

[70240413 Statistical Machine Learning, Spring, 2017]

# **Deep Generative Models**

**Jun Zhu**

dcszj@mail.tsinghua.edu.cn

<http://bigml.cs.tsinghua.edu.cn/~jun>

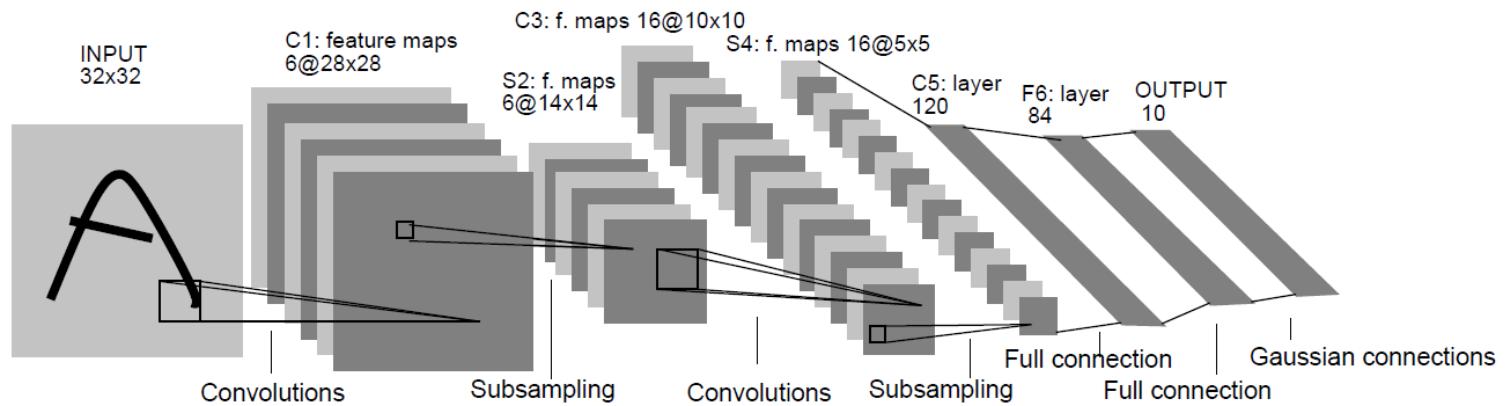
State Key Lab of Intelligent Technology & Systems

Tsinghua University

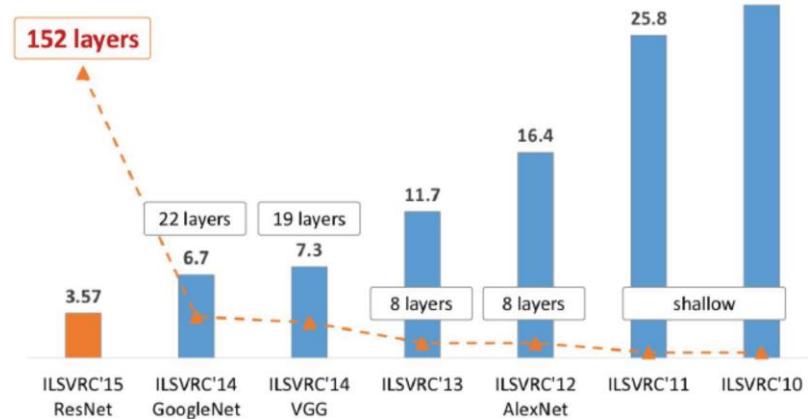
April 18, 2017

# Discriminative Deep Learning

- ◆ Learn a deep NN to map an input to output



- Gradient back-propagation
- Dropout
- Activation functions:
  - rectified linear



# Generative Modeling

- ◆ Have training examples

$$\mathbf{x} \sim p_{\text{data}}(\mathbf{x})$$

- ◆ Want a model that can draw samples:

$$\mathbf{x}' \sim p_{\text{model}}(\mathbf{x})$$

- where  $p_{\text{model}}(\mathbf{x}) \approx p_{\text{data}}(\mathbf{x})$

7	2	1	0	4	1	4	9	5	9
0	6	9	0	1	5	9	7	3	4
9	6	4	5	4	0	7	4	0	1
3	1	3	4	7	2	7	1	2	1
1	7	4	2	3	5	1	2	4	4
6	3	5	5	6	0	4	1	9	5
7	8	9	3	7	4	0	4	3	0
7	0	2	9	1	7	3	2	9	7
7	6	2	7	8	4	7	3	6	1
3	6	9	3	1	4	1	7	6	9



7	2	1	0	4	9	4	9	5	9
0	6	9	0	1	5	9	7	3	4
9	6	4	5	4	0	7	4	0	1
3	9	3	4	7	2	7	1	2	1
1	7	4	2	3	8	1	2	9	4
6	3	5	5	6	0	4	1	9	9
7	8	9	3	7	9	0	4	8	0
7	0	2	9	1	7	3	2	9	7
9	6	2	9	5	4	7	3	6	1
8	6	7	3	1	4	1	7	6	9

$$\mathbf{x} \sim p_{\text{data}}(\mathbf{x})$$

$$\mathbf{x}' \sim p_{\text{model}}(\mathbf{x})$$

# Why generative models?

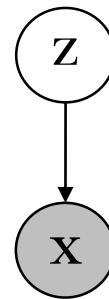
- ◆ Leverage unlabeled datasets, which are often much larger than labeled ones
  - Unsupervised learning
  - Semi-supervised learning
- ◆ Conditional generative models
  - Speech synthesis: Text  $\Rightarrow$  Speech
  - Machine Translation: French  $\Rightarrow$  English
  - Image captioning: Image  $\Rightarrow$  Text
- ◆ Many recent advances on VAE, GAN, etc

# Probabilistic Generative Models

- ◆ **Assumption:** data is described by some factors, which are often hidden

$$z \sim p(z|\alpha)$$

$$x \sim p(x|z, \beta)$$

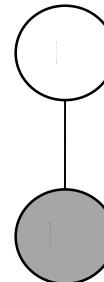


- **Inference:** infer the posterior distribution

$$p(z|x; \theta) \propto p(z|\alpha)p(x|z, \beta)$$

- ◆ Undirected random fields

$$z, x \sim p(z, x|\theta)$$

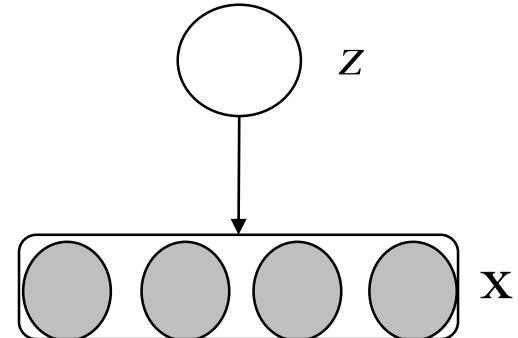


# Examples

- ◆ Latent-class mixture models

$$p(z|\alpha) = \text{Mult}(\alpha)$$

$$p(x|z, \beta) = N(\mu_z, \Sigma_z)$$

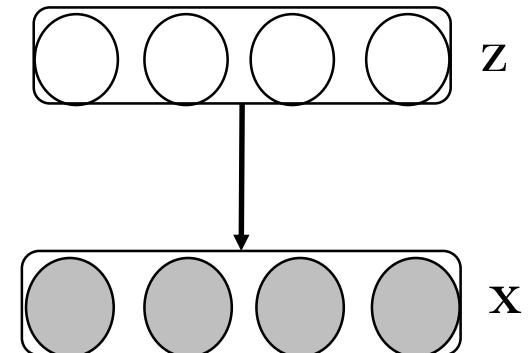


- ◆ Latent-feature factor analysis models

$$p(z|\alpha) = N(\mu_0, \Sigma_0)$$

$$p(x|z, \beta) = N(Wz, \Sigma)$$

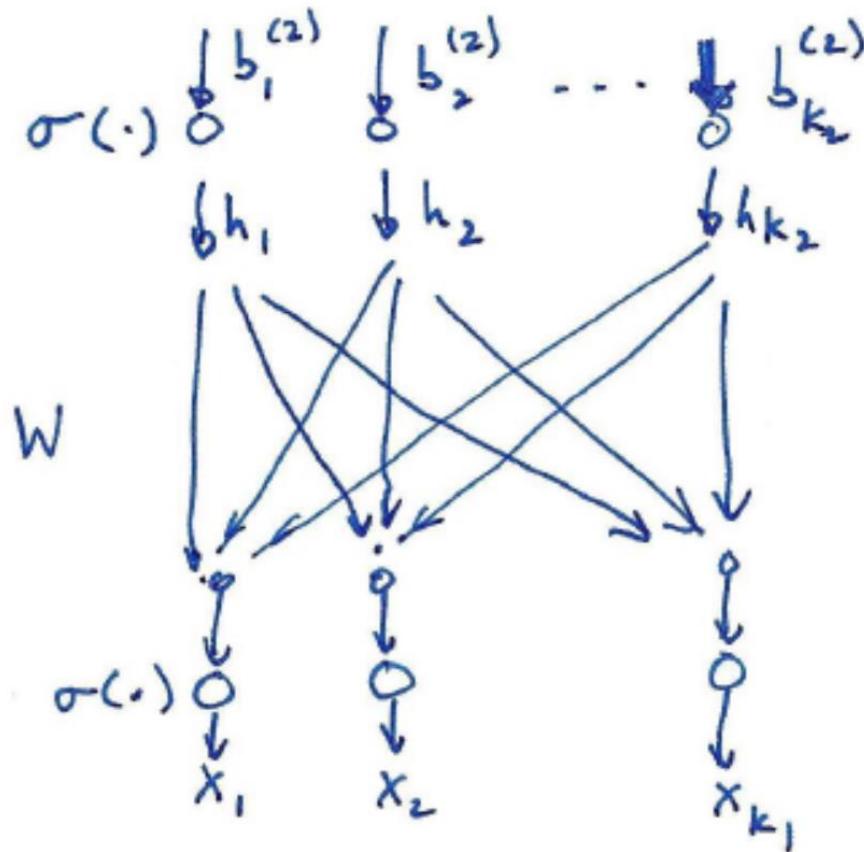
- $W$  is a loading matrix



- ◆ Building blocks for deep models

# Examples

- ◆ Sigmoid belief network



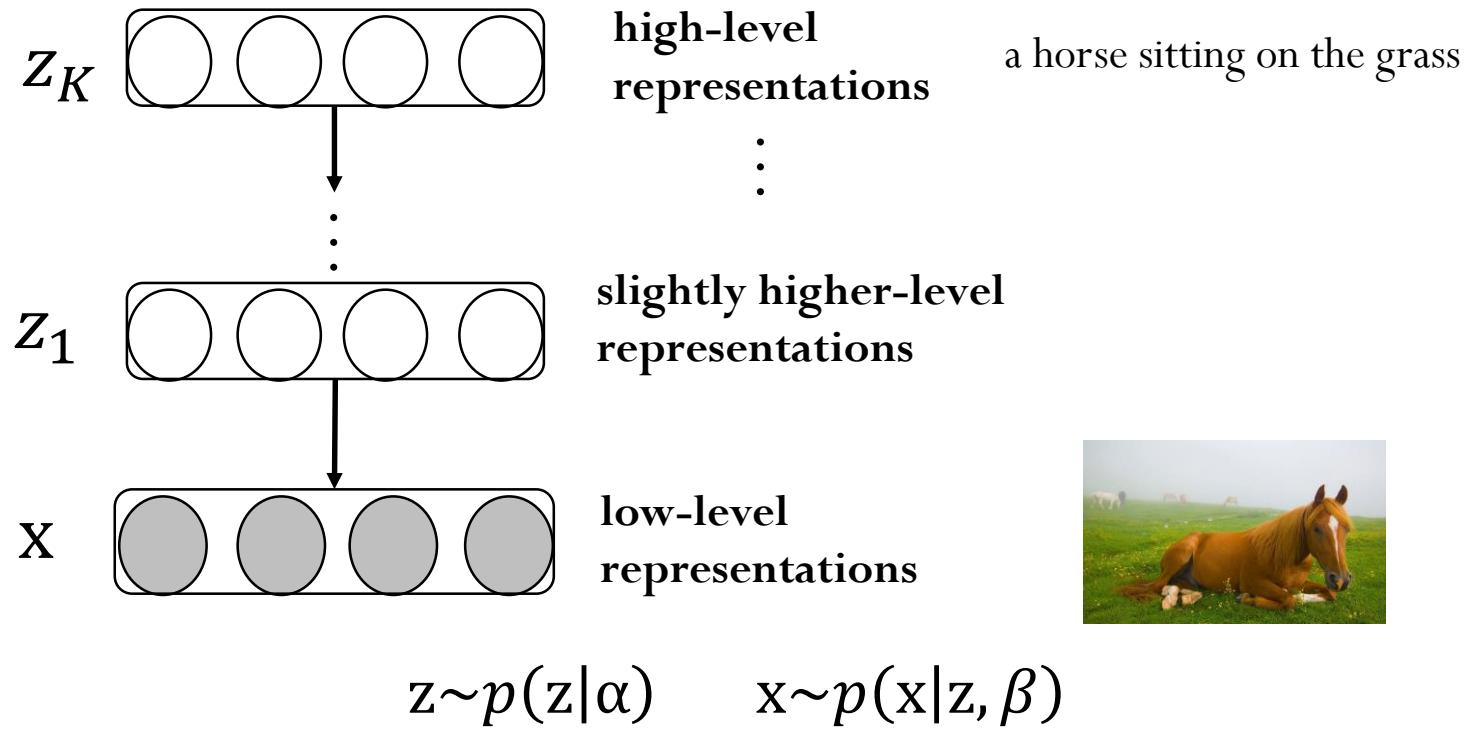
$$h_k \sim \text{Bern}(\sigma(b_k^{(2)})), \quad k = 1, \dots, K_2$$

