

# Probabilistic Graphical Models

Deep generative models: overview of the theoretical basis and connections

Eric Xing

Lecture 12, February 25, 2020

**Reading: see class homepage**





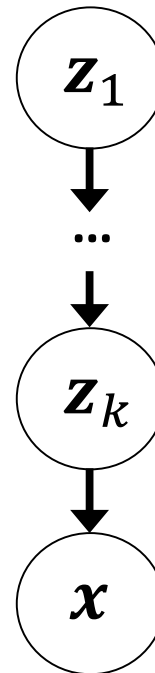
# Deep generative models





# Deep generative models

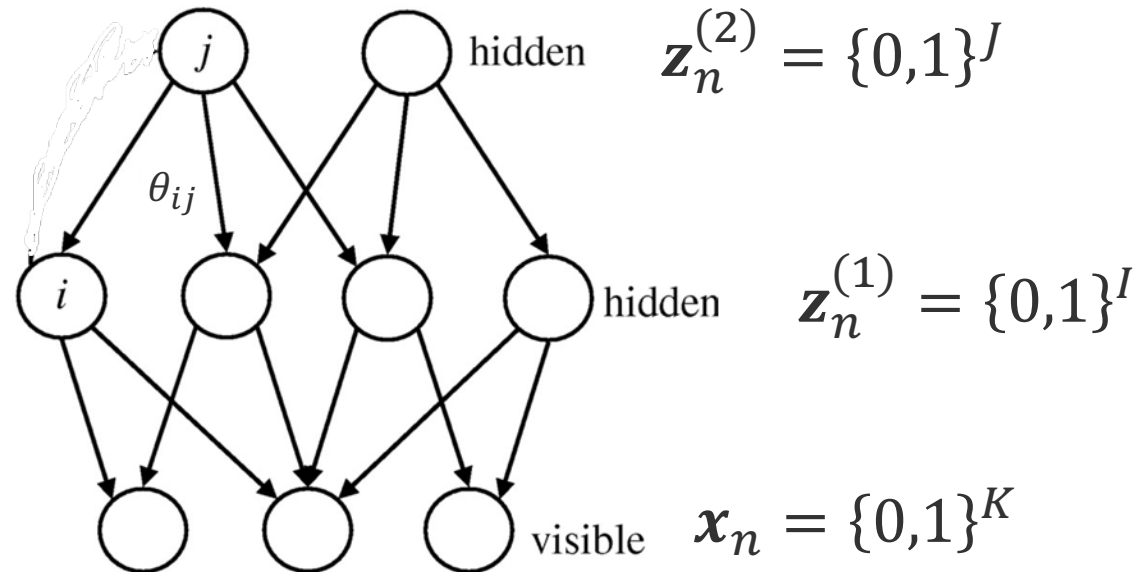
- Define probabilistic distributions over a set of variables
- "Deep" means multiple layers of hidden variables!





# Early forms of deep generative models

- Hierarchical Bayesian models
  - Sigmoid belief nets [Neal 1992]



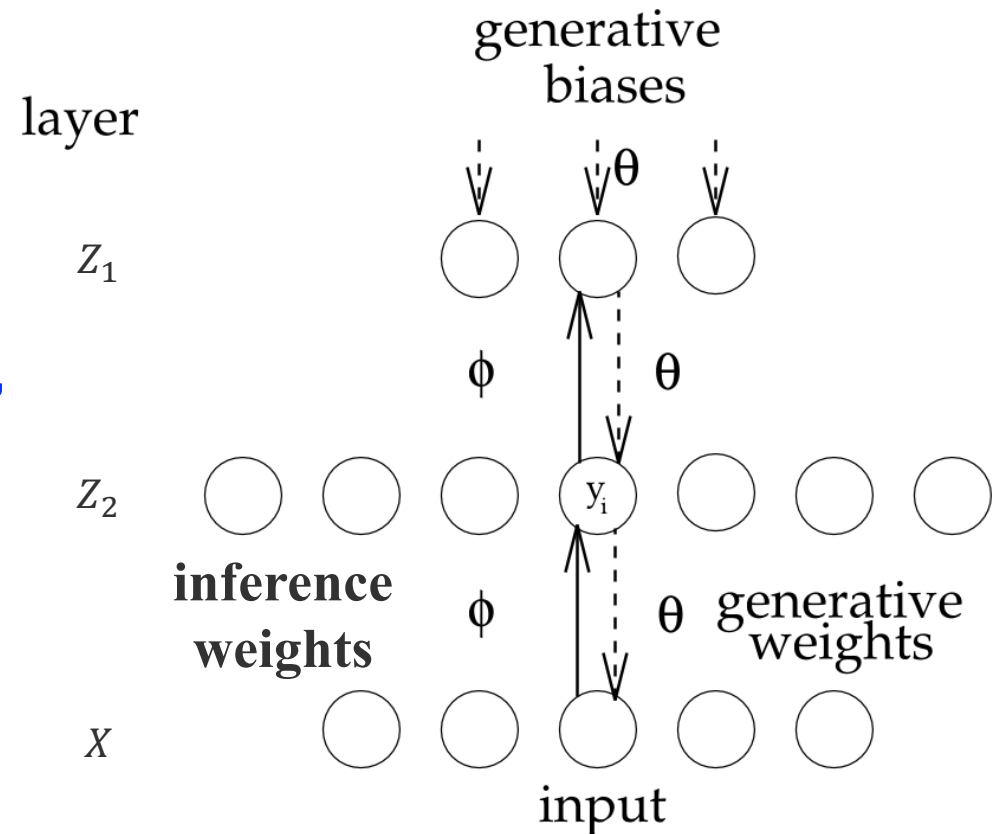
$$p\left(x_{kn} = 1 \mid \boldsymbol{\theta}_k, \mathbf{z}_n^{(1)}\right) = \sigma\left(\boldsymbol{\theta}_k^T \mathbf{z}_n^{(1)}\right)$$
$$p\left(z_{in}^{(1)} = 1 \mid \boldsymbol{\theta}_i, \mathbf{z}_n^{(2)}\right) = \sigma\left(\boldsymbol{\theta}_i^T \mathbf{z}_n^{(2)}\right)$$





# Early forms of deep generative models

- Hierarchical Bayesian models
  - Sigmoid belief nets [Neal 1992]
- Neural network models
  - Helmholtz machines [Dayan et al., 1995]
    - alternative inference/learning methods
    - the process, not a global math expression, defines the model



[Dayan et al. 1995]





# Early forms of deep generative models

- Hierarchical Bayesian models
  - Sigmoid belief nets [Neal 1992]
- Neural network models
  - Helmholtz machines [Dayan et al., 1995]
    - alternative inference/learning methods
    - the process, not a global math expression, defines the model
  - Predictability minimization [Schmidhuber 1995]
    - alternative loss-functions
    - again, no explicit global math expression. the training procedure defines it implicitly

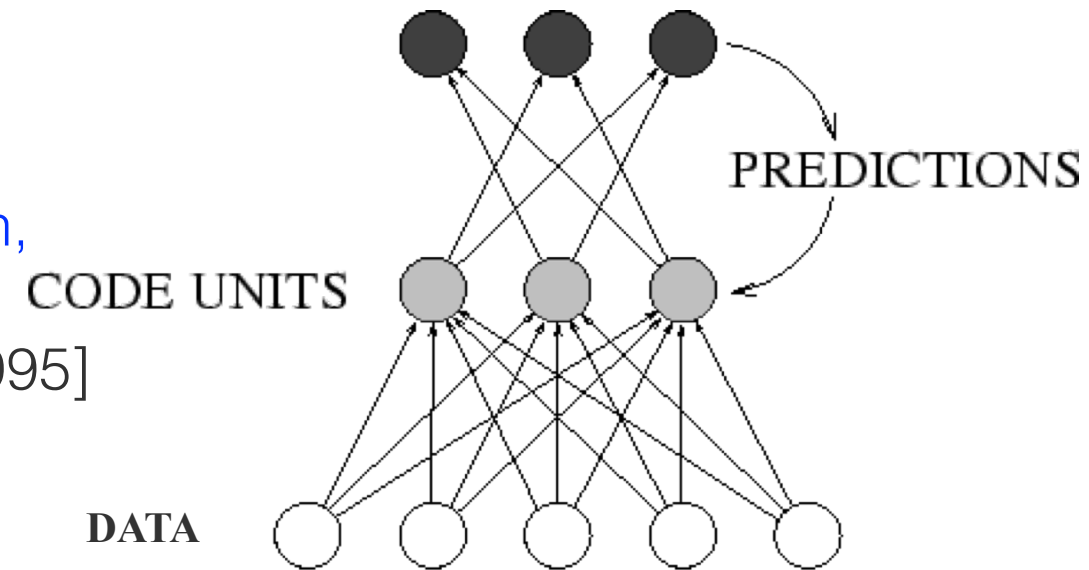


Figure courtesy: Schmidhuber 1996

**The word “model” is here not very rigorous anymore!**







# Early forms of deep generative models

- Training of DGMs via an EM style framework
  - Sampling / data augmentation

$$\mathbf{z} = \{\mathbf{z}_1, \mathbf{z}_2\}$$

$$\mathbf{z}_1^{new} \sim p(\mathbf{z}_1 | \mathbf{z}_2, \mathbf{x})$$

$$\mathbf{z}_2^{new} \sim p(\mathbf{z}_2 | \mathbf{z}_1^{new}, \mathbf{x})$$

- Variational inference

$$\log p(\mathbf{x}) \geq \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}, \mathbf{z})] - \text{KL}(q_\phi(\mathbf{z}|\mathbf{x}) || p(\mathbf{z})) := \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x})$$

$$\max_{\boldsymbol{\theta}, \boldsymbol{\phi}} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x})$$

- Wake sleep

$$\text{Wake: } \min_{\boldsymbol{\theta}} \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})]$$

