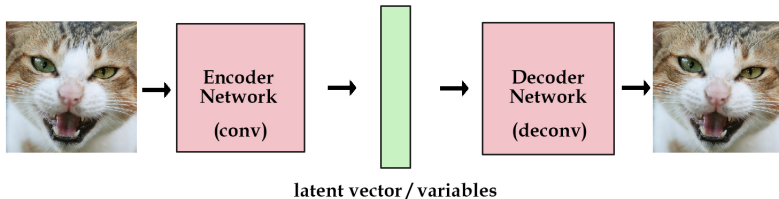


# Variational Autoencoder

# Introduction - Autoencoders



- ▶ Attempt to learn identity function
- ▶ Constrained in some way (e.g., small latent vector representation)
- ▶ Can generate new images by giving different latent vectors to trained network
- ▶ Variational: use probabilistic latent encoding

# Deep Learning Perspective

# Deep Learning Perspective

- ▶ Goal: Build a neural network that generates MNIST digits from random (Gaussian) noise

# Deep Learning Perspective

- ▶ Goal: Build a neural network that generates MNIST digits from random (Gaussian) noise
- ▶ Define two sub-networks: Encoder and Decoder

# Deep Learning Perspective

- ▶ Goal: Build a neural network that generates MNIST digits from random (Gaussian) noise
- ▶ Define two sub-networks: Encoder and Decoder
- ▶ Define a Loss Function

# Encoder

- ▶ A neural network  $q_{\theta}(z|x)$

# Encoder

- ▶ A neural network  $q_{\theta}(z|x)$
- ▶ Input: datapoint  $x$  (e.g.  $28 \times 28$ -pixel MNIST digit)



# Encoder

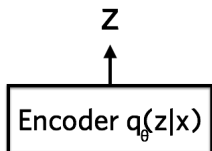
- ▶ A neural network  $q_{\theta}(z|x)$
- ▶ Input: datapoint  $x$  (e.g.  $28 \times 28$ -pixel MNIST digit)
- ▶ Output: encoding  $z$ , drawn from Gaussian density with parameters  $\theta$

# Encoder

- ▶ A neural network  $q_{\theta}(z|x)$
- ▶ Input: datapoint  $x$  (e.g.  $28 \times 28$ -pixel MNIST digit)
- ▶ Output: encoding  $z$ , drawn from Gaussian density with parameters  $\theta$
- ▶  $|z| \ll |x|$

# Encoder

- ▶ A neural network  $q_{\theta}(z|x)$
- ▶ Input: datapoint  $x$  (e.g.  $28 \times 28$ -pixel MNIST digit)
- ▶ Output: encoding  $z$ , drawn from Gaussian density with parameters  $\theta$
- ▶  $|z| \ll |x|$



▶ Data:  $x$

# Decoder

- ▶ A neural network  $p_{\phi}(x|z)$ , parameterized by  $\phi$

# Decoder

- ▶ A neural network  $p_{\phi}(x|z)$ , parameterized by  $\phi$
- ▶ Input: encoding  $z$ , output from encoder

# Decoder

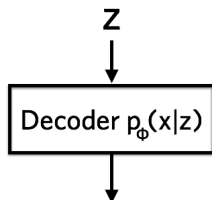
- ▶ A neural network  $p_{\phi}(x|z)$ , parameterized by  $\phi$
- ▶ Input: encoding  $z$ , output from encoder
- ▶ Output: reconstruction  $\tilde{x}$ , drawn from distribution of the data

# Decoder

- ▶ A neural network  $p_{\phi}(x|z)$ , parameterized by  $\phi$
- ▶ Input: encoding  $z$ , output from encoder
- ▶ Output: reconstruction  $\tilde{x}$ , drawn from distribution of the data
- ▶ E.g., output parameters for  $28 \times 28$  Bernoulli variables

# Decoder

- ▶ A neural network  $p_{\phi}(x|z)$ , parameterized by  $\phi$
- ▶ Input: encoding  $z$ , output from encoder
- ▶ Output: reconstruction  $\tilde{x}$ , drawn from distribution of the data
- ▶ E.g., output parameters for  $28 \times 28$  Bernoulli variables



- ▶ Reconstruction:  $\tilde{x}$



# Loss Function

- ▶  $\tilde{x}$  is reconstructed from  $z$  where  $|z| \ll |\tilde{x}|$

# Loss Function

- ▶  $\tilde{x}$  is reconstructed from  $z$  where  $|z| \ll |\tilde{x}|$
- ▶ How much information is lost when we go from  $x$  to  $z$  to  $\tilde{x}$ ?

