

# Lecture 12: Deep Generative Models

André Martins, Francisco Melo, Mário Figueiredo



Deep Learning Course, Winter 2023-2024

# Announcements

- Thursday (January 4): guest lecture by Chrysoula Zerva.

Grande Auditório, Centro de Congressos, 13:30

# Today's Roadmap

Most of the course was about supervised learning.

Today we'll talk about **deep generative models** and **unsupervised learning**.

- Deep auto-regressive models
- Evidence lower bound (ELBO) and variational inference
- Variational auto-encoders
- Generative adversarial networks

# Which of these people is real?



(<http://www.whichfaceisreal.com>)

# Which of these people is real?



(<http://www.whichfaceisreal.com>)



# Which of these people is real?



(<http://www.whichfaceisreal.com>)

# Which of these people is real?



[<http://www.whichfaceisreal.com>)

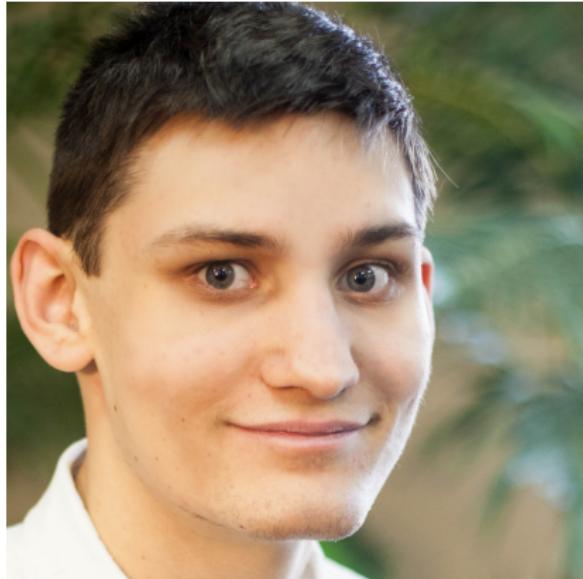


# Which of these people is real?



(<http://www.whichfaceisreal.com>)

# Which of these people is real?



[<http://www.whichfaceisreal.com>]



# Which of these people is real?



(<http://www.whichfaceisreal.com>)

# Which of these people is real?



(<http://www.whichfaceisreal.com>)



# Generative Modeling

- Modelling **high-dimensional** data (e.g., images) is a hard problem.

# Generative Modeling

- Modelling **high-dimensional** data (e.g., images) is a hard problem.
- **Deep generative models** are currently making remarkable progress.

# Generative Modeling

- Modelling **high-dimensional** data (e.g., images) is a hard problem.
- **Deep generative models** are currently making remarkable progress.
- **Unsupervised learning**: modelling  $\mathbb{P}(x)$ .
- **Supervised learning**: modelling  $\mathbb{P}(x, y)$ .

# Generative Modeling

- Modelling **high-dimensional** data (e.g., images) is a hard problem.
- **Deep generative models** are currently making remarkable progress.
- **Unsupervised learning**: modelling  $\mathbb{P}(\mathbf{x})$ .
- **Supervised learning**: modelling  $\mathbb{P}(\mathbf{x}, \mathbf{y})$ .
- Often, generative models use latent variables  $\mathbf{h}$ , such that

$$\mathbb{P}(\mathbf{x}) = \sum_{\mathbf{h}} \mathbb{P}(\mathbf{x}, \mathbf{h}) \quad \text{or} \quad \mathbb{P}(\mathbf{x}, \mathbf{y}) = \sum_{\mathbf{h}} \mathbb{P}(\mathbf{x}, \mathbf{y}, \mathbf{h})$$

where  $\mathbb{P}(\mathbf{x}, \mathbf{h})$  and  $\mathbb{P}(\mathbf{x}, \mathbf{y}, \mathbf{h})$  are “simple”.

# Examples of Deep Generative Models

- Auto-regressive networks
- Restricted Boltzmann machines
- Deep belief networks
- Deep Boltzmann machines
- Gaussian-Bernoulli RBMs
- Convolutional Boltzmann machines
- Variational auto-encoders (VAE)
- Generative adversarial networks (GAN)
- Denoising diffusion models

# Examples of Deep Generative Models

- Auto-regressive networks
- Restricted Boltzmann machines
- Deep belief networks
- Deep Boltzmann machines
- Gaussian-Bernoulli RBMs
- Convolutional Boltzmann machines
- Variational auto-encoders (VAE)
- Generative adversarial networks (GAN)
- Denoising diffusion models

# Outline

## ① Boltzmann Machines

## ② Variational Auto-Encoders

Variational Inference and ELBO

Gradients and Reparameterization Trick

## ③ Generative Adversarial Networks

## ④ Denoising Diffusion Models

## ⑤ Conclusions

# Energy Based Models

- A probability distribution (mixing observed and latent variables) via an **energy function**  $E(x, h; \theta)$ :

$$\mathbb{P}_{\theta}(x, h) = \frac{\exp(-E(x, h; \theta))}{Z(\theta)}, \quad \text{where } Z(\theta) = \sum_{x, h} \exp(-E(x, h; \theta))$$

# Energy Based Models

- A probability distribution (mixing observed and latent variables) via an **energy function**  $E(x, h; \theta)$ :

$$\mathbb{P}_{\theta}(x, h) = \frac{\exp(-E(x, h; \theta))}{Z(\theta)}, \quad \text{where } Z(\theta) = \sum_{x, h} \exp(-E(x, h; \theta))$$

- Maximizing probability corresponds to minimizing the energy.

# Energy Based Models

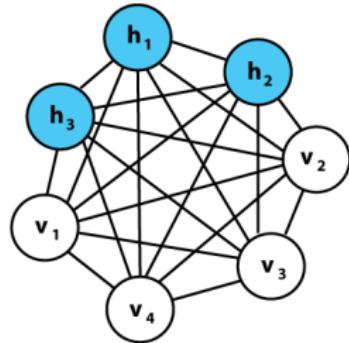
- A probability distribution (mixing observed and latent variables) via an **energy function**  $E(x, h; \theta)$ :

$$\mathbb{P}_{\theta}(x, h) = \frac{\exp(-E(x, h; \theta))}{Z(\theta)}, \quad \text{where } Z(\theta) = \sum_{x, h} \exp(-E(x, h; \theta))$$

- Maximizing probability corresponds to **minimizing** the energy.
- Challenges:
  - Computing the **partition function**  $Z(\theta)$
  - Computing the **evidence**  $\mathbb{P}_{\theta}(x) = \sum_h \mathbb{P}_{\theta}(x, h)$
  - Computing the **posterior**  $\mathbb{P}_{\theta}(h | x) = \mathbb{P}_{\theta}(x, h) / \mathbb{P}_{\theta}(x)$
  - Sampling from  $\mathbb{P}_{\theta}(x, h)$  or from  $\mathbb{P}_{\theta}(x)$

# Boltzmann Machine (Ackley et al., 1985)

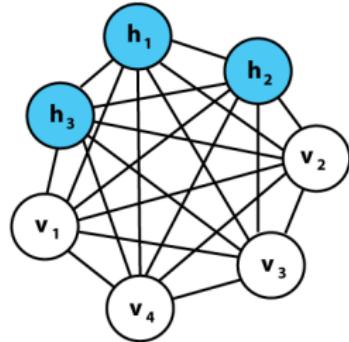
- Energy-based model over binary vectors
- Some variables are observed ( $v$ ), others are latent/hidden ( $h$ )



# Boltzmann Machine (Ackley et al., 1985)

- Energy-based model over binary vectors
- Some variables are observed ( $v$ ), others are latent/hidden ( $h$ )
- Probability distribution over  $(v, h) \in \{0, 1\}^{N+M}$ :

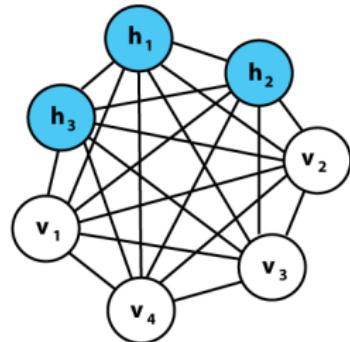
$$\mathbb{P}_{\theta}(v, h) = \frac{\exp(-E(v, h; \theta))}{Z(\theta)}$$



# Boltzmann Machine (Ackley et al., 1985)

- Energy-based model over binary vectors
- Some variables are observed ( $v$ ), others are latent/hidden ( $h$ )
- Probability distribution over  $(v, h) \in \{0, 1\}^{N+M}$ :

$$\mathbb{P}_{\theta}(v, h) = \frac{\exp(-E(v, h; \theta))}{Z(\theta)}$$



- Energy function, with  $\theta = (\mathbf{R}, \mathbf{W}, \mathbf{S}, \mathbf{b}, \mathbf{c})$ ,

$$E(v, h; \theta) = -v^T \mathbf{R} v - v^T \mathbf{W} h - h^T \mathbf{S} h - b^T v - c^T h$$

# Boltzmann Machine

- Boltzmann machine (BM): universal approximator of probability mass functions over discrete variables (Le Roux and Bengio, 2008)

# Boltzmann Machine

- Boltzmann machine (BM): universal approximator of probability mass functions over discrete variables (Le Roux and Bengio, 2008)
- Emulates the idea in Hebbian learning: “neurons that fire together wire together.”

# Boltzmann Machine

- Boltzmann machine (BM): universal approximator of probability mass functions over discrete variables (Le Roux and Bengio, 2008)
- Emulates the idea in Hebbian learning: “neurons that fire together wire together.”
- However, in general,
  - Sampling is hard,
  - Inference is hard,
  - Learning is hard.

# Boltzmann Machine

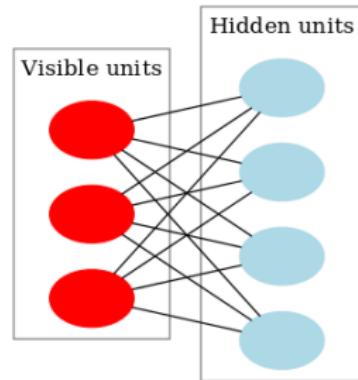
- Boltzmann machine (BM): universal approximator of probability mass functions over discrete variables (Le Roux and Bengio, 2008)
- Emulates the idea in Hebbian learning: “neurons that fire together wire together.”
- However, in general,
  - Sampling is hard,
  - Inference is hard,
  - Learning is hard.
- In summary: learning a general BM is usually very challenging, so we typically use a particular version.

# Boltzmann Machine

- Boltzmann machine (BM): universal approximator of probability mass functions over discrete variables (Le Roux and Bengio, 2008)
- Emulates the idea in Hebbian learning: “neurons that fire together wire together.”
- However, in general,
  - Sampling is hard,
  - Inference is hard,
  - Learning is hard.
- In summary: learning a general BM is usually very challenging, so we typically use a particular version.
- Not covered here, but check Goodfellow et al. (2016, Chapter 18).

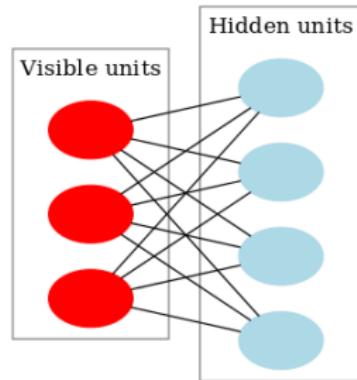
# Particular Case: Restricted Boltzmann Machines

- Also called [harmonium](#) (Smolensky, 1986)
- A layer of observable variables
- A single layer of latent variables.



# Particular Case: Restricted Boltzmann Machines

- Also called **harmonium** (Smolensky, 1986)
- A layer of observable variables
- A single layer of latent variables.

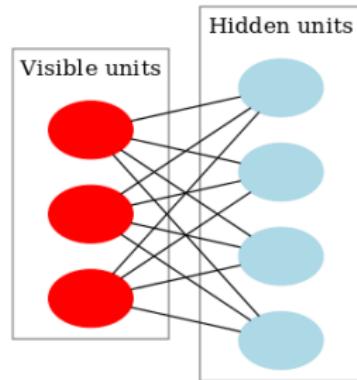


- Bipartite graph, **no intra-layer connections**:  $\mathbf{R} = 0$ ;  $\mathbf{S} = 0$ .
- The energy function becomes:

$$E(\mathbf{v}, \mathbf{h}; \theta) = -\mathbf{v}^\top \mathbf{W} \mathbf{h} - \mathbf{b}^\top \mathbf{v} - \mathbf{c}^\top \mathbf{h}$$

# Particular Case: Restricted Boltzmann Machines

- Also called **harmonium** (Smolensky, 1986)
- A layer of observable variables
- A single layer of latent variables.



- Bipartite graph, **no intra-layer connections**:  $\mathbf{R} = 0$ ;  $\mathbf{S} = 0$ .
- The energy function becomes:

$$E(\mathbf{v}, \mathbf{h}; \theta) = -\mathbf{v}^\top \mathbf{W} \mathbf{h} - \mathbf{b}^\top \mathbf{v} - \mathbf{c}^\top \mathbf{h}$$

- What is the advantage?

# Restricted Boltzmann Machines

- Unfortunately, the partition function  $Z(\theta)$  is still intractable

# Restricted Boltzmann Machines

- Unfortunately, the partition function  $Z(\theta)$  is still intractable
- ... but the **conditionals**  $\mathbb{P}_\theta(h \mid v)$  and  $\mathbb{P}_\theta(v \mid h)$  are now tractable!
  - easy to compute!
  - easy to sample!
  - using Markov-Chain Monte Carlo (MCMC) with Gibbs sampling.

# Restricted Boltzmann Machines

- Unfortunately, the partition function  $Z(\theta)$  is still intractable
- ... but the **conditionals**  $\mathbb{P}_\theta(h \mid v)$  and  $\mathbb{P}_\theta(v \mid h)$  are now tractable!
  - easy to compute!
  - easy to sample!
  - using Markov-Chain Monte Carlo (MCMC) with Gibbs sampling.
- Why are these easy? **Conditional independence** (next slide)

# Restricted Boltzmann Machines

- Without intra-layer connections,  $h_1, \dots, h_N$  are **conditionally independent** given  $v$ :

$$\mathbb{P}_{\theta}(h \mid v) = \prod_{j=1}^M \mathbb{P}(h_j \mid v)$$

where

$$\mathbb{P}_{\theta}(h_j = 1 \mid v) = \sigma(c_j + (\mathbf{W}v)_j), \quad j = 1, \dots, M.$$

# Restricted Boltzmann Machines

- Without intra-layer connections,  $h_1, \dots, h_N$  are **conditionally independent** given  $v$ :

$$\mathbb{P}_{\theta}(h | v) = \prod_{j=1}^M \mathbb{P}(h_j | v)$$

where

$$\mathbb{P}_{\theta}(h_j = 1 | v) = \sigma(c_j + (\mathbf{W}v)_j), \quad j = 1, \dots, M.$$

- Reciprocally for  $\mathbb{P}_{\theta}(v | h)$ .