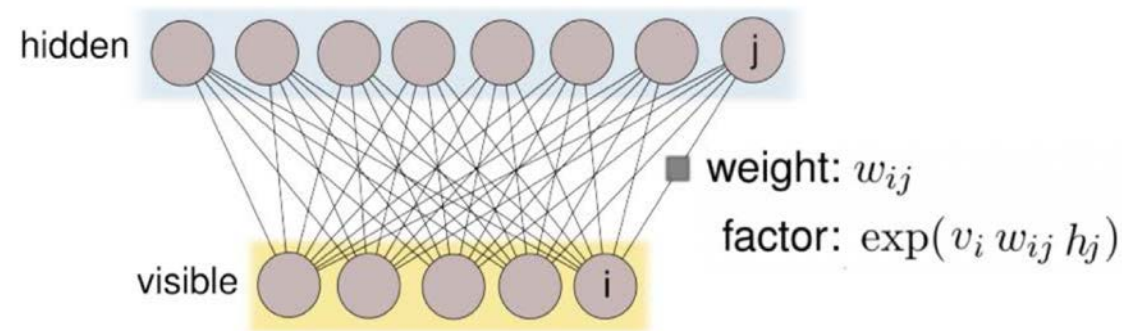
$$\text{Sleep: } \min_{\boldsymbol{\phi}} \mathbb{E}_{p_\theta(\mathbf{x}|\mathbf{z})}[\log q_\phi(\mathbf{z}|\mathbf{x})]$$





# Resurgence of deep generative models

- Restricted Boltzmann machines (RBMs) [Smolensky, 1986]
  - Building blocks of deep probabilistic models



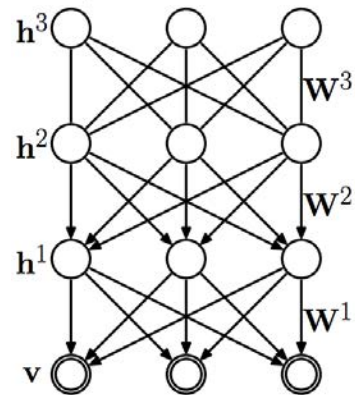




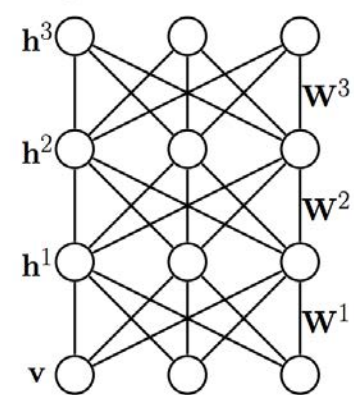
# Resurgence of deep generative models

- ❑ Restricted Boltzmann machines (RBMs) [Smolensky, 1986]
  - ❑ Building blocks of deep probabilistic models
- ❑ Deep belief networks (DBNs) [Hinton et al., 2006]
  - ❑ Hybrid graphical model
  - ❑ Inference in DBNs is problematic due to explaining away
- ❑ Deep Boltzmann Machines (DBMs) [Salakhutdinov & Hinton, 2009]
  - ❑ Undirected model

Deep Belief Network



Deep Boltzmann Machine





# Resurgence of deep generative models

- ▣ Variational autoencoders (VAEs) [Kingma & Welling, 2014]  
/ Neural Variational Inference and Learning (NVIL) [Mnih & Gregor, 2014]

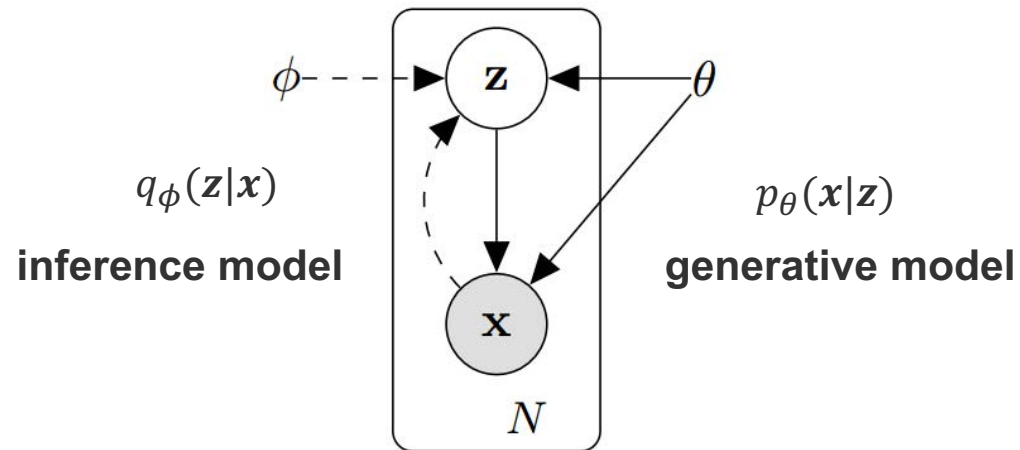


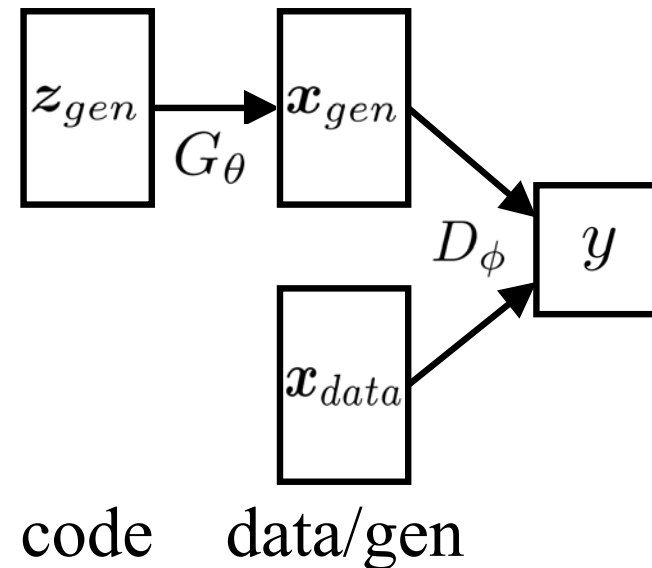
Figure courtesy: Kingma & Welling, 2014





# Resurgence of deep generative models

- Variational autoencoders (VAEs) [Kingma & Welling, 2014]  
/ Neural Variational Inference and Learning (NVIL) [Mnih & Gregor, 2014]
- Generative adversarial networks (GANs) [Goodfellow et al., 2014]



$G_\theta$ : generative model ?

$D_\phi$ : discriminator





# Resurgence of deep generative models

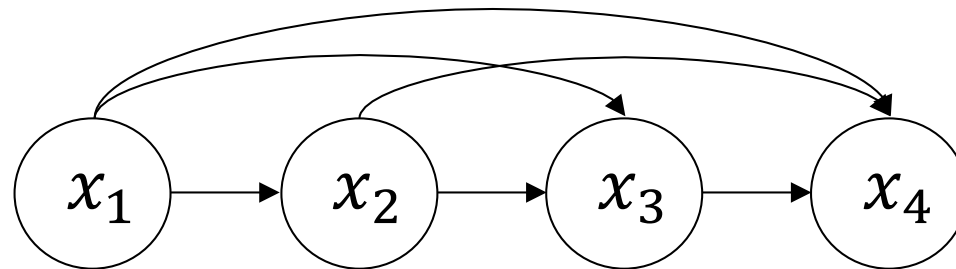
- Variational autoencoders (VAEs) [Kingma & Welling, 2014]  
/ Neural Variational Inference and Learning (NVIL) [Mnih & Gregor, 2014]
- Generative adversarial networks (GANs) [Goodfellow et al., 2014]
- Generative moment matching networks (GMMNs) [Li et al., 2015; Dziugaite et al., 2015]





# Resurgence of deep generative models

- Variational autoencoders (VAEs) [Kingma & Welling, 2014]  
/ Neural Variational Inference and Learning (NVIL) [Mnih & Gregor, 2014]
- Generative adversarial networks (GANs) [Goodfellow et al., 2014]
- Generative moment matching networks (GMMNs) [Li et al., 2015; Dziugaite et al., 2015]
- Autoregressive neural networks





# Outline

- Theoretical Basis of deep generative models
  - Wake sleep algorithm
  - Variational autoencoders
  - Generative adversarial networks
- A unified view of deep generative models
  - new formulations of deep generative models
  - Symmetric modeling of latent and visible variables





# Synonyms in the literature

- ❑ Posterior Distribution -> Inference model
  - ❑ Variational approximation
  - ❑ Recognition model
  - ❑ Inference network (if parameterized as neural networks)
  - ❑ Recognition network (if parameterized as neural networks)
  - ❑ (Probabilistic) encoder
- ❑ "The Model" (prior + conditional, or joint) -> Generative model
  - ❑ The (data) likelihood model
  - ❑ Generative network (if parameterized as neural networks)
  - ❑ Generator
  - ❑ (Probabilistic) decoder







# Recap: Variational Inference

- Consider a generative model  $p_{\theta}(\mathbf{x}|\mathbf{z})$ , and prior  $p(\mathbf{z})$ 
  - Joint distribution:  $p_{\theta}(\mathbf{x}, \mathbf{z}) = p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})$
- Assume **variational distribution**  $q_{\phi}(\mathbf{z}|\mathbf{x})$
- Objective: Maximize **lower bound** for log likelihood

$$\begin{aligned} \log p(\mathbf{x}) \\ &= KL(q_{\phi}(\mathbf{z}|\mathbf{x}) || p_{\theta}(\mathbf{z}|\mathbf{x})) + \int_{\mathbf{z}} q_{\phi}(\mathbf{z}|\mathbf{x}) \log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \\ &\geq \int_{\mathbf{z}} q_{\phi}(\mathbf{z}|\mathbf{x}) \log \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \\ &:= \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}) \end{aligned}$$

- Equivalently, minimize **free energy**

$$F(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}) = -\log p(\mathbf{x}) + KL(q_{\phi}(\mathbf{z}|\mathbf{x}) || p_{\theta}(\mathbf{z}|\mathbf{x}))$$





