{f} \mathbf{z}$ 

• Extracted features might work better than raw data in supervised learning tasks (especially with little labeled data):

$$\mathbf{x} \xrightarrow{f} \mathbf{z} 
ightarrow \mathbf{y}$$

- Problem: we do not know for which downstream tasks we need to prepare.
- We can extract patterns (features) that appear frequently in the data and hope that those features will be useful later.

### Bottleneck autoencoders

#### Principal component analysis (PCA)

- One of the simplest methods of unsupervised learning is PCA.
- Assume that the data have been centered by subtracting its mean:  $\mathbf{x} \leftarrow \mathbf{x} \mathbb{E}\{\mathbf{x}\}$ .
- The first principal component is found by maximizing the variance of the data multiplied by a unit-length vector **w**<sub>1</sub>:

$$y_1 = \mathbf{w}_1^{\top} \mathbf{x}, \qquad \|\mathbf{w}_1\| = 1$$
$$\mathbb{E}\{y_1^2\} = \mathbb{E}\{\mathbf{x}\mathbf{x}^{\top}\} = \mathbf{w}_1^{\top} \mathbb{E}\{\mathbf{x}\mathbf{x}^{\top}\}\mathbf{w}_1 = \mathbf{w}_1^{\top} \mathbf{C}_{\mathbf{x}}\mathbf{w}_1$$

$$\mathbf{w}_1^* = rg\max_{\mathbf{w}_1} \mathbf{w}_1^{ op} \mathbf{C}_{\mathbf{x}} \mathbf{w}_1, \quad ext{s.t.} \|\mathbf{w}_1\| = 1$$



The solution is given by the first dominant eigenvector of the covariance matrix  $\boldsymbol{C}_{\boldsymbol{x}}.$ 

• The second principal component is found by maximizing the variance in the subspace orthogonal to the first eigenvector of  $C_x$  (and so on).

#### PCA as minimum mean-square error compression

• We find an *m*-dimensional subspace spanned by orthonormal basis

$$\mathbf{W}_{n \times m} = \begin{bmatrix} \mathbf{w}_1 & \dots & \mathbf{w}_m \end{bmatrix} \qquad \mathbf{W}^\top \mathbf{W} = \mathbf{I}$$

- We project *n*-dimensional data vectors **x** onto the subspace:  $\mathbf{z} = \mathbf{W}^{\top} \mathbf{x}$ .
- The reconstruction of x that stays within the *m*-dimensional subspace defined by W:

$$\hat{\mathbf{x}} = \mathbf{W}\mathbf{z} = \sum_{i=1}^m (\mathbf{w}_i^{ op}\mathbf{x})\mathbf{w}_i = \mathbf{W}\mathbf{W}^{ op}\mathbf{x}$$

• We find **W** such that the mean-square error between original data and reconstruction is minimized:

$$\mathbf{W}_{\mathsf{PCA}} = \underset{\mathbf{W}}{\operatorname{arg\,min}} \mathbb{E}\{ \left\| \mathbf{x} - \underbrace{\mathbf{W}\mathbf{W}^{\top}\mathbf{x}}_{\hat{\mathbf{x}}} \right\|^2 \}, \qquad \text{s.t. } \mathbf{W}^{\top}\mathbf{W} = \mathbf{I}$$



#### PCA as a bootleneck autoencoder

• PCA as an autoencoder: We learn a mapping from x to x:

 $\hat{\mathbf{x}} = g(f(\mathbf{x}))$ encoder:  $f(\mathbf{x}) = \mathbf{W}_f \mathbf{x} + \mathbf{b}_f$ decoder:  $\hat{\mathbf{x}} = g(\mathbf{z}) = \mathbf{W}_g \mathbf{z} + \mathbf{b}_g$   $\mathcal{L} = \mathbb{E}\{\|\mathbf{x} - g(f(\mathbf{x}))\|^2\}$ 



• If we do not restrict f and g, we can learn a trivial identity mapping:

$$\hat{\mathbf{x}} = g(f(\mathbf{x})) = (\mathbf{W}_g \mathbf{W}_f) \mathbf{x} + (\mathbf{W}_g \mathbf{b}_f + \mathbf{b}_g) = \mathbf{x}, \quad \text{if } \mathbf{W}_g = \mathbf{W}_f^{-1} \text{ and } \mathbf{b}_g = -\mathbf{W}_g \mathbf{b}_f$$

- If the dimensionality of z is smaller than the dimensionality of x, autoencoding is useful: we compress the data.
  - z is often called a bottleneck.
  - Thus PCA can be implemented with a bottleneck autoencoder.