# Restricted Boltzmann Machines

- Without intra-layer connections,  $h_1, \dots, h_N$  are **conditionally independent** given  $v$ :

$$\mathbb{P}_{\theta}(h | v) = \prod_{j=1}^M \mathbb{P}(h_j | v)$$

where

$$\mathbb{P}_{\theta}(h_j = 1 | v) = \sigma(c_j + (\mathbf{W}v)_j), \quad j = 1, \dots, M.$$

- Reciprocally for  $\mathbb{P}_{\theta}(v | h)$ .
- RBM $s$  are relatively easy to train by approximating  $Z(\theta)$  (see Goodfellow et al. (2016, Chapter 18)).

# Restricted Boltzmann Machines

- Without intra-layer connections,  $h_1, \dots, h_N$  are **conditionally independent** given  $v$ :

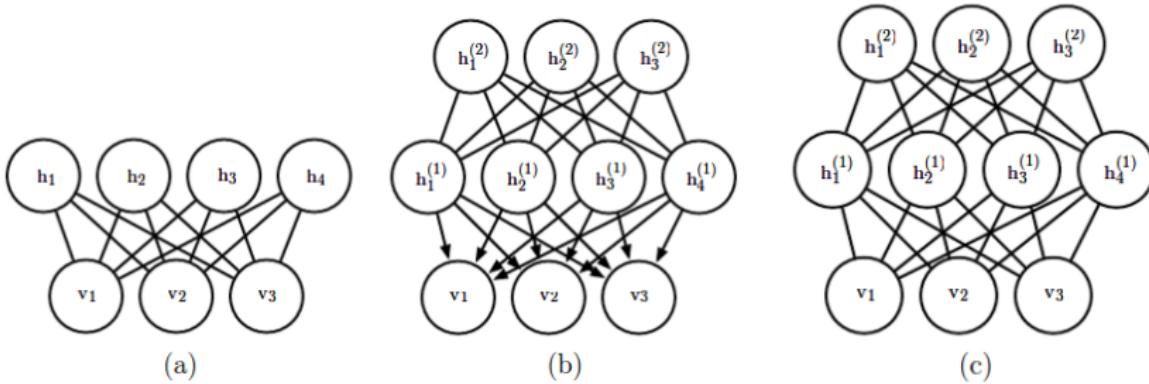
$$\mathbb{P}_{\theta}(h | v) = \prod_{j=1}^M \mathbb{P}(h_j | v)$$

where

$$\mathbb{P}_{\theta}(h_j = 1 | v) = \sigma(c_j + (\mathbf{W}v)_j), \quad j = 1, \dots, M.$$

- Reciprocally for  $\mathbb{P}_{\theta}(v | h)$ .
- RBM $s$  are relatively easy to train by approximating  $Z(\theta)$  (see Goodfellow et al. (2016, Chapter 18)).
- RBM $s$  may be stacked to form deeper models.

# Some RBM's Friends



(Image from Goodfellow et al. (2016))

- (a) **Restricted Boltzmann machine (RBM)**
- (b) **Deep belief network (DBN)**: hybrid directed/undirected GM with multiple latent layers
- (c) **Deep Boltzmann machine (DBM)**: undirected GM with several layers of latent variables.

# Examples of Deep Generative Models

- Restricted Boltzmann Machines
- Deep Belief Networks
- Deep Boltzmann Machines
- Gaussian-Bernoulli RBMs
- Convolutional Boltzmann Machines
- Sigmoid Belief Nets
- Variational Auto-Encoders
- Generative Adversarial Networks
- Denoising diffusion models

# Examples of Deep Generative Models

- Restricted Boltzmann Machines ✓
- Deep Belief Networks
- Deep Boltzmann Machines
- Gaussian-Bernoulli RBMs
- Convolutional Boltzmann Machines
- Sigmoid Belief Nets
- Variational Auto-Encoders
- Generative Adversarial Networks
- Denoising diffusion models

# Outline

① Boltzmann Machines

② Variational Auto-Encoders

Variational Inference and ELBO

Gradients and Reparameterization Trick

③ Generative Adversarial Networks

④ Denoising Diffusion Models

⑤ Conclusions

# Differentiable Generator Networks

- Several recent models are based on **differentiable generator networks**.

# Differentiable Generator Networks

- Several recent models are based on **differentiable generator networks**.
- This is a differentiable function  $G(\mathbf{h}; \boldsymbol{\theta})$  that maps latent variables  $\mathbf{h}$  into sample reconstructions  $\mathbf{x}$  (or distributions  $\mathbb{P}_{\boldsymbol{\theta}}(\mathbf{x} | \mathbf{h})$ ).

# Differentiable Generator Networks

- Several recent models are based on **differentiable generator networks**.
- This is a differentiable function  $G(\mathbf{h}; \theta)$  that maps latent variables  $\mathbf{h}$  into sample reconstructions  $\mathbf{x}$  (or distributions  $\mathbb{P}_\theta(\mathbf{x} | \mathbf{h})$ ).
- This idea underlies
  - **Variational auto-encoders** (VAE)
  - **Generative adversarial networks** (GAN)

# Variational Auto-Encoders

- Many latent variable models have:
  - intractable evidence  $\mathbb{P}(x)$
  - intractable posterior  $\mathbb{P}(h \mid x)$ .

# Variational Auto-Encoders

- Many latent variable models have:
  - intractable evidence  $\mathbb{P}(x)$
  - intractable posterior  $\mathbb{P}(h | x)$ .
- **Variational inference** is used to approximate these quantities.

# Variational Auto-Encoders

- Many latent variable models have:
  - intractable evidence  $\mathbb{P}(x)$
  - intractable posterior  $\mathbb{P}(h | x)$ .
- **Variational inference** is used to approximate these quantities.
- Widely used in Bayesian inference, topic models, etc...

# Variational Auto-Encoders

- Many latent variable models have:
  - intractable evidence  $\mathbb{P}(x)$
  - intractable posterior  $\mathbb{P}(h | x)$ .
- **Variational inference** is used to approximate these quantities.
- Widely used in Bayesian inference, topic models, etc...
- **Auto-encoders**: effective to learn data representations or codes, i.e.,

$$x \longrightarrow h \longrightarrow \hat{x}$$

# Variational Auto-Encoders

- Many latent variable models have:
  - intractable evidence  $\mathbb{P}(x)$
  - intractable posterior  $\mathbb{P}(h | x)$ .
- **Variational inference** is used to approximate these quantities.
- Widely used in Bayesian inference, topic models, etc...
- **Auto-encoders**: effective to learn data representations or codes, i.e.,

$$x \longrightarrow h \longrightarrow \hat{x}$$

- **Key idea:** combine auto-encoders with variational inference.

# Assumptions

- We assume that:
  - prior  $\mathbb{P}_\theta(h)$  is **tractable** (e.g.. standard Gaussian, Bernoulli, ...)

# Assumptions

- We assume that:
  - prior  $\mathbb{P}_\theta(h)$  is **tractable** (e.g.. standard Gaussian, Bernoulli, ...)
  - conditional  $\mathbb{P}_\theta(x | h)$  is **tractable** (e.g., feed-forward NN or an RNN).

# Assumptions

- We assume that:
  - prior  $\mathbb{P}_\theta(h)$  is **tractable** (e.g.. standard Gaussian, Bernoulli, ...)
  - conditional  $\mathbb{P}_\theta(x | h)$  is **tractable** (e.g., feed-forward NN or an RNN).
- However,
  - evidence  $\mathbb{P}_\theta(x)$  (i.e. marginalizing out  $h$ ) is still **intractable**
  - posterior  $\mathbb{P}_\theta(h | x)$  is still **intractable**.

# Assumptions

- We assume that:
  - prior  $\mathbb{P}_\theta(h)$  is **tractable** (e.g.. standard Gaussian, Bernoulli, ...)
  - conditional  $\mathbb{P}_\theta(x | h)$  is **tractable** (e.g., feed-forward NN or an RNN).
- However,
  - evidence  $\mathbb{P}_\theta(x)$  (i.e. marginalizing out  $h$ ) is still **intractable**
  - posterior  $\mathbb{P}_\theta(h | x)$  is still **intractable**.
- Next: using **variational inference** to approximate these computations

# Outline

① Boltzmann Machines

② Variational Auto-Encoders

Variational Inference and ELBO

Gradients and Reparameterization Trick

③ Generative Adversarial Networks

④ Denoising Diffusion Models

⑤ Conclusions

## Recap: Shannon's Entropy

- Let  $\mathbb{P}$  be a distribution over  $\mathcal{X}$ . The **entropy** of  $\mathbb{P}$  is

$$H(\mathbb{P}) = - \sum_{x \in \mathcal{X}} \mathbb{P}(x) \log \mathbb{P}(x)$$

## Recap: Shannon's Entropy

- Let  $\mathbb{P}$  be a distribution over  $\mathcal{X}$ . The **entropy** of  $\mathbb{P}$  is

$$H(\mathbb{P}) = - \sum_{x \in \mathcal{X}} \mathbb{P}(x) \log \mathbb{P}(x) = \mathbb{E}_{\mathbb{P}(x)}[-\log \mathbb{P}(x)]$$

## Recap: Shannon's Entropy

- Let  $\mathbb{P}$  be a distribution over  $\mathcal{X}$ . The **entropy** of  $\mathbb{P}$  is

$$H(\mathbb{P}) = - \sum_{x \in \mathcal{X}} \mathbb{P}(x) \log \mathbb{P}(x) = \mathbb{E}_{\mathbb{P}(x)}[-\log \mathbb{P}(x)]$$

- Always **non-negative**:  $H(\mathbb{P}) \geq 0$
- $H(\mathbb{P}) = 0$  iff  $\mathbb{P}(x) = 1$ , for some  $x$  and  $\mathbb{P}(x') = 0$ , for any  $x' \neq x$ .

## Recap: Shannon's Entropy

- Let  $\mathbb{P}$  be a distribution over  $\mathcal{X}$ . The **entropy** of  $\mathbb{P}$  is

$$H(\mathbb{P}) = - \sum_{x \in \mathcal{X}} \mathbb{P}(x) \log \mathbb{P}(x) = \mathbb{E}_{\mathbb{P}(x)}[-\log \mathbb{P}(x)]$$

- Always **non-negative**:  $H(\mathbb{P}) \geq 0$
- $H(\mathbb{P}) = 0$  iff  $\mathbb{P}(x) = 1$ , for some  $x$  and  $\mathbb{P}(x') = 0$ , for any  $x' \neq x$ .
- Upper bounded**:  $H(\mathbb{P}) \leq \log |\mathcal{X}|$
- $H(\mathbb{P}) = \log |\mathcal{X}|$  iff  $\mathbb{P}(x) = 1/|\mathcal{X}|$  (uniform distribution)

## Recap: Shannon's Entropy

- Let  $\mathbb{P}$  be a distribution over  $\mathcal{X}$ . The **entropy** of  $\mathbb{P}$  is

$$H(\mathbb{P}) = - \sum_{x \in \mathcal{X}} \mathbb{P}(x) \log \mathbb{P}(x) = \mathbb{E}_{\mathbb{P}(x)}[-\log \mathbb{P}(x)]$$

- Always **non-negative**:  $H(\mathbb{P}) \geq 0$
- $H(\mathbb{P}) = 0$  iff  $\mathbb{P}(x) = 1$ , for some  $x$  and  $\mathbb{P}(x') = 0$ , for any  $x' \neq x$ .
- Upper bounded**:  $H(\mathbb{P}) \leq \log |\mathcal{X}|$
- $H(\mathbb{P}) = \log |\mathcal{X}|$  iff  $\mathbb{P}(x) = 1/|\mathcal{X}|$  (uniform distribution)
- Intuition**:  $H(\mathbb{P})$  measures how close to uniform the distribution is

## Recap: Shannon's Entropy

- Let  $\mathbb{P}$  be a distribution over  $\mathcal{X}$ . The **entropy** of  $\mathbb{P}$  is

$$H(\mathbb{P}) = - \sum_{x \in \mathcal{X}} \mathbb{P}(x) \log \mathbb{P}(x) = \mathbb{E}_{\mathbb{P}(x)}[-\log \mathbb{P}(x)]$$

- Always **non-negative**:  $H(\mathbb{P}) \geq 0$
- $H(\mathbb{P}) = 0$  iff  $\mathbb{P}(x) = 1$ , for some  $x$  and  $\mathbb{P}(x') = 0$ , for any  $x' \neq x$ .
- Upper bounded**:  $H(\mathbb{P}) \leq \log |\mathcal{X}|$
- $H(\mathbb{P}) = \log |\mathcal{X}|$  iff  $\mathbb{P}(x) = 1/|\mathcal{X}|$  (uniform distribution)
- Intuition**:  $H(\mathbb{P})$  measures how close to uniform the distribution is
- Coding perspective**: expected number of bits (using  $\log_2$ ) to optimally encode  $x \sim \mathbb{P}(x)$

## Recap: Kullback-Leibler (KL) Divergence

- Let  $\mathbb{P}$  and  $\mathbb{Q}$  be two distributions over  $\mathcal{X}$ .

$$KL(\mathbb{P} \parallel \mathbb{Q}) = \sum_x \mathbb{P}(x) \log \frac{\mathbb{P}(x)}{\mathbb{Q}(x)}$$

## Recap: Kullback-Leibler (KL) Divergence

- Let  $\mathbb{P}$  and  $\mathbb{Q}$  be two distributions over  $\mathcal{X}$ .

$$\begin{aligned} KL(\mathbb{P}\|\mathbb{Q}) &= \sum_x \mathbb{P}(x) \log \frac{\mathbb{P}(x)}{\mathbb{Q}(x)} \\ &= \mathbb{E}_{\mathbb{P}(x)}[-\log \mathbb{Q}(x)] + \mathbb{E}_{\mathbb{P}(x)}[\log \mathbb{P}(x)] \end{aligned}$$

## Recap: Kullback-Leibler (KL) Divergence

- Let  $\mathbb{P}$  and  $\mathbb{Q}$  be two distributions over  $\mathcal{X}$ .

$$\begin{aligned} KL(\mathbb{P}||\mathbb{Q}) &= \sum_x \mathbb{P}(x) \log \frac{\mathbb{P}(x)}{\mathbb{Q}(x)} \\ &= \mathbb{E}_{\mathbb{P}(x)}[-\log \mathbb{Q}(x)] + \mathbb{E}_{\mathbb{P}(x)}[\log \mathbb{P}(x)] \\ &= \mathbb{E}_{\mathbb{P}(x)}[-\log \mathbb{Q}(x)] - H(\mathbb{P}) \end{aligned}$$

## Recap: Kullback-Leibler (KL) Divergence

- Let  $\mathbb{P}$  and  $\mathbb{Q}$  be two distributions over  $\mathcal{X}$ .

$$\begin{aligned} KL(\mathbb{P}\|\mathbb{Q}) &= \sum_x \mathbb{P}(x) \log \frac{\mathbb{P}(x)}{\mathbb{Q}(x)} \\ &= \mathbb{E}_{\mathbb{P}(x)}[-\log \mathbb{Q}(x)] + \mathbb{E}_{\mathbb{P}(x)}[\log \mathbb{P}(x)] \\ &= \mathbb{E}_{\mathbb{P}(x)}[-\log \mathbb{Q}(x)] - H(\mathbb{P}) \end{aligned}$$

- Always **non-negative**:  $KL(\mathbb{P}\|\mathbb{Q}) \geq 0$

## Recap: Kullback-Leibler (KL) Divergence

- Let  $\mathbb{P}$  and  $\mathbb{Q}$  be two distributions over  $\mathcal{X}$ .