# Recap: Variational Inference

Maximize the variational lower bound:

$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}) &= \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})}[\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})] + KL(q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z})) \\ &= \log p(\mathbf{x}) - KL(q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}) || p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}))\end{aligned}$$

- **E-step:** maximize  $\mathcal{L}$  wrt.  $\boldsymbol{\phi}$ , with  $\boldsymbol{\theta}$  fixed

$$\max_{\boldsymbol{\phi}} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x})$$

- If closed form solutions exist:

$$q_{\boldsymbol{\phi}}^*(\mathbf{z}|\mathbf{x}) \propto \exp[\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})]$$

- **M-step:** maximize  $\mathcal{L}$  wrt.  $\boldsymbol{\theta}$ , with  $\boldsymbol{\phi}$  fixed

$$\max_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x})$$





# Wake Sleep Algorithm [Hinton et al., Science 1995]

- Train a separate inference model along with the generative model
  - Generally applicable to a wide range of generative models, e.g., Helmholtz machines
- Consider a generative model  $p_{\theta}(\mathbf{x}|\mathbf{z})$  and prior  $p(\mathbf{z})$ 
  - Joint distribution  $p_{\theta}(\mathbf{x}, \mathbf{z}) = p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})$
  - E.g., multi-layer belief nets
- Inference model  $q_{\phi}(\mathbf{z}|\mathbf{x})$
- Maximize data log-likelihood with **two steps of loss relaxation**:
  - Maximize the **variational lower bound** of log-likelihood, or equivalently, minimize the free energy

$$F(\theta, \phi; \mathbf{x}) = -\log p(\mathbf{x}) + KL(q_{\phi}(\mathbf{z}|\mathbf{x}) || p_{\theta}(\mathbf{z}|\mathbf{x}))$$

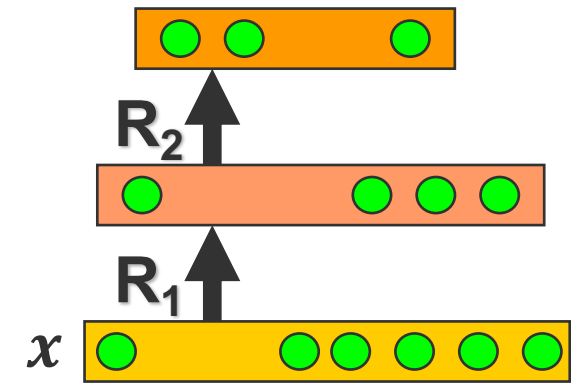
- Minimize a different objective (**reversed KLD**) wrt  $\phi$  to ease the optimization
  - Disconnect to the original variational lower bound loss

$$F'(\theta, \phi; \mathbf{x}) = -\log p(\mathbf{x}) + KL(p_{\theta}(\mathbf{z}|\mathbf{x}) || q_{\phi}(\mathbf{z}|\mathbf{x}))$$





# Wake Sleep Algorithm



- Free energy:

$$F(\theta, \phi; \mathbf{x}) = -\log p(\mathbf{x}) + KL(q_{\phi}(\mathbf{z}|\mathbf{x}) || p_{\theta}(\mathbf{z}|\mathbf{x}))$$

- Minimize the free energy wrt.  $\theta$  of  $p_{\theta}$   $\rightarrow$  *wake phase*

$$\max_{\theta} E_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}, \mathbf{z})]$$

- Get samples from  $q_{\phi}(\mathbf{z}|\mathbf{x})$  through inference on hidden variables
- Use the samples as targets for updating the generative model  $p_{\theta}(\mathbf{z}|\mathbf{x})$
- Correspond to the variational *M step*





# Wake Sleep Algorithm

- Free energy:

$$F(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}) = -\log p(\mathbf{x}) + KL(q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}) || p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}))$$

- Minimize the free energy wrt.  $\boldsymbol{\phi}$  of  $q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})$

- Correspond to the variational E step

$$\max_{\boldsymbol{\phi}} \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})} [\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})]$$

- Difficulties:

- Optimal  $q_{\boldsymbol{\phi}}^*(\mathbf{z}|\mathbf{x}) = \frac{p_{\boldsymbol{\theta}}(\mathbf{z}, \mathbf{x})}{\int p_{\boldsymbol{\theta}}(\mathbf{z}, \mathbf{x}) d\mathbf{z}}$  intractable

- High variance of direct gradient estimate  $\nabla_{\boldsymbol{\phi}} F(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}) = \dots + \nabla_{\boldsymbol{\phi}} \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})} [\log p_{\boldsymbol{\theta}}(\mathbf{z}, \mathbf{x})] + \dots$

- Gradient estimate with the log-derivative trick:

$$\nabla_{\boldsymbol{\phi}} \mathbb{E}_{q_{\boldsymbol{\phi}}} [\log p_{\boldsymbol{\theta}}] = \int \nabla_{\boldsymbol{\phi}} q_{\boldsymbol{\phi}} \log p_{\boldsymbol{\theta}} = \int q_{\boldsymbol{\phi}} \log p_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\phi}} \log q_{\boldsymbol{\phi}} = \mathbb{E}_{q_{\boldsymbol{\phi}}} [\log p_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\phi}} \log q_{\boldsymbol{\phi}}]$$

- Monte Carlo estimation:

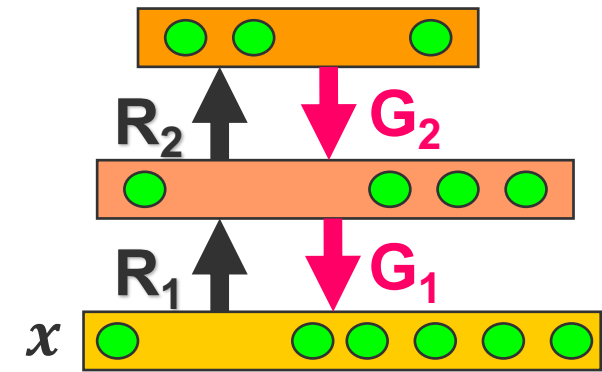
$$\nabla_{\boldsymbol{\phi}} \mathbb{E}_{q_{\boldsymbol{\phi}}} [\log p_{\boldsymbol{\theta}}] \approx \mathbb{E}_{\mathbf{z}_i \sim q_{\boldsymbol{\phi}}} [\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}_i) \nabla_{\boldsymbol{\phi}} \log q_{\boldsymbol{\phi}}(\mathbf{z}_i|\mathbf{x})]$$

- The scale factor  $\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}_i)$  can have arbitrary large magnitude





# Wake Sleep Algorithm



- Free energy:

$$F(\theta, \phi; \mathbf{x}) = -\log p(\mathbf{x}) + KL(q_{\phi}(\mathbf{z}|\mathbf{x}) || p_{\theta}(\mathbf{z}|\mathbf{x}))$$

- WS works around the difficulties with the **sleep phase approximation**
- Minimize the following objective  $\rightarrow$  *sleep* phase

$$F'(\theta, \phi; \mathbf{x}) = -\log p(\mathbf{x}) + KL(p_{\theta}(\mathbf{z}|\mathbf{x}) || q_{\phi}(\mathbf{z}|\mathbf{x}))$$

$$\max_{\phi} E_{p_{\theta}(\mathbf{z}, \mathbf{x})} [\log q_{\phi}(\mathbf{z}|\mathbf{x})]$$

$$\max_{\phi} E_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}, \mathbf{z})]$$

- “Dreaming” up samples from  $p_{\theta}(\mathbf{x}|\mathbf{z})$  through top-down pass
- Use the samples as targets for updating the inference model
- (Recent approaches other than sleep phase are developed to reduce the variance of gradient estimate: slides later)





# Wake Sleep Algorithm

## Wake sleep

- Parametrized inference model  $q_\phi(\mathbf{z}|\mathbf{x})$
- Wake phase:
  - minimize  $KL(q_\phi(\mathbf{z}|\mathbf{x}) || p_\theta(\mathbf{z}|\mathbf{x}))$  wrt.  $\theta$
  - $E_{q_\phi(\mathbf{z}|\mathbf{x})} [\nabla_\theta \log p_\theta(\mathbf{x}|\mathbf{z})]$
- Sleep phase:
  - minimize  $KL(p_\theta(\mathbf{z}|\mathbf{x}) || q_\phi(\mathbf{z}|\mathbf{x}))$  wrt.  $\phi$
  - $E_{p_\theta(\mathbf{z},\mathbf{x})} [\nabla_\phi \log q_\phi(\mathbf{z}, \mathbf{x})]$
  - low variance
  - Learning with generated samples of  $\mathbf{x}$
- Two objective, not guaranteed to converge

## Variational EM

- Variational distribution  $q_\phi(\mathbf{z}|\mathbf{x})$
- Variational M step:
  - minimize  $KL(q_\phi(\mathbf{z}|\mathbf{x}) || p_\theta(\mathbf{z}|\mathbf{x}))$  wrt.  $\theta$
  - $E_{q_\phi(\mathbf{z}|\mathbf{x})} [\nabla_\theta \log p_\theta(\mathbf{x}|\mathbf{z})]$
- Variational E step:
  - minimize  $KL(q_\phi(\mathbf{z}|\mathbf{x}) || p_\theta(\mathbf{z}|\mathbf{x}))$  wrt.  $\phi$
  - $q_\phi^* \propto \exp[\log p_\theta]$  if with closed-form
  - $\nabla_\phi E_{q_\phi} [\log p_\theta(\mathbf{z}, \mathbf{x})]$ 
    - need variance-reduce in practice
  - Learning with real data  $\mathbf{x}$
- Single objective, guaranteed to converge







# Variational Autoencoders (VAEs)

- [Kingma & Welling, 2014]
- Use variational inference with an inference model
  - Enjoy similar applicability with wake-sleep algorithm
- Generative model  $p_{\theta}(\mathbf{x}|\mathbf{z})$ , and prior  $p(\mathbf{z})$ 
  - Joint distribution  $p_{\theta}(\mathbf{x}, \mathbf{z}) = p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})$
- Inference model  $q_{\phi}(\mathbf{z}|\mathbf{x})$