• We have a linear autoencoder:

$$\hat{\mathbf{x}} = g(f(\mathbf{x}))$$
  
 $f(\mathbf{x}) = \mathbf{W}_f \mathbf{x} + \mathbf{b}_f$   
 $g(\mathbf{z}) = \mathbf{W}_g \mathbf{z} + \mathbf{b}_g$ 

• How can we improve compression so that we get a smaller reconstruction error

$$\mathbb{E}\{\|\mathbf{x} - g(f(\mathbf{x}))\|^2\}$$

with a bottleneck layer of the same size?



• We have a linear autoencoder:

$$\hat{\mathbf{x}} = g(f(\mathbf{x}))$$
  
 $f(\mathbf{x}) = \mathbf{W}_f \mathbf{x} + \mathbf{b}_f$   
 $g(\mathbf{z}) = \mathbf{W}_g \mathbf{z} + \mathbf{b}_g$ 

• How can we improve compression so that we get a smaller reconstruction error

$$\mathbb{E}\{\|\mathbf{x} - g(f(\mathbf{x}))\|^2\}$$

with a bottleneck layer of the same size?

• We can use nonlinear encoder f and decoder g.



- Deep autoencoder: both the encoder and the decoder are deep neural networks.
- The optimization criterion is the mean-squared reconstruction error:

$$\boldsymbol{ heta}_{f}, \boldsymbol{ heta}_{g} = \operatorname*{arg\,min}_{\boldsymbol{ heta}_{f}, \boldsymbol{ heta}_{g}} \mathbb{E}\{\|\mathbf{x} - g(\mathbf{z}, \boldsymbol{ heta}_{g})\|^{2}\}, \qquad \mathbf{z} = f(\mathbf{x}, \boldsymbol{ heta}_{f})$$

- Bottleneck autoencoder: To prevent learning a trivial (identity) function, we use z with fewer dimensions (a bottleneck layer).
- Bottleneck autoencoders were proposed by Bourlard and Kamp (1988), Oja (1991).



- In this hypothetical example, the data lie on one-dimensional manifold.
- Principal component analysis is not be able to learn the one-dimensional manifold because it is a linear model.



A one-dimensional data manifold in the two-dimensional space.

- In this hypothetical example, the data lie on one-dimensional manifold.
- Principal component analysis is not be able to learn the one-dimensional manifold because it is a linear model.
- With a nonlinear autoencoder, we can learn a curved data manifold.
- In our example, colors represents the values of the latent code *z* that may be found by an autoencoder.



A one-dimensional data manifold in the two-dimensional space.

#### Deep autoencoder as feature extractor

- The most popular use case for deep autoencoders is data compression.
- Consider, for example, reinforcement learning (RL) tasks such as playing Doom (Ha and Schmidhuber, 2018).
- Learning from raw images (pixels) is likely to require a huge number of training episodes because the amount of input data is large.
- The authors first compress the data using an autoencoder and then train the agent using as observations compressed representations z.



#### Deep bottleneck autoencoder: MNIST example

• In the home assignment, you will train a bottleneck autoencoder for the MNIST dataset.

 Visualization of the z-space using t-SNE:





# Denoising autoencoders

• Feed inputs corrupted with noise (for example, Gaussian):

$$ilde{\mathbf{x}} = \mathbf{x} + \boldsymbol{\epsilon} \quad ext{with} \quad \epsilon_i \sim \mathcal{N}(\mathbf{0}, \sigma^2)$$

• Train  $\hat{\mathbf{x}} = g(f(\tilde{\mathbf{x}}))$  to minimize the reconstruction error:

$$c = \mathbb{E}\{\|\hat{\mathbf{x}} - \mathbf{x}\|^2\}$$

• One can view adding noise to inputs as a way to regularize the autoencoder (regularization by noise injection) but there is more theory behind denoising autoencoders.



For Gaussian corruption ε<sub>i</sub> ~ N(0, σ<sup>2</sup>), the optimal denoising is

 $d(\mathbf{\tilde{x}}) = \mathbf{\tilde{x}} + \sigma^2 
abla_{\mathbf{\tilde{x}}} \log p(\mathbf{\tilde{x}})$ 

(see Alain and Bengio, 2014, Raphan and Simoncelli, 2011)

- $d(\cdot)$  learns to point towards higher probability density.
- Thus, by learning the optimal denoising function  $d(\mathbf{x})$ , we implicitly model the data distribution  $p(\mathbf{x})$ .



Image from (Alain and Bengio, 2014)

#### Denoising autoencoder: variance MNIST example

• In the home assignment, we create a synthetic dataset (which we call variance MNIST).

• For this dataset, a vanilla bottleneck autoencoder with mean-squared error reconstruction loss cannot extract high-level features **z** that would capture the shapes of the digits.