$$z = \mathbf{W}h + b^{(1)}$$

$$x_j \sim \text{Bern}(\sigma(z_j)), \quad j = 1, \dots, K_1$$

# Deep Generative Models

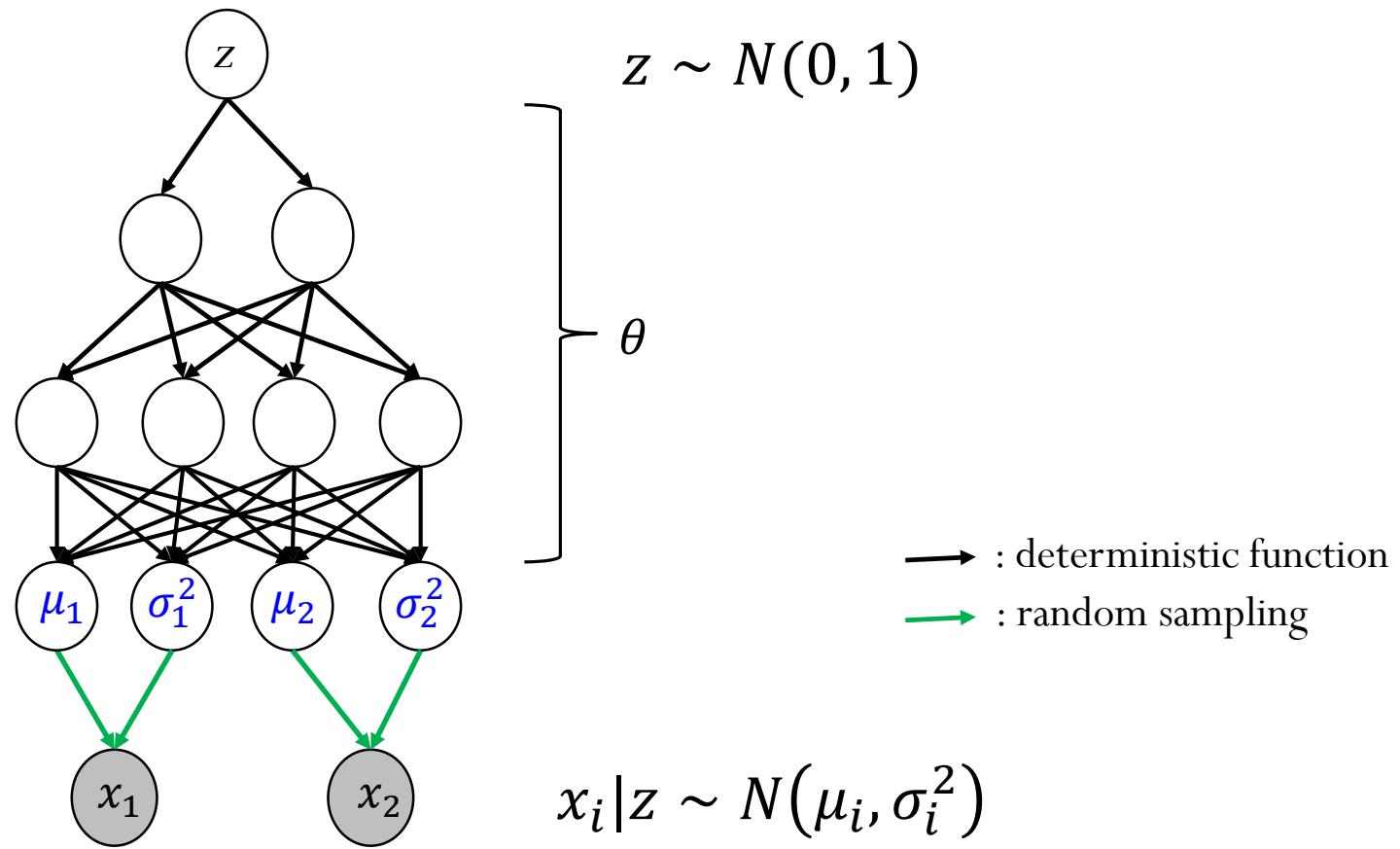
- ◆ Multi-layer *latent-feature* representations with nonlinear transformations



- ◆ Many variants by combining different building blocks

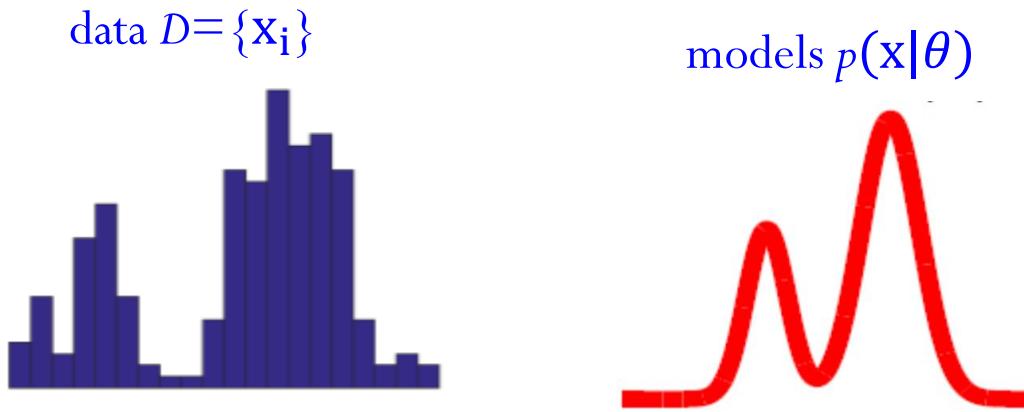
# An example with MLP

- ◆ 1D latent variable  $z$ ; 2D observation  $x$
- ◆ **Idea:** NN + Gaussian (or Bernoulli) with a diagonal covariance



# Learning Deep Generative Models

- ◆ Given a set  $D$  of unlabeled samples, learn the unknown parameters (or a distribution)



- ◆ Find a model that minimizes

$$\mathbb{D}(\text{data } \{x_i\}_{i=1}^n, \text{model } p)$$

# Learning Deep Generative Models

- ◆ Maximum likelihood estimation (MLE):

$$\hat{\theta} = \operatorname{argmax} p(D|\theta)$$

- ◆ Moment-matching:

- draw samples from  $p$ :  $\hat{D} = \{y_i\}_{i=1}^M$ , where  $y_i \sim p(x|\theta)$

- Kernel MMD: 
$$\mathcal{L}_{MMD^2} = \left\| \frac{1}{N} \sum_{i=1}^N \phi(x_i) - \frac{1}{M} \sum_{j=1}^M \phi(y_j) \right\|_{\mathcal{H}}^2$$

- rich enough to distinguish any two distributions in certain RKHS

- ◆ Minimax objective (e.g., GAN)

- ◆ Bayesian inference:

- infer the posterior distribution of parameters

$$p(\theta|D) \propto p_0(\theta)p(D|\theta)$$

# Variational Bayes

- ◆ Consider the log-likelihood of *a single example*

$$\log p(\mathbf{x}; \theta) = \log \int p(\mathbf{z}, \mathbf{x}; \theta) d\mathbf{z}$$

- ◆ Log-integral/sum is annoying to handle directly
- ◆ Derive a variational lower bound  $L(\theta, \phi, \mathbf{x})$

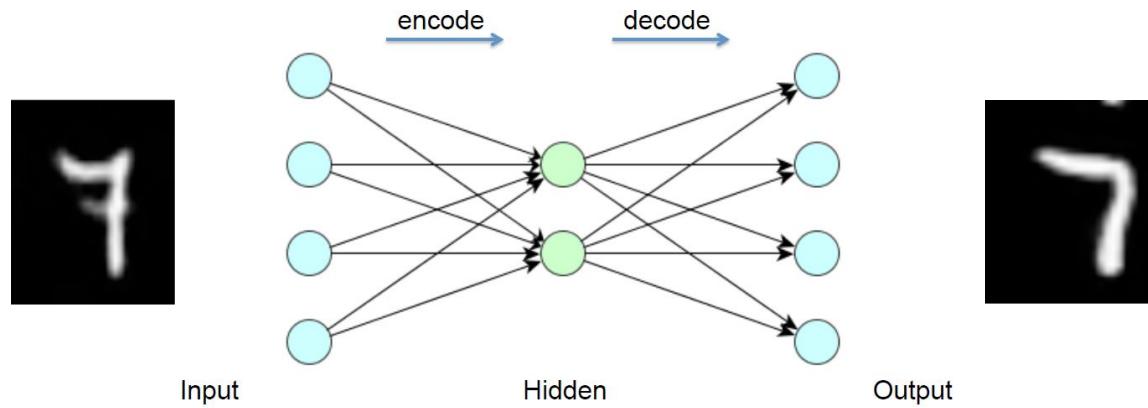
$$\log p(\mathbf{x}; \theta) = L(\theta, \phi, \mathbf{x}) + \text{KL}(q(\mathbf{z}|\mathbf{x}; \phi) \| p(\mathbf{z}|\mathbf{x}; \theta))$$

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

reconstruction term

prior regularization

# Recap: Auto-Encoder

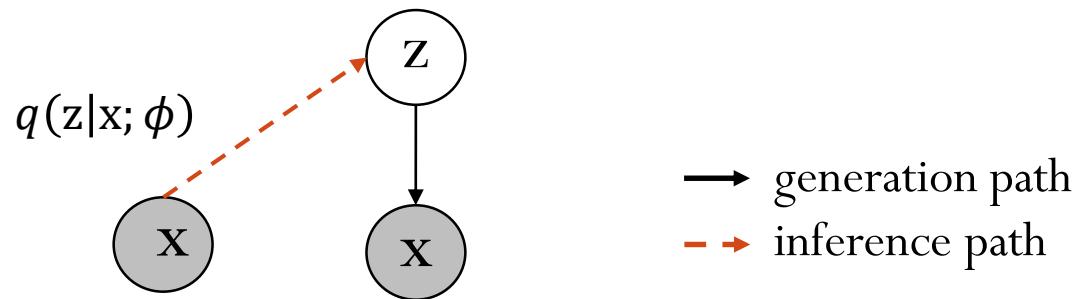


- ◆ Encoder:  $\mathbf{h} = s(W\mathbf{x} + b)$
- ◆ Decoder:  $\mathbf{x}' = s(W'\mathbf{h} + b')$
- ◆ Training: minimize the reconstruction error (e.g., square loss, cross-entropy loss)
- ◆ Denoising AE: randomly corrupted inputs are restored to learn more robust features

# Auto-Encoding Variational Bayes (AEVB)

- ◆ What's unique in AEVB is that *the variational distribution is parameterized by a deep neural network*

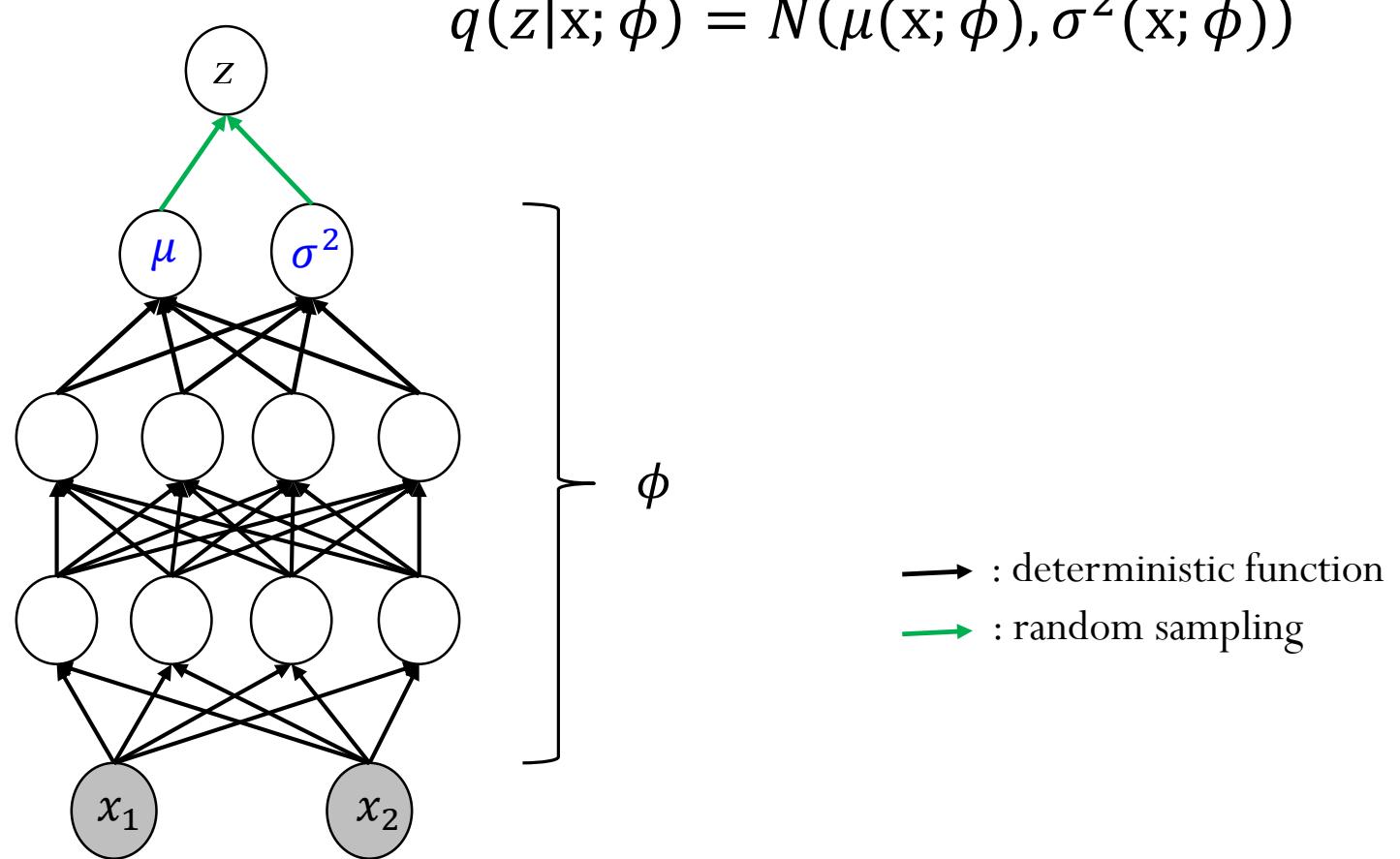
$$q(z|x; \phi) \approx p(z|x; \theta)$$



- We call it an **inference (recognition, encoder) network** or a **Q-network**
- All the parameters are learned jointly via SGD with variance reduction