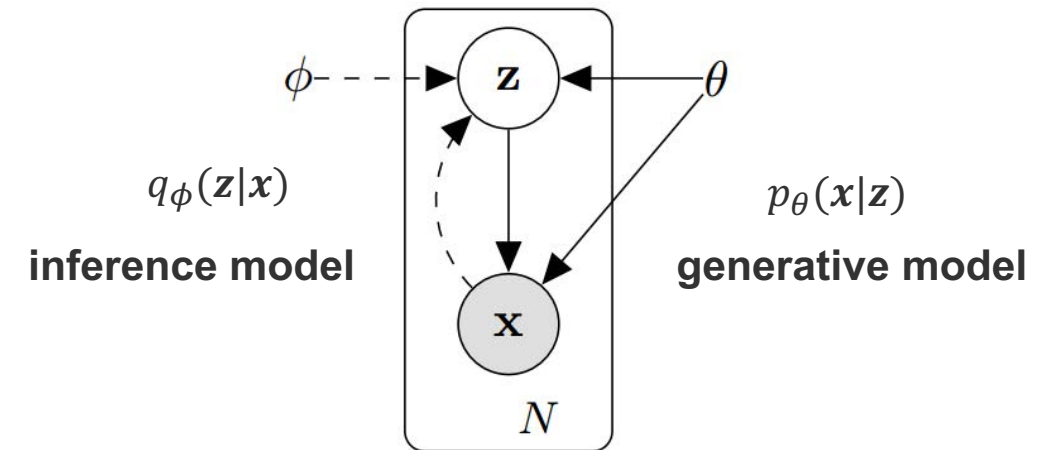


Figure courtesy: Kingma & Welling, 2014





# Variational Autoencoders (VAEs)

- Variational lower bound

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}) = \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})}[\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})] - \text{KL}(q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))$$

- Optimize  $\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x})$  wrt.  $\boldsymbol{\theta}$  of  $p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})$ 
  - The same with the wake phase
- Optimize  $\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x})$  wrt.  $\boldsymbol{\phi}$  of  $q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})$

$$\nabla_{\boldsymbol{\phi}} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}) = \cdots + \nabla_{\boldsymbol{\phi}} \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})}[\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})] + \cdots$$

- Use *reparameterization trick* to reduce variance
- Alternatives: use control variates as in reinforcement learning [Mnih & Gregor, 2014; Paisley et al., 2012]





# Reparametrized gradient

- Optimize  $\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x})$  wrt.  $\boldsymbol{\phi}$  of  $q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})$ 
  - Recap: gradient estimate with log-derivative trick:

$$\nabla_{\boldsymbol{\phi}} \mathbb{E}_{q_{\boldsymbol{\phi}}} [\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})] = \mathbb{E}_{q_{\boldsymbol{\phi}}} [\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}) \nabla_{\boldsymbol{\phi}} \log q_{\boldsymbol{\phi}}]$$

- High variance:  $\nabla_{\boldsymbol{\phi}} \mathbb{E}_{q_{\boldsymbol{\phi}}} [\log p_{\boldsymbol{\theta}}] \approx \mathbb{E}_{\mathbf{z}_i \sim q_{\boldsymbol{\phi}}} [\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}_i) \nabla_{\boldsymbol{\phi}} \log q_{\boldsymbol{\phi}}(\mathbf{z}_i|\mathbf{x})]$ 
    - The scale factor  $\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}_i)$  can have arbitrary large magnitude
- gradient estimate with *reparameterization trick*

$$\mathbf{z} \sim q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}) \iff \mathbf{z} = \mathbf{g}_{\boldsymbol{\phi}}(\boldsymbol{\epsilon}, \mathbf{x}), \quad \boldsymbol{\epsilon} \sim p(\boldsymbol{\epsilon})$$

$$\nabla_{\boldsymbol{\phi}} \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})} [\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})] = \mathbb{E}_{\boldsymbol{\epsilon} \sim p(\boldsymbol{\epsilon})} \left[ \nabla_{\boldsymbol{\phi}} \log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}_{\boldsymbol{\phi}}(\boldsymbol{\epsilon})) \right]$$

- (Empirically) lower variance of the gradient estimate
  - E.g.,  $\mathbf{z} \sim N(\boldsymbol{\mu}(\mathbf{x}), L(\mathbf{x})L(\mathbf{x})^T) \iff \boldsymbol{\epsilon} \sim N(0,1), \mathbf{z} = \boldsymbol{\mu}(\mathbf{x}) + L(\mathbf{x})\boldsymbol{\epsilon}$





# VAEs: algorithm

---

**Algorithm 1** Minibatch version of the Auto-Encoding VB (AEVB) algorithm. Either of the two SGVB estimators in section 2.3 can be used. We use settings  $M = 100$  and  $L = 1$  in experiments.

---

$\theta, \phi \leftarrow$  Initialize parameters

**repeat**

$\mathbf{X}^M \leftarrow$  Random minibatch of  $M$  datapoints (drawn from full dataset)

$\epsilon \leftarrow$  Random samples from noise distribution  $p(\epsilon)$

$\mathbf{g} \leftarrow \nabla_{\theta, \phi} \tilde{\mathcal{L}}^M(\theta, \phi; \mathbf{X}^M, \epsilon)$  (Gradients of minibatch estimator (8))

$\theta, \phi \leftarrow$  Update parameters using gradients  $\mathbf{g}$  (e.g. SGD or Adagrad [DHS10])

**until** convergence of parameters  $(\theta, \phi)$

**return**  $\theta, \phi$

---

[Kingma & Welling, 2014]





# VAEs: example results

- VAEs tend to generate **blurred** images due to the mode covering behavior (more later)



Celebrity faces [Radford 2015]

- Latent code interpolation and sentences generation from VAEs [Bowman et al., 2015].

---

“ i want to talk to you . ”  
“*i want to be with you .*”  
“*i do n’t want to be with you .*”  
*i do n’t want to be with you .*  
she did n’t want to be with him .

---





# Generative Adversarial Nets (GANs)

- [Goodfellow et al., 2014]
- Generative model  $\mathbf{x} = G_{\theta}(\mathbf{z})$ ,  $\mathbf{z} \sim p(\mathbf{z})$ 
  - Map noise variable  $\mathbf{z}$  to data space  $\mathbf{x}$
  - Define an **implicit distribution** over  $\mathbf{x}$ :  $p_{g_{\theta}}(\mathbf{x})$ 
    - a stochastic process to simulate data  $\mathbf{x}$
    - Intractable to evaluate likelihood
- Discriminator  $D_{\phi}(\mathbf{x})$ 
  - Output the probability that  $\mathbf{x}$  came from the data rather than the generator
- No explicit inference model
- No obvious connection to previous models with inference networks like VAEs
  - We will build formal connections between GANs and VAEs later

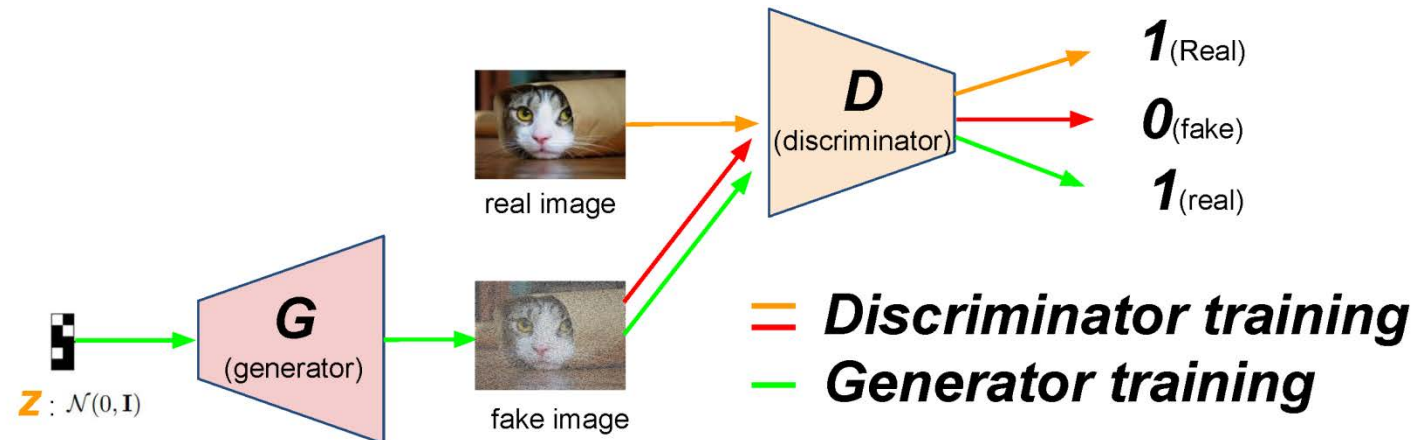




# Generative Adversarial Nets (GANs)

- Learning
  - A minimax game between the generator and the discriminator
  - Train  $D$  to maximize the probability of assigning the correct label to both training examples and generated samples
  - Train  $G$  to fool the discriminator

$$\max_D \mathcal{L}_D = \mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim G(\mathbf{z}), \mathbf{z} \sim p(\mathbf{z})} [\log(1 - D(\mathbf{x}))]$$
$$\min_G \mathcal{L}_G = \mathbb{E}_{\mathbf{x} \sim G(\mathbf{z}), \mathbf{z} \sim p(\mathbf{z})} [\log(1 - D(\mathbf{x}))].$$



