

# Variational Autoencoders (VAEs) (chapter 20)

University of Padova, A.A. 2022/23

## Generative Models

**Aim: learn a model so that  $p_{\text{model}}(x)$  is as close as possible to  $p_{\text{data}}(x)$**

Two main families of approaches

- Explicit density estimation: explicitly define and solve for  $p_{\text{model}}(x)$ 
  - up to now we have studied models from the Boltzmann Machine family:
    - they use an **approximation** of the density via **Markov Chains (Gibbs sampling)**
    - Variational Autoencoders (VAEs) use an **approximation** of the density via a **variational approach (ELBO)**
- Implicit density estimation: learn a model that can sample from  $p_{\text{model}}(x)$  without explicitly defining it
  - Generative Adversarial Networks (GANs) use a **direct** approach based on **Game Theory**

Both VAEs and GANs exploit **differentiable generator networks**

# Evidence Lower Bound (ELBO)

$$\mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, q) = \log p(\mathbf{v}; \boldsymbol{\theta}) - D_{\text{KL}}(q(\mathbf{h} | \mathbf{v}) \| p(\mathbf{h} | \mathbf{v}; \boldsymbol{\theta})) \geq 0$$

where  $q$  is an arbitrary probability distribution over  $\mathbf{h}$  Kullback–Leibler divergence

$$\begin{aligned} \mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, q) &= \log p(\mathbf{v}; \boldsymbol{\theta}) - D_{\text{KL}}(q(\mathbf{h} | \mathbf{v}) \| p(\mathbf{h} | \mathbf{v}; \boldsymbol{\theta})) \\ &= \log p(\mathbf{v}; \boldsymbol{\theta}) - \mathbb{E}_{\mathbf{h} \sim q} \log \frac{q(\mathbf{h} | \mathbf{v})}{p(\mathbf{h} | \mathbf{v})} \\ &= \log p(\mathbf{v}; \boldsymbol{\theta}) - \mathbb{E}_{\mathbf{h} \sim q} \log \frac{q(\mathbf{h} | \mathbf{v})}{\frac{p(\mathbf{h}, \mathbf{v}; \boldsymbol{\theta})}{p(\mathbf{v}; \boldsymbol{\theta})}} \\ &= \log p(\mathbf{v}; \boldsymbol{\theta}) - \mathbb{E}_{\mathbf{h} \sim q} [\log q(\mathbf{h} | \mathbf{v}) - \log p(\mathbf{h}, \mathbf{v}; \boldsymbol{\theta}) + \log p(\mathbf{v}; \boldsymbol{\theta})] \\ &= -\mathbb{E}_{\mathbf{h} \sim q} [\log q(\mathbf{h} | \mathbf{v}) - \log p(\mathbf{h}, \mathbf{v}; \boldsymbol{\theta})] \end{aligned}$$

maximize  $\longrightarrow \boxed{\mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, q) = \mathbb{E}_{\mathbf{h} \sim q} [\log p(\mathbf{h}, \mathbf{v})] + H(q)}$

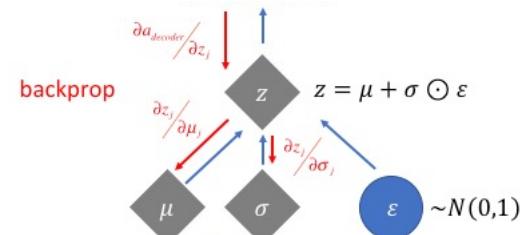
A.A. 2022/23: Deep Learning

## Differentiable Generator Networks

**Basic idea: parametrized functions for generating samples**

A differentiable **generator network** is

- a differentiable function  $g(z; \boldsymbol{\theta}^{(g)})$  transforming samples of latent variables  $z$  to
  - samples  $x$  (**direct**) or
  - to distributions over samples  $x$  (**indirect**)
- typically represented by a (deep) neural network
  - the architecture provides the family of possible distributions to sample from
  - the parameters select a distribution from within that family
- they allow the definition of a distribution  $p_g(x)$  and SGD optimization of criteria defined on  $p_g(x)$  via the **reparametrization trick**:
  - $\nabla$  of stochastic transformations of  $x$  ?
  - augment the neural network with extra inputs  $\varepsilon$  that are sampled from some simple probability distribution