# Loss Function

- ▶  $\tilde{x}$  is reconstructed from  $z$  where  $|z| \ll |\tilde{x}|$
- ▶ How much information is lost when we go from  $x$  to  $z$  to  $\tilde{x}$ ?
- ▶ Measure this with reconstruction log-likelihood:  $\log p_{\phi}(x|z)$

# Loss Function

- ▶  $\tilde{x}$  is reconstructed from  $z$  where  $|z| \ll |\tilde{x}|$
- ▶ How much information is lost when we go from  $x$  to  $z$  to  $\tilde{x}$ ?
- ▶ Measure this with reconstruction log-likelihood:  $\log p_{\phi}(x|z)$
- ▶ Measures how effectively the decoder has learned to reconstruct  $x$  given the latent representation  $z$

# Loss Function

- ▶ Loss function is negative reconstruction log-likelihood + regularizer

# Loss Function

- ▶ Loss function is negative reconstruction log-likelihood + regularizer
- ▶ Loss decomposes into term for each datapoint:

$$L(\theta, \phi) = \sum_{i=1}^N l_i(\theta, \phi)$$

# Loss Function

- ▶ Loss function is negative reconstruction log-likelihood + regularizer
- ▶ Loss decomposes into term for each datapoint:

$$L(\theta, \phi) = \sum_{i=1}^N l_i(\theta, \phi)$$

- ▶ Loss for datapoint  $x_i$ :

$$l_i(\theta, \phi) = -\mathbb{E}_{z \sim q_\theta(z|x_i)} [\log p_\phi(x_i|z)] + KL(q_\theta(z|x_i) || p(z))$$

# Loss Function

- ▶ Negative reconstruction log-likelihood:

$$-\mathbb{E}_{z \sim q_{\theta}(z|x_i)} [\log p_{\phi}(x_i|z)]$$



# Loss Function

- ▶ Negative reconstruction log-likelihood:

$$-\mathbb{E}_{z \sim q_{\theta}(z|x_i)} [\log p_{\phi}(x_i|z)]$$

- ▶ Encourages decoder to learn to reconstruct the data

# Loss Function

- ▶ Negative reconstruction log-likelihood:

$$-\mathbb{E}_{z \sim q_{\theta}(z|x_i)} [\log p_{\phi}(x_i|z)]$$

- ▶ Encourages decoder to learn to reconstruct the data
- ▶ Expectation taken over distribution of latent representations

# Loss Function

- ▶ KL Divergence as regularizer:

$$KL(q_{\theta}(z|x_i)||p(z)) = \mathbb{E}_{z \sim q_{\theta}(z|x_i)} [\log q_{\theta}(z|x_i) - \log p(z)]$$

# Loss Function

- ▶ KL Divergence as regularizer:

$$KL(q_{\theta}(z|x_i)||p(z)) = \mathbb{E}_{z \sim q_{\theta}(z|x_i)} [\log q_{\theta}(z|x_i) - \log p(z)]$$

- ▶ Measures information lost when using  $q_{\theta}$  to represent  $p$

# Loss Function

- ▶ KL Divergence as regularizer:

$$KL(q_{\theta}(z|x_i)||p(z)) = \mathbb{E}_{z \sim q_{\theta}(z|x_i)} [\log q_{\theta}(z|x_i) - \log p(z)]$$

- ▶ Measures information lost when using  $q_{\theta}$  to represent  $p$
- ▶ We will use  $p(z) = \mathcal{N}(\mathbf{0}, \mathbf{I})$

# Loss Function

- ▶ KL Divergence as regularizer:

$$KL(q_{\theta}(z|x_i)||p(z)) = \mathbb{E}_{z \sim q_{\theta}(z|x_i)} [\log q_{\theta}(z|x_i) - \log p(z)]$$

- ▶ Measures information lost when using  $q_{\theta}$  to represent  $p$
- ▶ We will use  $p(z) = \mathcal{N}(\mathbf{0}, \mathbf{I})$
- ▶ Encourages encoder to produce  $z$ 's that are close to standard normal distribution

# Loss Function

- ▶ KL Divergence as regularizer:

$$KL(q_{\theta}(z|x_i)||p(z)) = \mathbb{E}_{z \sim q_{\theta}(z|x_i)} [\log q_{\theta}(z|x_i) - \log p(z)]$$

- ▶ Measures information lost when using  $q_{\theta}$  to represent  $p$
- ▶ We will use  $p(z) = \mathcal{N}(\mathbf{0}, \mathbf{I})$
- ▶ Encourages encoder to produce  $z$ 's that are close to standard normal distribution
- ▶ Encoder learns a meaningful representation of MNIST digits

# Loss Function

- ▶ KL Divergence as regularizer:

$$KL(q_{\theta}(z|x_i)||p(z)) = \mathbb{E}_{z \sim q_{\theta}(z|x_i)} [\log q_{\theta}(z|x_i) - \log p(z)]$$