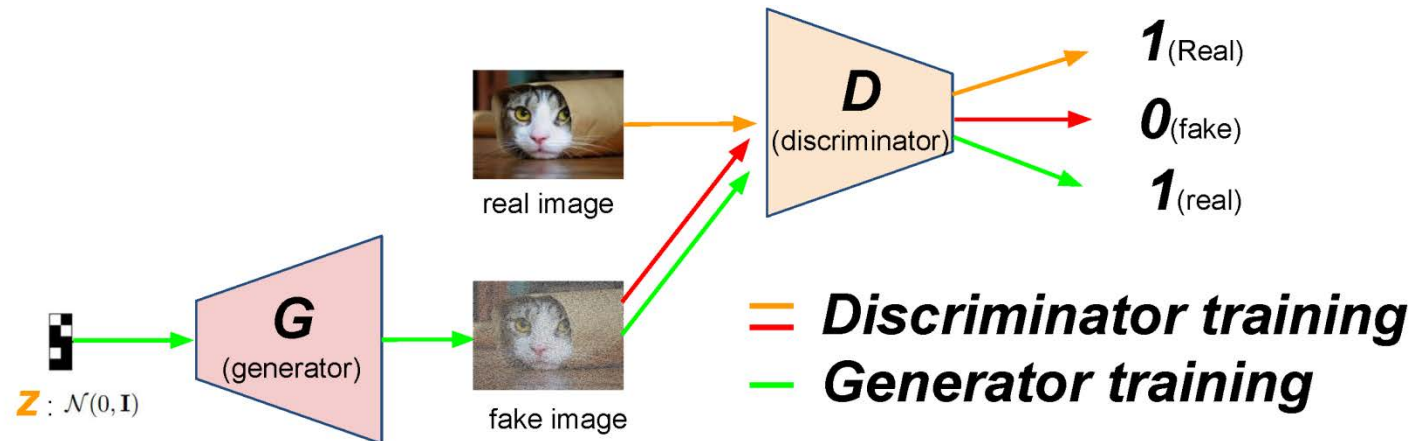




# Generative Adversarial Nets (GANs)

- Learning
  - Train  $G$  to fool the discriminator
    - The original loss suffers from vanishing gradients when  $D$  is too strong
    - Instead use the following in practice

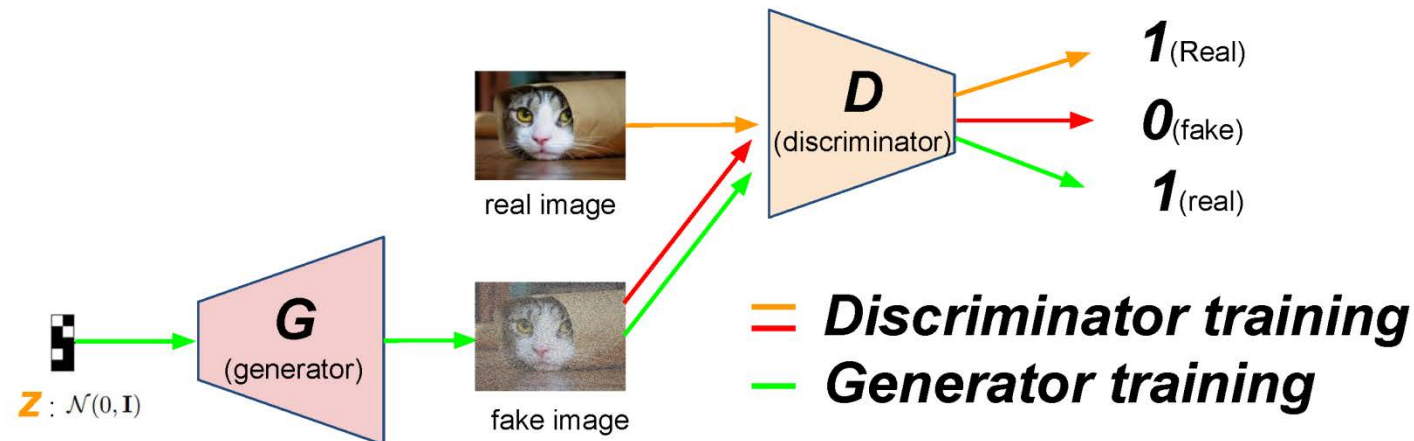
$$\max_G \mathcal{L}_G = \mathbb{E}_{\mathbf{x} \sim G(\mathbf{z}), \mathbf{z} \sim p(\mathbf{z})} [\log D(\mathbf{x})]$$





# Generative Adversarial Nets (GANs)

- Learning
  - Aim to achieve equilibrium of the game
  - Optimal state:
    - $p_g(\mathbf{x}) = p_{data}(\mathbf{x})$
    - $D(\mathbf{x}) = \frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})} = \frac{1}{2}$





# GANs: example results



Generated bedrooms [Radford et al., 2016]





# The Zoo of DGMs

- ❑ Variational autoencoders (VAEs) [Kingma & Welling, 2014]
  - ❑ Adversarial autoencoder [Makhzani et al., 2015]
  - ❑ Importance weighted autoencoder [Burda et al., 2015]
  - ❑ Implicit variational autoencoder [Mescheder., 2017]
- ❑ Generative adversarial networks (GANs) [Goodfelloos et al., 2014]
  - ❑ InfoGAN [Chen et al., 2016]
  - ❑ CycleGAN [Zhu et al., 2017]
  - ❑ Wasserstein GAN [Arjovsky et al., 2017]
- ❑ Autoregressive neural networks
  - ❑ PixelRNN / PixelCNN [Oord et al., 2016]
  - ❑ RNN (e.g., for language modeling)
- ❑ Generative moment matching networks (GMMNs) [Li et al., 2015; Dziugaite et al., 2015]
- ❑ Restricted Boltzmann Machines (RBMs) [Smolensky, 1986]



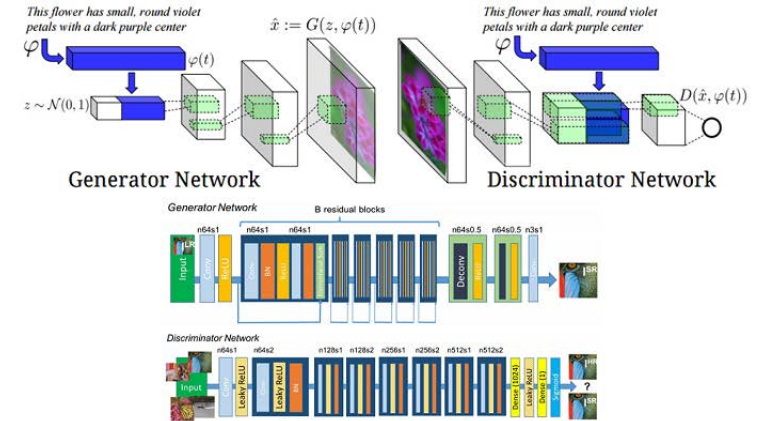
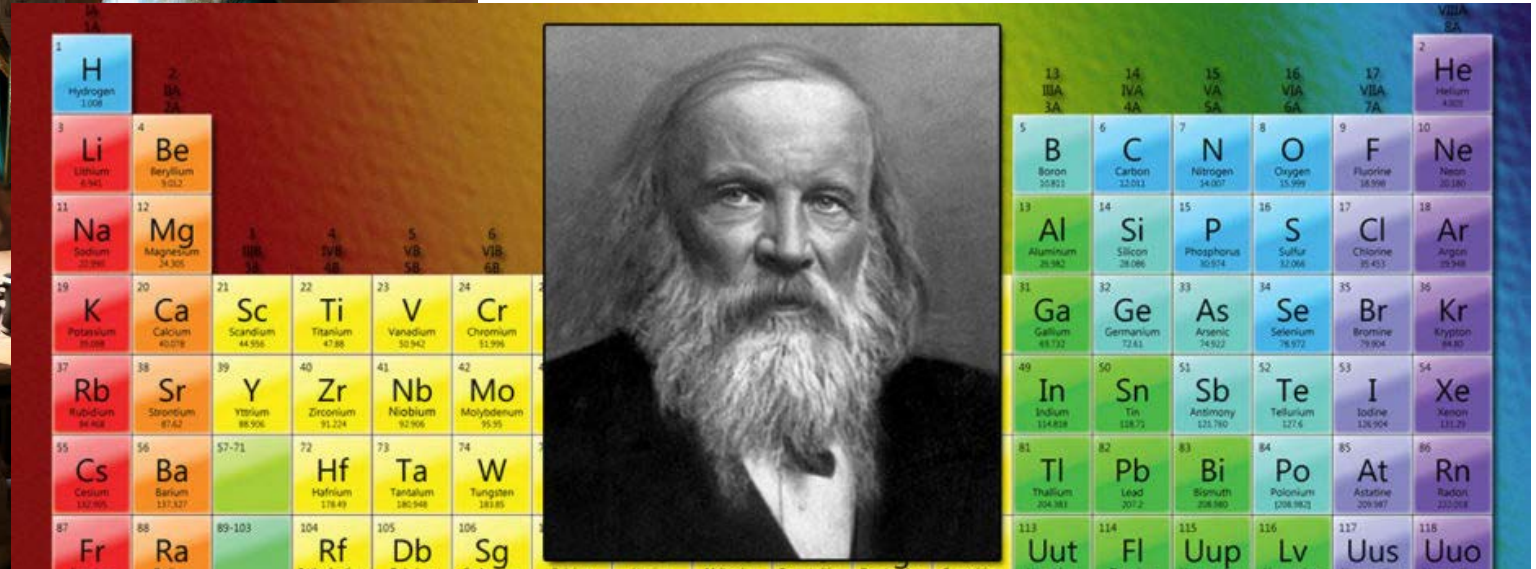




# Alchemy Vs Modern Chemistry



© Mind Juice Media Inc. All rights reserved





# Outline

- Overview of advances in deep generative models
- Theoretical backgrounds of deep generative models
  - Wake sleep algorithm
  - Variational autoencoders
  - Generative adversarial networks
- A unified view of deep generative models
  - new formulations of deep generative models
  - Symmetric modeling of latent and visible variables