A.A. 2022/23: Deep Learning

# Differentiable Generator Networks: Examples

- Example of **direct** sample generation
  - Drawing samples  $\mathbf{x}$  from a normal distribution with mean  $\mu$  and covariance  $\Sigma$ 
    - draw samples  $\mathbf{z}$  from a normal distribution with zero mean and *identity* covariance
    - feed  $\mathbf{z}$  to a simple generator network  $g(\cdot)$  made of a single affine layer:
- $\mathbf{x} = g(\mathbf{z}) = \mu + L\mathbf{z}$ ,  
where  $L\mathbf{L}^T = \Sigma$  (Cholesky decomposition of  $\Sigma$ )
- Example of **indirect** sample generation (define distribution over samples)
  - use  $g(\cdot)$  with sigmoid outputs to provide the mean parameters of Bernoulli distributions:

$$p(x_i = 1 | \mathbf{z}) = g(\mathbf{z})_i$$

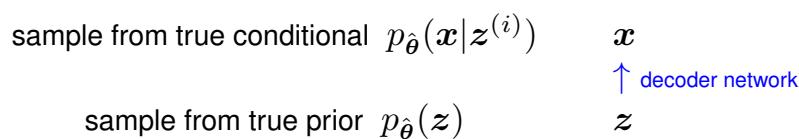
- when using  $g(\cdot)$  to define  $p(\mathbf{x}|\mathbf{z})$ , it is possible to impose a distribution over  $\mathbf{x}$  by marginalizing  $\mathbf{z}$ :

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}) = \sum_{\mathbf{z}} \underbrace{p(\mathbf{z})p(\mathbf{x}|\mathbf{z})}_{\text{product rule}} = \mathbb{E}_{\mathbf{z}} p(\mathbf{x}|\mathbf{z})$$

A.A. 2022/23: Deep Learning

## Variational Autoencoders

Basic idea:  
assume training data  $\mathcal{T} = \{\mathbf{x}^{(i)}\}_{i=1}^n$  is generated by latent variables  $\mathbf{z}$



Goal: estimate the true parameters  $\hat{\theta}$  of this generative model

Solution:

- pick a simple distribution for prior  $p_\theta(\mathbf{z})$ , e.g. Gaussian
- use a generator network for the posterior  $p_\theta(\mathbf{x}|\mathbf{z})$
- ...but need to maximize likelihood of  $p_\theta(\mathbf{x})$ :

$$p_\theta(\mathbf{x}) = \int p_\theta(\mathbf{z})p_\theta(\mathbf{x}|\mathbf{z})d\mathbf{z} \longrightarrow \text{intractable!!}$$

A.A. 2022/23: Deep Learning

# Variational Autoencoders

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{z}) p_{\theta}(\mathbf{x}|\mathbf{z}) d\mathbf{z} \rightarrow \text{intractable!!}$$

Also the posterior  $p_{\theta}(\mathbf{z}|\mathbf{x}) = p_{\theta}(\mathbf{x}|\mathbf{z})p_{\theta}(\mathbf{z})/p_{\theta}(\mathbf{x})$  is intractable due to  $p_{\theta}(\mathbf{x})$

**Solution:**

In addition to decoder network modeling  $p_{\theta}(\mathbf{x}|\mathbf{z})$ , define additional encoder network  $q_{\phi}(\mathbf{z}|\mathbf{x})$  that approximates  $p_{\theta}(\mathbf{z}|\mathbf{x})$  using... ELBO!

$$\begin{aligned} \mathcal{L}(q) &= \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} \log p_{\text{model}}(\mathbf{z}, \mathbf{x}) + \mathcal{H}(q(\mathbf{z} \mid \mathbf{x})) \\ &= \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} \log p_{\text{model}}(\mathbf{x} \mid \mathbf{z}) - D_{\text{KL}}(q(\mathbf{z} \mid \mathbf{x}) \parallel p_{\text{model}}(\mathbf{z})) \\ &\leq \log p_{\text{model}}(\mathbf{x}) \end{aligned}$$

decoder      encoder

forcing the encoder to be a Gaussian distribution  
(has closed-form solution)