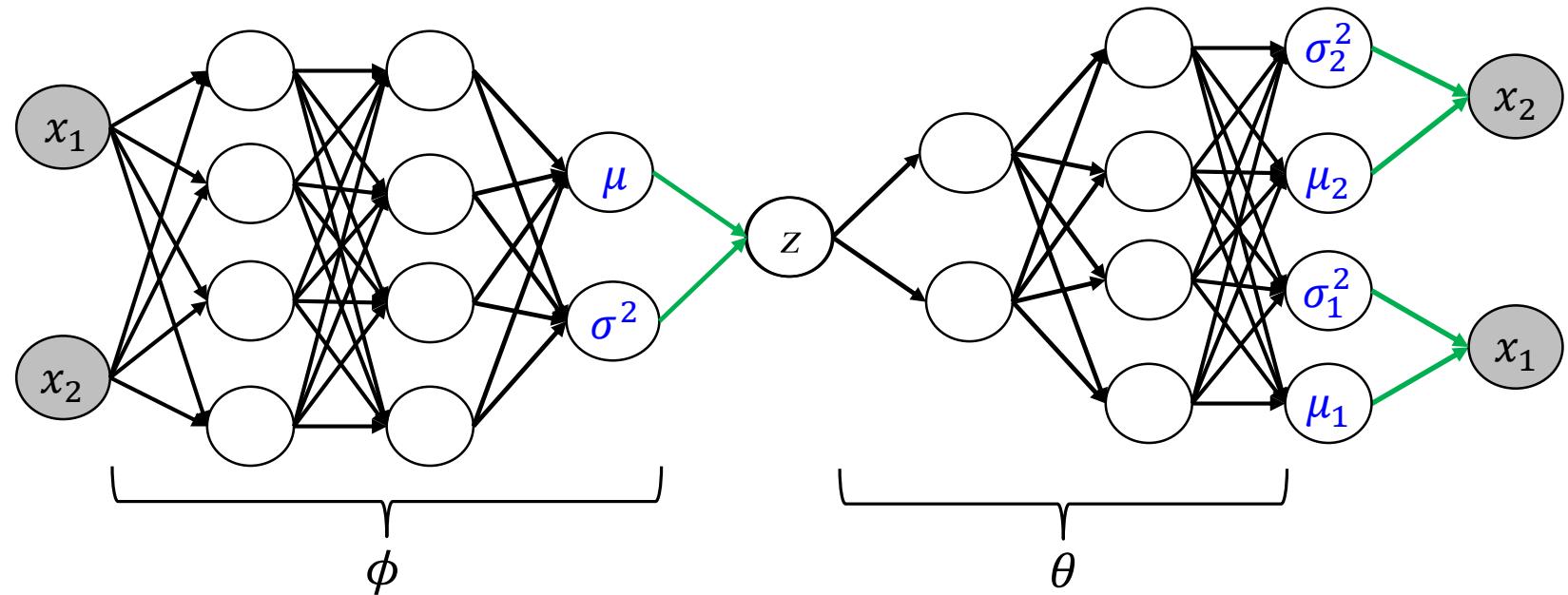
# The Encoder Network

- ◆ A feedforward NN + Gaussian



# The Complete Auto-encoder

- ◆ The Q-P network architecture:



→ : deterministic function  
→ : random sampling

# Stochastic Variational Inference

- ◆ Variational lower-bound for *a single example*

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

- ◆ Variational lower-bound for *a set of examples*

$$L(\theta, \phi, D) = \sum_i \mathbf{E}_{q(\mathbf{z}_i|\mathbf{x}_i;\phi)}[\log p(\mathbf{x}_i|\mathbf{z}_i;\theta)] - \text{KL}(q(\mathbf{z}_i|\mathbf{x}_i;\phi) \| p(\mathbf{z}_i;\theta))$$

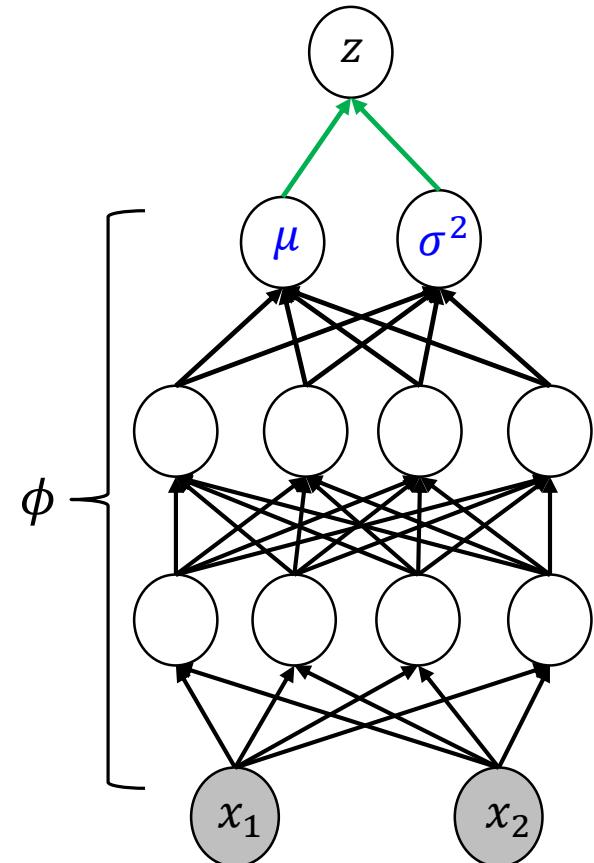
- Use stochastic gradient methods to handle large datasets
- Random mini-batch
  - for each  $i$ , infer the posterior  $q(\mathbf{z}_i|\mathbf{x}_i;\phi)$  ; As we parameterize as a neural network, this in fact optimizes  $\phi$
- ◆ *However, calculating the expectation and its gradients is non-trivial, often intractable*

# Example with Gaussian Distributions

- ◆ Use  $N(0, 1)$  as prior for  $z$ ;  $q(z|x; \phi)$  is Gaussian with parameters  $(\mu(x; \phi), \sigma^2(x; \phi))$  determined by NN
  - The KL-divergence

$$-\text{KL}(q(z|x; \phi) \| p(z; \theta)) = \frac{1}{2} (1 + \log \sigma^2 - \mu^2 - \sigma^2)$$

- **Exercise:** finish the derivation



# Example with Gaussian Distributions

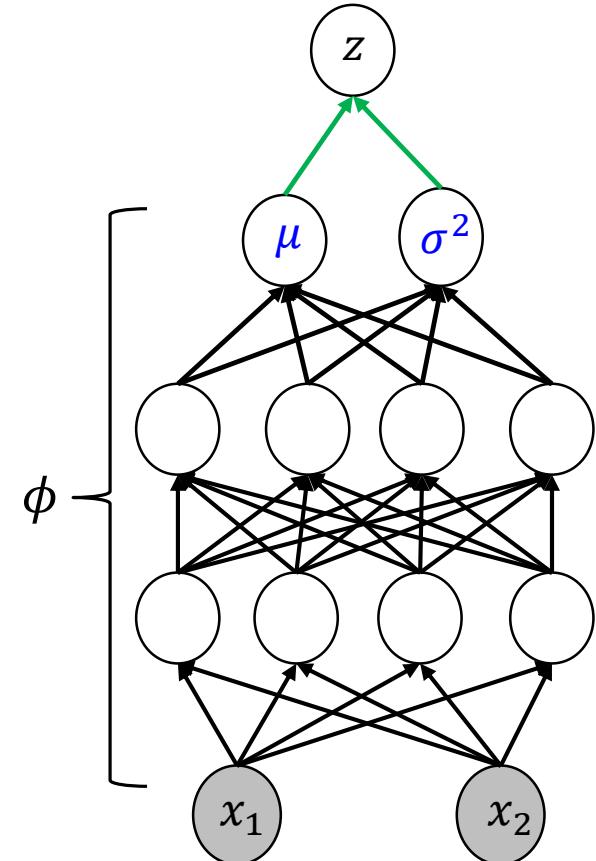
- ◆ Use  $N(0, 1)$  as prior for  $z$ ;  $q(z|x; \phi)$  is Gaussian with parameters  $(\mu(x; \phi), \sigma^2(x; \phi))$  determined by NN
  - The expected log-likelihood

$$\mathbf{E}_{q(z|x;\phi)} [\log p(x|z; \theta)]$$

- If the likelihood is Gaussian

$$-\log p(x_i|z_i) = \sum_j \frac{1}{2} \log \sigma_j^2 + \frac{(x_{ij} - \mu_{xi})^2}{2\sigma_j^2}$$

- *The expectation is still hard to compute because of nonlinearity functions*



# Example with Gaussian Distributions

- ◆ Use  $N(0, 1)$  as prior for  $z$ ;  $q(z|x; \phi)$  is Gaussian with parameters  $(\mu(x; \phi), \sigma^2(x; \phi))$  determined by NN

- The expected log-likelihood

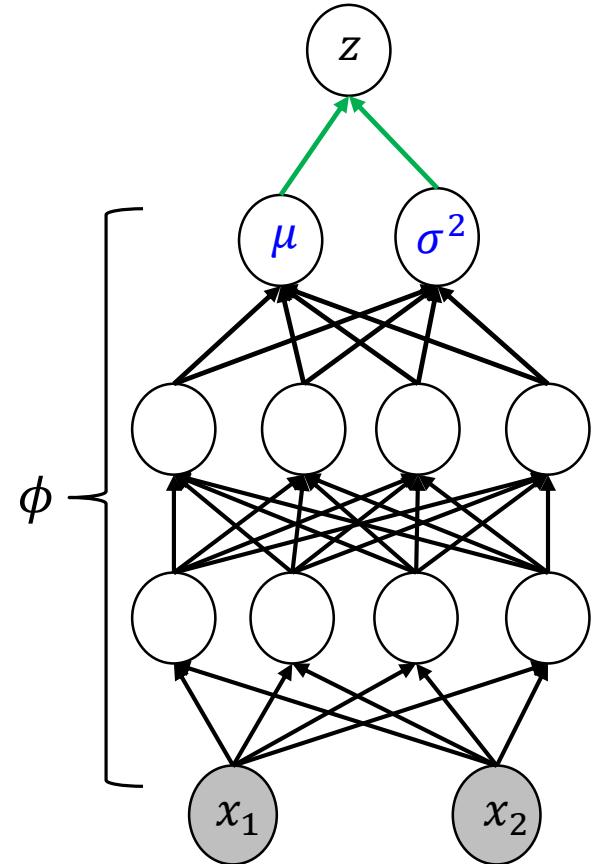
$$\mathbf{E}_{q(z|x;\phi)} [\log p(x|z; \theta)]$$

- Approximate via Monte Carlo methods

$$\mathbf{E}_{q(z|x;\phi)} [\log p(x|z; \theta)] \approx \frac{1}{L} \sum_k \log p(x|z^{(k)})$$

$$z^{(k)} \sim q(z|x; \phi)$$

- An unbiased estimator



# Example with Gaussian Distributions

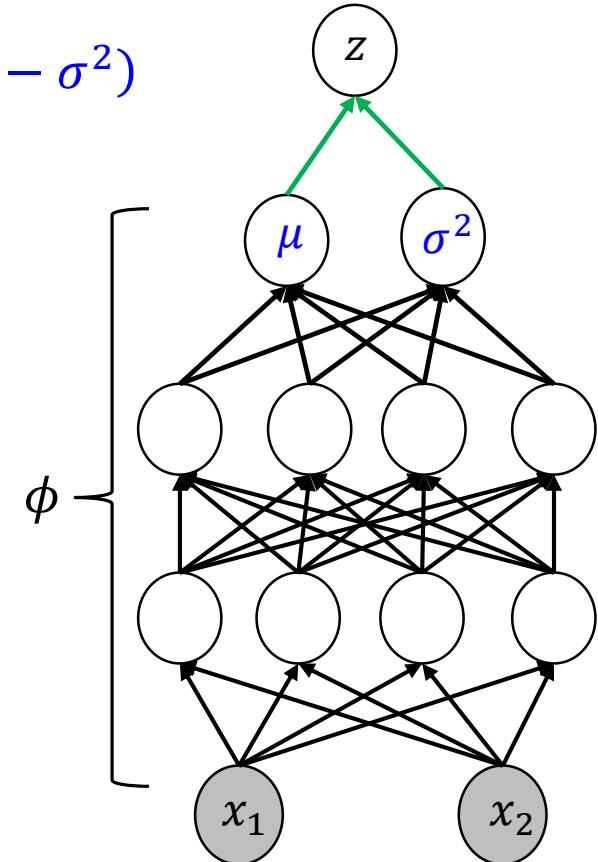
- ◆ The KL-regularization term (closed-form):

$$-\text{KL}(q(z|x; \phi) \| p(z; \theta)) = \frac{1}{2} (1 + \log \sigma^2 - \mu^2 - \sigma^2)$$