• A denoising autoencoder can extract meaningful features. Visualization of the z-space using t-SNE:



Converting autoencoders into generative models with latent variables

- Generative models:
  - learn to represent the data distribution  $p(\mathbf{x})$
  - can be used to generate new examples from  $p(\mathbf{x})$ .
- An example: a mixture-of-Gaussians model

$$p(x \mid \boldsymbol{\theta}) = w_1 \mathcal{N}(x \mid \mu_1, \sigma_1^2) + w_2 \mathcal{N}(x \mid \mu_2, \sigma_2^2)$$

Parameters  $\theta = \{w_1, \mu_1, \sigma_1, w_2, \mu_2, \sigma_2\}$  can be estimated by maximum likelihood.

This model is an example of an explicit density model:
 p(x | θ) has an explicit parametric form.



#### Converting autoencoders into generative models

- Vanilla autoencoders are not generative models.
  - We cannot generate new samples from  $p(\mathbf{x})$ .
  - We cannot compute the probability that a new sample x comes from the same distribution (e.g., for novelty detection).



#### Converting autoencoders into generative models

- Vanilla autoencoders are not generative models.
  - We cannot generate new samples from  $p(\mathbf{x})$ .
  - We cannot compute the probability that a new sample x comes from the same distribution (e.g., for novelty detection).
- We can build a generative model, for example, in this way:
  - Assume that variables z are normally distributed:

 $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 

• Data samples x are nonlinear transformations of latent variables z:

 $\mathbf{x} = g(\mathbf{z}, \boldsymbol{ heta}) + \boldsymbol{arepsilon}$ 

with possibly noise added:  $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ 

- Function  $g(\mathbf{z}, \theta)$  can be modeled as a neural network.
- Now we can draw samples from the model.



• Our model contains latent (unobserved) variables z:

 $egin{aligned} \mathbf{z} &\sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \ \mathbf{x} &= g(\mathbf{z}, oldsymbol{ heta}) + arepsilon \ arepsilon &\sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}) \end{aligned}$ 

- A simple example to illustrate the idea: We model one-dimensional data x as a Gaussian variable z transformed with nonlinearity g with some noise added.
- We need to learn the latent variable model from training data {x<sub>i</sub>}. We should tune parameters θ, σ<sup>2</sup> so that the training examples are likely to be produced by the model.



#### Learning the parameters of the latent variable model

• We can tune parameters  $\theta$ ,  $\sigma^2$  by maximizing the probability of the training data (maximum likelihood estimate):

$$\theta_{\mathsf{ML}} = \operatorname*{arg\,max}_{\theta} \log p(\mathbf{x}_1, ..., \mathbf{x}_N \mid \theta)$$
$$\log p(\mathbf{x}_1, ..., \mathbf{x}_N \mid \theta) = \sum_{i=1}^N \log p(\mathbf{x}_i \mid \theta) = \sum_{i=1}^N \log \int p(\mathbf{x}_i \mid \mathbf{z}_i, \theta) p(\mathbf{z}_i) d\mathbf{z}$$

• The probability density functions are defined by our model:

$$p(\mathbf{x}_i \mid \mathbf{z}_i, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}_i \mid g(\mathbf{z}_i, \boldsymbol{\theta}), \sigma^2 \mathbf{I})$$
$$p(\mathbf{z}_i) = \mathcal{N}(\mathbf{z}_i \mid \mathbf{0}, \mathbf{I})$$

 Direct optimization of log p(x<sub>1</sub>,..., x<sub>N</sub> | θ) is difficult because the above integrals are intractable.



#### ML estimation with the EM algorithm

• The classical way to estimate parameters  $\theta$  of a latent variable model

$$p(\mathbf{x}_1,...,\mathbf{x}_N,\mathbf{z}_1,...,\mathbf{z}_N \mid \boldsymbol{ heta}) = \prod_{i=1}^N p(\mathbf{x}_i \mid \mathbf{z}_i,\boldsymbol{ heta}) p(\mathbf{z}_i)$$

is the expectation-maximization (EM) algorithm.

- The EM-algorithm iterates between two steps: E-step and M-step.
  - E-step: Compute posterior probabilities  $p(\mathbf{z}_i | \mathbf{x}_i, \theta)$  given current values of  $\theta$ .
  - M-step: Update the values of  $\theta$  using computed  $p(\mathbf{z}_i \mid \mathbf{x}_i, \theta)$ .



Consider our simple example. We initialize  $\theta$  with values that give us **g** of the form shown in the figure.

$$egin{aligned} \mathbf{z} &\sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \ \mathbf{x} &= g(\mathbf{z}, m{ heta}) + m{arepsilon} \ m{arepsilon} &\sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}) \end{aligned}$$

 The E-step: Compute the posterior probabilities of the unobserved latent variables z<sub>i</sub> given the data and the current estimates of the model parameters θ:

1

$$egin{aligned} q(\mathbf{z}_1,...,\mathbf{z}_N) &= q(\mathbf{z}_1)...q(\mathbf{z}_N) \ q(\mathbf{z}_i) &= p(\mathbf{z}_i \mid \mathbf{x}_i, oldsymbol{ heta}) \end{aligned}$$