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

Training: Maximize lower bound

# Variational Autoencoders

$$\log p_{\text{model}}(\mathbf{z}, \mathbf{x}) = \log p_{\text{model}}(\mathbf{x} \mid \mathbf{z}) + \log p_{\text{model}}(\mathbf{z}) \quad \mathcal{H}(q(\mathbf{z} \mid \mathbf{x})) = -\mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} \log q(\mathbf{z} \mid \mathbf{x})$$

$$\mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} \log p_{\text{model}}(\mathbf{z}, \mathbf{x}) + \mathcal{H}(q(\mathbf{z} \mid \mathbf{x}))$$

=

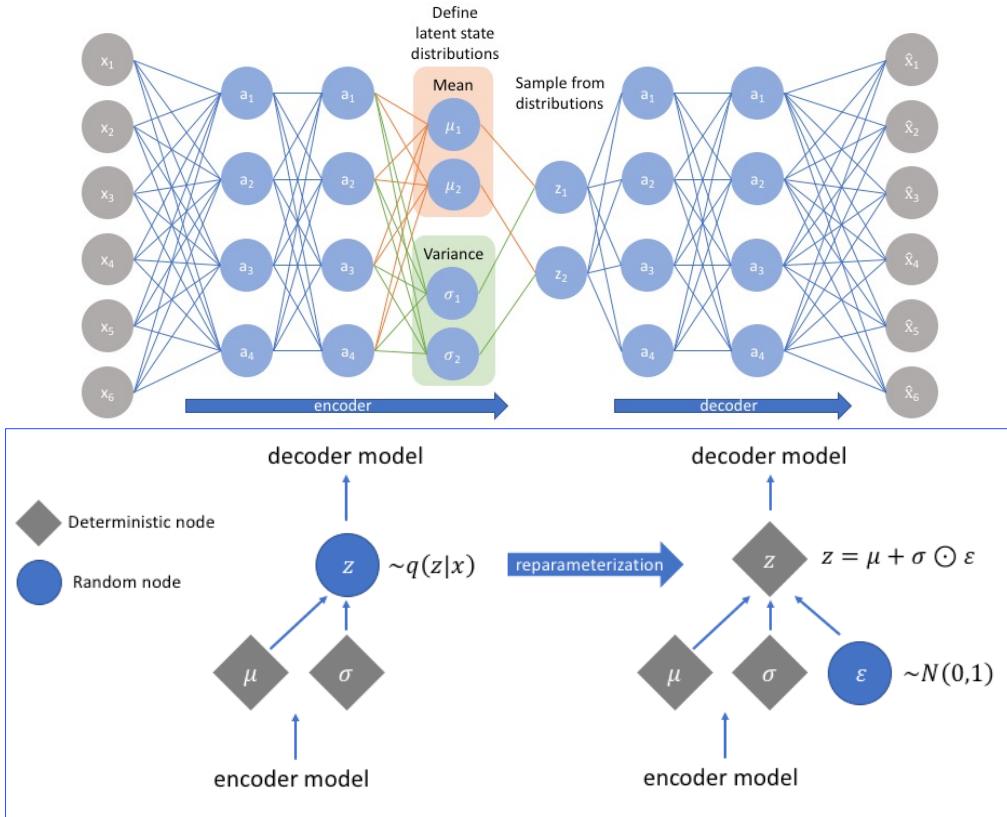
$$\mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} [\log p_{\text{model}}(\mathbf{x} \mid \mathbf{z}) + \log p_{\text{model}}(\mathbf{z})] - \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} \log q(\mathbf{z} \mid \mathbf{x})$$