**Z Hu, Z YANG, R Salakhutdinov, E Xing,**  
**“On Unifying Deep Generative Models”, arxiv 1706.00550**





# A unified view of deep generative models

- Literatures have viewed these DGM approaches as distinct model training paradigms
  - GANs: achieve an equilibrium between generator and discriminator
  - VAEs: maximize lower bound of the data likelihood
- Let's study a new formulation for DGMs
  - Connects GANs, VAEs, and other variants, under a unified view
  - Links them back to inference and learning of Graphical Models, and the wake-sleep heuristic that approximates this
  - Provides a tool to analyze many GAN-/VAE-based algorithms
  - Encourages mutual exchange of ideas from each individual class of models







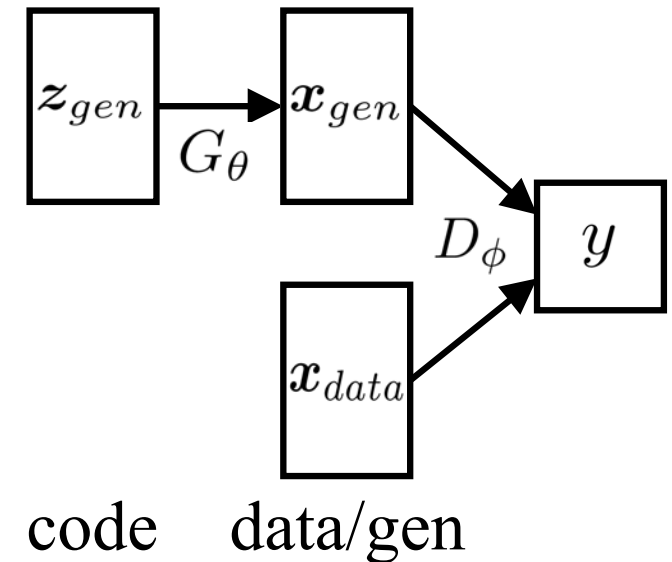
# Generative Adversarial Nets (GANs):

- Implicit distribution over  $\mathbf{x} \sim p_{\theta}(\mathbf{x}|y)$

$$p_{\theta}(\mathbf{x}|y) = \begin{cases} p_{g_{\theta}}(\mathbf{x}) & y = 0 \\ p_{data}(\mathbf{x}) & y = 1. \end{cases}$$

(distribution of generated images)  
(distribution of real images)

- $\mathbf{x} \sim p_{g_{\theta}}(\mathbf{x}) \Leftrightarrow \mathbf{x} = G_{\theta}(\mathbf{z}), \mathbf{z} \sim p(\mathbf{z}|y = 0)$
- $\mathbf{x} \sim p_{data}(\mathbf{x})$ 
  - the code space of  $\mathbf{z}$  is degenerated
  - sample directly from data





# A new formulation

- Rewrite GAN objectives in the "variational-EM" format

- Recap: conventional formulation:

$$\max_{\phi} \mathcal{L}_{\phi} = \mathbb{E}_{\mathbf{x}=G_{\theta}(\mathbf{z}), \mathbf{z} \sim p(\mathbf{z}|y=0)} [\log(1 - D_{\phi}(\mathbf{x}))] + \mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})} [\log D_{\phi}(\mathbf{x})]$$

$$\begin{aligned} \max_{\theta} \mathcal{L}_{\theta} &= \mathbb{E}_{\mathbf{x}=G_{\theta}(\mathbf{z}), \mathbf{z} \sim p(\mathbf{z}|y=0)} [\log D_{\phi}(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})} [\log(1 - D_{\phi}(\mathbf{x}))] \\ &= \mathbb{E}_{\mathbf{x}=G_{\theta}(\mathbf{z}), \mathbf{z} \sim p(\mathbf{z}|y=0)} [\log D_{\phi}(\mathbf{x})] \end{aligned}$$

- Rewrite in the new form

- Implicit distribution over  $\mathbf{x} \sim p_{\theta}(\mathbf{x}|y)$

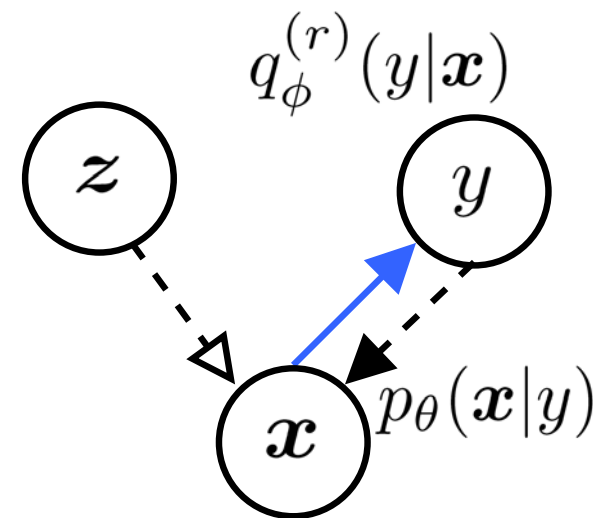
$$\mathbf{x} = G_{\theta}(\mathbf{z}), \mathbf{z} \sim p(\mathbf{z}|y)$$

- Discriminator distribution  $q_{\phi}(y|\mathbf{x})$

$$q_{\phi}^r(y|\mathbf{x}) = q_{\phi}(1 - y|\mathbf{x}) \quad (\text{reverse})$$

$$\max_{\phi} \mathcal{L}_{\phi} = \mathbb{E}_{p_{\theta}(\mathbf{x}|y)p(y)} [\log q_{\phi}(y|\mathbf{x})]$$

$$\max_{\theta} \mathcal{L}_{\theta} = \mathbb{E}_{p_{\theta}(\mathbf{x}|y)p(y)} [\log q_{\phi}^r(y|\mathbf{x})]$$





# GANs vs. Variational EM

## Variational EM

- Objectives

$$\max_{\phi} \mathcal{L}_{\phi, \theta} = \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z)] + KL(q_{\phi}(z|x) || p(z))$$

$$\max_{\theta} \mathcal{L}_{\phi, \theta} = \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z)] + KL(q_{\phi}(z|x) || p(z))$$

- Single objective for both  $\theta$  and  $\phi$
- Extra prior regularization by  $p(z)$
- The reconstruction term: maximize the conditional log-likelihood of  $x$  with the generative distribution  $p_{\theta}(x|z)$  conditioning on the latent code  $z$  inferred by  $q_{\phi}(z|x)$



- $p_{\theta}(x|z)$  is the generative model
- $q_{\phi}(z|x)$  is the inference model

## GAN

- Objectives

$$\max_{\phi} \mathcal{L}_{\phi} = \mathbb{E}_{p_{\theta}(x|y)p(y)} [\log q_{\phi}(y|x)]$$

$$\max_{\theta} \mathcal{L}_{\theta} = \mathbb{E}_{p_{\theta}(x|y)p(y)} [\log q_{\phi}^r(y|x)]$$

- Two objectives
- Have global optimal state in the game theoretic view
- The objectives: maximize the conditional log-likelihood of  $y$  (or  $1 - y$ ) with the distribution  $q_{\phi}(y|x)$  conditioning on data/generation  $x$  inferred by  $p_{\theta}(x|y)$



- Interpret  $q_{\phi}(y|x)$  as the generative model
- Interpret  $p_{\theta}(x|y)$  as the inference model





# GANs vs. Variational EM

- Interpret  $x$  as latent variables
- Interpret generation of  $x$  as performing inference over latent

## Variational EM

### Objectives

$$\max_{\phi} \mathcal{L}_{\phi, \theta} = \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z)] + KL(q_{\phi}(z|x) || p(z))$$

$$\max_{\theta} \mathcal{L}_{\phi, \theta} = \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z)] + KL(q_{\phi}(z|x) || p(z))$$

- Single objective for both  $\theta$  and  $\phi$
- Extra prior regularization by  $p(z)$
- **The reconstruction term:** maximize the conditional log-likelihood of  $x$  with the generative distribution  $p_{\theta}(x|z)$  conditioning on the latent code  $z$  inferred by  $q_{\phi}(z|x)$



- $p_{\theta}(x|z)$  is the generative model
- $q_{\phi}(z|x)$  is the inference model

## GAN

### Objectives

$$\max_{\phi} \mathcal{L}_{\phi} = \mathbb{E}_{p_{\theta}(x|y)p(y)} [\log q_{\phi}(y|x)]$$

$$\max_{\theta} \mathcal{L}_{\theta} = \mathbb{E}_{p_{\theta}(x|y)p(y)} [\log q_{\phi}^r(y|x)]$$

- Two objectives
- Have global optimal state in the game theoretic view
- The objectives: maximize the conditional log-likelihood of  $y$  (or  $1 - y$ ) with the distribution  $q_{\phi}(y|x)$  conditioning on data/generation  $x$  inferred by  $p_{\theta}(x|y)$



- Interpret  $q_{\phi}(y|x)$  as the generative model
- Interpret  $p_{\theta}(x|y)$  as the inference model





# GANs: minimizing KLD

- As in Variational EM, we can further rewrite in the form of **minimizing KLD** to reveal more insights into the optimization problem
- For each optimization step of  $p_{\theta}(\mathbf{x}|y)$  at point  $(\theta = \theta_0, \phi = \phi_0)$ , let
  - $p(y)$ : uniform prior distribution
  - $p_{\theta=\theta_0}(\mathbf{x}) = \mathbb{E}_{p(y)}[p_{\theta=\theta_0}(\mathbf{x}|y)]$
  - $q^r(\mathbf{x}|y) \propto q_{\phi=\phi_0}^r(y|\mathbf{x})p_{\theta=\theta_0}(\mathbf{x})$

- **Lemma 1:** The updates of  $\theta$  at  $\theta_0$  have

$$\nabla_{\theta} \left[ - \mathbb{E}_{p_{\theta}(\mathbf{x}|y)p(y)} [\log q_{\phi=\phi_0}^r(y|\mathbf{x})] \right] \Big|_{\theta=\theta_0} =$$
$$\nabla_{\theta} \left[ \mathbb{E}_{p(y)} [KL(p_{\theta}(\mathbf{x}|y) \| q^r(\mathbf{x}|y))] - JSD(p_{\theta}(\mathbf{x}|y=0) \| p_{\theta}(\mathbf{x}|y=1)) \right] \Big|_{\theta=\theta_0},$$