- ▶ Measures information lost when using  $q_{\theta}$  to represent  $p$
- ▶ We will use  $p(z) = \mathcal{N}(\mathbf{0}, \mathbf{I})$
- ▶ Encourages encoder to produce  $z$ 's that are close to standard normal distribution
- ▶ Encoder learns a meaningful representation of MNIST digits
- ▶ Representation for images of the same digit are close together in latent space



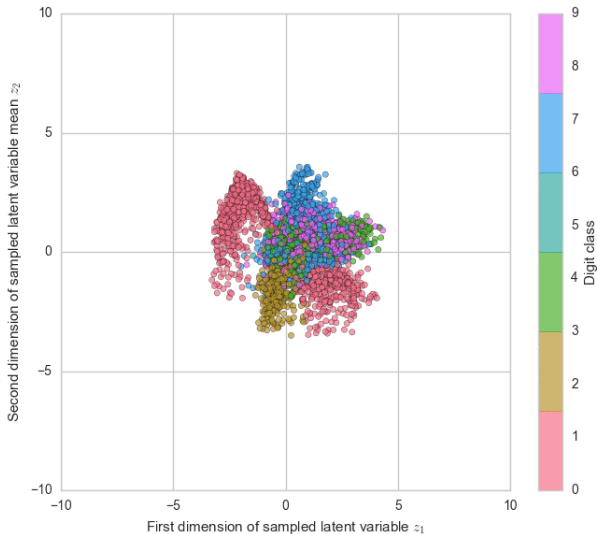
# Loss Function

- ▶ KL Divergence as regularizer:

$$KL(q_{\theta}(z|x_i)||p(z)) = \mathbb{E}_{z \sim q_{\theta}(z|x_i)} [\log q_{\theta}(z|x_i) - \log p(z)]$$

- ▶ Measures information lost when using  $q_{\theta}$  to represent  $p$
- ▶ We will use  $p(z) = \mathcal{N}(\mathbf{0}, \mathbf{I})$
- ▶ Encourages encoder to produce  $z$ 's that are close to standard normal distribution
- ▶ Encoder learns a meaningful representation of MNIST digits
- ▶ Representation for images of the same digit are close together in latent space
- ▶ Otherwise could “memorize” the data and map each observed datapoint to a distinct region of space

# MNIST latent variable space



# Reparameterization trick

- ▶ We want to use gradient descent to learn the model's parameters

# Reparameterization trick

- ▶ We want to use gradient descent to learn the model's parameters
- ▶ Given  $z$  drawn from  $q_{\theta}(z|x)$ , how do we take derivatives of (a function of)  $z$  w.r.t.  $\theta$ ?

# Reparameterization trick

- ▶ We want to use gradient descent to learn the model's parameters
- ▶ Given  $z$  drawn from  $q_\theta(z|x)$ , how do we take derivatives of (a function of)  $z$  w.r.t.  $\theta$ ?
- ▶ We can reparameterize:  $z = \mu + \sigma \odot \epsilon$

# Reparameterization trick

- ▶ We want to use gradient descent to learn the model's parameters
- ▶ Given  $z$  drawn from  $q_\theta(z|x)$ , how do we take derivatives of (a function of)  $z$  w.r.t.  $\theta$ ?
- ▶ We can reparameterize:  $z = \mu + \sigma \odot \epsilon$
- ▶  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , and  $\odot$  is element-wise product

# Reparameterization trick

- ▶ We want to use gradient descent to learn the model's parameters
- ▶ Given  $z$  drawn from  $q_{\theta}(z|x)$ , how do we take derivatives of (a function of)  $z$  w.r.t.  $\theta$ ?
- ▶ We can reparameterize:  $z = \mu + \sigma \odot \epsilon$
- ▶  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , and  $\odot$  is element-wise product
- ▶ Can take derivatives of (functions of)  $z$  w.r.t.  $\mu$  and  $\sigma$

# Reparameterization trick

- ▶ We want to use gradient descent to learn the model's parameters
- ▶ Given  $z$  drawn from  $q_\theta(z|x)$ , how do we take derivatives of (a function of)  $z$  w.r.t.  $\theta$ ?
- ▶ We can reparameterize:  $z = \mu + \sigma \odot \epsilon$
- ▶  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , and  $\odot$  is element-wise product
- ▶ Can take derivatives of (functions of)  $z$  w.r.t.  $\mu$  and  $\sigma$
- ▶ Output of  $q_\theta(z|x)$  is vector of  $\mu$ 's and vector of  $\sigma$ 's



# Summary

- ▶ Deep Learning objective is to minimize the loss function:

$$L(\theta, \phi) = \sum_{i=1}^N \left( -\mathbb{E}_{z \sim q_{\theta}(z|x_i)} [\log p_{\phi}(x_i|z)] + KL(q_{\theta}(z|x_i) || p(z)) \right)$$