=

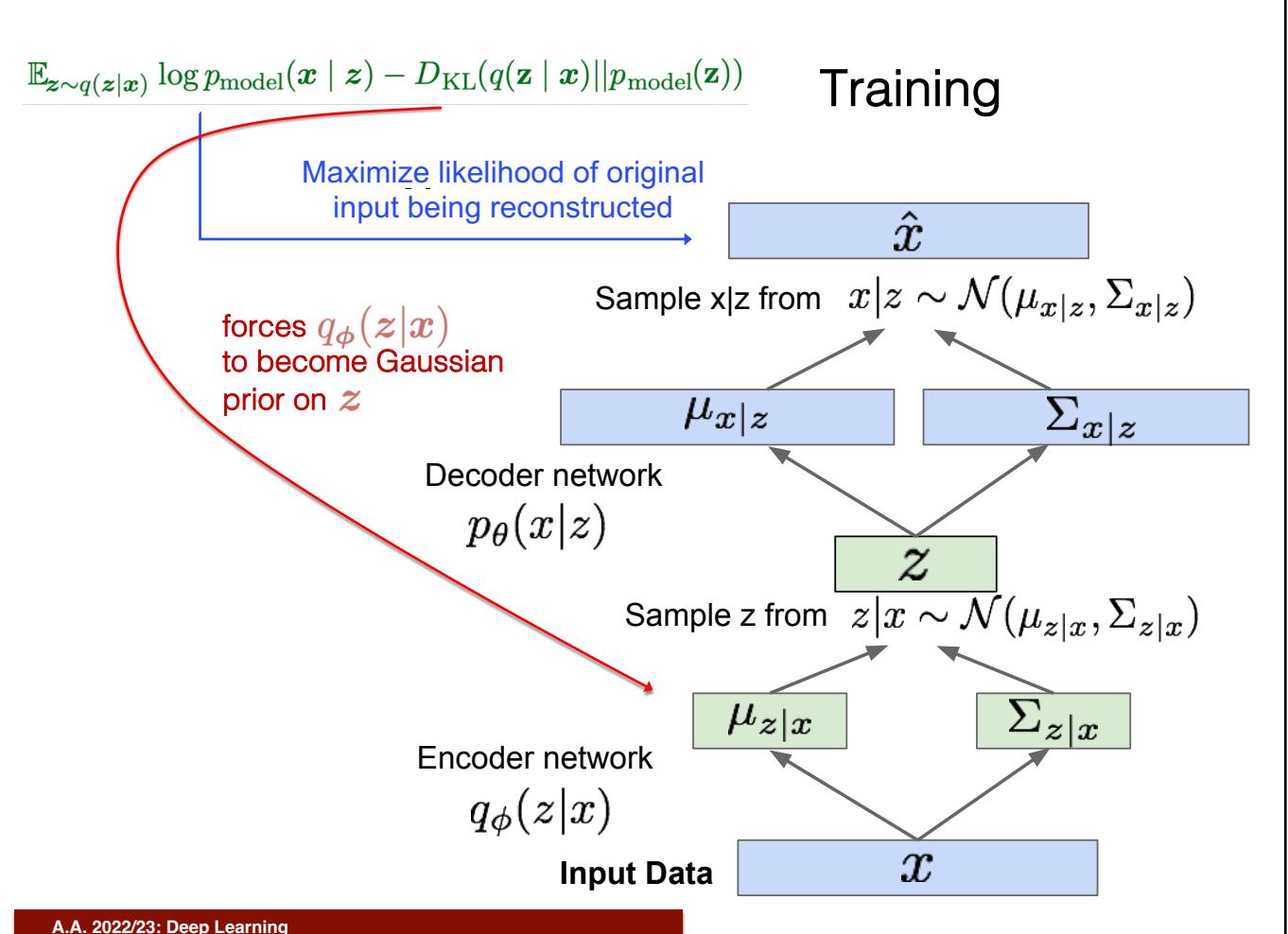
$$\mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} \log p_{\text{model}}(\mathbf{x} \mid \mathbf{z}) - \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} [\log q(\mathbf{z} \mid \mathbf{x}) - \log p_{\text{model}}(\mathbf{z})]$$

$D_{\text{KL}}(q(\mathbf{z} \mid \mathbf{x}) \parallel p_{\text{model}}(\mathbf{z}))$