$$\begin{aligned} KL(\mathbb{P}\|\mathbb{Q}) &= \sum_x \mathbb{P}(x) \log \frac{\mathbb{P}(x)}{\mathbb{Q}(x)} \\ &= \mathbb{E}_{\mathbb{P}(x)}[-\log \mathbb{Q}(x)] + \mathbb{E}_{\mathbb{P}(x)}[\log \mathbb{P}(x)] \\ &= \mathbb{E}_{\mathbb{P}(x)}[-\log \mathbb{Q}(x)] - H(\mathbb{P}) \end{aligned}$$

- Always **non-negative**:  $KL(\mathbb{P}\|\mathbb{Q}) \geq 0$
- $KL(\mathbb{P}\|\mathbb{Q}) = 0$  iff  $\mathbb{P}(x) = \mathbb{Q}(x)$ , for all  $x \in \mathcal{X}$

## Recap: Kullback-Leibler (KL) Divergence

- Let  $\mathbb{P}$  and  $\mathbb{Q}$  be two distributions over  $\mathcal{X}$ .

$$\begin{aligned} KL(\mathbb{P}\|\mathbb{Q}) &= \sum_x \mathbb{P}(x) \log \frac{\mathbb{P}(x)}{\mathbb{Q}(x)} \\ &= \mathbb{E}_{\mathbb{P}(x)}[-\log \mathbb{Q}(x)] + \mathbb{E}_{\mathbb{P}(x)}[\log \mathbb{P}(x)] \\ &= \mathbb{E}_{\mathbb{P}(x)}[-\log \mathbb{Q}(x)] - H(\mathbb{P}) \end{aligned}$$

- Always **non-negative**:  $KL(\mathbb{P}\|\mathbb{Q}) \geq 0$
- $KL(\mathbb{P}\|\mathbb{Q}) = 0$  iff  $\mathbb{P}(x) = \mathbb{Q}(x)$ , for all  $x \in \mathcal{X}$
- Not symmetric: in general,  $KL(\mathbb{P}\|\mathbb{Q}) \neq KL(\mathbb{Q}\|\mathbb{P})$

## Recap: Kullback-Leibler (KL) Divergence

- Let  $\mathbb{P}$  and  $\mathbb{Q}$  be two distributions over  $\mathcal{X}$ .

$$\begin{aligned} KL(\mathbb{P}||\mathbb{Q}) &= \sum_x \mathbb{P}(x) \log \frac{\mathbb{P}(x)}{\mathbb{Q}(x)} \\ &= \mathbb{E}_{\mathbb{P}(x)}[-\log \mathbb{Q}(x)] + \mathbb{E}_{\mathbb{P}(x)}[\log \mathbb{P}(x)] \\ &= \mathbb{E}_{\mathbb{P}(x)}[-\log \mathbb{Q}(x)] - H(\mathbb{P}) \end{aligned}$$

- Always **non-negative**:  $KL(\mathbb{P}||\mathbb{Q}) \geq 0$
- $KL(\mathbb{P}||\mathbb{Q}) = 0$  iff  $\mathbb{P}(x) = \mathbb{Q}(x)$ , for all  $x \in \mathcal{X}$
- Not symmetric: in general,  $KL(\mathbb{P}||\mathbb{Q}) \neq KL(\mathbb{Q}||\mathbb{P})$
- Intuition**:  $KL(\mathbb{P}||\mathbb{Q})$  measures how different  $\mathbb{Q}$  is from  $\mathbb{P}$

## Recap: Kullback-Leibler (KL) Divergence

- Let  $\mathbb{P}$  and  $\mathbb{Q}$  be two distributions over  $\mathcal{X}$ .

$$\begin{aligned} KL(\mathbb{P}||\mathbb{Q}) &= \sum_x \mathbb{P}(x) \log \frac{\mathbb{P}(x)}{\mathbb{Q}(x)} \\ &= \mathbb{E}_{\mathbb{P}(x)}[-\log \mathbb{Q}(x)] + \mathbb{E}_{\mathbb{P}(x)}[\log \mathbb{P}(x)] \\ &= \mathbb{E}_{\mathbb{P}(x)}[-\log \mathbb{Q}(x)] - H(\mathbb{P}) \end{aligned}$$

- Always **non-negative**:  $KL(\mathbb{P}||\mathbb{Q}) \geq 0$
- $KL(\mathbb{P}||\mathbb{Q}) = 0$  iff  $\mathbb{P}(x) = \mathbb{Q}(x)$ , for all  $x \in \mathcal{X}$
- Not symmetric: in general,  $KL(\mathbb{P}||\mathbb{Q}) \neq KL(\mathbb{Q}||\mathbb{P})$
- Intuition**:  $KL(\mathbb{P}||\mathbb{Q})$  measures how different  $\mathbb{Q}$  is from  $\mathbb{P}$
- Coding perspective**: expected number of extra bits needed to encode  $x \sim \mathbb{P}(x)$  using a code that is optimal for  $\mathbb{Q}(x)$ .

## Recap: Entropy and KL Divergence in the Continuous Case

- Let  $\mathbb{P}$  be a probability density over  $\mathcal{X}$ . The **differential entropy** of  $\mathbb{P}$  is

$$H(\mathbb{P}) = - \int_{\mathcal{X}} \mathbb{P}(x) \log \mathbb{P}(x) dx$$

...no longer guaranteed to be non-negative.

## Recap: Entropy and KL Divergence in the Continuous Case

- Let  $\mathbb{P}$  be a probability density over  $\mathcal{X}$ . The **differential entropy** of  $\mathbb{P}$  is

$$H(\mathbb{P}) = - \int_{\mathcal{X}} \mathbb{P}(x) \log \mathbb{P}(x) dx = \mathbb{E}_{\mathbb{P}(x)} [-\log \mathbb{P}(x)]$$

...no longer guaranteed to be non-negative.

## Recap: Entropy and KL Divergence in the Continuous Case

- Let  $\mathbb{P}$  be a probability density over  $\mathcal{X}$ . The **differential entropy** of  $\mathbb{P}$  is

$$H(\mathbb{P}) = - \int_{\mathcal{X}} \mathbb{P}(x) \log \mathbb{P}(x) dx = \mathbb{E}_{\mathbb{P}(x)}[-\log \mathbb{P}(x)]$$

...no longer guaranteed to be non-negative.

- Let  $\mathbb{P}$  and  $\mathbb{Q}$  be two probability densities over  $\mathcal{X}$ .

$$\begin{aligned} KL(\mathbb{P}||\mathbb{Q}) &= \int_{\mathcal{X}} \mathbb{P}(x) \log \frac{\mathbb{P}(x)}{\mathbb{Q}(x)} dx \\ &= \mathbb{E}_{\mathbb{P}(x)}[-\log \mathbb{Q}(x)] + \mathbb{E}_{\mathbb{P}(x)}[\log \mathbb{P}(x)] \\ &= \mathbb{E}_{\mathbb{P}(x)}[-\log \mathbb{Q}(x)] - H(\mathbb{P}) \end{aligned}$$

## Recap: Entropy and KL Divergence in the Continuous Case

- Let  $\mathbb{P}$  be a probability density over  $\mathcal{X}$ . The **differential entropy** of  $\mathbb{P}$  is

$$H(\mathbb{P}) = - \int_{\mathcal{X}} \mathbb{P}(x) \log \mathbb{P}(x) dx = \mathbb{E}_{\mathbb{P}(x)}[-\log \mathbb{P}(x)]$$

...no longer guaranteed to be non-negative.

- Let  $\mathbb{P}$  and  $\mathbb{Q}$  be two probability densities over  $\mathcal{X}$ .

$$\begin{aligned} KL(\mathbb{P}||\mathbb{Q}) &= \int_{\mathcal{X}} \mathbb{P}(x) \log \frac{\mathbb{P}(x)}{\mathbb{Q}(x)} dx \\ &= \mathbb{E}_{\mathbb{P}(x)}[-\log \mathbb{Q}(x)] + \mathbb{E}_{\mathbb{P}(x)}[\log \mathbb{P}(x)] \\ &= \mathbb{E}_{\mathbb{P}(x)}[-\log \mathbb{Q}(x)] - H(\mathbb{P}) \end{aligned}$$

- Also, always **non-negative**:  $KL(\mathbb{P}||\mathbb{Q}) \geq 0$

## Recap: Entropy and KL Divergence in the Continuous Case

- Let  $\mathbb{P}$  be a probability density over  $\mathcal{X}$ . The **differential entropy** of  $\mathbb{P}$  is

$$H(\mathbb{P}) = - \int_{\mathcal{X}} \mathbb{P}(x) \log \mathbb{P}(x) dx = \mathbb{E}_{\mathbb{P}(x)}[-\log \mathbb{P}(x)]$$

...no longer guaranteed to be non-negative.

- Let  $\mathbb{P}$  and  $\mathbb{Q}$  be two probability densities over  $\mathcal{X}$ .

$$\begin{aligned} KL(\mathbb{P}||\mathbb{Q}) &= \int_{\mathcal{X}} \mathbb{P}(x) \log \frac{\mathbb{P}(x)}{\mathbb{Q}(x)} dx \\ &= \mathbb{E}_{\mathbb{P}(x)}[-\log \mathbb{Q}(x)] + \mathbb{E}_{\mathbb{P}(x)}[\log \mathbb{P}(x)] \\ &= \mathbb{E}_{\mathbb{P}(x)}[-\log \mathbb{Q}(x)] - H(\mathbb{P}) \end{aligned}$$

- Also, always **non-negative**:  $KL(\mathbb{P}||\mathbb{Q}) \geq 0$
- $KL(\mathbb{P}||\mathbb{Q}) = 0$  iff  $\mathbb{P}(x) = \mathbb{Q}(x)$ , almost everywhere

## Recap: Entropy and KL Divergence in the Continuous Case

- Let  $\mathbb{P}$  be a probability density over  $\mathcal{X}$ . The **differential entropy** of  $\mathbb{P}$  is

$$H(\mathbb{P}) = - \int_{\mathcal{X}} \mathbb{P}(x) \log \mathbb{P}(x) dx = \mathbb{E}_{\mathbb{P}(x)}[-\log \mathbb{P}(x)]$$

...no longer guaranteed to be non-negative.

- Let  $\mathbb{P}$  and  $\mathbb{Q}$  be two probability densities over  $\mathcal{X}$ .

$$\begin{aligned} KL(\mathbb{P}||\mathbb{Q}) &= \int_{\mathcal{X}} \mathbb{P}(x) \log \frac{\mathbb{P}(x)}{\mathbb{Q}(x)} dx \\ &= \mathbb{E}_{\mathbb{P}(x)}[-\log \mathbb{Q}(x)] + \mathbb{E}_{\mathbb{P}(x)}[\log \mathbb{P}(x)] \\ &= \mathbb{E}_{\mathbb{P}(x)}[-\log \mathbb{Q}(x)] - H(\mathbb{P}) \end{aligned}$$

- Also, always **non-negative**:  $KL(\mathbb{P}||\mathbb{Q}) \geq 0$
- $KL(\mathbb{P}||\mathbb{Q}) = 0$  iff  $\mathbb{P}(x) = \mathbb{Q}(x)$ , almost everywhere
- Intuition**:  $KL(\mathbb{P}||\mathbb{Q})$  still measure how different  $\mathbb{Q}$  is from  $\mathbb{P}$

# Evidence Lower Bound (ELBO)

- ELBO is a central concept in variational inference.
- True posterior and **evidence**:  $\mathbb{P}_{\theta}(h | x)$  and  $\mathbb{P}_{\theta}(x)$

## Evidence Lower Bound (ELBO)

- ELBO is a central concept in variational inference.
- True posterior and **evidence**:  $\mathbb{P}_\theta(\mathbf{h} \mid \mathbf{x})$  and  $\mathbb{P}_\theta(\mathbf{x})$
- For any distribution  $\mathbb{Q}(\mathbf{h})$ ,

$$\begin{aligned} 0 &\geq -KL(\mathbb{Q}(\mathbf{h}) \parallel \mathbb{P}_\theta(\mathbf{h} \mid \mathbf{x})) \\ &= \mathbb{E}_{\mathbb{Q}(\mathbf{h})}[\log \mathbb{P}_\theta(\mathbf{h} \mid \mathbf{x})] - \overbrace{\mathbb{E}_{\mathbb{Q}(\mathbf{h})}[\log \mathbb{Q}(\mathbf{h})]}^{H(\mathbb{Q})} \end{aligned}$$

# Evidence Lower Bound (ELBO)

- ELBO is a central concept in variational inference.
- True posterior and **evidence**:  $\mathbb{P}_\theta(\mathbf{h} | \mathbf{x})$  and  $\mathbb{P}_\theta(\mathbf{x})$
- For any distribution  $\mathbb{Q}(\mathbf{h})$ ,

$$\begin{aligned} 0 &\geq -KL(\mathbb{Q}(\mathbf{h}) \| \mathbb{P}_\theta(\mathbf{h} | \mathbf{x})) \\ &= \mathbb{E}_{\mathbb{Q}(\mathbf{h})}[\log \mathbb{P}_\theta(\mathbf{h} | \mathbf{x})] - \overbrace{\mathbb{E}_{\mathbb{Q}(\mathbf{h})}[\log \mathbb{Q}(\mathbf{h})]}^{H(\mathbb{Q})} \\ &\stackrel{(a)}{=} \underbrace{\mathbb{E}_{\mathbb{Q}(\mathbf{h})}[\log \mathbb{P}_\theta(\mathbf{x}, \mathbf{h})] - \mathbb{E}_{\mathbb{Q}(\mathbf{h})}[\log \mathbb{Q}(\mathbf{h})]}_{\text{ELBO}(\mathbb{Q})} - \log \mathbb{P}_\theta(\mathbf{x}). \end{aligned}$$

where (a) uses Bayes law:  $\log \mathbb{P}_\theta(\mathbf{h} | \mathbf{x}) = \log \mathbb{P}_\theta(\mathbf{x}, \mathbf{h}) - \log \mathbb{P}_\theta(\mathbf{x})$

# Evidence Lower Bound (ELBO)

- ELBO is a central concept in variational inference.
- True posterior and **evidence**:  $\mathbb{P}_\theta(\mathbf{h} | \mathbf{x})$  and  $\mathbb{P}_\theta(\mathbf{x})$
- For any distribution  $\mathbb{Q}(\mathbf{h})$ ,

$$\begin{aligned} 0 &\geq -KL(\mathbb{Q}(\mathbf{h}) \| \mathbb{P}_\theta(\mathbf{h} | \mathbf{x})) \\ &= \mathbb{E}_{\mathbb{Q}(\mathbf{h})}[\log \mathbb{P}_\theta(\mathbf{h} | \mathbf{x})] - \overbrace{\mathbb{E}_{\mathbb{Q}(\mathbf{h})}[\log \mathbb{Q}(\mathbf{h})]}^{H(\mathbb{Q})} \\ &\stackrel{(a)}{=} \underbrace{\mathbb{E}_{\mathbb{Q}(\mathbf{h})}[\log \mathbb{P}_\theta(\mathbf{x}, \mathbf{h})] - \mathbb{E}_{\mathbb{Q}(\mathbf{h})}[\log \mathbb{Q}(\mathbf{h})]}_{\text{ELBO}(\mathbb{Q})} - \log \mathbb{P}_\theta(\mathbf{x}). \end{aligned}$$

where (a) uses Bayes law:  $\log \mathbb{P}_\theta(\mathbf{h} | \mathbf{x}) = \log \mathbb{P}_\theta(\mathbf{x}, \mathbf{h}) - \log \mathbb{P}_\theta(\mathbf{x})$

- Moving  $\log \mathbb{P}_\theta(\mathbf{x})$  to the l.h.s.,

$$\log \mathbb{P}_\theta(\mathbf{x}) \geq \text{ELBO}(\mathbb{Q})$$

# Variational Inference

- Evidence lower bound (ELBO):

$$\log \mathbb{P}_{\theta}(x) = \text{ELBO}(\mathbb{Q}) + KL(\mathbb{Q}(\mathbf{h}) \| \mathbb{P}_{\theta}(\mathbf{h} | x)) \geq \text{ELBO}(\mathbb{Q}).$$

# Variational Inference

- Evidence lower bound (ELBO):

$$\log \mathbb{P}_{\theta}(x) = \text{ELBO}(\mathbb{Q}) + KL(\mathbb{Q}(\mathbf{h}) \| \mathbb{P}_{\theta}(\mathbf{h} | x)) \geq \text{ELBO}(\mathbb{Q}).$$

- Equality achieved for  $\mathbb{Q}(\mathbf{h}) = \mathbb{P}_{\theta}(\mathbf{h} | x)$ , but intractable in general

# Variational Inference

- Evidence lower bound (ELBO):

$$\log \mathbb{P}_{\theta}(x) = \text{ELBO}(\mathbb{Q}) + KL(\mathbb{Q}(\mathbf{h}) \| \mathbb{P}_{\theta}(\mathbf{h} | x)) \geq \text{ELBO}(\mathbb{Q}).$$

- Equality achieved for  $\mathbb{Q}(\mathbf{h}) = \mathbb{P}_{\theta}(\mathbf{h} | x)$ , but intractable in general
- Key idea:
  - ① constrain  $\mathbb{Q}(\mathbf{h})$  to a chosen tractable family;
  - ② look for the  $\mathbb{Q}(\mathbf{h})$  in this family that maximizes the ELBO.