- Easy to calculate gradient
- ◆ The expected log-likelihood term (MC estimate)

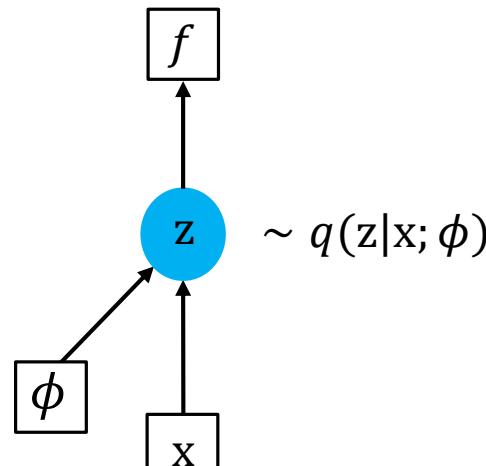
$$\mathbf{E}_{q(z|x;\phi)}[\log p(x|z; \theta)] \approx \frac{1}{L} \sum_k \log p(x|z^{(k)})$$
$$z^{(k)} \sim q(z|x; \phi)$$

- Gradient needs back-propagation!
- *However,  $Z^{(k)}$  is a random variable, we can't take gradient over a randomly drawn number*



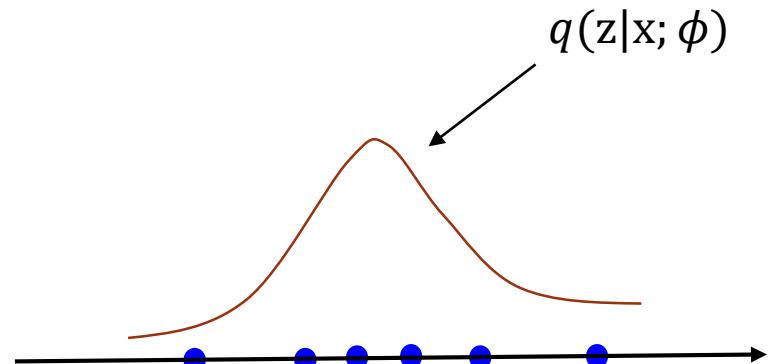
# Reparameterization Trick

- ◆ Backpropagation not possible through random sampling



◻ : deterministic node

● : random node



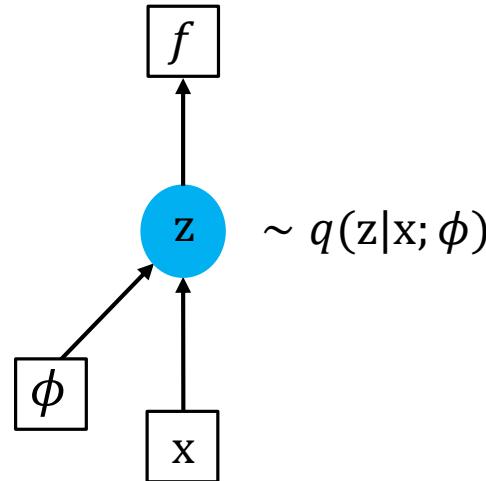
$$z^{(k)} \sim N(\mu(x, \phi), \sigma^2(x, \phi))$$

$$\{-1.5, -0.5, 0.3, 0.6, 1.5, \dots\}$$

Cannot back-propagate through a randomly drawn number

# Reparameterization Trick

- ◆ Backpropagation not possible through random sampling

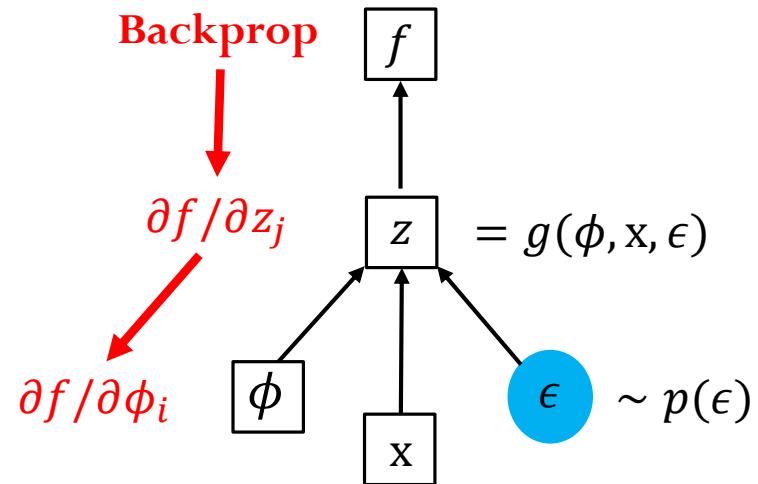


◻ : deterministic node

● : random node

$$z^{(k)} \sim N(\mu(x, \phi), \sigma^2(x, \phi))$$

Cannot back-propagate through a randomly drawn number



$$\epsilon^{(k)} \sim N(0,1)$$

$$z^{(k)} = \mu(x, \phi) + \sigma(x, \phi) \cdot \epsilon^{(k)}$$

$Z$  has the same distribution, but now can back-prop  
Separate into a deterministic part and noise

# The General Form

- ◆ The VAE bound

$$\begin{aligned} L(\theta, \phi, \mathbf{x}) &= \mathbf{E}_{q(z|x;\phi)} [\log p(z, \mathbf{x}; \theta) - \log q(z|x; \phi)] \\ &= \mathbf{E}_{q(z|x;\phi)} \left[ \log \frac{p(z, \mathbf{x}; \theta)}{q(z|x; \phi)} \right] \end{aligned}$$

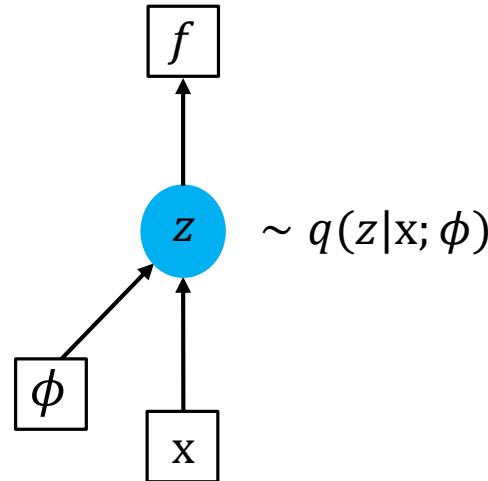
- ◆ Monte Carlo estimate:

$$\begin{aligned} L(\theta, \phi, \mathbf{x}) &\approx \frac{1}{L} \sum_k \log \frac{p(z^{(k)}, \mathbf{x}; \theta)}{q(z^{(k)}|x; \phi)} \\ z^{(k)} &\sim q(z|x; \phi) \end{aligned}$$

- Again, we cannot back-prop through the randomly drawn numbers

# Reparameterization Trick

- ◆ Backpropagation not possible through random sampling



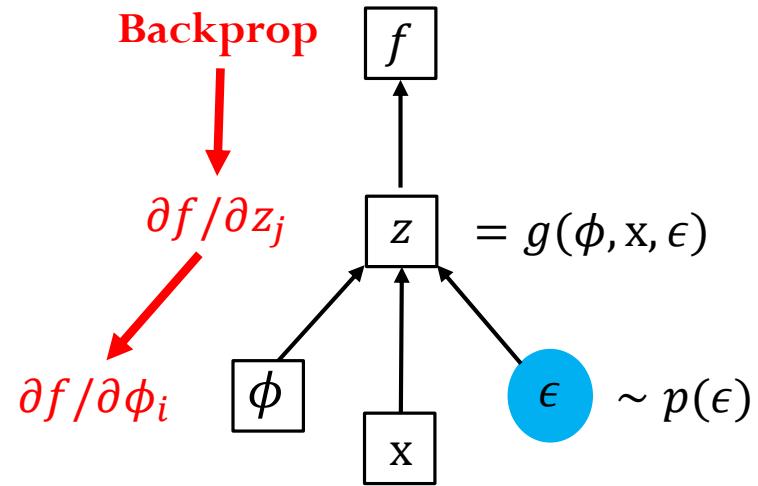
$$\sim q(z|x; \phi)$$

◻ : deterministic node

● : random node

$$z^{(k)} \sim q(z|x; \phi)$$

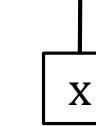
Cannot back-propagate through a randomly drawn number



$$\text{Backprop}$$

$$\partial f / \partial z_j$$

$$\partial f / \partial \phi_i$$



$$= g(\phi, x, \epsilon)$$

$$\epsilon^{(k)} \sim p(\epsilon)$$

$$z^{(k)} = g(\phi, x, \epsilon^{(k)})$$

$Z$  has the same distribution, but now can back-prop  
Separate into a deterministic part and noise

# Reparam-Trick Summary

## ◆ The VAE bound

$$L(\theta, \phi, x) = \mathbf{E}_{q(z|x; \phi)} \left[ \log \frac{p(z, x; \theta)}{q(z|x; \phi)} \right]$$

□ Reparameterized as

$$L(\theta, \phi, x) = \mathbf{E}_{p(\epsilon)} \left[ \log \frac{p(g(x, \epsilon, \phi), x; \theta)}{q(g(x, \epsilon, \phi)|x; \phi)} \right]$$

- where  $\epsilon$  is a simple distribution (e.g., standard normal) and  $g$  is a deep NN

## ◆ The gradients are

$$\nabla_{\theta} L(\theta, \phi, x) = \mathbf{E}_{p(\epsilon)} \left[ \nabla_{\theta} \log \frac{p(g(x, \epsilon, \phi), x; \theta)}{q(g(x, \epsilon, \phi)|x; \phi)} \right]$$

- Back-prop is applied over the deep NN
- Similar for  $\phi$

# Importance Weighted Auto-Encoder (IWAE)

- ◆ The VAE lower bound of log-likelihood

$$L(\theta, \phi, \mathbf{x}) = \mathbf{E}_{q(z|x;\phi)} \left[ \log \frac{p(z, \mathbf{x}; \theta)}{q(z|\mathbf{x}; \phi)} \right]$$

- ◆ A better variational lower bound (IWAE)

$$L_K(\theta, \phi, \mathbf{x}) = \mathbf{E}_{q(z|x;\phi)} \left[ \log \left( \frac{1}{K} \sum_{k=1:K} \frac{p(z^{(k)}, \mathbf{x}; \theta)}{q(z^{(k)}|\mathbf{x}; \phi)} \right) \right]$$

where  $z^{(k)} \sim q(z|\mathbf{x}; \phi)$

- This is a lower-bound of the log-likelihood
- When  $K=1$ , recovers the VAE bound
- When  $K = \infty$ , recovers the log-likelihood
- A monotonic sequence:

$$L_K(\theta, \phi, \mathbf{x}) \leq L_{K+1}(\theta, \phi, \mathbf{x}), \quad \forall \theta, \phi, \mathbf{x}$$

**Homework (\*):**  
prove these properties

# Reparametrization Trick

- ◆ The IWAE bound:

$$L_K(\theta, \phi, \mathbf{x}) = \mathbf{E}_{q(z|x;\phi)} \left[ \log \left( \frac{1}{K} \sum_{k=1:K} w(z^{(k)}, \mathbf{x}; \theta) \right) \right]$$

$$\text{where } z^{(k)} \sim q(z|x; \phi) \quad w(z^{(k)}, \mathbf{x}; \theta, \phi) = \frac{p(z^{(k)}, \mathbf{x}; \theta)}{q(z^{(k)}|x; \phi)}$$