E-step: For each training data point, find the distribution over the latent variables that could have produced that data point according to the model.

$$egin{aligned} \mathbf{z} &\sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \ \mathbf{x} &= g(\mathbf{z}, m{ heta}) + m{arepsilon} \ m{arepsilon} &\sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}) \end{aligned}$$

 The E-step: Compute the posterior probabilities of the unobserved latent variables z<sub>i</sub> given the data and the current estimates of the model parameters θ:

1

$$egin{aligned} q(\mathsf{z}_1,...,\mathsf{z}_N) &= q(\mathsf{z}_1)...q(\mathsf{z}_N) \ q(\mathsf{z}_i) &= p(\mathsf{z}_i \mid \mathsf{x}_i, heta) \end{aligned}$$



E-step: For each training data point, find the distribution over the latent variables that could have produced that data point according to the model.

$$egin{aligned} \mathbf{z} &\sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \ \mathbf{x} &= g(\mathbf{z}, m{ heta}) + m{arepsilon} \ m{arepsilon} &\sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}) \end{aligned}$$

 The E-step: Compute the posterior probabilities of the unobserved latent variables z<sub>i</sub> given the data and the current estimates of the model parameters θ:

1

$$egin{aligned} q(\mathsf{z}_1,...,\mathsf{z}_N) &= q(\mathsf{z}_1)...q(\mathsf{z}_N) \ q(\mathsf{z}_i) &= p(\mathsf{z}_i \mid \mathsf{x}_i, heta) \end{aligned}$$



E-step: For each training data point, find the distribution over the latent variables that could have produced that data point according to the model.

 In the M-step, we use the computed distributions q(z<sub>i</sub>) to form the following objective function:

$$egin{aligned} \mathcal{F}(oldsymbol{ heta}) &= \langle \log p(\mathbf{x}_1,...,\mathbf{x}_N,\mathbf{z}_1,...,\mathbf{z}_N \mid oldsymbol{ heta}) 
angle_{q(\mathbf{z}_1,...,\mathbf{z}_N)} \ &= \sum_{i=1}^N \left\langle \log p(\mathbf{x}_i,\mathbf{z}_i \mid oldsymbol{ heta}) 
angle_{q(\mathbf{z}_i)} \ &= \sum_{i=1}^N \int q(\mathbf{z}_i) \log p(\mathbf{x}_i,\mathbf{z}_i \mid oldsymbol{ heta}) d\mathbf{z}_i \end{aligned}$$

and maximize it wrt model parameters  $\boldsymbol{\theta}$ .

• We are guaranteed to improve the likelihood

$$\log p(\mathbf{x}_1, ..., \mathbf{x}_N \mid \boldsymbol{\theta})$$



Iteration 1

 In the M-step, we use the computed distributions q(z<sub>i</sub>) to form the following objective function:

$$egin{aligned} \mathcal{F}(oldsymbol{ heta}) &= \langle \log p(\mathbf{x}_1,...,\mathbf{x}_N,\mathbf{z}_1,...,\mathbf{z}_N \mid oldsymbol{ heta}) 
angle_{q(\mathbf{z}_1,...,\mathbf{z}_N)} \ &= \sum_{i=1}^N \langle \log p(\mathbf{x}_i,\mathbf{z}_i \mid oldsymbol{ heta}) 
angle_{q(\mathbf{z}_i)} \ &= \sum_{i=1}^N \int q(\mathbf{z}_i) \log p(\mathbf{x}_i,\mathbf{z}_i \mid oldsymbol{ heta}) d\mathbf{z}_i \end{aligned}$$

and maximize it wrt model parameters  $\boldsymbol{\theta}$ .

• We are guaranteed to improve the likelihood

$$\log p(\mathbf{x}_1, ..., \mathbf{x}_N \mid \boldsymbol{\theta})$$



Iteration 2

 In the M-step, we use the computed distributions q(z<sub>i</sub>) to form the following objective function:

$$egin{aligned} \mathcal{F}(oldsymbol{ heta}) &= \langle \log p(\mathbf{x}_1,...,\mathbf{x}_N,\mathbf{z}_1,...,\mathbf{z}_N \mid oldsymbol{ heta}) 
angle_{q(\mathbf{z}_1,...,\mathbf{z}_N)} \ &= \sum_{i=1}^N \langle \log p(\mathbf{x}_i,\mathbf{z}_i \mid oldsymbol{ heta}) 
angle_{q(\mathbf{z}_i)} \ &= \sum_{i=1}^N \int q(\mathbf{z}_i) \log p(\mathbf{x}_i,\mathbf{z}_i \mid oldsymbol{ heta}) d\mathbf{z}_i \end{aligned}$$

and maximize it wrt model parameters  $\boldsymbol{\theta}$ .