# Variational Inference

- Evidence lower bound (ELBO):

$$\log \mathbb{P}_{\theta}(x) = \text{ELBO}(\mathbb{Q}) + KL(\mathbb{Q}(\mathbf{h}) \parallel \mathbb{P}_{\theta}(\mathbf{h} \mid x)) \geq \text{ELBO}(\mathbb{Q}).$$

- Equality achieved for  $\mathbb{Q}(\mathbf{h}) = \mathbb{P}_{\theta}(\mathbf{h} \mid x)$ , but intractable in general
- Key idea:
  - ① constrain  $\mathbb{Q}(\mathbf{h})$  to a chosen tractable family;
  - ② look for the  $\mathbb{Q}(\mathbf{h})$  in this family that maximizes the ELBO.
- Since  $\log \mathbb{P}_{\theta}(x)$  fixed,

$$\text{maximizing } \text{ELBO}(\mathbb{Q}) \Leftrightarrow \text{minimizing } KL(\mathbb{Q}(\mathbf{h}) \parallel \mathbb{P}_{\theta}(\mathbf{h} \mid x))$$

# Variational Inference

- Evidence lower bound (ELBO):

$$\log \mathbb{P}_{\theta}(x) = \text{ELBO}(\mathbb{Q}) + KL(\mathbb{Q}(\mathbf{h}) \parallel \mathbb{P}_{\theta}(\mathbf{h} \mid x)) \geq \text{ELBO}(\mathbb{Q}).$$

- Equality achieved for  $\mathbb{Q}(\mathbf{h}) = \mathbb{P}_{\theta}(\mathbf{h} \mid x)$ , but intractable in general
- Key idea:
  - ① constrain  $\mathbb{Q}(\mathbf{h})$  to a chosen tractable family;
  - ② look for the  $\mathbb{Q}(\mathbf{h})$  in this family that maximizes the ELBO.
- Since  $\log \mathbb{P}_{\theta}(x)$  fixed,  
$$\text{maximizing } \text{ELBO}(\mathbb{Q}) \Leftrightarrow \text{minimizing } KL(\mathbb{Q}(\mathbf{h}) \parallel \mathbb{P}_{\theta}(\mathbf{h} \mid x))$$
- Old roots in calculus of variations (Euler, Lagrange, ...)

# Evidence Lower Bound

- The ELBO can be written in different ways:

$$\begin{aligned}\text{ELBO}(\mathbb{Q}) &= \mathbb{E}_{\mathbb{Q}(h)}[\log \mathbb{P}_{\theta}(x, h)] - \mathbb{E}_{\mathbb{Q}(h)}[\log \mathbb{Q}(h)] \\ &= \mathbb{E}_{\mathbb{Q}(h)}[\log \mathbb{P}_{\theta}(x | h)] - \mathbb{E}_{\mathbb{Q}(h)}\left[\log \frac{\mathbb{Q}(h)}{\mathbb{P}_{\theta}(h)}\right] \\ &= \mathbb{E}_{\mathbb{Q}(h)}[\log \mathbb{P}_{\theta}(x | h)] - \textcolor{red}{KL}(\mathbb{Q}(h) \| \mathbb{P}_{\theta}(h)).\end{aligned}$$

# Evidence Lower Bound

- The ELBO can be written in different ways:

$$\begin{aligned}\text{ELBO}(\mathbb{Q}) &= \mathbb{E}_{\mathbb{Q}(h)}[\log \mathbb{P}_{\theta}(x, h)] - \mathbb{E}_{\mathbb{Q}(h)}[\log \mathbb{Q}(h)] \\ &= \mathbb{E}_{\mathbb{Q}(h)}[\log \mathbb{P}_{\theta}(x | h)] - \mathbb{E}_{\mathbb{Q}(h)}\left[\log \frac{\mathbb{Q}(h)}{\mathbb{P}_{\theta}(h)}\right] \\ &= \mathbb{E}_{\mathbb{Q}(h)}[\log \mathbb{P}_{\theta}(x | h)] - \text{KL}(\mathbb{Q}(h) \| \mathbb{P}_{\theta}(h)).\end{aligned}$$

- Which values of  $h$  is  $\mathbb{Q}(h)$  encouraged to place its mass on?
  - First term: **expected likelihood** encourages placing mass on latent variables  $h$  that explain the observed data  $x$ .
  - Second term: **negative KLD between  $\mathbb{Q}(h)$  and the prior** encourages staying close to the prior.

# Evidence Lower Bound

- The ELBO can be written in different ways:

$$\begin{aligned}\text{ELBO}(\mathbb{Q}) &= \mathbb{E}_{\mathbb{Q}(h)}[\log \mathbb{P}_{\theta}(x, h)] - \mathbb{E}_{\mathbb{Q}(h)}[\log \mathbb{Q}(h)] \\ &= \mathbb{E}_{\mathbb{Q}(h)}[\log \mathbb{P}_{\theta}(x | h)] - \mathbb{E}_{\mathbb{Q}(h)}\left[\log \frac{\mathbb{Q}(h)}{\mathbb{P}_{\theta}(h)}\right] \\ &= \mathbb{E}_{\mathbb{Q}(h)}[\log \mathbb{P}_{\theta}(x | h)] - \text{KL}(\mathbb{Q}(h) \| \mathbb{P}_{\theta}(h)).\end{aligned}$$

- Which values of  $h$  is  $\mathbb{Q}(h)$  encouraged to place its mass on?
  - First term: **expected likelihood** encourages placing mass on latent variables  $h$  that explain the observed data  $x$ .
  - Second term: **negative KLD between  $\mathbb{Q}(h)$  and the prior** encourages staying close to the prior.
- ELBO mirrors the **Bayesian** trade-off between **likelihood** and **prior**.

# Mean Field Approximation

- Which **tractable family** to use for  $\mathbb{Q}(\mathbf{h})$ ?

# Mean Field Approximation

- Which tractable family to use for  $\mathbb{Q}(\mathbf{h})$ ?
- Mean field approximation (MFA):

$$\mathbb{Q}(\mathbf{h}) = \prod_i \mathbb{Q}(h_i).$$

i.e., model the  $h_i$  are independent.

# Mean Field Approximation

- Which tractable family to use for  $\mathbb{Q}(\mathbf{h})$ ?
- Mean field approximation (MFA):

$$\mathbb{Q}(\mathbf{h}) = \prod_i \mathbb{Q}(h_i).$$

i.e., model the  $h_i$  are independent.

- Structured mean field: imposes a graphical model on  $\mathbb{Q}$  capturing (some) interactions among the  $h_i$ , but still tractable.  
(Wainwright and Jordan, 2008).

# Amortized Variational Inference

- As seen so far, the variational distribution  $\mathbb{Q}(\mathbf{h})$  has to be optimized for every example  $\mathbf{x}$ .

# Amortized Variational Inference

- As seen so far, the variational distribution  $\mathbb{Q}(\boldsymbol{h})$  has to be optimized for every example  $\boldsymbol{x}$ .
- This can be **expensive**: requires several gradient/coordinate ascent iterations per example.

# Amortized Variational Inference

- As seen so far, the variational distribution  $\mathbb{Q}(\mathbf{h})$  has to be optimized for every example  $\mathbf{x}$ .
- This can be **expensive**: requires several gradient/coordinate ascent iterations per example.
- Alternative: **amortized variational inference!**
- **Key idea:** instead of optimizing  $\mathbb{Q}(\mathbf{h})$  for every example, use an **encoder** with shared parameters  $\phi$  and define  $\mathbb{Q}_\phi(\mathbf{h} \mid \mathbf{x})$ .

For each example:

- make a forward pass on the encoder to obtain  $\mathbb{Q}_\phi(\mathbf{h} \mid \mathbf{x})$
- backpropagate through the encoder to update  $\phi$ .

## Example: Multivariate Bernoulli with Continuous Latent Variables (Kingma and Welling, 2013)

- **Prior:** standard Gaussian  $\mathbb{P}_{\theta}(\mathbf{h}) = \mathcal{N}(\mathbf{h}; \mathbf{0}, \mathbf{I})$

## Example: Multivariate Bernoulli with Continuous Latent Variables (Kingma and Welling, 2013)

- **Prior:** standard Gaussian  $\mathbb{P}_{\theta}(\mathbf{h}) = \mathcal{N}(\mathbf{h}; \mathbf{0}, \mathbf{I})$
- **Conditional:** multivariate Bernoulli

$$\mathbb{P}_{\theta}(x \mid h) = \prod_{i=1}^D \sigma(f_i(h; \theta))^{x_i} (1 - \sigma(f_i(h; \theta)))^{1-x_i},$$

where  $f(h; \theta)$  is an MLP, with parameters  $\theta$  and input  $h$

## Example: Multivariate Bernoulli with Continuous Latent Variables (Kingma and Welling, 2013)

- **Prior:** standard Gaussian  $\mathbb{P}_{\theta}(\mathbf{h}) = \mathcal{N}(\mathbf{h}; \mathbf{0}, \mathbf{I})$
- **Conditional:** multivariate Bernoulli

$$\mathbb{P}_{\theta}(x \mid h) = \prod_{i=1}^D \sigma(f_i(h; \theta))^{x_i} (1 - \sigma(f_i(h; \theta)))^{1-x_i},$$

where  $f(\mathbf{h}; \theta)$  is an MLP, with parameters  $\theta$  and input  $\mathbf{h}$

- **True posterior**  $\mathbb{P}_{\theta}(\mathbf{h} \mid x)$ : intractable

## Example: Multivariate Bernoulli with Continuous Latent Variables (Kingma and Welling, 2013)

- **Prior:** standard Gaussian  $\mathbb{P}_{\theta}(\mathbf{h}) = \mathcal{N}(\mathbf{h}; \mathbf{0}, \mathbf{I})$
- **Conditional:** multivariate Bernoulli

$$\mathbb{P}_{\theta}(\mathbf{x} \mid \mathbf{h}) = \prod_{i=1}^D \sigma(f_i(\mathbf{h}; \theta))^{x_i} (1 - \sigma(f_i(\mathbf{h}; \theta)))^{1-x_i},$$

where  $f(\mathbf{h}; \theta)$  is an MLP, with parameters  $\theta$  and input  $\mathbf{h}$

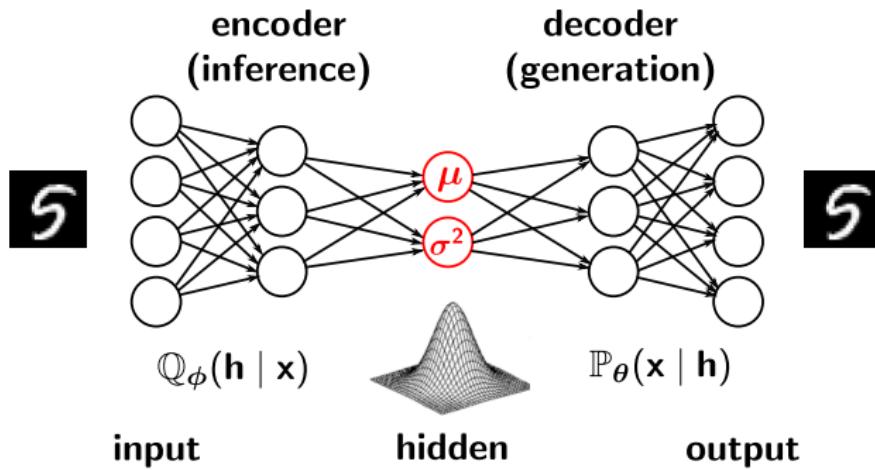
- **True posterior**  $\mathbb{P}_{\theta}(\mathbf{h} \mid \mathbf{x})$ : intractable
- Approximate the posterior with a **variational distribution**

$$\mathbb{Q}_{\phi}(\mathbf{h} \mid \mathbf{x}) = \mathcal{N}(\mathbf{h}; \boldsymbol{\mu}(\mathbf{x}; \phi), \boldsymbol{\sigma}^2(\mathbf{x}; \phi))$$

where  $\boldsymbol{\mu}(\mathbf{x}; \phi)$  and  $\boldsymbol{\sigma}^2(\mathbf{x}; \phi)$  are (learned) MLPs

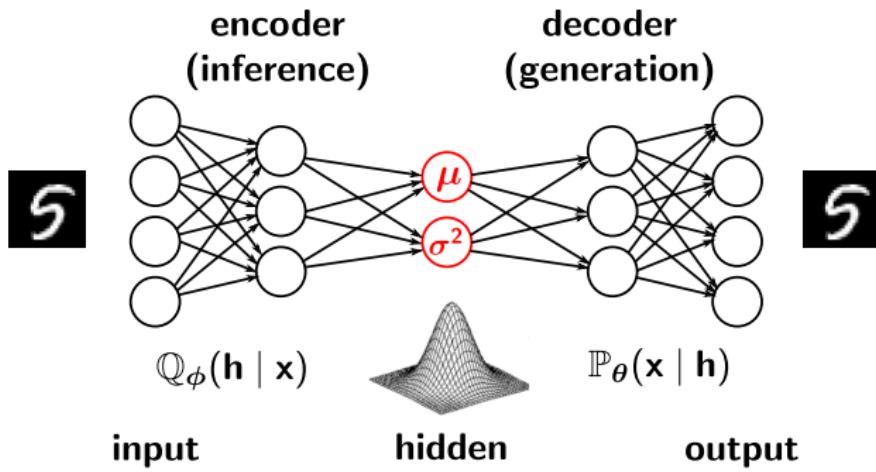
# Example: Multivariate Bernoulli with Continuous Latent Variables (Kingma and Welling, 2013)

- This leads to **variational auto-encoders**:



# Example: Multivariate Bernoulli with Continuous Latent Variables (Kingma and Welling, 2013)

- This leads to **variational auto-encoders**:



- ... we will come back to this!