- KL: KL divergence
- JSD: Jensen-shannon divergence





# GANs: minimizing KLD

- *Lemma 1:* The updates of  $\theta$  at  $\theta_0$  have

$$\nabla_{\theta} \left[ -\mathbb{E}_{p_{\theta}(\mathbf{x}|y)p(y)} [\log q_{\phi=\phi_0}^r(y|\mathbf{x})] \right] \Big|_{\theta=\theta_0} =$$
$$\nabla_{\theta} \left[ \mathbb{E}_{p(y)} [\text{KL}(p_{\theta}(\mathbf{x}|y) \| q^r(\mathbf{x}|y))] - \text{JSD}(p_{\theta}(\mathbf{x}|y=0) \| p_{\theta}(\mathbf{x}|y=1)) \right] \Big|_{\theta=\theta_0}$$

- Connection to variational inference
  - See  $\mathbf{x}$  as latent variables,  $y$  as visible
  - $p_{\theta=\theta_0}(\mathbf{x})$ : prior distribution
  - $q^r(\mathbf{x}|y) \propto q_{\phi=\phi_0}^r(y|\mathbf{x})p_{\theta=\theta_0}(\mathbf{x})$ : posterior distribution
  - $p_{\theta}(\mathbf{x}|y)$ : variational distribution
    - Amortized inference: updates model parameter  $\theta$
- Suggests relations to VAEs, as we will explore shortly

In VEM, we minimize the following:

$$F(\theta, \phi; \mathbf{x}) = -\log p(\mathbf{x}) + \text{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) \| p_{\theta}(\mathbf{z}|\mathbf{x}))$$

⇒ KL(inference model | posterior)

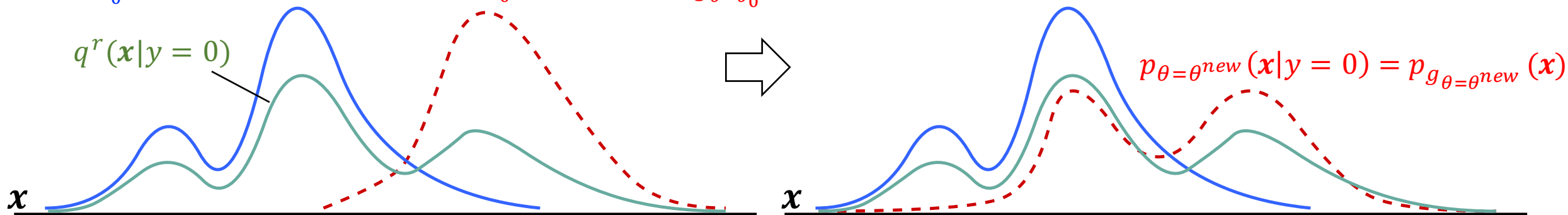




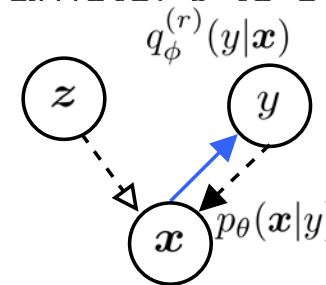
# GANs: minimizing KLD

$$p_{\theta=\theta_0}(\mathbf{x}|y=1) = p_{data}(\mathbf{x}) \quad p_{\theta=\theta_0}(\mathbf{x}|y=0) = p_{g_{\theta=\theta_0}}(\mathbf{x})$$

$$q^r(\mathbf{x}|y=0)$$

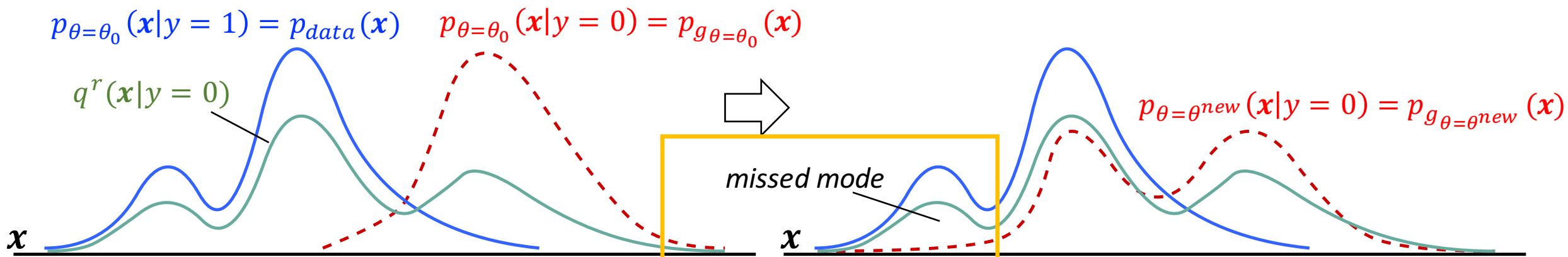


- Minimizing the KLD drives  $p_{g_{\theta}}(\mathbf{x})$  to  $p_{data}(\mathbf{x})$ 
  - By definition:  $p_{\theta=\theta_0}(\mathbf{x}) = E_{p(y)}[p_{\theta=\theta_0}(\mathbf{x}|y)] = (p_{g_{\theta=\theta_0}}(\mathbf{x}) + p_{data}(\mathbf{x})) / 2$
  - $KL(p_{\theta}(\mathbf{x}|y=1) || q^r(\mathbf{x}|y=1)) = KL(p_{data}(\mathbf{x}) || q^r(\mathbf{x}|y=1))$  : constant, no free parameters
  - $KL(p_{\theta}(\mathbf{x}|y=0) || q^r(\mathbf{x}|y=0)) = KL(p_{g_{\theta}}(\mathbf{x}) || q^r(\mathbf{x}|y=0))$  : parameter  $\theta$  to optimize
    - $q^r(\mathbf{x}|y=0) \propto q_{\phi=\phi_0}^r(y=0|\mathbf{x})p_{\theta=\theta_0}(\mathbf{x})$ 
      - seen as a mixture of  $p_{g_{\theta=\theta_0}}(\mathbf{x})$  and  $p_{data}(\mathbf{x})$
      - mixing weights induced from  $q_{\phi=\phi_0}^r(y=0|\mathbf{x})$
    - Drives  $p_{g_{\theta}}(\mathbf{x}|y)$  to mixture of  $p_{g_{\theta=\theta_0}}(\mathbf{x})$  and  $p_{data}(\mathbf{x})$   
 $\Rightarrow$  Drives  $p_{g_{\theta}}(\mathbf{x})$  to  $p_{data}(\mathbf{x})$





# GANs: minimizing KLD



- Missing mode phenomena of GANs
  - Asymmetry of KLD
    - Concentrates  $p_{\theta}(x|y=0)$  to large modes of  $q^r(x|y)$   
 $\Rightarrow p_{g_{\theta}}(x)$  misses modes of  $p_{data}(x)$
  - Symmetry of JSD
    - Does not affect the behavior of mode missing

$$\text{KL}(p_{g_{\theta}}(x) || q^r(x|y=0)) \\ = \int p_{g_{\theta}}(x) \log \frac{p_{g_{\theta}}(x)}{q^r(x|y=0)} dx$$

- Large positive contribution to the KLD in the regions of  $x$  space where  $q^r(x|y=0)$  is small, unless  $p_{g_{\theta}}(x)$  is also small
- $\Rightarrow p_{g_{\theta}}(x)$  tends to avoid regions where  $q^r(x|y=0)$  is small







# Recap: conventional formulation of VAEs

- Objective:

$$\max_{\theta, \eta} \mathcal{L}_{\theta, \eta}^{\text{vae}} = \mathbb{E}_{p_{\text{data}}(\mathbf{x})} \left[ \mathbb{E}_{\tilde{q}_{\eta}(\mathbf{z}|\mathbf{x})} [\log \tilde{p}_{\theta}(\mathbf{x}|\mathbf{z})] - \text{KL}(\tilde{q}_{\eta}(\mathbf{z}|\mathbf{x}) \parallel \tilde{p}(\mathbf{z})) \right]$$

- $\tilde{p}(\mathbf{z})$ : prior over  $\mathbf{z}$
  - $\tilde{p}_{\theta}(\mathbf{x}|\mathbf{z})$ : generative model
  - $\tilde{q}_{\eta}(\mathbf{z}|\mathbf{x})$ : inference model
  - Only uses real examples from  $p_{\text{data}}(\mathbf{x})$ , lacks adversarial mechanism
- To align with GANs, let's introduce the real/fake indicator  $y$  and adversarial discriminator





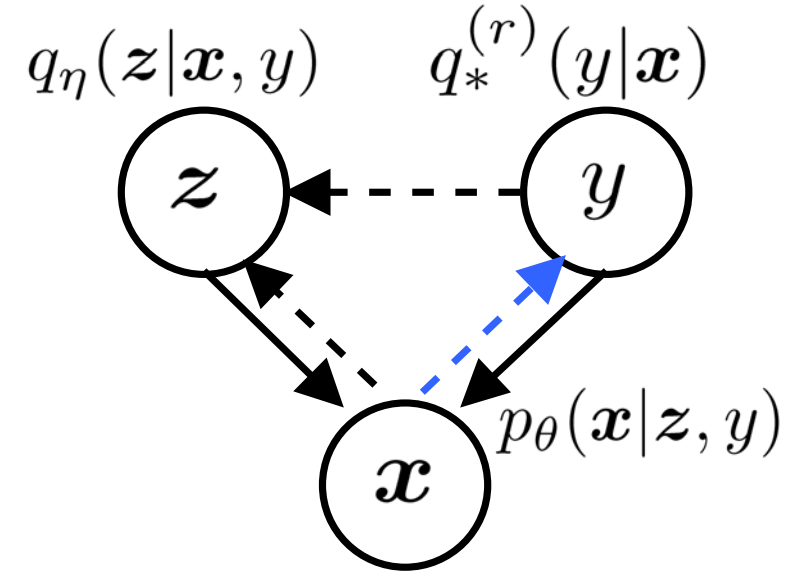
# VAEs: new formulation

- Assume a *perfect* discriminator  $q_*(y|\mathbf{x})$ 
  - $q_*(y = 1|\mathbf{x}) = 1$  if  $\mathbf{x}$  is real examples
  - $q_*(y = 0|\mathbf{x}) = 1$  if  $\mathbf{x}$  is generated samples
  - $q_*^r(y|\mathbf{x}) := q_*(1 - y|\mathbf{x})$
- Generative distribution

$$p_\theta(\mathbf{x}|\mathbf{z}, y) = \begin{cases} p_\theta(\mathbf{x}|\mathbf{z}) & y = 0 \\ p_{data}(\mathbf{x}) & y = 1. \end{cases}$$