- ◆ Reparameterization form:

$$L_K(\theta, \phi, \mathbf{x}) = \mathbf{E}_{p(\epsilon)} \left[ \log \left( \frac{1}{K} \sum_{k=1:K} w(g(\epsilon^{(k)}, \mathbf{x}, \phi), \mathbf{x}; \theta) \right) \right]$$

$$\text{where } \epsilon^{(k)} \sim p(\epsilon)$$

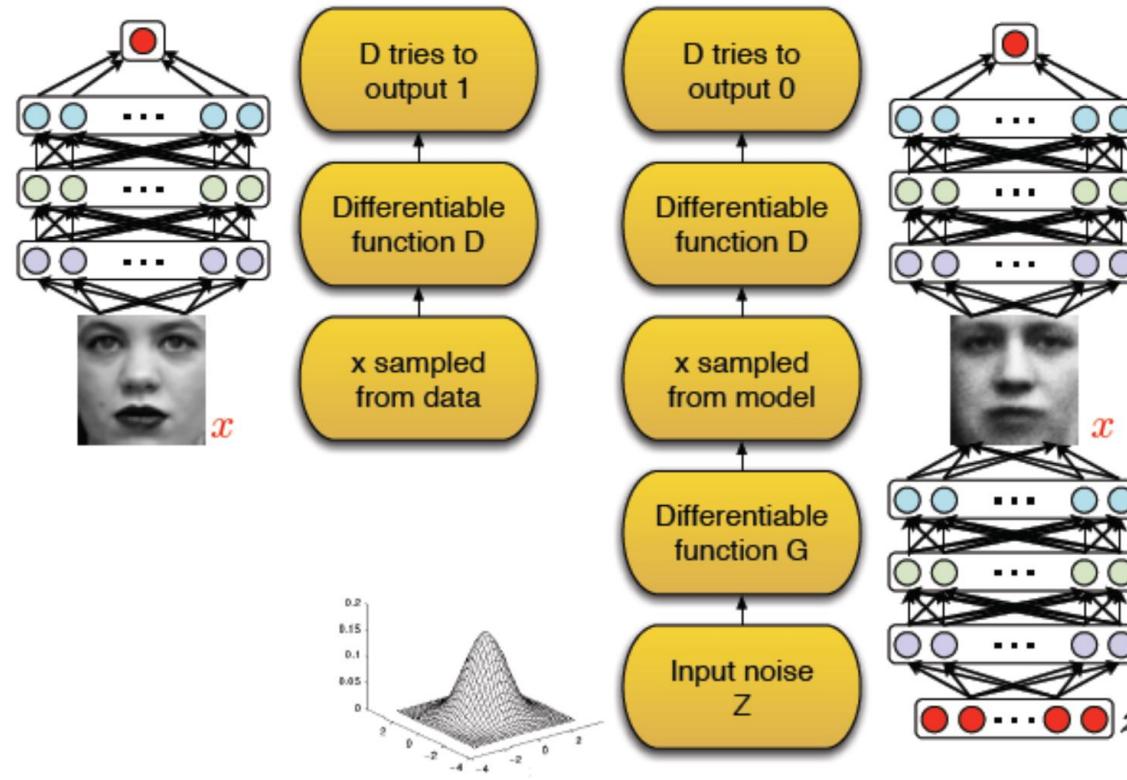
- The gradient can be calculated as in VAE

# Generative Adversarial Networks (GAN)

- ◆ A game between two players:
  - A discriminator  $D$
  - A generator  $G$
- ◆  $D$  tries to discriminate between:
  - A sample from the data distribution.
  - And a sample from the generator  $G$ .
- ◆  $G$  tries to “trick”  $D$  by generating samples that are hard for  $D$  to distinguish from data.

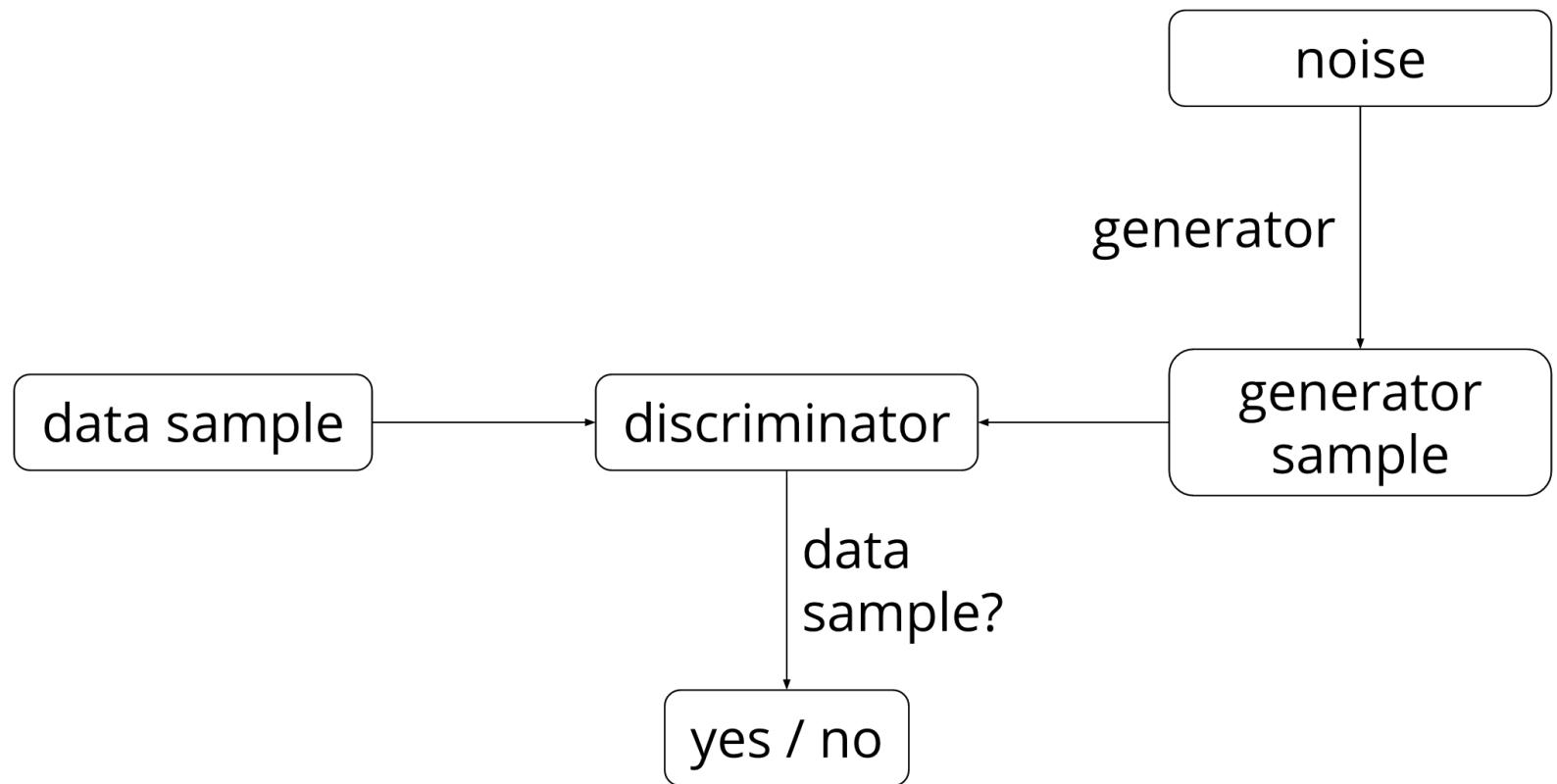
# GAN – architecture

- ◆ The discriminator-network  $D$  is a binary classifier
  - It aims to assign the correct label to both training samples and the samples from  $G$



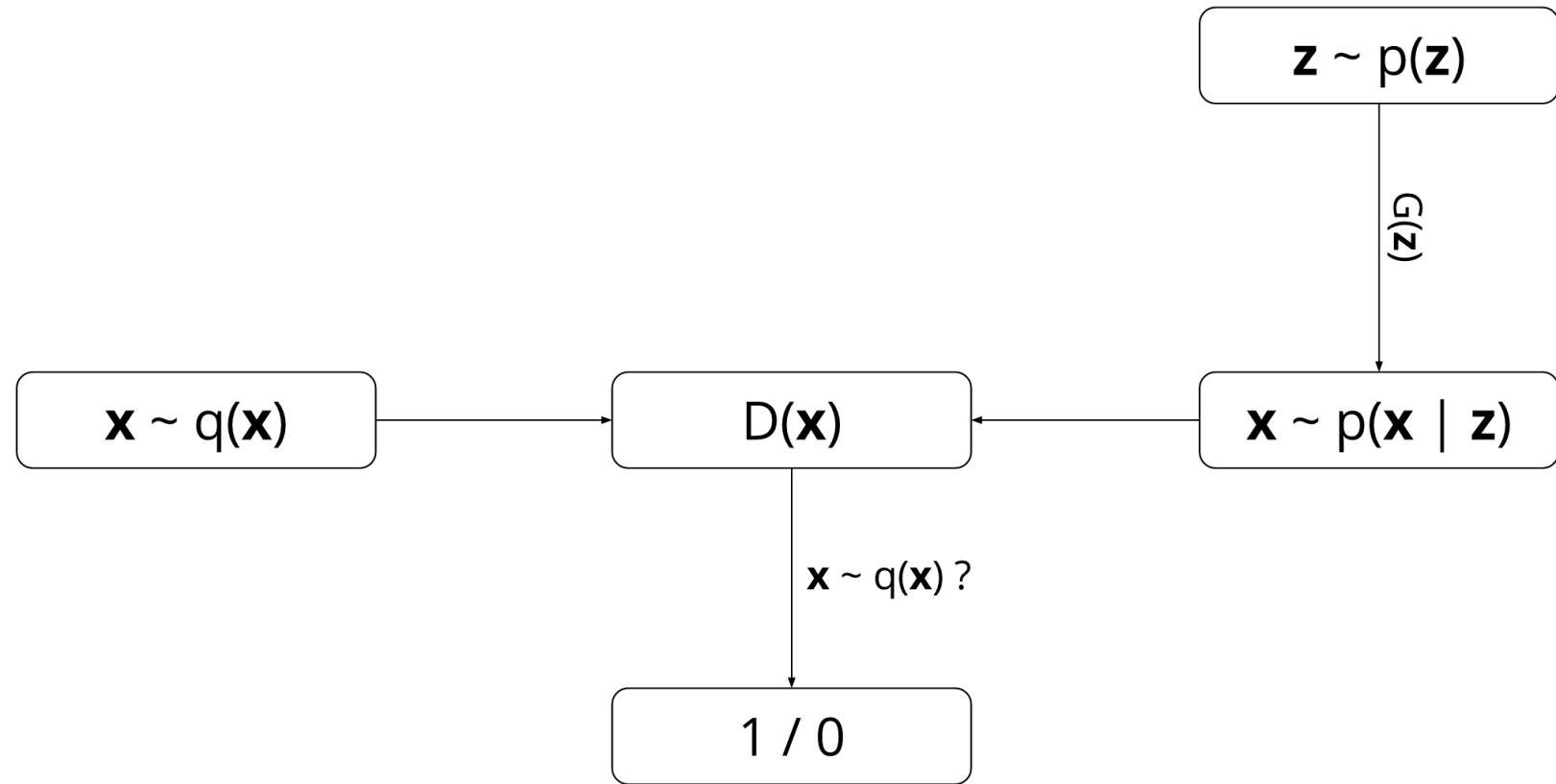
[figure courtesy: Ian Goodfellow]

# GAN – architecture



# GAN – architecture

- ◆ Probabilistic view



# GAN – objective

- ◆ Minimax objective function

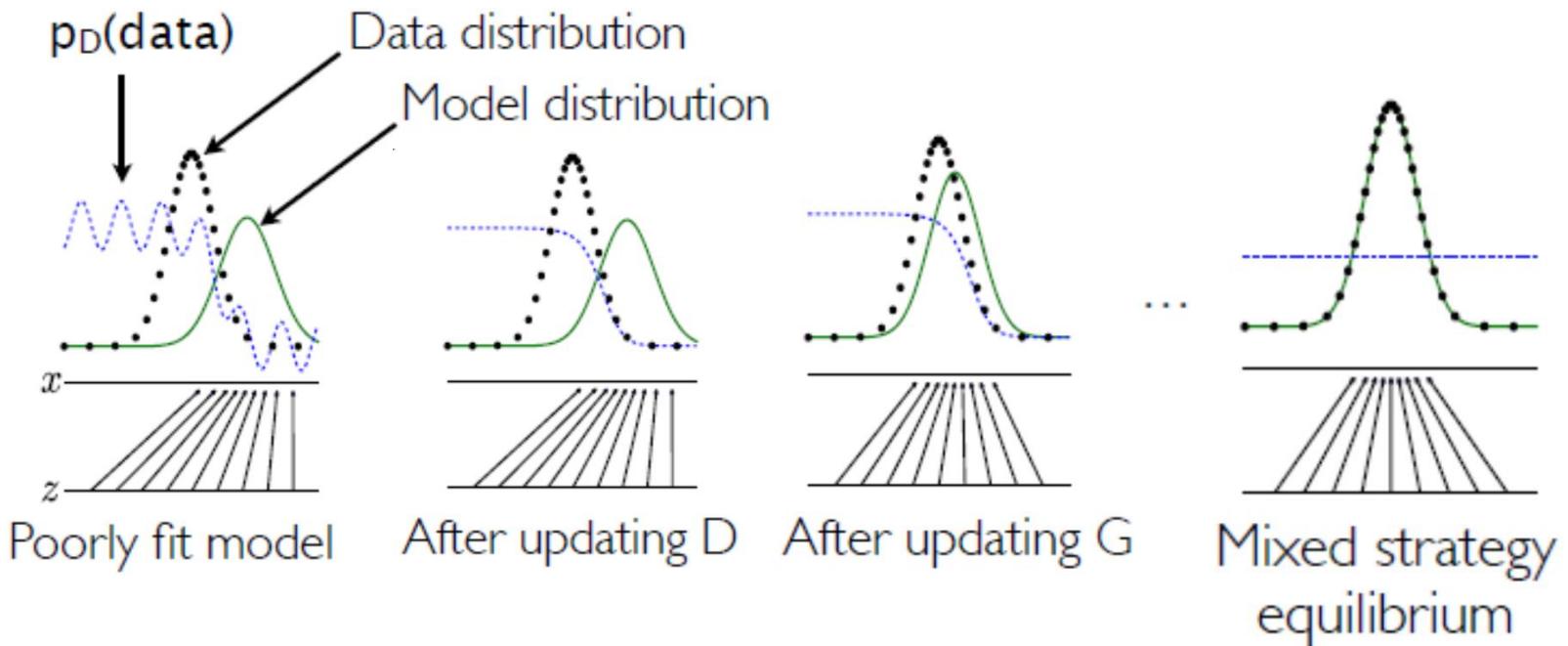
$$\min_G \max_D \mathbf{E}_{p_{\text{data}}(x)}[\log D(x)] + \mathbf{E}_{p(z)} [\log (1 - D(G(z)))]$$