• We are guaranteed to improve the likelihood

$$\log p(\mathbf{x}_1, ..., \mathbf{x}_N \mid \boldsymbol{\theta})$$



Iteration 3

 In the M-step, we use the computed distributions q(z<sub>i</sub>) to form the following objective function:

$$egin{aligned} \mathcal{F}(oldsymbol{ heta}) &= \langle \log p(\mathbf{x}_1,...,\mathbf{x}_N,\mathbf{z}_1,...,\mathbf{z}_N \mid oldsymbol{ heta}) 
angle_{q(\mathbf{z}_1,...,\mathbf{z}_N)} \ &= \sum_{i=1}^N \left\langle \log p(\mathbf{x}_i,\mathbf{z}_i \mid oldsymbol{ heta}) 
ight
angle_{q(\mathbf{z}_i)} \ &= \sum_{i=1}^N \int q(\mathbf{z}_i) \log p(\mathbf{x}_i,\mathbf{z}_i \mid oldsymbol{ heta}) d\mathbf{z}_i \end{aligned}$$

and maximize it wrt model parameters  $\boldsymbol{\theta}$ .

• We are guaranteed to improve the likelihood

$$\log p(\mathbf{x}_1, ..., \mathbf{x}_N \mid \boldsymbol{\theta})$$



Iteration 4

Learning latent variable models with variational approximations

- There are a few problems with the direct application of the EM-algorithm in nonlinear latent variable models.
- One problem is the intractability of the true conditional distributions q(z<sub>i</sub>) = p(z<sub>i</sub> | x<sub>i</sub>, θ) that we need to compute on the E-step.
- The true distributions can be very complex (for example, a multi-modal distribution in our simple example).



Example of multi-modal  $p(\mathbf{z}_i | \mathbf{x}_i, \boldsymbol{\theta})$ 

#### E-step: Variational approximations

- Solution: Instead of using true conditional distributions, use their approximations q(z<sub>i</sub>) ≈ p(z<sub>i</sub> | x<sub>i</sub>, θ).
- q(z<sub>i</sub>) is selected to have a simple form, most often a Gaussian:

$$q(\mathbf{z}_i) = \mathcal{N}(\mu_{\mathbf{z}_i}, \sigma_{\mathbf{z}_i}^2)$$

Note: we have two parameters  $\mu_{z_i}$  and  $\sigma_{z_i}^2$  describing  $q(z_i)$  for *each training sample*.

• Parameters describing the posterior distributions of the latent variables  $\theta_q = \{\mu_{z_i}, \sigma_{z_i}^2\}_{i=1}^N$  are called variational parameters.



• A popular way to find the approximation is by minimizing the Kullback-Leibler divergence between  $q(\mathbf{z}_i)$  and  $p(\mathbf{z}_i \mid \mathbf{x}_i, \boldsymbol{\theta})$ .

#### **E-step: Variational approximations**

- We can minimize the KL divergence between  $q(\mathbf{z}_i)$  and  $p(\mathbf{z}_i | \mathbf{x}_i, \theta)$  using the following trick:
  - Add to the objective function used in the M-step the entropies of the approximate distributions:

$$\mathcal{F}(\theta, \theta_q) = \sum_{i=1}^{N} \underbrace{\int q(\mathbf{z}_i) \log p(\mathbf{x}_i, \mathbf{z}_i \mid \theta) d\mathbf{z}_i}_{\text{what we had in the M-step}} \underbrace{- \int q(\mathbf{z}_i) \log q(\mathbf{z}_i) d\mathbf{z}_i}_{\text{entropy}}$$

$$= \sum_{i=1}^{N} \int q(\mathbf{z}_i) \log \frac{p(\mathbf{x}_i, \mathbf{z}_i \mid \theta)}{q(\mathbf{z}_i)} dz_i = \sum_{i=1}^{N} \int q(\mathbf{z}_i) \log \frac{p(\mathbf{z}_i \mid \mathbf{x}_i, \theta) p(\mathbf{x}_i \mid \theta)}{q(\mathbf{z}_i)} dz_i$$

$$= \sum_{i=1}^{N} -D_{\mathsf{KL}}(q(\mathbf{z}_i) \parallel p(\mathbf{z}_i \mid \mathbf{x}_i, \theta)) + \log p(\mathbf{x}_i \mid \theta)$$

 One can see that maximizing *F*(θ, θ<sub>q</sub>) wrt variational parameters θ<sub>q</sub> is equivalent to minimizing the KL divergence between q(z<sub>i</sub>) and p(z<sub>i</sub> | x<sub>i</sub>, θ). We can now maximize a single function *F* wrt *θ* and *θ<sub>q</sub>* jointly without the need to alternate between the E- and M-steps:

$$egin{aligned} \mathcal{F}(m{ heta},m{ heta}_q) &= \sum_{i=1}^N \int q(\mathbf{z}_i)\log p(\mathbf{x}_i,\mathbf{z}_i\midm{ heta})d\mathbf{z}_i - \int q(\mathbf{z}_i)\log q(\mathbf{z}_i)d\mathbf{z}_i \ &= \sum_{i=1}^N - D_{\mathsf{KL}}(q(\mathbf{z}_i)\parallel p(\mathbf{z}_i\mid\mathbf{x}_i,m{ heta})) + \log p(\mathbf{x}_i\midm{ heta}) \end{aligned}$$