# Outline

① Boltzmann Machines

② Variational Auto-Encoders

Variational Inference and ELBO

Gradients and Reparameterization Trick

③ Generative Adversarial Networks

④ Denoising Diffusion Models

⑤ Conclusions

# Parameter Gradients

- Recall that:

$$\text{ELBO}(\phi; \theta) = \mathbb{E}_{\mathbb{Q}_\phi(h|x)} [\log \mathbb{P}_\theta(x, h) - \log \mathbb{Q}_\phi(h | x)].$$

## Parameter Gradients

- Recall that:

$$\text{ELBO}(\phi; \theta) = \mathbb{E}_{\mathbb{Q}_\phi(h|x)} [\log \mathbb{P}_\theta(x, h) - \log \mathbb{Q}_\phi(h | x)].$$

- We need to compute gradients with respect to  $\theta$  and  $\phi$ .

# Parameter Gradients

- Recall that:

$$\text{ELBO}(\phi; \theta) = \mathbb{E}_{\mathbb{Q}_\phi(h|x)} [\log \mathbb{P}_\theta(x, h) - \log \mathbb{Q}_\phi(h | x)].$$

- We need to compute gradients with respect to  $\theta$  and  $\phi$ .
- Gradient w.r.t.  $\theta$ , parameters of the generation (decoder) network:

$$\nabla_\theta \text{ELBO}(\phi; \theta) = \mathbb{E}_{\mathbb{Q}_\phi(h|x)} [\nabla_\theta \log \mathbb{P}_\theta(x, h)]$$

- Follows from linearity of the expectation.
- This is simple and can be well approximated with Monte Carlo samples.

# Parameter Gradients

- Recall that:

$$\text{ELBO}(\phi; \theta) = \mathbb{E}_{\mathbb{Q}_\phi(\mathbf{h}|\mathbf{x})}[\log \mathbb{P}_\theta(\mathbf{x}, \mathbf{h}) - \log \mathbb{Q}_\phi(\mathbf{h} \mid \mathbf{x})].$$

# Parameter Gradients

- Recall that:

$$\text{ELBO}(\phi; \theta) = \mathbb{E}_{\mathbb{Q}_\phi(h|x)}[\log \mathbb{P}_\theta(x, h) - \log \mathbb{Q}_\phi(h | x)].$$

- Gradient w.r.t.  $\phi$ , parameters of the **inference network**:

$$\begin{aligned} & \nabla_\phi \text{ELBO}(\phi; \theta) \\ &= \mathbb{E}_{\mathbb{Q}_\phi(h|x)} \underbrace{[(\log \mathbb{P}_\theta(x, h) - \log \mathbb{Q}_\phi(h | x))]}_{\text{"reward" } R_{\theta, \phi}(h)} \nabla_\phi \log \mathbb{Q}_\phi(h | x). \end{aligned}$$

(derivation in the next slide)

# Parameter Gradients

- Recall that:

$$\text{ELBO}(\phi; \theta) = \mathbb{E}_{\mathbb{Q}_\phi(h|x)} [\log \mathbb{P}_\theta(x, h) - \log \mathbb{Q}_\phi(h | x)].$$

- Gradient w.r.t.  $\phi$ , parameters of the **inference network**:

$$\begin{aligned} & \nabla_\phi \text{ELBO}(\phi; \theta) \\ &= \mathbb{E}_{\mathbb{Q}_\phi(h|x)} \underbrace{[(\log \mathbb{P}_\theta(x, h) - \log \mathbb{Q}_\phi(h | x))]}_{\text{"reward" } R_{\theta, \phi}(h)} \nabla_\phi \log \mathbb{Q}_\phi(h | x). \end{aligned}$$

(derivation in the next slide)

- High variance of Monte Carlo estimates due to the left part!
- Requires variance reduction techniques.

# Derivation of the Inference Network Gradient

$$\begin{aligned}\nabla_{\phi} \text{ELBO}(\phi; \theta) &= \nabla_{\phi} \mathbb{E}_{\mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x})} \left[ \log \mathbb{P}_{\theta}(\mathbf{x}, \mathbf{h}) - \log \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x}) \right] \\ &= \nabla_{\phi} \sum_{\mathbf{h}} \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x}) \log \mathbb{P}_{\theta}(\mathbf{x}, \mathbf{h}) - \nabla_{\phi} \sum_{\mathbf{h}} \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x}) \log \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x}) \\ &= \sum_{\mathbf{h}} \log \mathbb{P}_{\theta}(\mathbf{x}, \mathbf{h}) \nabla_{\phi} \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x}) - \sum_{\mathbf{h}} (1 + \log \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x})) \nabla_{\phi} \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x}) \\ &= \sum_{\mathbf{h}} \left[ \log \mathbb{P}_{\theta}(\mathbf{x}, \mathbf{h}) - \log \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x}) \right] \nabla_{\phi} \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x}) \\ &= \mathbb{E}_{\mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x})} \left[ (\log \mathbb{P}_{\theta}(\mathbf{x}, \mathbf{h}) - \log \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x})) \nabla_{\phi} \log \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x}) \right],\end{aligned}$$

where we used the facts:

$$\sum_{\mathbf{h}} \nabla_{\phi} \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x}) = \nabla_{\phi} \sum_{\mathbf{h}} \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x}) = \nabla_{\phi} 1 = 0.$$

$$\nabla_{\phi} \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x}) = \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x}) \nabla_{\phi} \log \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x}).$$

- Summing up, the difficulty is the gradient w.r.t.  $\phi$ , the parameters of the inference (encoder) network
- ... the Monte Carlo approximation has large variance.

- Summing up, the difficulty is the gradient w.r.t.  $\phi$ , the parameters of the inference (encoder) network
- ... the Monte Carlo approximation has large variance.
- Is there a better strategy?

- Summing up, the difficulty is the gradient w.r.t.  $\phi$ , the parameters of the inference (encoder) network
- ... the Monte Carlo approximation has large variance.
- Is there a better strategy?
- Yes: **the reparameterization trick.**

# Reparameterization Trick (Kingma and Welling, 2013)

- How to draw samples from  $\mathbb{Q}_\phi(\mathbf{h} \mid \mathbf{x})$ ?

## Reparameterization Trick (Kingma and Welling, 2013)

- How to draw samples from  $\mathbb{Q}_\phi(\mathbf{h} \mid \mathbf{x})$ ?
- Trick:
  - Use an auxiliary random variable  $\epsilon$  with fixed distribution  $\mathbb{P}(\epsilon)$
  - Sample  $\epsilon \sim \mathbb{P}(\epsilon)$ , and obtain  $\mathbf{h}$  as a **deterministic function** of  $\epsilon$  and  $\mathbf{x}$

$$\mathbf{h} = \mathbf{g}_\phi(\epsilon, \mathbf{x})$$

## Reparameterization Trick (Kingma and Welling, 2013)

- How to draw samples from  $\mathbb{Q}_\phi(\mathbf{h} \mid \mathbf{x})$ ?
- Trick:
  - Use an auxiliary random variable  $\epsilon$  with fixed distribution  $\mathbb{P}(\epsilon)$
  - Sample  $\epsilon \sim \mathbb{P}(\epsilon)$ , and obtain  $\mathbf{h}$  as a **deterministic function** of  $\epsilon$  and  $\mathbf{x}$

$$\mathbf{h} = \mathbf{g}_\phi(\epsilon, \mathbf{x})$$

- Consequently, for any function  $f$ ,

$$\mathbb{E}_{\mathbb{Q}_\phi(\mathbf{h}|\mathbf{x})}[f(\mathbf{h})] \approx \frac{1}{N} \sum_{i=1}^N f(\mathbf{g}_\phi(\mathbf{x}, \epsilon^{(i)}))$$

## Reparameterization Trick (Kingma and Welling, 2013)

- How to draw samples from  $\mathbb{Q}_\phi(\mathbf{h} \mid \mathbf{x})$ ?
- Trick:
  - Use an auxiliary random variable  $\epsilon$  with fixed distribution  $\mathbb{P}(\epsilon)$
  - Sample  $\epsilon \sim \mathbb{P}(\epsilon)$ , and obtain  $\mathbf{h}$  as a **deterministic function** of  $\epsilon$  and  $\mathbf{x}$

$$\mathbf{h} = \mathbf{g}_\phi(\epsilon, \mathbf{x})$$

- Consequently, for any function  $f$ ,

$$\mathbb{E}_{\mathbb{Q}_\phi(\mathbf{h}|\mathbf{x})}[f(\mathbf{h})] \approx \frac{1}{N} \sum_{i=1}^N f(\mathbf{g}_\phi(\mathbf{x}, \epsilon^{(i)}))$$

- $\nabla_\phi$  can be estimated with regular backpropagation over  $\mathbf{g}_\phi$ .

# Reparameterization Trick

- This is possible in many cases for continuous latent variables:
  - exponential
  - Gaussian
  - location-scale families
  - log-normal
  - etc...

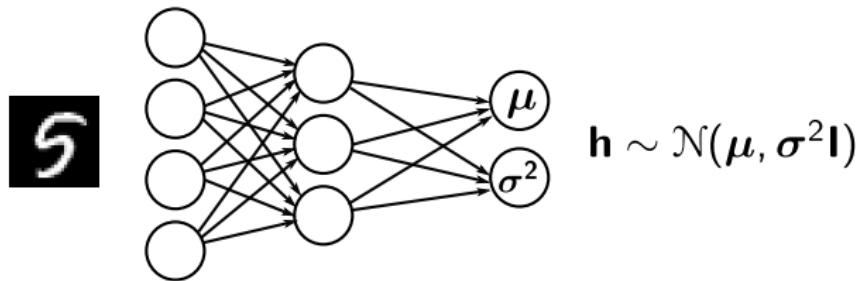
# Reparameterization Trick

- This is possible in many cases for continuous latent variables:
  - exponential
  - Gaussian
  - location-scale families
  - log-normal
  - etc...
- For discrete latent variables, it is still possible via the **Gumbel-softmax trick** (not covered in the course).

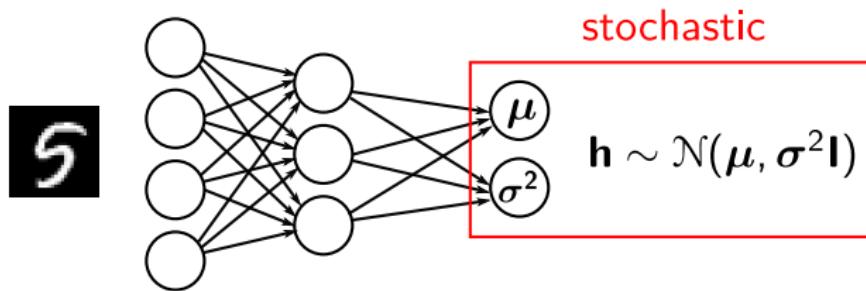
## Example: Gaussian

- ① Sample  $d$ -variate standard Gaussian  $\epsilon \sim \mathcal{N}(\epsilon; 0, I)$
- ② Use inference network  $g_\phi$  with input  $x$  to output:
  - mean  $\mu(x)$
  - component-wise variance  $\sigma^2(x) = \text{diag}[\sigma_1^2(x), \dots, \sigma_d^2(x)]$
- ③ Set  $h = \mu(x) + \epsilon \odot \sqrt{\sigma^2(x)}$ .
- ④ Thus,  $h | x \sim \mathcal{N}(\mu(x), \sigma^2(x))$