# Encoding/Decoding implementation

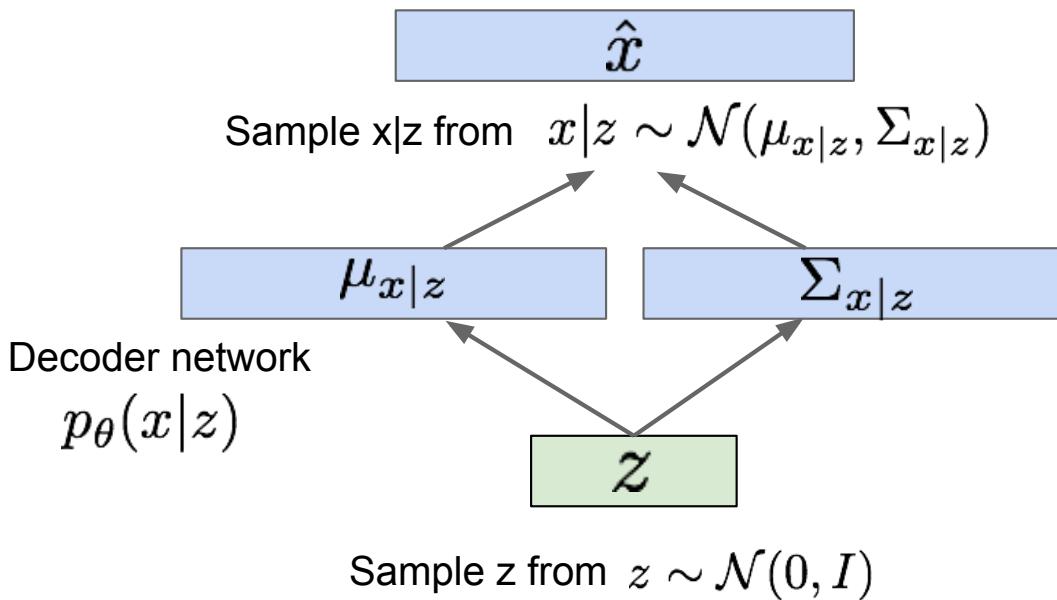


A.A. 2022/23: Deep Learning



A.A. 2022/23: Deep Learning

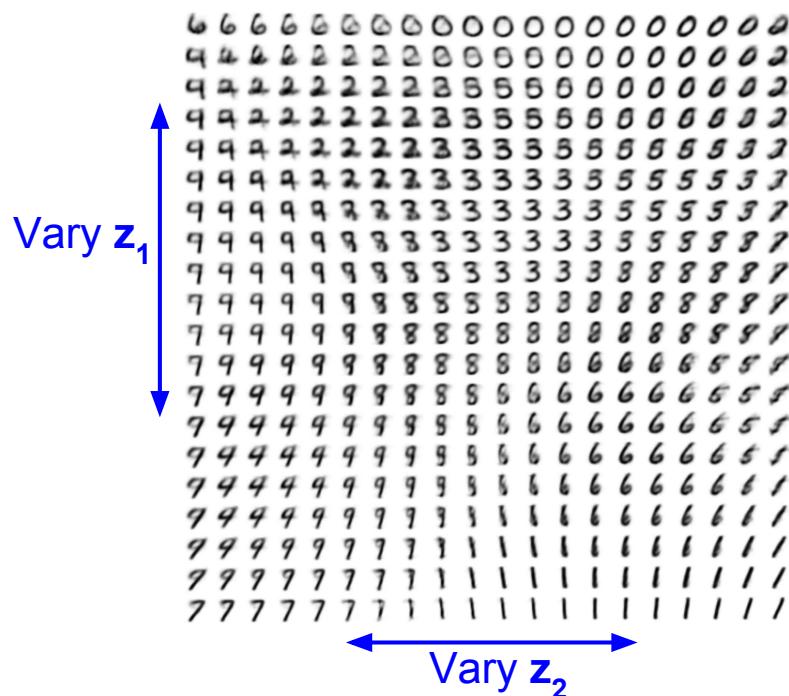
# Data Generation



A.A. 2022/23: Deep Learning

## Example of Data Generation

Data manifold for 2-d  $z$

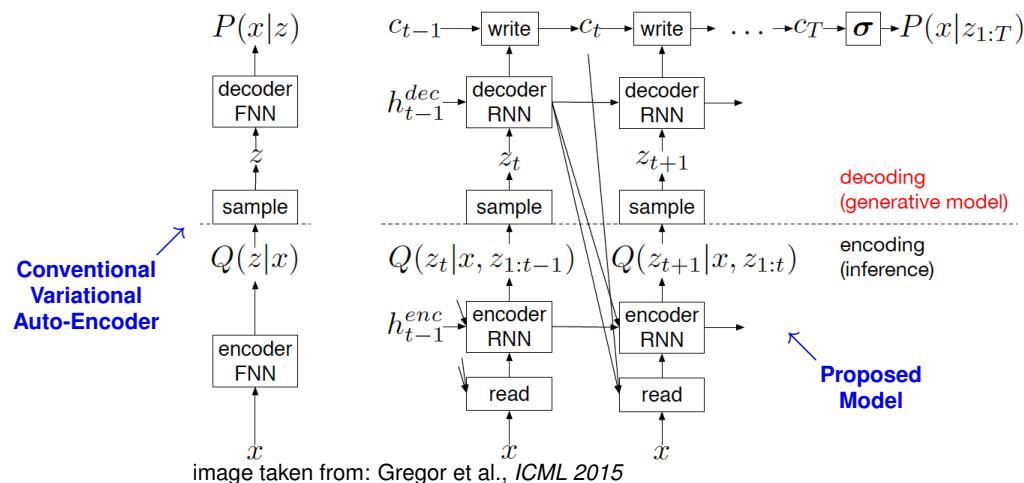


A.A. 2022/23: Deep Learning

# Application to Generation of Sequences: Draw

## Iterative construction of images by a spatial attention mechanism

- sequential auto-encoding framework where both the encoder and decoder are RNN
- strictly coupled encoder-decoder → encoder adapts to decoder output
- image generated incrementally by the decoder (as opposed to emitting it in a single step)
- dynamically updated attention mechanism: at each time-step “where to read”, “where to write”



A.A. 2022/23: Deep Learning

# Application to Generation of Sequences: Draw

Reading MNIST

A.A. 2022/23: Deep Learning

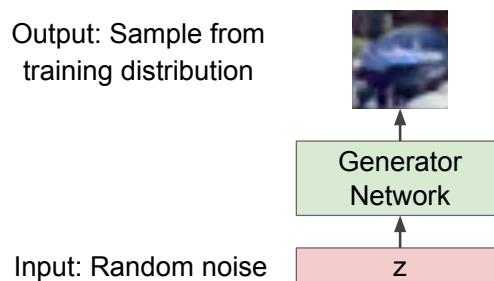
# Generative Adversarial Networks (GANs) (chapter 20)