- Let  $p_\theta(\mathbf{z}, y|\mathbf{x}) \propto p_\theta(\mathbf{x}|\mathbf{z}, y)p(\mathbf{z}|y)p(y)$
- Lemma 2*

$$\begin{aligned} \mathcal{L}_{\theta, \eta}^{vae} &= 2 \cdot \mathbb{E}_{p_{\theta_0}(\mathbf{x})} \left[ \mathbb{E}_{q_\eta(\mathbf{z}|\mathbf{x}, y) q_*^r(y|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z}, y)] - KL(q_\eta(\mathbf{z}|\mathbf{x}, y) q_*^r(y|\mathbf{x}) \| p(\mathbf{z}|y)p(y)) \right] \\ &= 2 \cdot \mathbb{E}_{p_{\theta_0}(\mathbf{x})} [-KL(q_\eta(\mathbf{z}|\mathbf{x}, y) q_*^r(y|\mathbf{x}) \| p_\theta(\mathbf{z}, y|\mathbf{x}))]. \end{aligned}$$





# GANs vs VAEs side by side

$$p_{\theta}(\mathbf{z}, y|\mathbf{x}) \propto p_{\theta}(\mathbf{x}|\mathbf{z}, y)p(\mathbf{z}|y)p(y)$$

	GANs (InfoGAN)	VAEs
Generative distribution	$p_{\theta}(\mathbf{x} y) = \begin{cases} p_{g_{\theta}}(\mathbf{x}) & y = 0 \\ p_{data}(\mathbf{x}) & y = 1. \end{cases}$	$p_{\theta}(\mathbf{x} \mathbf{z}, y) = \begin{cases} p_{\theta}(\mathbf{x} \mathbf{z}) & y = 0 \\ p_{data}(\mathbf{x}) & y = 1. \end{cases}$
Discriminator distribution	$q_{\phi}(y \mathbf{x})$	$q_{*}(y \mathbf{x}), \text{ perfect, degenerated}$
$\mathbf{z}$ -inference model	$q_{\eta}(\mathbf{z} \mathbf{x}, y) \text{ of InfoGAN}$	$q_{\eta}(\mathbf{z} \mathbf{x}, y)$
KLD to minimize	$\min_{\theta} \text{KL}(p_{\theta}(\mathbf{x} y)    q^r(\mathbf{x} \mathbf{z}, y))$ $\sim \min_{\theta} \text{KL}(P_{\theta}    Q)$	$\min_{\theta} \text{KL}(q_{\eta}(\mathbf{z} \mathbf{x}, y)q_{*}^r(y \mathbf{x})    p_{\theta}(\mathbf{z}, y \mathbf{x}))$ $\sim \min_{\theta} \text{KL}(Q    P_{\theta})$

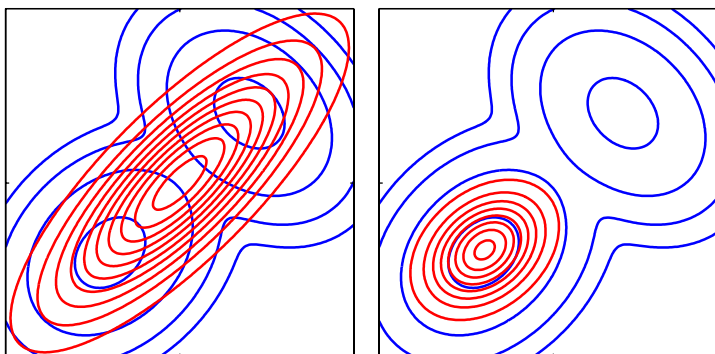




# GANs vs VAEs side by side

	GANs (InfoGAN)	VAEs
KLD to minimize	$\min_{\theta} \text{KL}(p_{\theta}(\mathbf{x} y) \parallel q^r(\mathbf{x} \mathbf{z}, y))$ $\sim \min_{\theta} \text{KL}(P_{\theta} \parallel Q)$	$\min_{\theta} \text{KL}(q_{\eta}(\mathbf{z} \mathbf{x}, y)q_{*}^r(y \mathbf{x}) \parallel p_{\theta}(\mathbf{z}, y \mathbf{x}))$ $\sim \min_{\theta} \text{KL}(Q \parallel P_{\theta})$

- Asymmetry of KLDs inspires combination of GANs and VAEs
  - GANs:  $\min_{\theta} \text{KL}(P_{\theta} \parallel Q)$  tends to missing mode
  - VAEs:  $\min_{\theta} \text{KL}(Q \parallel P_{\theta})$  tends to cover regions with small values of  $p_{data}$



Mode covering

Mode missing





# Link back to wake sleep algorithm

- Denote
  - Latent variables  $\mathbf{h}$
  - Parameters  $\lambda$
- Recap: wake sleep algorithm

$$\text{Wake : } \max_{\theta} \mathbb{E}_{q_{\lambda}(\mathbf{h}|\mathbf{x})p_{data}(\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{h})]$$

$$\text{Sleep : } \max_{\lambda} \mathbb{E}_{p_{\theta}(\mathbf{x}|\mathbf{h})p(\mathbf{h})} [\log q_{\lambda}(\mathbf{h}|\mathbf{x})]$$





# VAEs vs. Wake-sleep

- Wake sleep algorithm

$$\text{Wake : } \max_{\boldsymbol{\theta}} \mathbb{E}_{q_{\lambda}(\mathbf{h}|\mathbf{x})p_{data}(\mathbf{x})} [\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{h})]$$

$$\text{Sleep : } \max_{\lambda} \mathbb{E}_{p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{h})p(\mathbf{h})} [\log q_{\lambda}(\mathbf{h}|\mathbf{x})]$$

- Let  $\mathbf{h}$  be  $\mathbf{z}$ , and  $\lambda$  be  $\boldsymbol{\eta}$

$$\Rightarrow \max_{\boldsymbol{\theta}} \mathbb{E}_{q_{\boldsymbol{\eta}}(\mathbf{z}|\mathbf{x})p_{data}(\mathbf{x})} [\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})], \quad \text{recovers VAE objective of optimizing } \boldsymbol{\theta}$$

- VAEs extend wake phase by also learning the inference model ( $\boldsymbol{\eta}$ )

$$\max_{\boldsymbol{\theta}, \boldsymbol{\eta}} \mathcal{L}_{\boldsymbol{\theta}, \boldsymbol{\eta}}^{\text{vae}} = \mathbb{E}_{q_{\boldsymbol{\eta}}(\mathbf{z}|\mathbf{x})p_{data}(\mathbf{x})} [\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})] - \mathbb{E}_{p_{data}(\mathbf{x})} [\text{KL}(q_{\boldsymbol{\eta}}(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))]$$

- Minimize the KLD in the original variational free energy wrt.  $\boldsymbol{\eta}$
- Stick to minimizing the wake-phase KLD wrt. both  $\boldsymbol{\theta}$  and  $\boldsymbol{\eta}$
- Do not involve sleep-phase objective
- Recall: sleep phase minimizes the *reverse* KLD in the variational free energy





# GANs vs. Wake-sleep

- Wake sleep algorithm

$$\text{Wake : } \max_{\theta} \mathbb{E}_{q_{\lambda}(\mathbf{h}|\mathbf{x})p_{data}(\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{h})]$$

$$\text{Sleep : } \max_{\lambda} \mathbb{E}_{p_{\theta}(\mathbf{x}|\mathbf{h})p(\mathbf{h})} [\log q_{\lambda}(\mathbf{h}|\mathbf{x})]$$

- Let  $\mathbf{h}$  be  $y$ , and  $\lambda$  be  $\phi$

$\Rightarrow \max_{\phi} \mathbb{E}_{p_{\theta}(\mathbf{x}|y)p(y)} [\log q_{\phi}(y|\mathbf{x})]$ , recovers GAN objective of optimizing  $\phi$

- GANs extend sleep phase by also learning the generative model ( $\theta$ )

- Directly extending sleep phase:  $\max_{\phi} \mathcal{L}_{\phi} = \mathbb{E}_{p_{\theta}(\mathbf{x}|y)p(y)} [\log q_{\phi}(y|\mathbf{x})]$
- GANs:  $\max_{\theta} \mathcal{L}_{\theta} = \mathbb{E}_{p_{\theta}(\mathbf{x}|y)p(y)} [\log q_{\phi}^r(y|\mathbf{x})]$
- The only difference is replacing  $q_{\phi}$  with  $q_{\phi}^r$
- This is where adversarial mechanism come about !
- GANs stick to minimizing the sleep-phase KLD
- Do not involve wake-phase objective





# Conclusions

Z Hu, Z YANG, R Salakhutdinov, E Xing,  
“On Unifying Deep Generative Models”, arxiv 1706.00550

- ❑ Deep generative models research have a long history
  - ❑ Deep belief nets / Helmholtz machines / Predictability Minimization / ...
- ❑ Unification of deep generative models
  - ❑ GANs and VAEs are essentially minimizing KLD in opposite directions
    - ❑ Extends two phases of classic wake sleep algorithm, respectively
  - ❑ A general formulation framework useful for
    - ❑ Analyzing broad class of existing DGM and variants: ADA/InfoGAN/Joint-models/...
    - ❑ Inspiring new models and algorithms by borrowing ideas across research fields