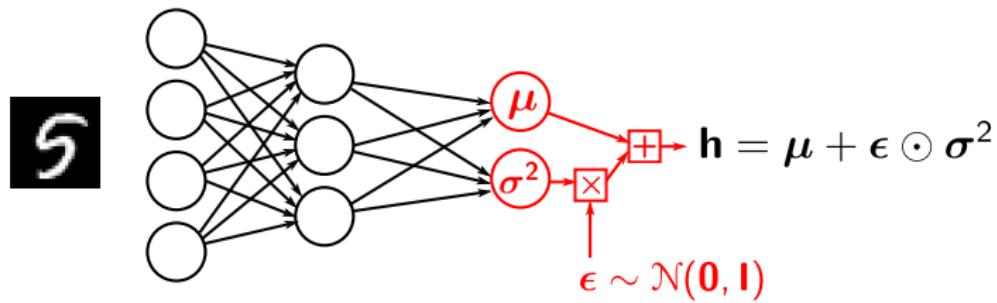
# Reparameterization Trick



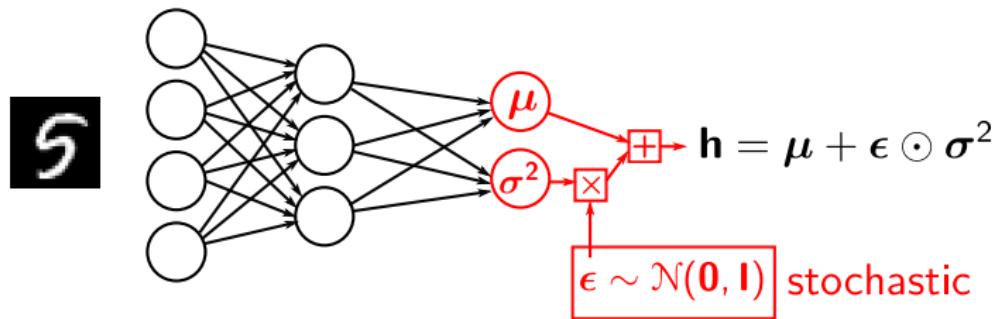
# Reparameterization Trick



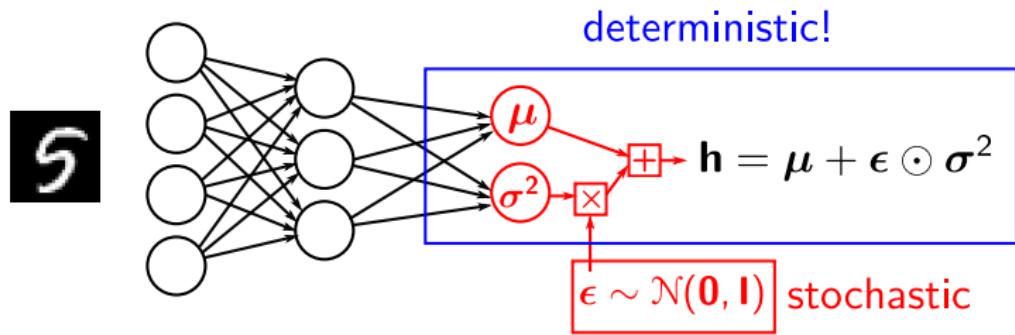
# Reparameterization Trick



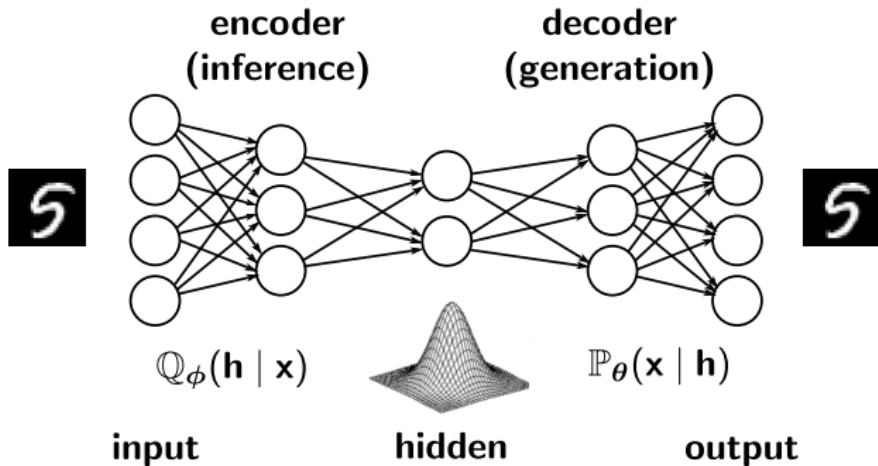
# Reparameterization Trick



# Reparameterization Trick

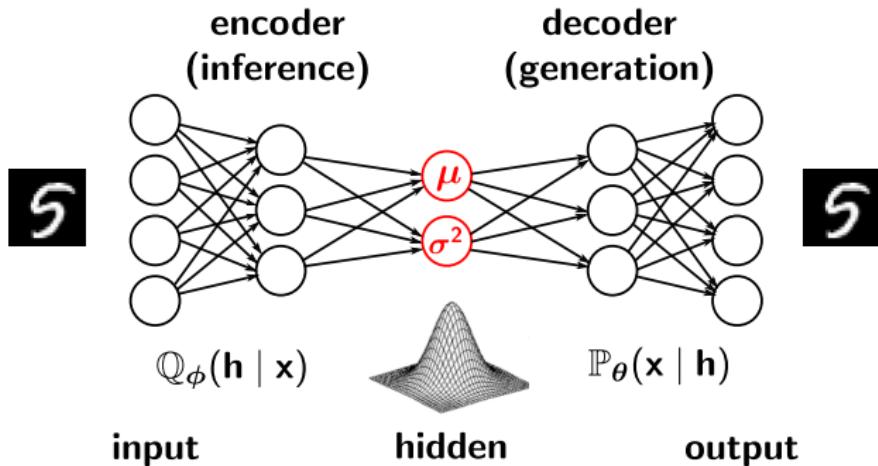


# Variational Auto-Encoders (Kingma and Welling, 2013)



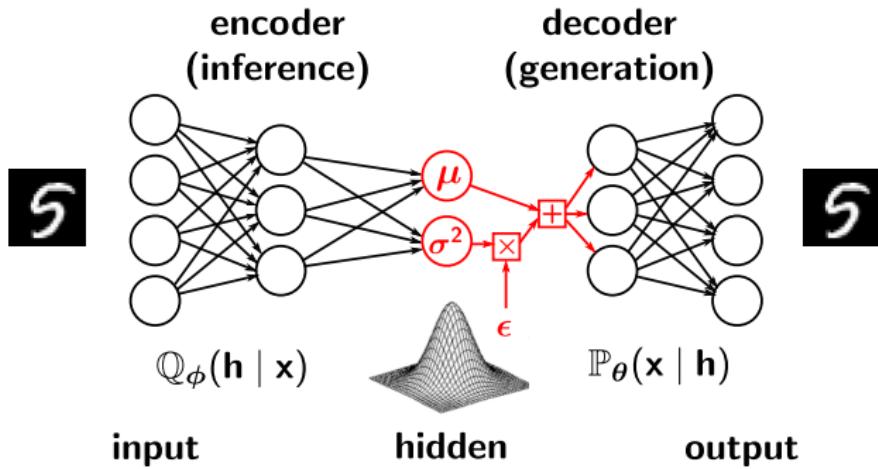
- Decoder computes  $\mathbb{P}_\theta(\mathbf{h})$  and  $\mathbb{P}_\theta(\mathbf{x} \mid \mathbf{h})$
- Encoder computes  $\mathbb{Q}_\phi(\mathbf{h} \mid \mathbf{x}) = \mathcal{N}(\mathbf{h}; \boldsymbol{\mu}_\phi(\mathbf{x}), \boldsymbol{\sigma}_\phi^2(\mathbf{x}))$
- Loss function: ELBO.

# Variational Auto-Encoders (Kingma and Welling, 2013)



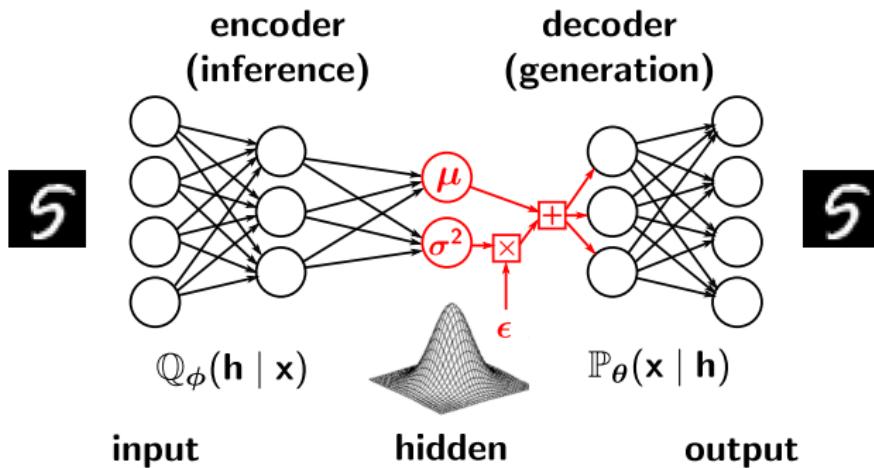
- Decoder computes  $P_\theta(h)$  and  $P_\theta(x | h)$
- Encoder computes  $Q_\phi(h | x) = \mathcal{N}(h; \mu_\phi(x), \sigma_\phi^2(x))$
- Loss function: ELBO.

# Variational Auto-Encoders (Kingma and Welling, 2013)

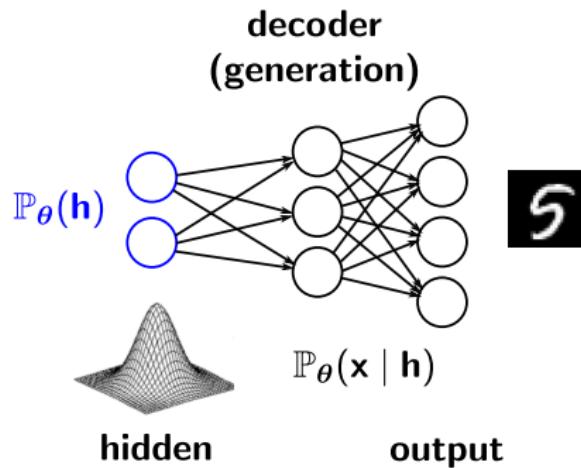


- Decoder computes  $\mathbb{P}_\theta(h)$  and  $\mathbb{P}_\theta(x | h)$
- Encoder computes  $\mathbb{Q}_\phi(h | x) = \mathcal{N}(h; \mu_\phi(x), \sigma_\phi^2(x))$
- Loss function: ELBO.

# Summing Up: VAEs at Training Time



# Summing Up: VAEs at Test Time



## What is the Latent Variable Representing?

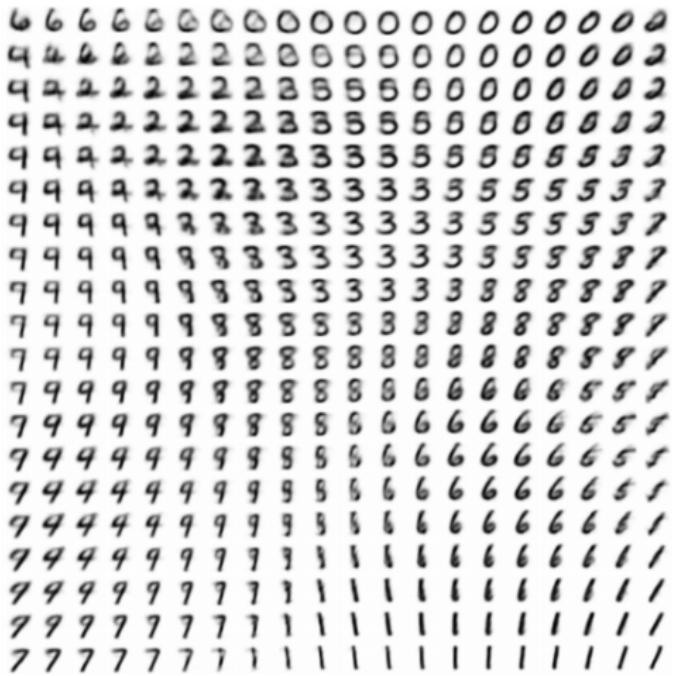
- Nice property of VAE: simultaneously training a **parametric encoder** in combination with a **generator network** forces the model to learn a predictable coordinate system that the encoder can capture.
- This makes it an excellent manifold learning algorithm.
- Example: the algorithm discovered two independent factors of variation present in images of faces: angle of rotation and emotional expression.

# What is the Latent Variable Representing?



From Kingma and Welling (2013).

# What is the Latent Variable Representing?



(b) Learned MNIST manifold

From Kingma and Welling (2013).

## Issues with VAEs

- **Posterior collapse:** if the generative part is strong, the model may learn to ignore the latent variables:

$$\begin{aligned}\mathbb{P}_{\theta}(x | h) &\approx \mathbb{P}(x) \\ \mathbb{Q}_{\phi}(h | x) &\approx \mathbb{P}_{\theta}(h).\end{aligned}$$

## Issues with VAEs

- **Posterior collapse:** if the generative part is strong, the model may learn to ignore the latent variables:

$$\begin{aligned}\mathbb{P}_{\theta}(x | h) &\approx \mathbb{P}(x) \\ \mathbb{Q}_{\phi}(h | x) &\approx \mathbb{P}_{\theta}(h).\end{aligned}$$

- Can be mitigated with a few tricks:
  - Decrease/anneal the weight of  $KL(\mathbb{Q}_{\phi}(h | x) || \mathbb{P}_{\theta}(h))$  in the ELBO
  - Use auxiliary losses
  - Combine stochastic and amortized inference.

## Issues with VAEs

- **Posterior collapse:** if the generative part is strong, the model may learn to ignore the latent variables:

$$\begin{aligned}\mathbb{P}_{\theta}(x | h) &\approx \mathbb{P}(x) \\ \mathbb{Q}_{\phi}(h | x) &\approx \mathbb{P}_{\theta}(h).\end{aligned}$$

- Can be mitigated with a few tricks:
  - Decrease/anneal the weight of  $KL(\mathbb{Q}_{\phi}(h | x) || \mathbb{P}_{\theta}(h))$  in the ELBO
  - Use auxiliary losses
  - Combine stochastic and amortized inference.
- In general, reporting both reconstruction loss and the KL term is needed to assess if the model makes use of the latent variables.

# Examples of Deep Generative Models

- Auto-Regressive Networks ✓
- Restricted Boltzmann Machines ✓
- Deep Belief Networks
- Deep Boltzmann Machines
- Gaussian-Bernoulli RBMs
- Convolutional Boltzmann Machines
- Sigmoid Belief Nets
- Variational Auto-Encoders ✓
- Generative Adversarial Networks
- Denoising diffusion models

# Examples of Deep Generative Models

- Auto-Regressive Networks ✓
- Restricted Boltzmann Machines ✓
- Deep Belief Networks
- Deep Boltzmann Machines
- Gaussian-Bernoulli RBMs
- Convolutional Boltzmann Machines
- Sigmoid Belief Nets
- Variational Auto-Encoders ✓
- **Generative Adversarial Networks**
- Denoising diffusion models

# Outline

① Boltzmann Machines

② Variational Auto-Encoders

Variational Inference and ELBO

Gradients and Reparameterization Trick

③ Generative Adversarial Networks

④ Denoising Diffusion Models

⑤ Conclusions

# Why Maximum Likelihood?

- All models discussed aim to maximize the likelihood (evidence)  $\mathbb{P}(x)$
- In fact, since this is intractable, they maximize a lower bound (ELBO)
- But if we want to build a generator, is this really the best criterion?
- Maximum likelihood tends to produce blurry images

# Generative Adversarial Networks (GANs)

(Goodfellow et al., 2014)

- **Key idea:**

- keep the **generation network**  $G = \{\mathbb{P}_{\theta}(h), \mathbb{P}_{\theta}(x | h)\}$
- drop the inference network: use a **discriminator network**  $D : \mathcal{X} \rightarrow \{0, 1\}$

# Generative Adversarial Networks (GANs)

(Goodfellow et al., 2014)

- Key idea:
  - keep the **generation network**  $G = \{\mathbb{P}_{\theta}(h), \mathbb{P}_{\theta}(x | h)\}$
  - drop the inference network: use a **discriminator network**  $D : \mathcal{X} \rightarrow \{0, 1\}$
- Formulate the learning problem as a game between two players:
  - the generator's job is to generate data that looks real
  - the discriminator's job is to distinguish between real data and fake data generated by the generator

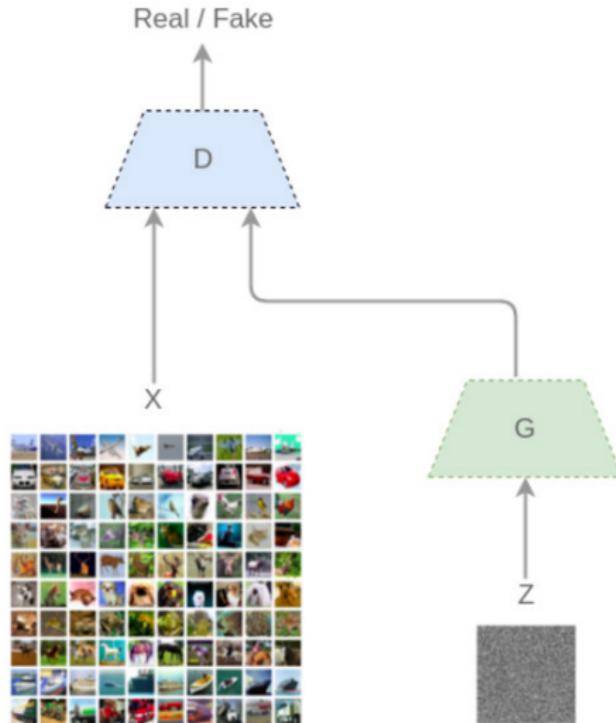
# Generative Adversarial Networks (GANs)

(Goodfellow et al., 2014)

- Key idea:
  - keep the **generation network**  $G = \{\mathbb{P}_{\theta}(h), \mathbb{P}_{\theta}(x | h)\}$
  - drop the inference network: use a **discriminator network**  $D : \mathcal{X} \rightarrow \{0, 1\}$
- Formulate the learning problem as a game between two players:
  - the generator's job is to generate data that looks real
  - the discriminator's job is to distinguish between real data and fake data generated by the generator
- Sort of like a **Turing test**: distinguish artificial from real.