- Optimal strategy of the discriminator for any  $p_{\text{model}}(x)$  is

$$D(x) = \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_{\text{model}}(x)}$$

- Assume infinite data, infinite model capacity, direct updating generator's distribution
  - Unique global optimum
  - Optimum corresponds to data distribution

# GAN – learning process

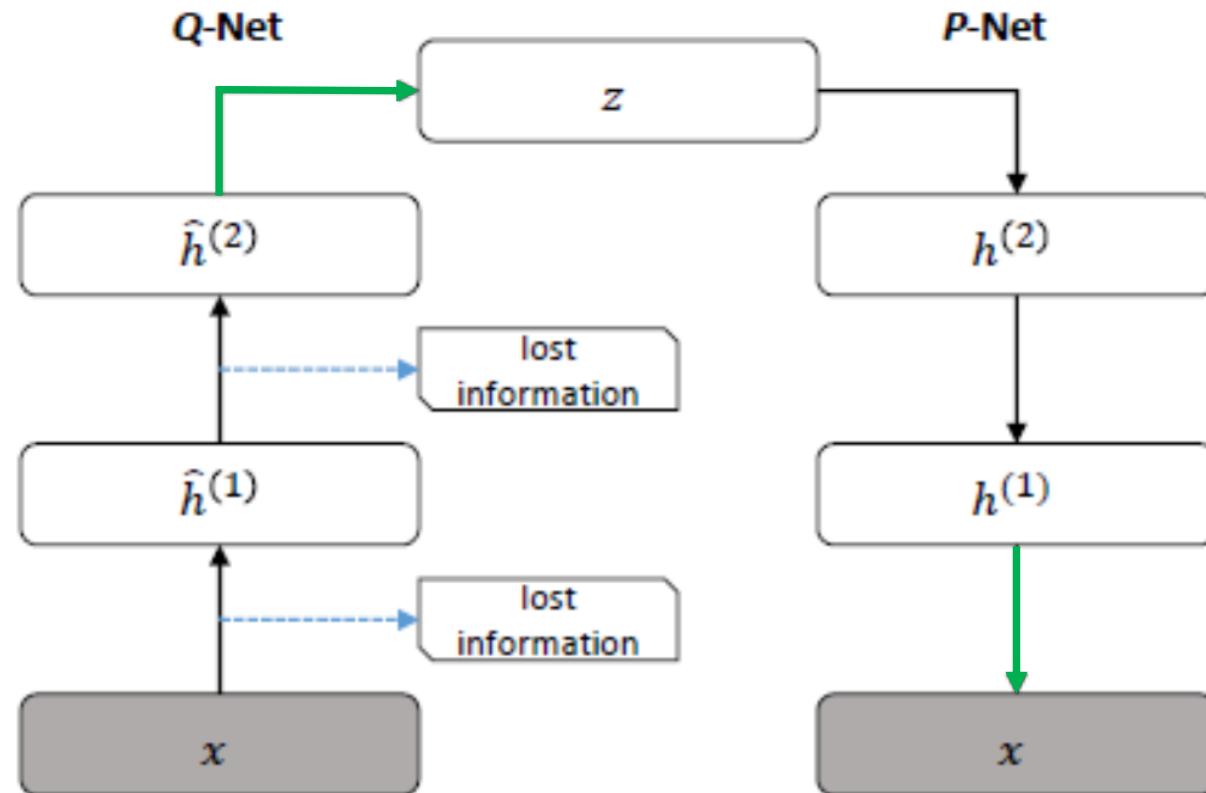


- An adversarial pair near convergence:  $D$  is partially accurate
- Update discriminator distribution  $D$  according to the optimal strategy
- Update generator distribution to be more similar to data distribution
- Iterate until convergence

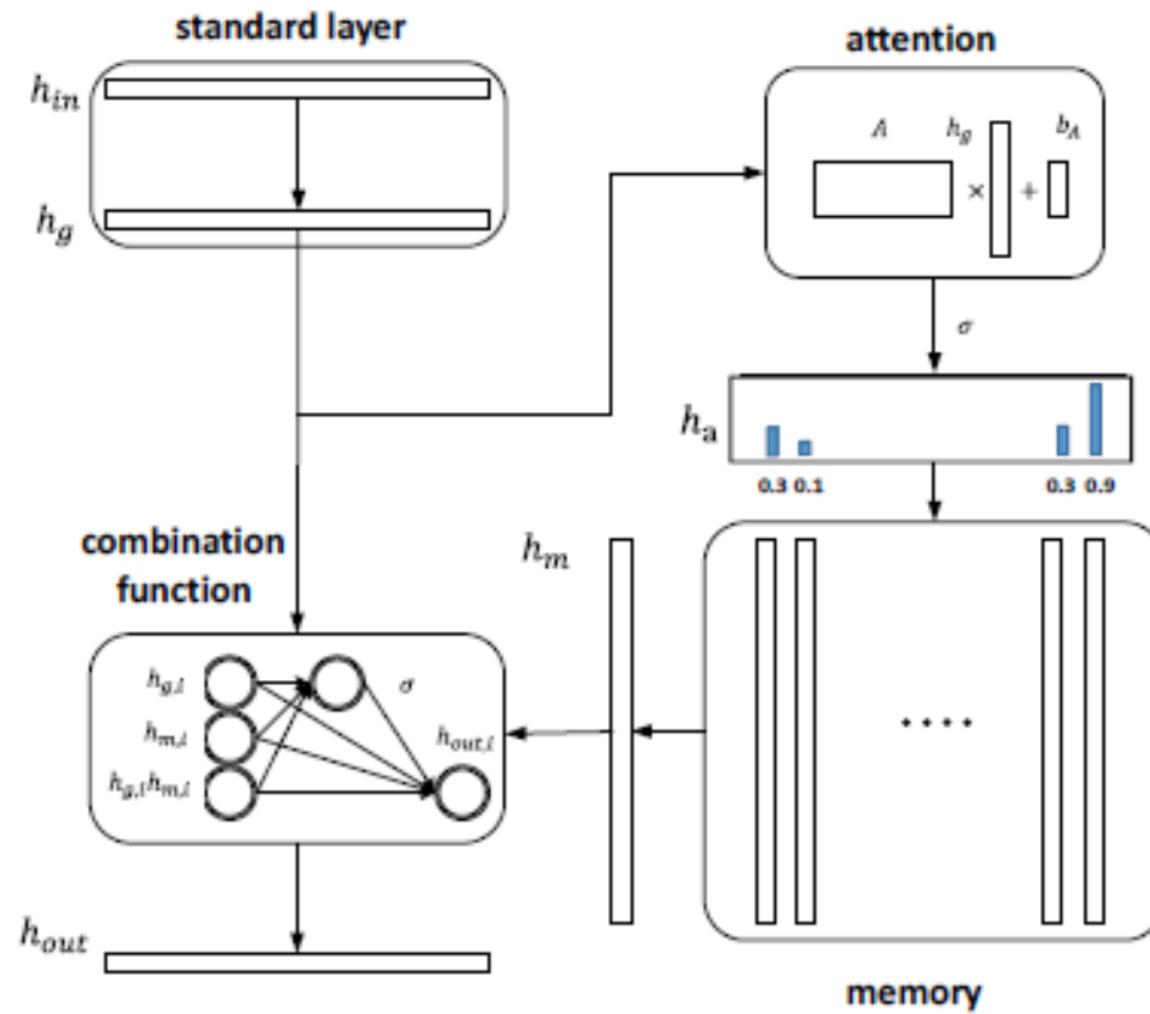
Some extensions

# Symmetric Q-P Network

- ◆ **Problem:** detail information is lost during abstraction



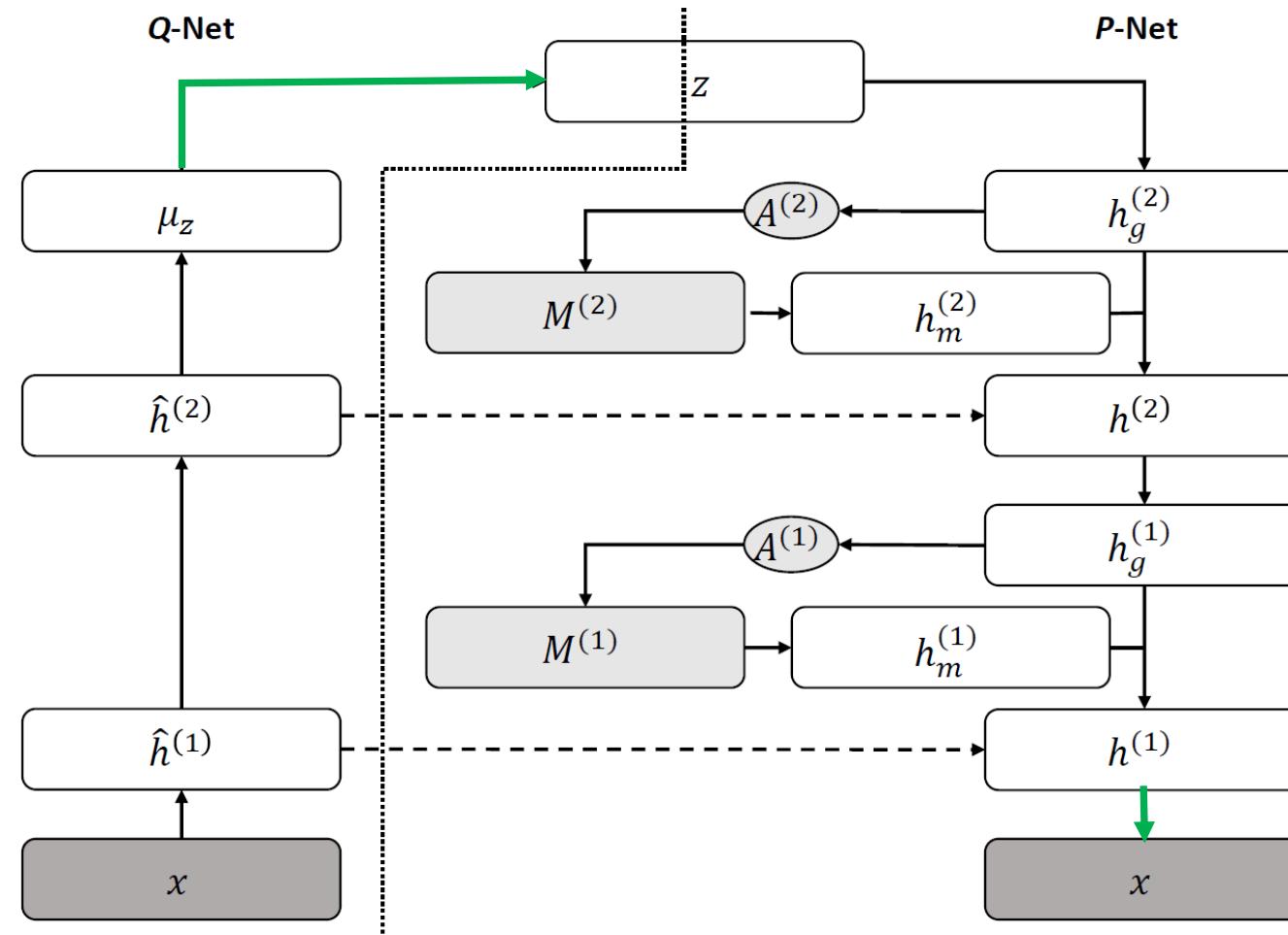
# A Layer with Memory and Attention



[Li, et al., ICML 2016]

# A Stacked Deep Model with Memory

- ◆ Asymmetric architecture



# Some Results

- ◆ Density estimation

MODELS	MNIST	OCR-LETTERS
VAE	-85.69	-30.09
<i>MEM-VAE(ours)</i>	<b>-84.41</b>	<b>-29.09</b>
IWAE-5	-84.43	-28.69
<i>MEM-IWAE-5(ours)</i>	<b>-83.26</b>	<b>-27.65</b>
IWAE-50	-83.58	-27.60
<i>MEM-IWAE-50(ours)</i>	<b>-82.84</b>	<b>-26.90</b>

- Better than symmetric VAE networks
- Comparable with state-of-the-art with much fewer parameters

# Missing Value Imputation

7	2	1	0	4	1	4	9
5	9	0	6	9	0	1	5
9	7	3	4	9	6	6	5
4	0	7	4	0	1	3	1
3	4	7	2	7	1	2	1
1	7	4	2	3	5	1	2
4	4	6	3	5	5	6	0
4	1	9	5	7	8	9	3

(a) Data

7	2	1	0	4	1	4	9
5	9	0	6	9	0	1	5
9	7	3	4	9	6	6	5
4	0	7	4	0	1	3	1
3	4	7	2	7	1	2	1
1	7	4	2	3	5	1	2
4	4	6	3	5	5	6	0
4	1	9	5	7	8	9	3

(b) Noisy data

7	2	1	0	4	1	4	9
5	9	0	6	9	0	1	5
9	7	3	4	9	6	6	5
4	0	7	4	0	1	3	1
3	4	7	2	7	1	2	1
1	7	4	2	3	5	1	2
4	4	6	3	5	5	6	0
4	1	9	5	7	8	9	3