- Maximizing  $\mathcal{F}(\theta, \theta_q)$  wrt  $\theta$  is equivalent to the M-step.
- Maximizing  $\mathcal{F}(\theta, \theta_q)$  wrt  $\theta_q$  is done in the E-step with variational approxiations.
- We can solve this optimization problem using any optimizer of our choice.

• The objective function

$$\mathcal{F}(oldsymbol{ heta},oldsymbol{ heta}_q) = \sum_{i=1}^N - D_{\mathsf{KL}}(oldsymbol{q}(\mathsf{z}_i) \parallel p(oldsymbol{z}_i \mid \mathsf{x}_i,oldsymbol{ heta})) + \log p(oldsymbol{x}_i \mid oldsymbol{ heta})$$

is the *lower bound* of the true likelihood that we want to optimize. Since  $D_{\mathsf{KL}}(q \parallel p) \ge 0$ :

$$\mathcal{F}(\boldsymbol{ heta}, \boldsymbol{ heta}_q) \leq \sum_{i=1}^N \log p(\mathbf{x}_i \mid \boldsymbol{ heta}) = \log p(\mathbf{x}_1, ..., \mathbf{x}_N \mid \boldsymbol{ heta})$$

- This function is often called evidence lower bound or ELBO.
- The closer our approximation  $q(\mathbf{z}_i)$  to the true posterior  $p(\mathbf{z}_i | \mathbf{x}_i, \theta)$ , the tighter the bound.

#### ELBO for our deep generative model

• ELBO can be re-written in the following form:

$$\mathcal{F}\left(\boldsymbol{\theta}, \boldsymbol{\theta}_{q}\right) = \sum_{i=1}^{N} \int q(\mathbf{z}_{i}) \log p(\mathbf{x}_{i} \mid \mathbf{z}_{i}, \boldsymbol{\theta}) dz_{i} - \int q(\mathbf{z}_{i}) \log \frac{q(\mathbf{z}_{i})}{p(\mathbf{z}_{i})} dz_{i}$$

- Recall our deep generative model:  $p(\mathbf{x}_i | \mathbf{z}_i, \theta) = \mathcal{N}(\mathbf{x}_i | g(\mathbf{z}_i, \theta), \sigma^2 \mathbf{I}),$
- The first term in equation (1) can be written as

$$\left\langle -rac{D}{2}\log 2\pi\sigma^2 - rac{1}{2\sigma^2}\sum_{d=1}^D (\mathbf{x}_i(d) - g_d(\mathbf{z}_i, oldsymbol{ heta}))^2 
ight
angle_{q(\mathbf{z}_i)}$$

where D is the number of dimensions in  $\mathbf{x}$ ,  $\mathbf{x}_i(d)$  is the d-th element of  $\mathbf{x}_i$  and  $g_d$  is the d-th element of the output of function g.

 The first term contains the expected mean-squared error between data sample x<sub>i</sub> and its reconstruction g<sub>d</sub>(z<sub>i</sub>, θ) from the latent code z<sub>i</sub>.



g

x

(1)

#### ELBO for our deep generative model

$$\mathcal{F}(\boldsymbol{\theta}, \boldsymbol{\theta}_{q}) = \sum_{i=1}^{N} \underbrace{\int q(\mathbf{z}_{i}) \log p(\mathbf{x}_{i} \mid \mathbf{z}_{i}, \boldsymbol{\theta}) d\mathbf{z}_{i}}_{\text{minus mean-square reconstruction error}} \underbrace{-\int q(\mathbf{z}_{i}) \log \frac{q(\mathbf{z}_{i})}{p(\mathbf{z}_{i})} d\mathbf{z}_{i}}_{\text{regularization term}}$$

• The second term is minus KL-divergence between  $q(\mathbf{z}_i)$  and the prior  $p(\mathbf{z}_i) = \mathcal{N}(0, \mathbf{I})$ :

$$-\int q(\mathbf{z}_i)\lograc{q(\mathbf{z}_i)}{p(\mathbf{z}_i)}dz_i = -D_{\mathsf{KL}}(q(\mathbf{z}_i)\parallel p(\mathbf{z}_i))$$

It is a kind of a regularization term: We want the conditional distributions q(z<sub>i</sub>) to be close to the prior p(z<sub>i</sub>) = N(0, I).