# Generative Adversarial Networks (GANs)

(Goodfellow et al., 2014)



# Minimax Game

- Saddle point problem:

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbb{P}_{\text{data}}(x)}[\log D(x)] + \mathbb{E}_{\mathbb{P}_{\theta}(h)}[\log(1 - D(G(h)))].$$

# Minimax Game

- Saddle point problem:

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbb{P}_{\text{data}}(x)}[\log D(x)] + \mathbb{E}_{\mathbb{P}_{\theta}(h)}[\log(1 - D(G(h)))].$$

- The optimal discriminator (intractable to compute) is:

$$D^*(x) = \frac{\mathbb{P}_{\text{data}}(x)}{\mathbb{P}_{\text{data}}(x) + \mathbb{P}_{\theta}(x)}, \quad \mathbb{P}_{\theta}(x) = \int \mathbb{P}_{\theta}(x | h) \mathbb{P}_{\theta}(h).$$

# Minimax Game

- Saddle point problem:

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbb{P}_{\text{data}}(x)}[\log D(x)] + \mathbb{E}_{\mathbb{P}_{\theta}(h)}[\log(1 - D(G(h)))].$$

- The optimal discriminator (intractable to compute) is:

$$D^*(x) = \frac{\mathbb{P}_{\text{data}}(x)}{\mathbb{P}_{\text{data}}(x) + \mathbb{P}_{\theta}(x)}, \quad \mathbb{P}_{\theta}(x) = \int \mathbb{P}_{\theta}(x | h) \mathbb{P}_{\theta}(h).$$

- The optimal generator  $G^*$  minimizes the **Jensen-Shannon divergence** between  $\mathbb{P}_{\text{data}}(x)$  and  $\mathbb{P}_{\theta}(x)$ :

$$JS(\mathbb{P}_{\text{data}}(x), \mathbb{P}_{\theta}(x)) = \frac{1}{2} KL(\mathbb{P}_{\text{data}}(x) \| \bar{\mathbb{P}}(x)) + \frac{1}{2} KL(\mathbb{P}_{\theta}(x) \| \bar{\mathbb{P}}(x)),$$

where  $\bar{\mathbb{P}}(x) = \frac{\mathbb{P}_{\text{data}}(x) + \mathbb{P}_{\theta}(x)}{2}$ .

# Wasserstein GANs (WGANS)

- Find the saddle point ( $\min_G \max_D$ ) of

$$V(D, G) = \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{\text{data}}(\mathbf{x})}[D(\mathbf{x})] - \mathbb{E}_{\mathbf{h} \sim \mathbb{P}_{\theta}(\mathbf{h})}[D(G(\mathbf{h}))].$$

- This is related to the **Wasserstein distance** (also called Earth mover's distance).
- A technical condition is that  $\nabla D$  is bounded; in practice this is ensured with gradient clipping.
- This improves stability and mitigates the mode collapse problem.

# Training GANs

- Use stochastic gradient descent! Alternate between:
  - Stochastic gradients updates of the generator parameters  $\theta$
  - Stochastic gradients updates of the discriminator  $D$ .

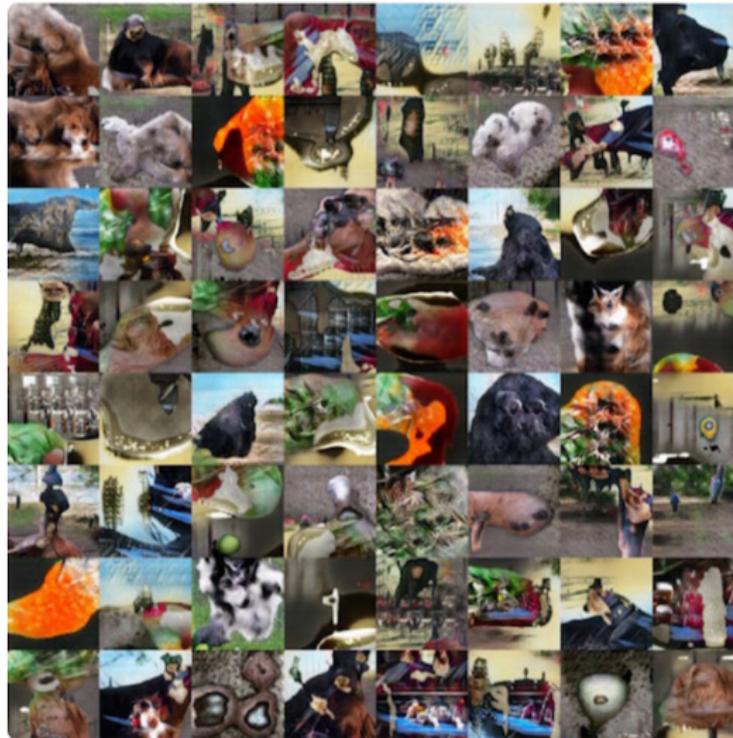
# Training GANs

- Use stochastic gradient descent! Alternate between:
  - Stochastic gradients updates of the generator parameters  $\theta$
  - Stochastic gradients updates of the discriminator  $D$ .
- Several variants and schedules have been proposed.

# Training GANs

- Use stochastic gradient descent! Alternate between:
  - Stochastic gradients updates of the generator parameters  $\theta$
  - Stochastic gradients updates of the discriminator  $D$ .
- Several variants and schedules have been proposed.
- Caveats:
  - no convergence guarantees
  - optimization in GANs is often difficult

# Images Generated by GANs



(<https://tryolabs.com/blog/2016/12/06/major-advancements-deep-learning-2016/>)

# Remarkable Improvement from 2014 to 2017



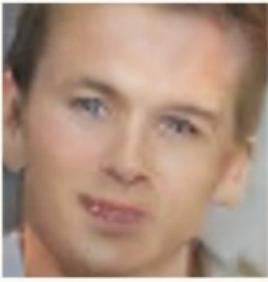
Ian Goodfellow  
@goodfellow\_ian

Follow

4 years of GAN progress (source:  
[eff.org/files/2018/02/ ...](http://eff.org/files/2018/02/))



2014



2015



2016



2017

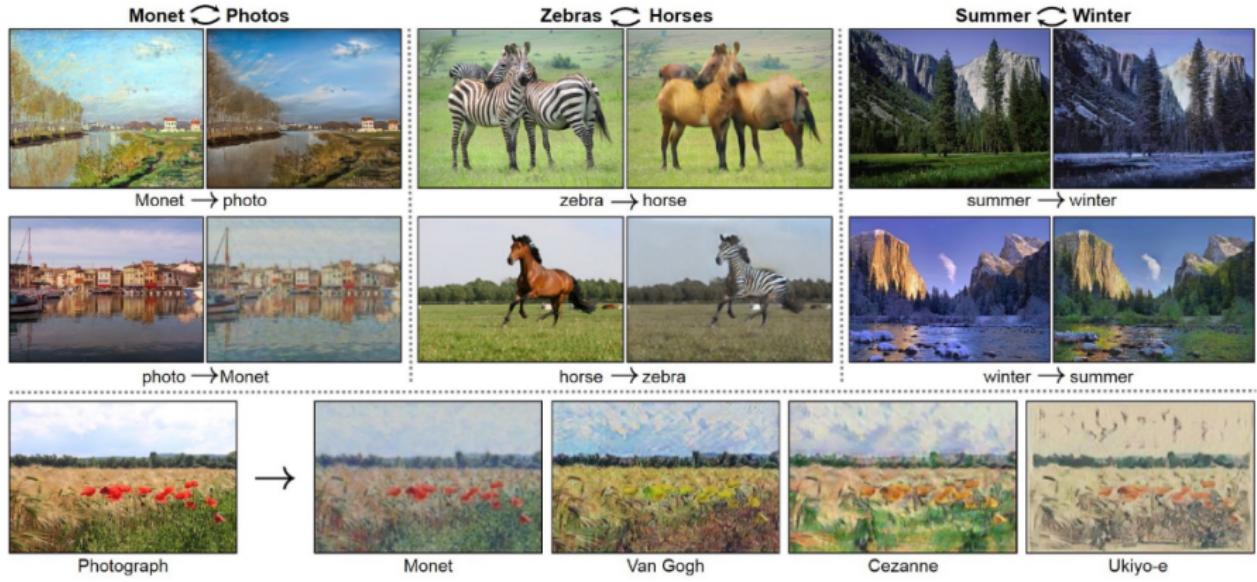
7:26 PM - 2 Mar 2018

# Some Extensions of GANs

- Augmenting GANs with an inference network (Dumoulin et al., 2016; Donahue et al., 2016)
- Domain adversarial training for domain adaptation (Ganin et al., 2016)
- Conditional GANs and semi-supervised GANs (Salimans et al., 2016)
- CycleGAN (Zhu et al., 2017): “translate” images from a source domain  $\mathcal{X}$  to a target domain  $\mathcal{Y}$  without paired examples. Use two generators  $G : \mathcal{Y} \rightarrow \mathcal{X}$  and  $F : \mathcal{X} \rightarrow \mathcal{Y}$  and introduce a **cycle consistency loss** to push  $F(G(y)) \approx y$  and  $G(F(x)) \approx x$ .

# Image-to-Image Translation with CycleGAN

## (Zhu et al., 2017)



(<https://junyanz.github.io/CycleGAN>)

# Failure Cases



(<https://junyanz.github.io/CycleGAN>)

# Examples of Deep Generative Models

- Auto-Regressive Networks ✓
- Restricted Boltzmann Machines ✓
- Deep Belief Networks
- Deep Boltzmann Machines
- Gaussian-Bernoulli RBMs
- Convolutional Boltzmann Machines
- Sigmoid Belief Nets
- Variational Auto-Encoders ✓
- Generative Adversarial Networks ✓
- Denoising diffusion models

# Examples of Deep Generative Models

- Auto-Regressive Networks ✓
- Restricted Boltzmann Machines ✓
- Deep Belief Networks
- Deep Boltzmann Machines
- Gaussian-Bernoulli RBMs
- Convolutional Boltzmann Machines
- Sigmoid Belief Nets
- Variational Auto-Encoders ✓
- Generative Adversarial Networks ✓
- Denoising diffusion models

# Outline

① Boltzmann Machines

② Variational Auto-Encoders

Variational Inference and ELBO

Gradients and Reparameterization Trick

③ Generative Adversarial Networks

④ Denoising Diffusion Models

⑤ Conclusions

# Denoising Diffusion Models

## The Landscape of Deep Generative Learning

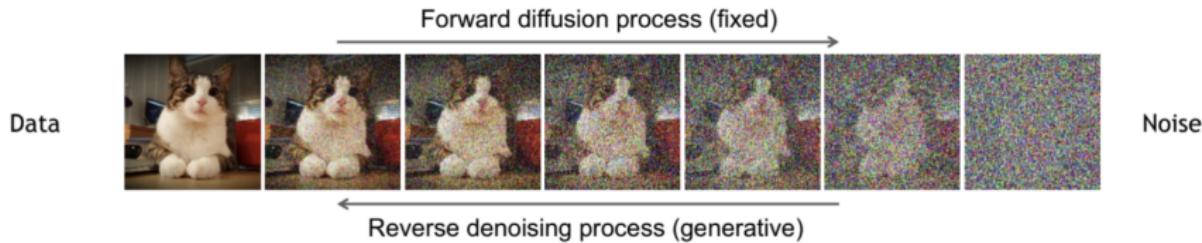


This and the next 10 slides: Diffusion Models Tutorial, CVPR2023, J. Song, C. Meng, A. Vahdat. <https://cvpr2023-tutorial-diffusion-models.github.io/>

# Denoising Diffusion Models: Basic Idea

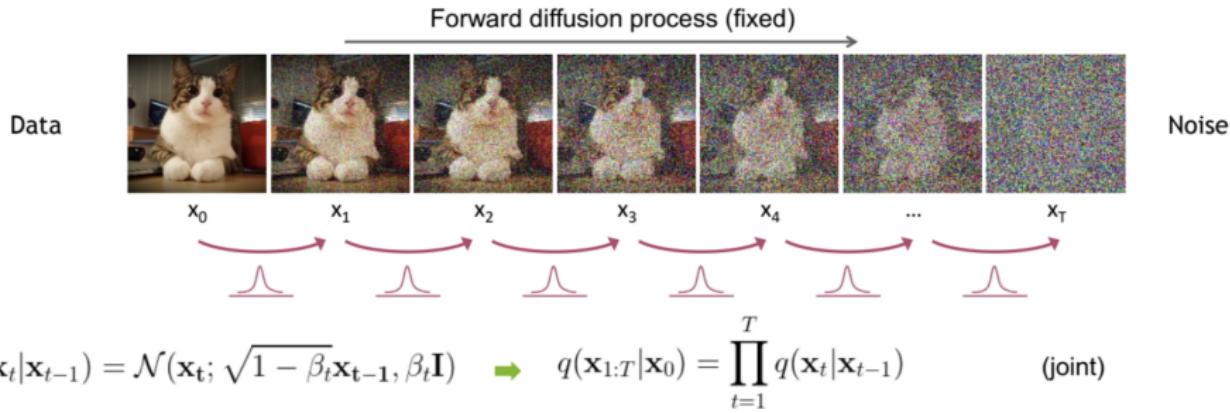
Denoising diffusion models consist of two processes:

- Forward diffusion process that gradually adds noise to input
- Reverse denoising process that learns to generate data by denoising

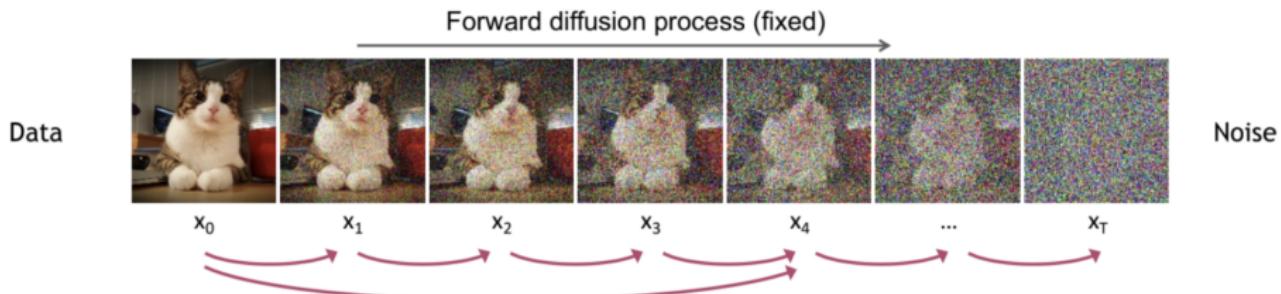


# Forward Diffusion

The formal definition of the forward process in T steps:



# Diffusion Kernel



Define  $\bar{\alpha}_t = \prod_{s=1}^t (1 - \beta_s)$   $\rightarrow q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$  (Diffusion Kernel)

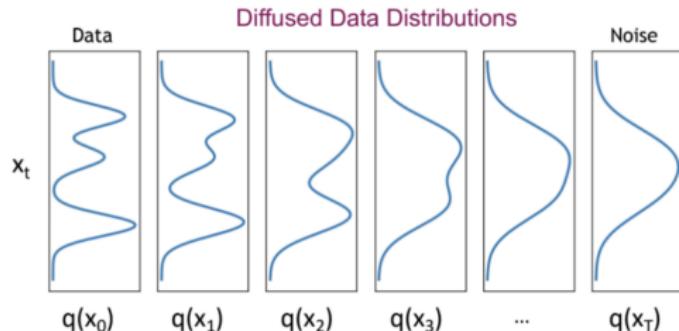
For sampling:  $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{(1 - \bar{\alpha}_t)} \boldsymbol{\epsilon}$  where  $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

$\beta_t$  values schedule (i.e., the noise schedule) is designed such that  $\bar{\alpha}_T \rightarrow 0$  and  $q(\mathbf{x}_T | \mathbf{x}_0) \approx \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$

# Diffusions are Convolutions

So far, we discussed the diffusion kernel  $q(\mathbf{x}_t|\mathbf{x}_0)$  but what about  $q(\mathbf{x}_t)$ ?

$$q(\mathbf{x}_t) = \int \underbrace{q(\mathbf{x}_0, \mathbf{x}_t)}_{\text{Diffused data dist.}} d\mathbf{x}_0 = \int \underbrace{q(\mathbf{x}_0)}_{\text{Input data dist.}} \underbrace{q(\mathbf{x}_t|\mathbf{x}_0)}_{\text{Diffusion kernel}} d\mathbf{x}_0$$



The diffusion kernel is Gaussian convolution.