(c) Results of VAE

7	2	1	0	4	1	4	9
5	9	0	6	9	0	1	5
9	7	3	4	9	5	6	5
4	0	7	4	0	1	3	1
3	4	7	2	7	1	2	1
1	2	4	2	3	5	1	2
4	4	6	3	5	5	6	0
4	1	9	5	7	8	9	8

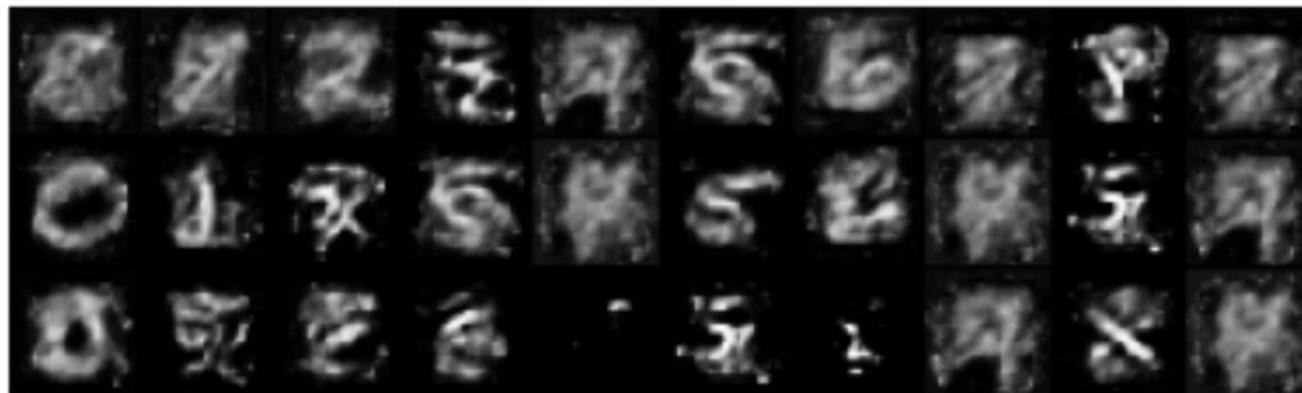
(d) Results of MEM-VAE

# Learnt Memory Slots

- ◆ Average preference over classes of the first 3 slots:

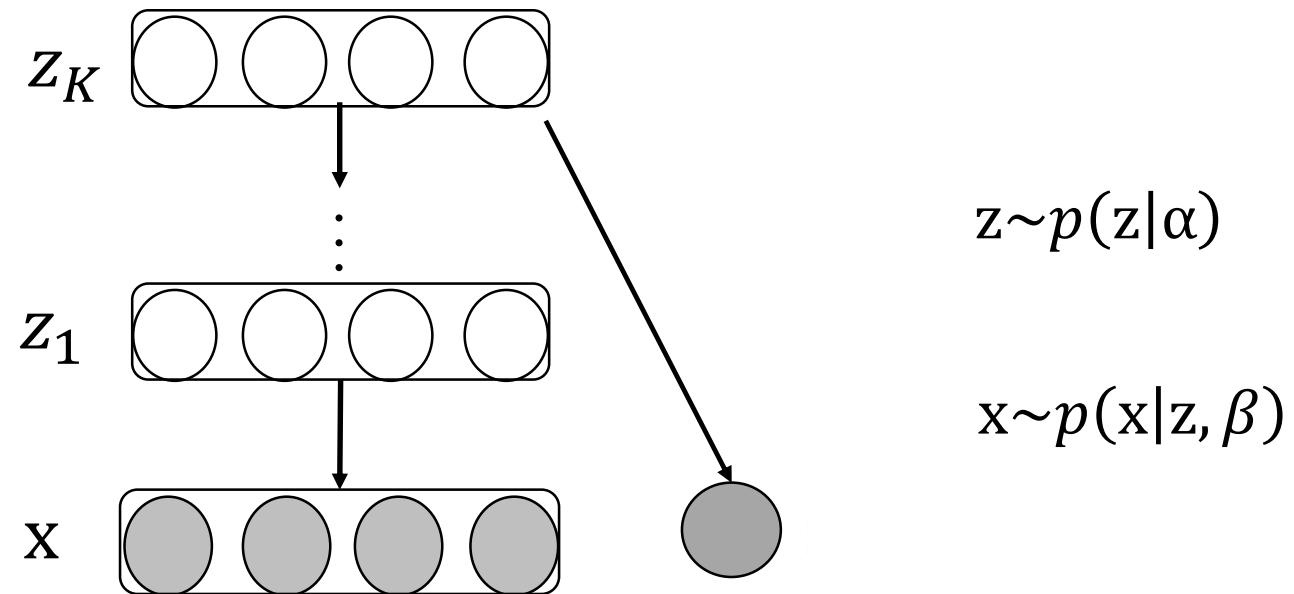
"0"	"1"	"2"	"3"	"4"	"5"	"6"	"7"	"8"	"9"
0.27	0.82	0.33	0.11	0.34	0.15	0.49	0.27	0.09	0.28
0.24	0.09	0.06	0.11	0.30	0.13	0.12	0.27	0.09	0.21
0.18	0.05	0.06	0.11	0.07	0.07	0.05	0.11	0.09	0.18

- ◆ Corresponding images:



# Learning Problem

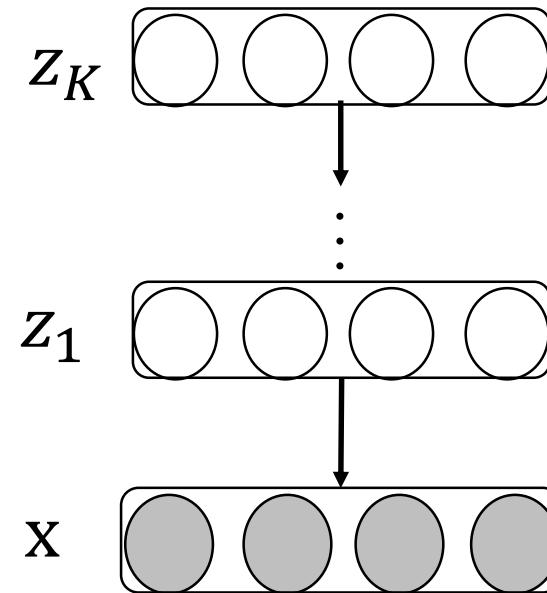
- ◆ A supervised deep generative model



$$\min_{\theta, q(\eta, \mathbf{Z})} \mathcal{L}(\theta, q(\eta, \mathbf{Z}); \mathbf{X}) + C\mathcal{R}(q(\eta, \mathbf{Z}; \mathbf{X}))$$

# Variational Objective

- ◆ Unsupervised data fitting

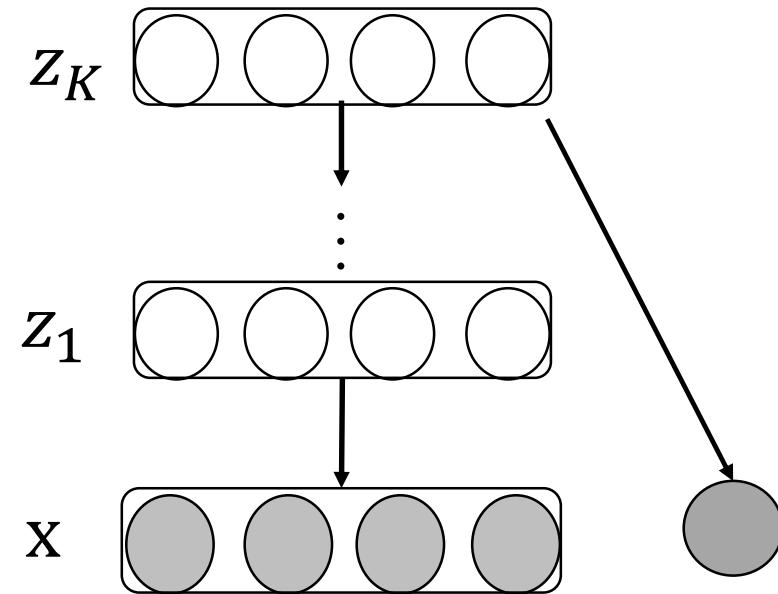


An upper bound of the per sample NLL –  $\log p(\mathbf{x}_n|\theta)$  :

$$\mathcal{L}(\theta, \phi; \mathbf{x}_n) \triangleq \text{KL}(q(\mathbf{z}_n|\phi) || p(\mathbf{z}_n|\theta)) - \mathbb{E}_{q(\mathbf{z}_n|\phi)} [\log p(\mathbf{x}_n|\mathbf{z}_n, \theta)]$$

# Classification Loss

- ◆ Hinge-loss of an averaging classifier



$$f(y; \mathbf{x}) = \mathbf{E}_q [\eta_y^T \mathbf{z}]$$

$$R(q(\eta, \mathbf{z}); \mathbf{X}) = \sum_{n=1}^N \max_y (\Delta \ell_n(y) + f(y; \mathbf{x}_n) - f(y_n; \mathbf{x}_n))$$

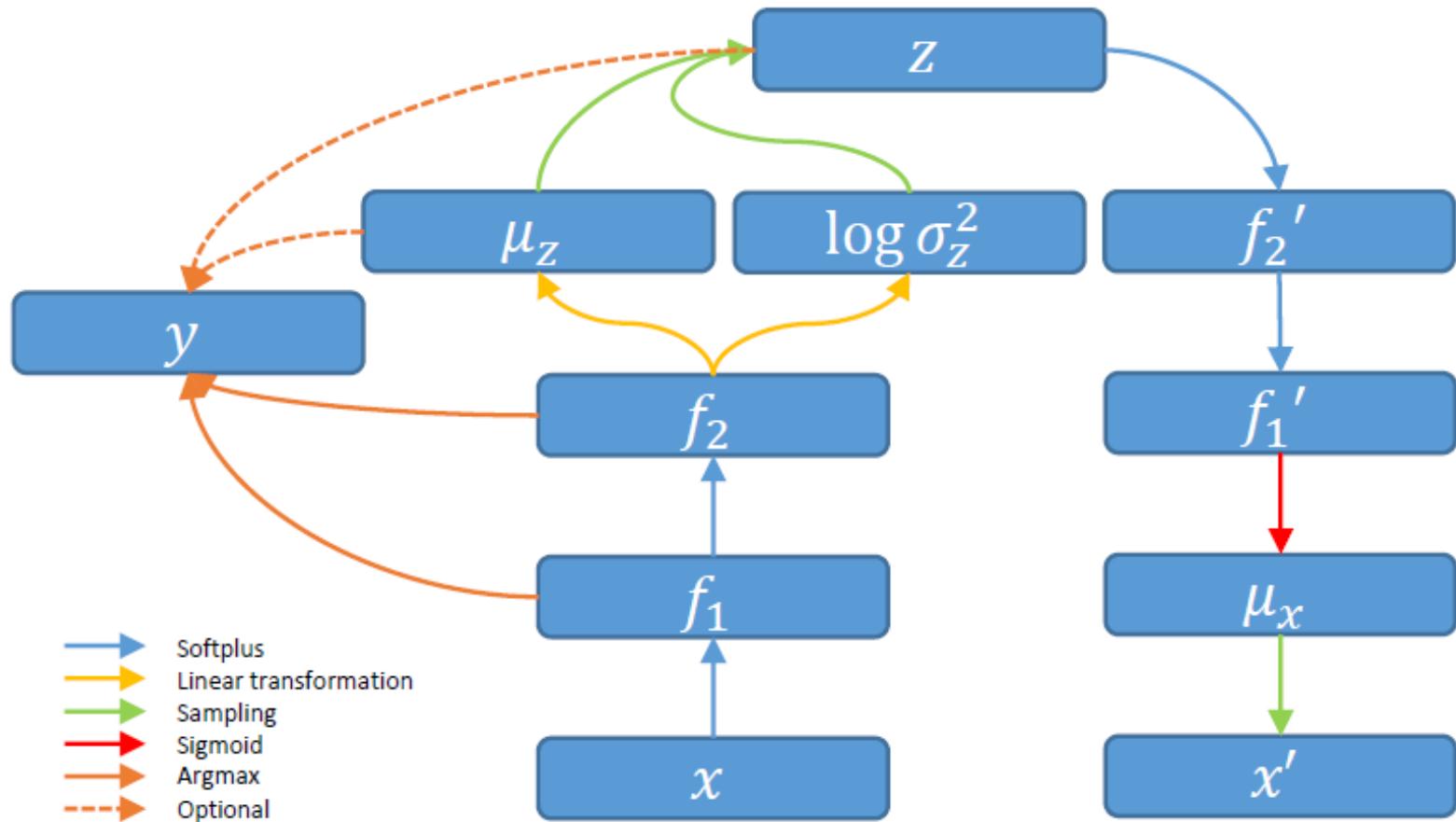
# Doubly stochastic subgradient descent

- ◆ A doubly stochastic generalization of Pegasos (Shalev-Shwartz et al., 2011)
- ◆ Stochastic **unbiased** estimate of the objective by random mini-batches:

$$\tilde{\mathcal{L}}_m := \frac{N}{m} \sum_{n=1}^m \mathcal{L}(\theta, \phi; \mathbf{x}_n) + \frac{||\boldsymbol{\lambda}||^2}{2\sigma^2} + \frac{NC}{m} \sum_{n=1}^m \ell(\boldsymbol{\lambda}, \phi; \mathbf{x}_n)$$