### Variational autoencoders

• The first algorithm for learning latent variable model

$$\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$
  $\mathbf{x} = g(\mathbf{z}, \boldsymbol{ heta}) + \boldsymbol{arepsilon}$   $\boldsymbol{arepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ 

using variational approximations was proposed in this university (Lappalainen and Honkela, 2001).

• The objective function was ELBO:

$$\mathcal{F}(\theta, \theta_q) = \sum_{i=1}^{N} \underbrace{\int q(\mathbf{z}_i) \log p(\mathbf{x}_i \mid \mathbf{z}_i, \theta) d\mathbf{z}_i}_{\text{needs approximations}} \underbrace{-\int q(\mathbf{z}_i) \log \frac{q(\mathbf{z}_i)}{p(\mathbf{z}_i)} d\mathbf{z}_i}_{\text{can be computed analytically}}$$

The posterior approximations were Gaussian q(z<sub>i</sub>) = N(μ<sub>z<sub>i</sub></sub>, σ<sup>2</sup><sub>z<sub>i</sub></sub>). The number of variational parameters θ<sub>q</sub> = {μ<sub>z<sub>i</sub></sub>, σ<sup>2</sup><sub>z<sub>i</sub></sub>}<sup>N</sup><sub>i=1</sub> was proportional to the number of training samples.

#### Adding encoder

- We want to get rid of the large number of variational parameters θ<sub>q</sub> = {μ<sub>z<sub>i</sub></sub>, σ<sup>2</sup><sub>z<sub>i</sub></sub>}<sup>N</sup><sub>i=1</sub>.
- For fixed model parameters θ, the optimal q(z) only depends on x. The inference procedure does the following mapping:

$$\mathsf{x} o q(\mathsf{z})$$

For Gaussian approximation:  $\mathbf{x} \rightarrow \mu_{\mathbf{z}}, \sigma_{\mathbf{z}}^2$ .



- In variational autoencoders (VAE) (Kingma and Welling, 2014), mapping x → q(z) is done using a neural network (encoder).
- The encoder performs so called *amortized inference*: When doing inference for a particular sample  $\mathbf{x}_i$ , we leverage the knowledge of the inference results for other samples. If two samples  $\mathbf{x}_i$  and  $\mathbf{x}_j$  are close to each other, the corresponding  $q(\mathbf{z}_i)$ ,  $q(\mathbf{z}_j)$  should be close as well.

• Our generative model is defined by the decoder.

$$\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$
  $\mathbf{x} = g(\mathbf{z}, \boldsymbol{ heta}) + \boldsymbol{\varepsilon}$   $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ 

• Encoder is a neural network that is trained to perform variational inference:

 ${\sf x} o q({\sf z})$ 

• For Gaussian approximation q(z), the neural network needs to produce:

$$\mathbf{x} 
ightarrow \mu_{\mathbf{z}}, \sigma_{\mathbf{z}}^2$$

- In practice, this is done using one neural network with two heads.
- The encoder is similar to the encoder in a bottleneck autoencoder but produces the mean and variance of the code z.
- The encoder and decoder are two components of the variational autoencoder.

#### Monte Carlo estimates of the objective function

• The first term of the objective function cannot computed analytically

$$\mathcal{F}(\boldsymbol{\theta}, \boldsymbol{\theta}_{q}) = \sum_{i=1}^{N} \underbrace{\int q(\mathbf{z}_{i}) \log p(\mathbf{x}_{i} \mid \mathbf{z}_{i}, \boldsymbol{\theta}) d\mathbf{z}_{i}}_{\text{needs approximations}} \underbrace{-\int q(\mathbf{z}_{i}) \log \frac{q(\mathbf{z}_{i})}{p(\mathbf{z}_{i})} d\mathbf{z}_{i}}_{\text{can be computed analytically}}$$

• Kingma and Welling (2014) proposed to use Monte Carlo estimates:

$$\int q(\mathbf{z}_i) \log \mathcal{N}(\mathbf{x}_i \mid g(\mathbf{z}_i, \boldsymbol{\theta}), \sigma^2 \mathbf{I}) d\mathbf{z}_i \approx \frac{1}{L} \sum_{l=1}^L \log \mathcal{N}(\mathbf{x}_i \mid g(\mathbf{z}_i^{(l)}, \boldsymbol{\theta}), \sigma^2 \mathbf{I})$$

where  $\mathbf{z}_{i}^{(l)}$  are drawn from  $q(\mathbf{z}_{i})$ . Using L = 1 works well in practice.