We can sample  $\mathbf{x}_t \sim q(\mathbf{x}_t)$  by first sampling  $\mathbf{x}_0 \sim q(\mathbf{x}_0)$  and then sampling  $\mathbf{x}_t \sim q(\mathbf{x}_t|\mathbf{x}_0)$  (i.e., ancestral sampling).

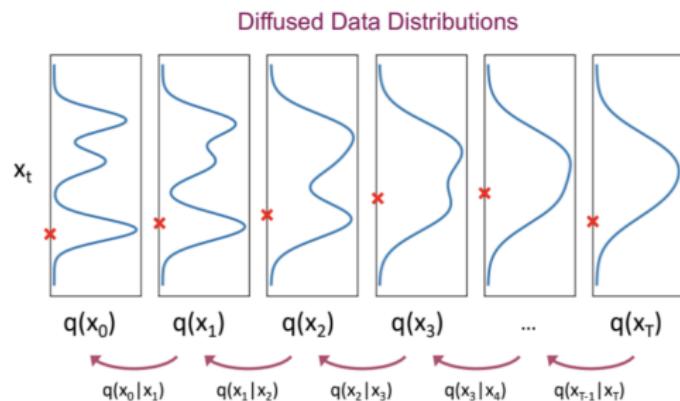
# Generative Learning by Learning to Denoise

Recall, that the diffusion parameters are designed such that  $q(\mathbf{x}_T) \approx \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$

**Generation:**

Sample  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$

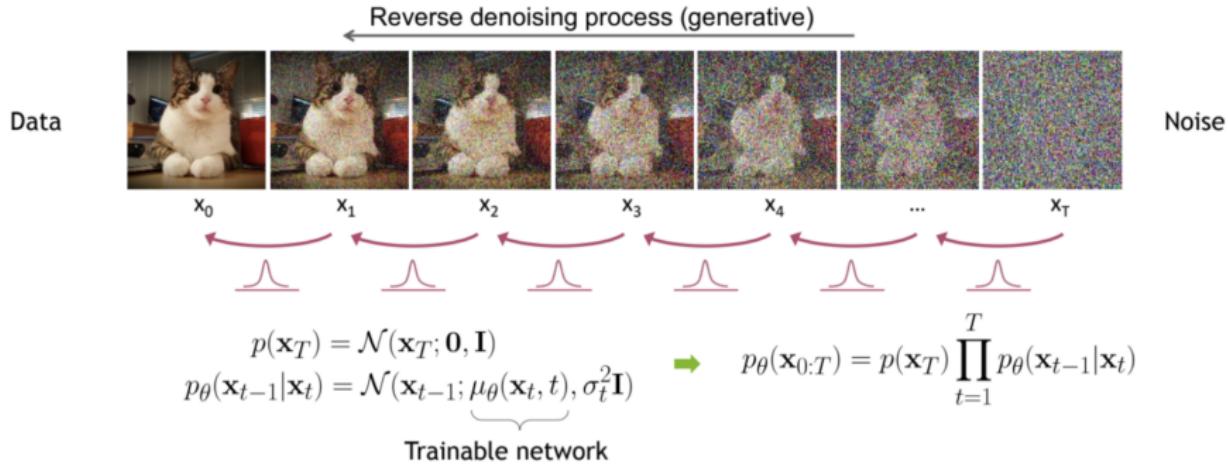
Iteratively sample  $\mathbf{x}_{t-1} \sim \underbrace{q(\mathbf{x}_{t-1} | \mathbf{x}_t)}_{\text{True Denoising Dist.}}$



Can we approximate  $q(\mathbf{x}_{t-1} | \mathbf{x}_t)$ ? Yes, we can use a **Normal distribution** if  $\beta_t$  is small in each forward diffusion step.

# Reverse Denoising Process

Formal definition of forward and reverse processes in T steps:



# Generative Learning by Learning to Denoise

For training, we can form variational upper bound that is commonly used for training variational autoencoders:

$$\mathbb{E}_{q(\mathbf{x}_0)} [-\log p_\theta(\mathbf{x}_0)] \leq \mathbb{E}_{q(\mathbf{x}_0)q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[ -\log \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] =: L$$

Recall that  $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{(1 - \bar{\alpha}_t)} \epsilon$ . [Ho et al. NeurIPS 2020](#) parameterized the mean of denoising model via:

$$\mu_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{1 - \beta_t}} \left( \mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right)$$

Using a few simple arithmetic operations, we can write down the variational objective as:

$$L = \mathbb{E}_{\mathbf{x}_0 \sim q(\mathbf{x}_0), t \sim \mathcal{U}\{1, T\}, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[ \lambda_t \|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)\|^2 \right]$$

[Ho et al. NeurIPS 2020](#) observe that simply setting  $\lambda_t$  to 1 for all  $t$  works best in practice.

---

**Algorithm 1** Training

---

```
1: repeat
2:    $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ 
3:    $t \sim \text{Uniform}(\{1, \dots, T\})$ 
4:    $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5:   Take gradient descent step on
      
$$\nabla_{\theta} \| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta} (\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t) \|^2$$

6: until converged
```

---

---

**Algorithm 2** Sampling

---

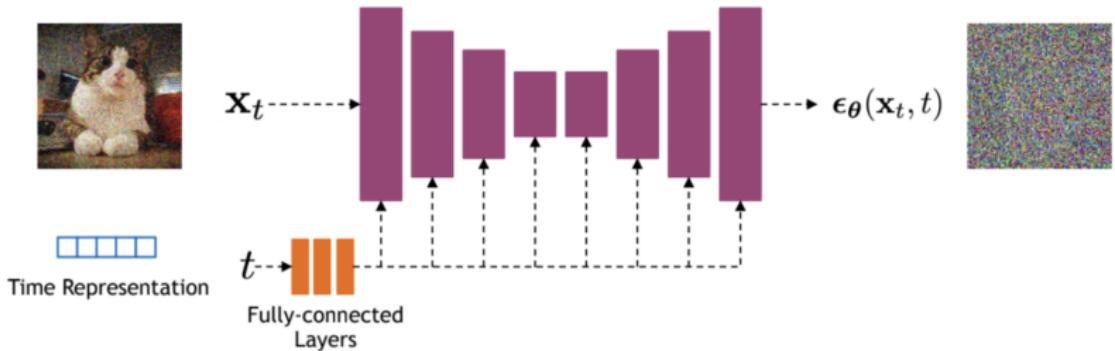
```
1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
4:   
$$\mathbf{x}_{t-1} = \frac{1}{\sqrt{\bar{\alpha}_t}} \left( \mathbf{x}_t - \frac{1 - \bar{\alpha}_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$$

5: end for
6: return  $\mathbf{x}_0$ 
```

---

# Implementation of the Denoisers

Diffusion models often use U-Net architectures with ResNet blocks and self-attention layers to represent  $\epsilon_\theta(\mathbf{x}_t, t)$



Time representation: sinusoidal positional embeddings or random Fourier features.

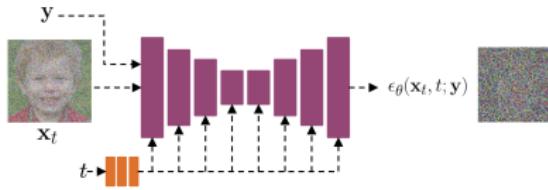
# Conditional Training and Generation

Conditional sampling can be considered as training  $p(\mathbf{x}|\mathbf{y})$  where  $\mathbf{y}$  is the input conditioning (e.g., text) and  $\mathbf{x}$  is generated output (e.g., image)

Train the score model for  $\mathbf{x}$  conditioned on  $\mathbf{y}$  using:

$$\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim p_{\text{data}}(\mathbf{x}, \mathbf{y})} \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \mathbb{E}_{t \sim \mathcal{U}[0, T]} \|\epsilon_\theta(\mathbf{x}_t, t; \mathbf{y}) - \epsilon\|_2^2$$

The conditional score is simply a U-Net with  $\mathbf{x}_t$  and  $\mathbf{y}$  together in the input.



# Text-to-Image Generation Examples

## DALL-E 2

*"a propaganda poster depicting a cat dressed as french emperor napoleon holding a piece of cheese"*



## IMAGEN

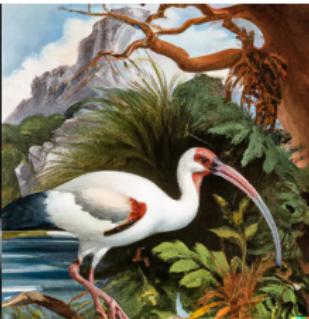
*"A photo of a raccoon wearing an astronaut helmet, looking out of the window at night."*



Imagen

## Artificial Imagination

OpenAI DALL-E 2



DALL-E

# Outline

① Boltzmann Machines

② Variational Auto-Encoders

Variational Inference and ELBO

Gradients and Reparameterization Trick

③ Generative Adversarial Networks

④ Denoising Diffusion Models

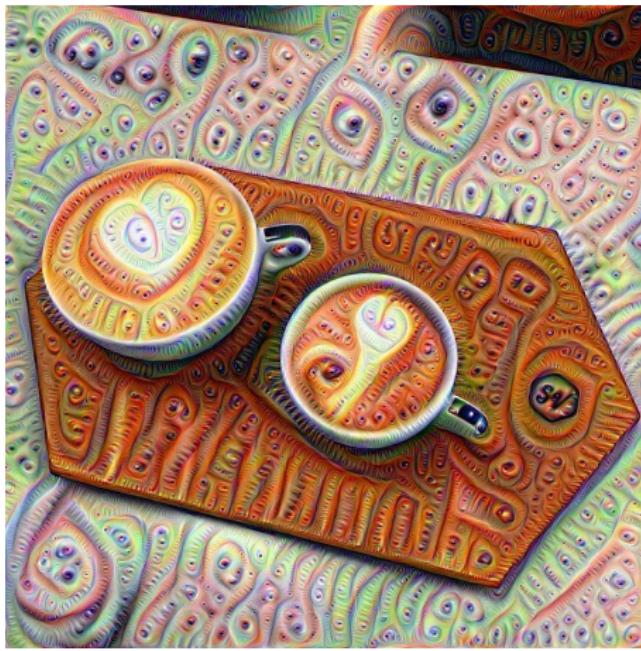
⑤ Conclusions

# Conclusions

- Generative models are useful to model high-dimensional data.
- Latent-variable generative models are appealing since they are more compact.
- Often, computing evidence and posterior distributions is intractable.
- A common surrogate for maximum likelihood is the evidence lower bound (ELBO).
- Variational auto-encoders optimize the ELBO with amortized VI.
- Generative adversarial networks (GANs) are formulated as a game between a generator and a discriminator.
- Denoising diffusion models have impressive results and can be conditioned by text.
- Open problem (for VAE, GAN, and diffusion models): discrete data.

# Thank you!

Questions?



# References I

- Ackley, D., Hinton, G., and Sejnowski, T. (1985). A learning algorithm for Boltzmann machines. *Cognitive science*, 9(1):147–169.
- Donahue, J., Krähenbühl, P., and Darrell, T. (2016). Adversarial feature learning. *arXiv preprint arXiv:1605.09782*.
- Dumoulin, V., Belghazi, I., Poole, B., Mastropietro, O., Lamb, A., Arjovsky, M., and Courville, A. (2016). Adversarially learned inference. *arXiv preprint arXiv:1606.00704*.
- Finn, C., Christiano, P., Abbeel, P., and Levine, S. (2016). A connection between generative adversarial networks, inverse reinforcement learning, and energy-based models. *arXiv preprint arXiv:1611.03852*.
- Ganin, Y., Ustinova, E., Ajakan, H., Germain, P., Larochelle, H., Laviolette, F., Marchand, M., and Lempitsky, V. (2016). Domain-adversarial training of neural networks. *The Journal of Machine Learning Research*, 17(1):2096–2030.
- Goodfellow, I., Bengio, Y., and Courville, A. (2016). Deep Learning. Book in preparation for MIT Press.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. (2014). Generative adversarial nets. In *Advances in Neural Information Processing Systems*, pages 2672–2680.
- Jang, E., Gu, S., and Poole, B. (2016). Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*.
- Kingma, D. P. and Welling, M. (2013). Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*.
- Le Roux, N. and Bengio, Y. (2008). Representational power of restricted boltzmann machines and deep belief networks. *Neural computation*, 20(6):1631–1649.
- Maddison, C. J., Mnih, A., and Teh, Y. W. (2016). The concrete distribution: A continuous relaxation of discrete random variables. *arXiv preprint arXiv:1611.00712*.
- Metz, L., Poole, B., Pfau, D., and Sohl-Dickstein, J. (2016). Unrolled generative adversarial networks. *arXiv preprint arXiv:1611.02163*.
- Oord, A. v. d., Kalchbrenner, N., and Kavukcuoglu, K. (2016). Pixel Recurrent Neural Networks. In *Proc. of the International Conference on Machine Learning*.
- Pfau, D. and Vinyals, O. (2016). Connecting generative adversarial networks and actor-critic methods. *arXiv preprint arXiv:1610.01945*.
- Salimans, T., Goodfellow, I., Zaremba, W., Cheung, V., Radford, A., and Chen, X. (2016). Improved techniques for training gans. In *Advances in Neural Information Processing Systems*, pages 2234–2242.

## References II

- Smolensky, P. (1986). Information processing in dynamical systems: Foundations of harmony theory. Technical report, DTIC Document.
- Theis, L., Oord, A. v. d., and Bethge, M. (2015). A note on the evaluation of generative models. *arXiv preprint arXiv:1511.01844*.
- Wainwright, M. and Jordan, M. (2008). *Graphical Models, Exponential Families, and Variational Inference*. Now Publishers.
- Williams, R. J. (1992). Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256.
- Zhu, J.-Y., Park, T., Isola, P., and Efros, A. A. (2017). Unpaired image-to-image translation using cycle-consistent adversarial networks. *arXiv preprint*.