University of Padova, A.A. 2022/23

## Generative Adversarial Networks

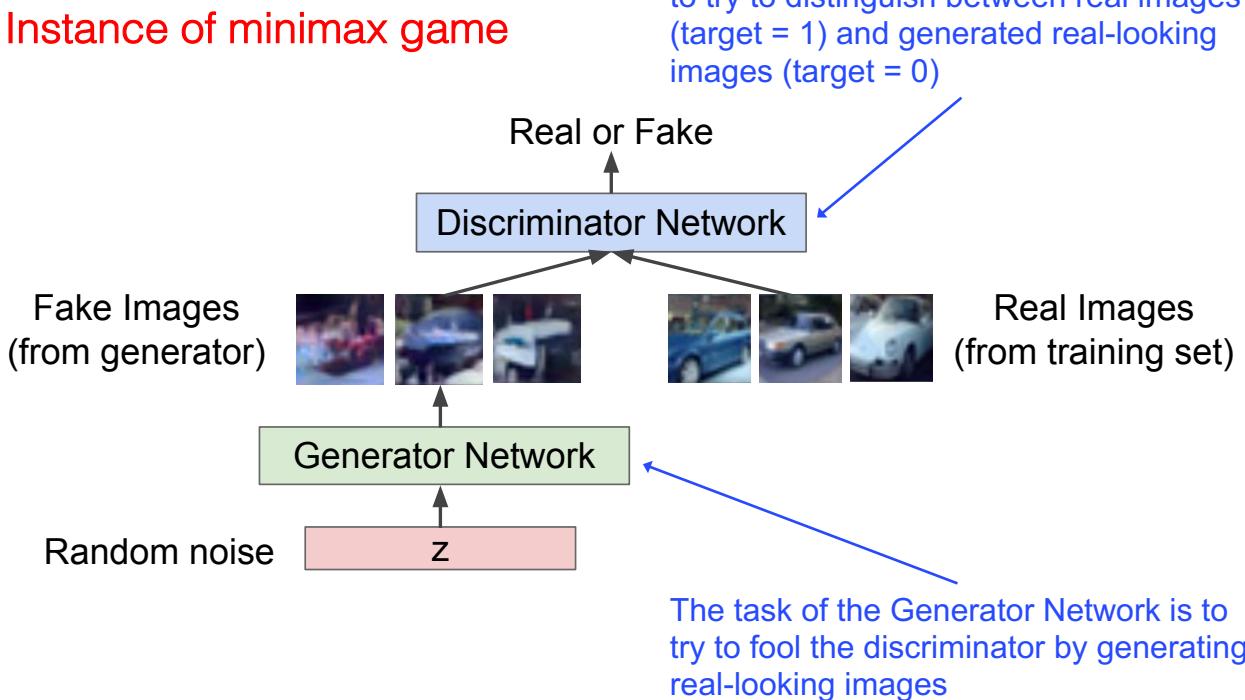
Unlike a VAE, a Generative Adversarial Network does not use an explicit density function, instead learns from the training data to draw samples

- Game-theoretic approach: learn to generate from training distribution through 2-player game
- Sample from a simple distribution, e.g. random noise, and learn transformation (generator network) to training distribution
- Transformation should be complex enough: **a (deep) neural network as generator network!**



# Generative Adversarial Networks: Basic Idea

Instance of minimax game



A.A. 2022/23: Deep Learning

## Objective Function

Train jointly in **minimax game**

Discriminator outputs likelihood in (0,1) of real image

Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[ \mathbb{E}_{x \sim p_{data}} \log \underbrace{D_{\theta_d}(x)}_{\text{Discriminator output for real data } x} + \mathbb{E}_{z \sim p(z)} \log \left( 1 - \underbrace{D_{\theta_d}(G_{\theta_g}(z))}_{\text{Discriminator output for generated fake data } G(z)} \right) \right]$$

For  $G$  fixed, the optimal discriminator  $D$  is

$$D_G^*(\mathbf{x}) = \frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})}$$



$$\mathbb{E}_{\mathbf{x} \sim p_{data}} \left[ \log \frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})} \right] + \mathbb{E}_{\mathbf{x} \sim p_g} \left[ \log \frac{p_g(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})} \right]$$

A.A. 2022/23: Deep Learning

# Training Procedure

## Training procedure

Alternate between:

### 1. Gradient ascent on discriminator

$$\max_{\theta_d} \left[ \mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

### 2. Gradient descent on generator

$$\min_{\theta_g} \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z)))$$

**Convergence** If  $G$  and  $D$  have enough capacity, and at each step the discriminator is allowed to reach its optimum given  $G$ , and  $p_g$  is updated so as to improve the criterion

$$\mathbb{E}_{x \sim p_{data}} [\log D_G^*(x)] + \mathbb{E}_{x \sim p_g} [\log(1 - D_G^*(x))]$$

then  $p_g$  converges to  $p_{data}$

# Training Procedure

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

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

**end for**

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

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

**end for**

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

# Training Procedure Modification

## Training procedure

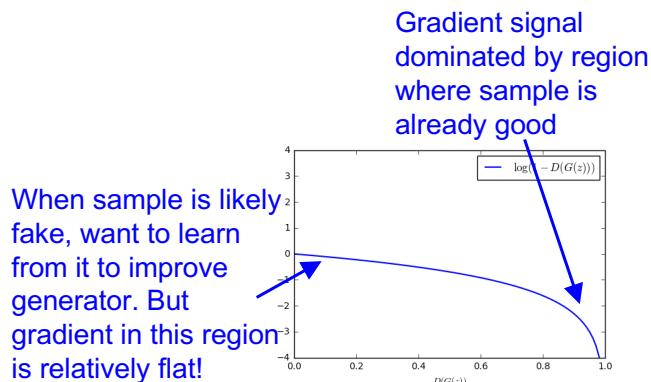
Alternate between:

### 1. Gradient ascent on discriminator

$$\max_{\theta_d} \left[ \mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

### 2. Gradient descent on generator Problem: does not work well in practice!

$$\min_{\theta_g} \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z)))$$



# Training Procedure Modification

## Training procedure

Alternate between:

### 1. Gradient ascent on discriminator

$$\max_{\theta_d} \left[ \mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

### 2. Gradient descent on generator

$$\min_{\theta_g} \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z)))$$

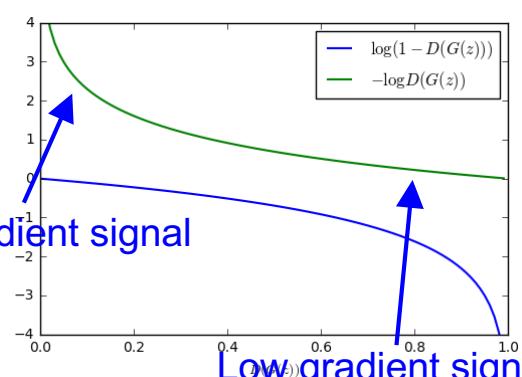
#### Gradient ascent

$$\max_{\theta_g} \mathbb{E}_{z \sim p(z)} -\log(D_{\theta_d}(G_{\theta_g}(z)))$$

maximize log-likelihood of discriminator being wrong

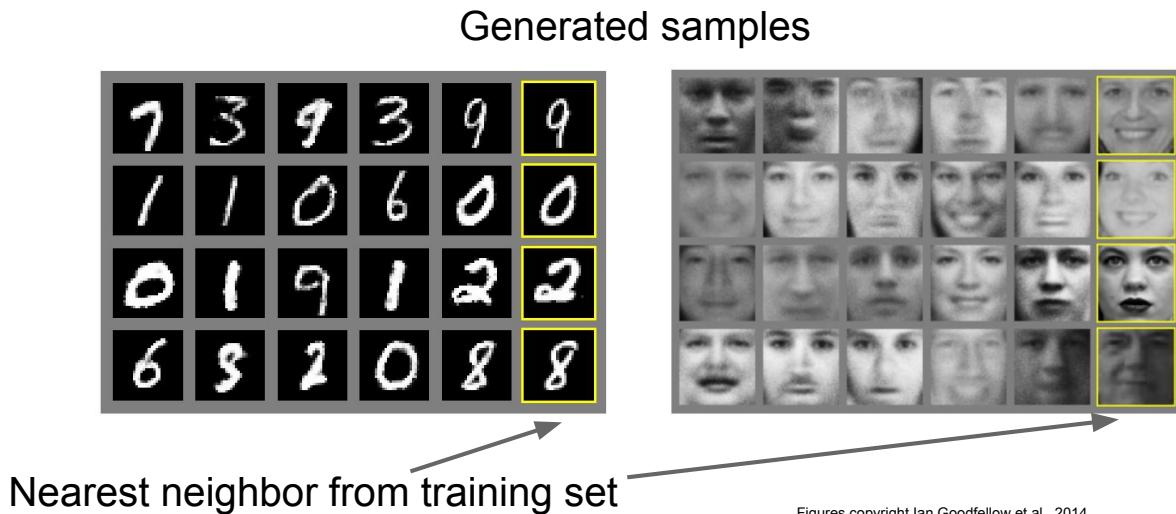
Change with

High gradient signal



# Generative Adversarial Networks: Examples

After training, use generator network to generate new images



## VAEs vs GANs

- Variational Autoencoders (VAE)
  - Optimize variational lower bound on likelihood
  - Useful latent representation that allows for inference queries
  - Usually sample quality is not the best
- Generative Adversarial Networks (GANs)
  - Game-theoretic approach
  - Best samples!
  - Can be tricky and unstable to train, no inference queries