#### Computation of the objective function

$$\mathcal{F}(\boldsymbol{\theta}, \boldsymbol{\theta}_{q}) = \sum_{i=1}^{N} \underbrace{\log \mathcal{N}(\mathbf{x}_{i} \mid g(\mathbf{z}_{i}^{(l)}, \boldsymbol{\theta}), \sigma^{2}\mathbf{I})}_{\text{Monte Carlo estimate}} \underbrace{-\int q(\mathbf{z}_{i}) \log \frac{q(\mathbf{z}_{i})}{p(\mathbf{z}_{i})} d\mathbf{z}_{i}}_{\text{can be computed analytically}}$$

- For each training example x<sub>i</sub>:
  - compute means  $\mu_{\mathbf{z}_i}$  and  $\sigma_{\mathbf{z}_i}$  using the encoder
  - compute the second term analytically
  - draw L = 1 samples  $\mathbf{z}_i^{(l)}$  from  $q(\mathbf{z}_i) = \mathcal{N}(\mu_{\mathbf{z}_i}, \sigma_{\mathbf{z}_i}^2)$
  - propagate  $\mathbf{z}_{i}^{(l)}$  through the decoder and compute the first term
- Problem: We can use backpropagation to compute the derivatives wrt the parameters of the decoder but we need an extra trick to propagate derivatives through the encoder.



- We need a computational block that would
  - take as inputs  $\mu_z$  and  $\sigma_z$
  - produce a sample from distribution  $\boldsymbol{\epsilon} \sim \mathcal{N}(\mu_{\mathbf{z}_i}, \sigma_{\mathbf{z}_i})$
  - would be differentiable wrt  $\mu_z$  and  $\sigma_z$



- We need a computational block that would
  - take as inputs  $\mu_z$  and  $\sigma_z$
  - produce a sample from distribution  $\epsilon \sim \mathcal{N}(\mu_{z_i}, \sigma_{z_i})$
  - would be differentiable wrt  $\mu_z$  and  $\sigma_z$
- We can obtain this with the reparameterization trick:
  - Sample  $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
  - Compute  $\mathbf{z}_i = \mu_{\mathbf{z}_i} + \sigma_{\mathbf{z}_i} \boldsymbol{\epsilon}_i$
- Now we can also backpropagate through the sampling block and then further through the encoder.



- VAE training algorithm:
  - Take a mini-batch  $\{\mathbf{x}_i\}$  of training samples.
  - Use the encoder to compute means μ<sub>zi</sub> and standard deviations σ<sub>zi</sub> for each sample x<sub>i</sub> in the mini-batch.
  - Draw  $\epsilon_i \sim \mathcal{N}(0, \mathbf{I})$  and compute samples  $\mathbf{z}_i = \mu_{\mathbf{z}_i} + \sigma_{\mathbf{z}_i} \epsilon_i$
  - Propagate samples  $z_i$  through the decoder to compute reconstructions  $\hat{x}_i$ .
  - Compute the loss which is the negative of

$$\mathcal{F}(\boldsymbol{\theta}, \boldsymbol{\theta}_{q}) = \frac{1}{n} \sum_{i=1}^{n} \underbrace{\log \mathcal{N}(\mathbf{x}_{i} \mid g(\mathbf{z}_{i}^{(I)}, \boldsymbol{\theta}), \sigma^{2}\mathbf{I})}_{\text{Monte Carlo estimate}} \underbrace{-\int q(\mathbf{z}_{i}) \log \frac{q(\mathbf{z}_{i})}{p(\mathbf{z}_{i})} dz_{i}}_{\text{can be computed analytically}}$$

 Perform backpropagation and update the parameters of the encoder and the decoder.



• In the home assignment, we train a variational autoencoder on a synthetic (variance MNIST) dataset.

• In order to extract meaningful features for this dataset, we need to use a generator (decoder) that models the variances of pixel intensities:





- VAE is more complex than a simple bottleneck autoencoder. Do we need these complications?
- As we will see in the home assignment, VAEs are more powerful. In some problems when vanilla autoencoders fail, VAEs can develop useful representations.
- The problem of the vanilla autoencoder is the mean-squared error loss, which makes too simplistic assumptions about the data distribution.
- One advantage of VAE is in greater flexibility in defining the generative model.
- Note that denoising autoencoders are more powerful than standard autoencoders even though they also use the mean-squared error loss.

- The main benefit of VAEs is that we can encode data into a lower-dimensional representation.
- But VAEs are generative models and we can draw samples using VAEs.
- So far, the quality of the VAE-generated samples have not been very impressive (the samples as well as reconstructions usually look blurry).

Reconstructions



Images from (Tolstikhin et al., 2017)

#### Generated samples

#### Nouveau VAE (NVAE; Vahdat and Kautz, 2020)

• Vahdat and Kautz (2020) presented a VAE model that is able to generate high-quality images.

-

• It is a hierarchical latent variable model, that is there are multiple levels of latent variables.



Generated samples





## Home assignment

- In the home assignment, you will have to implement three types of autoencoders:
  - 1. Vanilla bottleneck autoencoder
  - 2. Denoising autoencoder
  - 3. Variational autoencoder

- Chapter 14 of the Deep Learning book
- Papers cited in the lecture slides