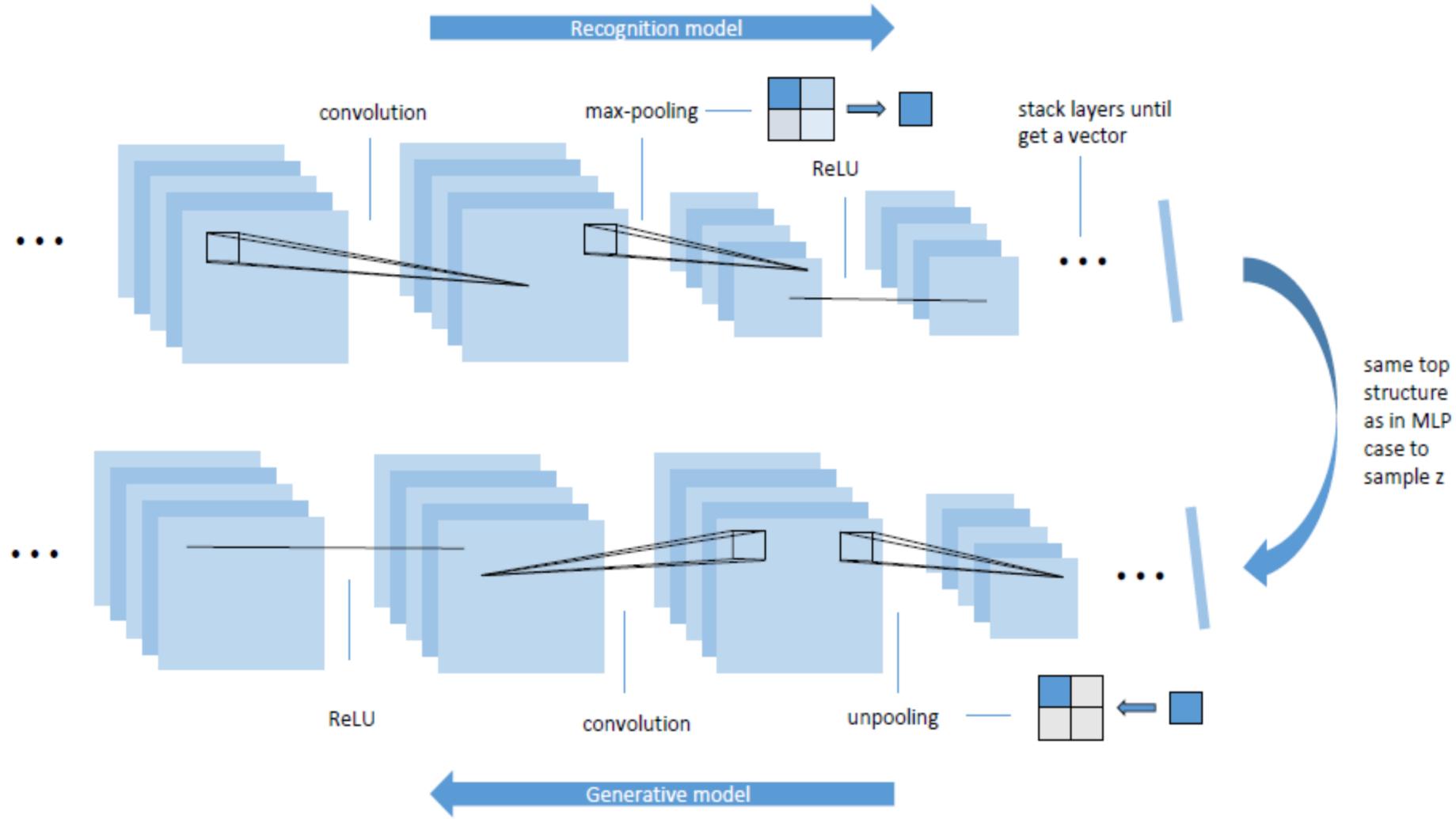
- ◆ Stochastic estimate of the per-sample variational bound with a Q-network and its subgradient (**unbiased**)

## 2-Layer MLP: Q-P network architecture



\*Same as in Auto-Encoding Variational Bayes (VA) [Kingma & Welling, 2014]

# 5-Layer CNN: Q-P network architecture



# Classification Performance

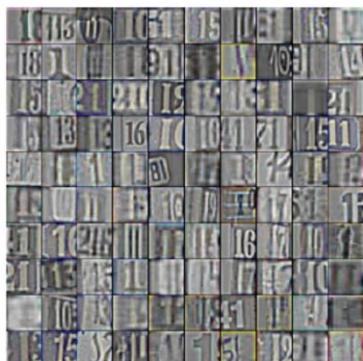
## ◆ MNIST



**Table:** Error rates (%) on MNIST dataset.

MODEL	ERROR RATE
VA+Pegasos	1.04
VA+Class-conditionVA	0.96
<b>MMVA</b>	<b>0.90</b>
CVA+Pegasos	1.35
<b>CMMVA</b>	<b>0.45</b>
<i>Stochastic Pooling (ZEILER AND FERGUS [2013])</i>	0.47
<i>Network in Network (LIN ET AL. [2014])</i>	0.47
<i>Maxout Network (GOODFELLOW ET AL. [2013])</i>	0.45
<i>DSN (LEE ET AL. [2015])</i>	0.39

## ◆ SVHN



MODEL	ERROR RATE
CVA+Pegasos	25.3
<b>CMMVA</b>	<b>3.09</b>
<i>CNN SERMANET ET AL. [2012]</i>	4.9
<i>Stochastic Pooling ZEILER AND FERGUS [2013]</i>	2.80
<i>Maxout Network GOODFELLOW ET AL. [2013]</i>	2.47
<i>Network in Network LIN ET AL. [2014]</i>	2.35
<i>DSN LEE ET AL. [2015]</i>	1.92

# Generative Capability

◆ MNIST & SVHN:

3 1 1 6 0 6 1 2 0 2 2 1
9 3 8 1 6 7 9 0 6 1 5 2
6 4 0 4 0 9 9 8 5 4 8 8
7 2 9 9 3 4 5 6 2 9 7 5
4 7 5 4 8 0 5 4 9 0 5 8
4 6 9 7 3 8 1 3 7 6 9 9
8 1 9 9 2 5 7 8 6 9 2 1
5 1 8 9 7 6 5 3 6 8 1 5
4 6 1 9 4 3 9 5 6 7 6 7
6 0 4 3 3 3 9 4 7 8 3 6
9 1 0 8 4 1 9 1 4 9 4 9
3 2 7 4 3 1 3 7 8 5 3 1

(a) VA

1 7 8 3 9 4 4 0 6 5 9 2
2 0 3 0 8 8 8 1 1 6 6 0
2 1 4 2 6 0 6 6 4 5 4 9
3 0 6 0 1 2 1 5 7 3 8 3
8 1 0 2 9 0 9 6 5 1 3 6
0 9 2 9 3 7 9 7 0 1 7 0
3 8 0 9 1 1 3 3 0 0 2 1
8 7 9 4 7 0 6 3 0 3 2 3
0 1 4 1 8 4 5 5 9 3 5 8
3 2 9 5 0 4 3 8 0 1 6 0
4 7 7 9 3 8 6 3 3 6 2 7
7 8 5 0 0 8 7 4 7 8 1 4

(b) MMVA

2 8 7 7 5 5 0 5 7 2 5 4
0 0 4 5 0 4 8 8 2 2 3 8 8
9 5 6 4 3 2 1 6 3 9 4 6
4 0 1 3 3 6 8 2 2 1 5 4
2 7 9 8 1 9 0 2 9 9 7 3
7 7 6 0 2 1 0 8 1 1 1 3
3 7 0 1 3 9 2 5 1 7 7 2
9 8 7 5 7 2 1 6 8 7 1 9
9 1 4 9 8 0 7 2 5 9 6
0 0 0 6 6 6 8 4 7 4 9 4
7 5 3 8 0 2 4 1 4 7 8 3
9 4 4 9 5 6 1 9 4 0 5 6
2 7 4 2 1 4 7 5 8 9 7 1

(c) CVA

9 3 5 5 8 7 1 0 6 7 8 8
5 5 9 7 2 9 9 5 0 3 7 9
4 3 3 2 1 1 1 2 6 2 3 2
1 5 6 4 3 8 7 1 8 6 8 0
8 8 3 4 7 9 1 9 3 6 9 0
3 8 8 8 9 0 1 9 5 6 4 9
5 3 4 2 9 6 7 0 8 9 0 0
6 6 3 4 8 0 8 7 8 8 1 2
0 0 0 6 6 6 8 4 7 4 9 4
4 1 9 3 2 2 1 0 2 9 8 6
7 0 6 0 4 2 7 6 3 6 8 8
3 8 1 7 4 1 3 9 3 6 2 7

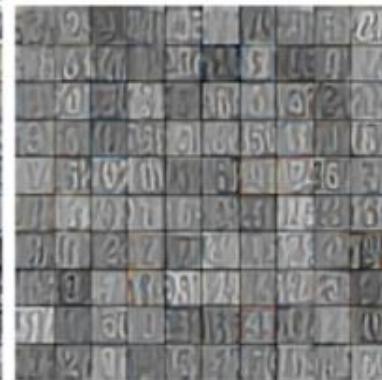
(d) CMMVA



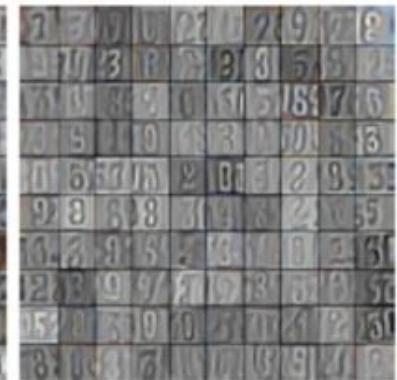
(a) Training data



(b) CVA



(c) CMMVA



(d) CMMVA

# Imputation examples

Considering  $\text{Rect}(12 \times 12)$  noise, CMMVA makes fewer mistakes and refines the images better, which accords with the MSE results.

7	2	1	0	4	1	4	9	5	9
0	6	9	0	1	5	9	7	3	4
9	6	4	5	4	0	7	4	0	1
3	1	3	4	7	2	7	1	2	1
1	7	4	2	3	5	1	2	4	4
6	3	5	5	6	0	4	1	9	5
7	8	9	3	7	4	0	4	3	0
7	0	2	9	1	7	3	2	9	7
1	6	2	7	8	4	7	3	6	1
3	6	9	3	1	4	1	7	6	9

(a) Original data



(b) Noisy data

7	2	4	0	4	1	4	9	5	9
0	6	9	0	4	5	9	7	3	4
9	6	4	5	4	0	7	4	0	1
3	4	3	4	7	2	7	1	2	4
1	7	4	2	3	8	1	2	9	4
6	3	5	5	6	0	4	1	9	9
7	3	9	3	7	4	0	4	3	0
7	0	2	9	4	7	3	2	1	7
9	6	2	9	8	4	7	3	6	4
8	6	9	3	4	4	1	7	6	9

(c) CVA

7	2	1	0	4	4	9	5	9	
0	6	9	0	1	5	9	7	3	4
9	6	4	5	4	0	7	4	0	1
3	9	3	4	7	2	7	1	2	1
1	7	4	2	3	8	1	2	9	4
6	3	5	5	6	0	4	1	9	9
7	3	9	3	7	4	0	4	3	0
7	0	2	9	1	7	3	2	9	7
9	6	2	9	5	4	7	3	6	1
8	6	7	3	1	4	1	7	6	9

(d) CMMVA

**Figure:** (a): original test data; (b) test data with missing value; (c-d): results inferred by CVA and CMMVA respectively for 100 epochs.

# Classification with Missing Values

CNN makes prediction on the incomplete data directly. CVA and CMMVA with default settings infer missing data for 100 iterations at first and then make prediction on the refined data.

**Table:** Error rates(%) with missing values on MNIST.

NOISE LEVEL	CNN	CVA	CMMVA
RECT (6 × 6)	7.5	2.5	1.9
RECT (8 × 8)	18.8	4.2	3.7
RECT (10 × 10)	30.3	8.4	7.7
RECT (12 × 12)	47.2	18.3	15.9

CMMVA outperforms both VA and CNN in this scenario.

# Moment-Matching DGMs

- ◆ Moment-matching:

- draw samples from  $p$  via DGM:

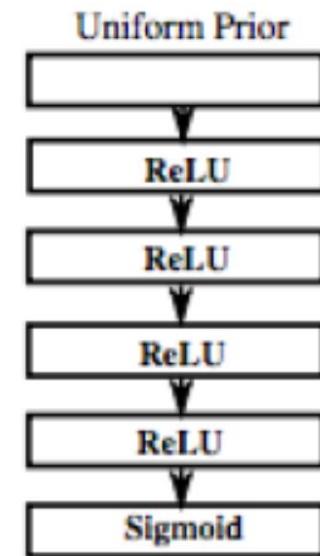
$$y_i = f(\epsilon_i; \theta), \quad \epsilon_i \sim \text{Uniform}(0,1)$$

$$\hat{D} = \{y_i\}_{i=1}^M$$

- Kernel MMD:

$$\mathcal{L}_{MMD^2} = \left\| \frac{1}{N} \sum_{i=1}^N \phi(x_i) - \frac{1}{M} \sum_{j=1}^M \phi(y_j) \right\|_{\mathcal{H}}^2$$

- Rich enough to distinguish any two distributions in certain RKHS
- Back-propagation to update the parameters



# Conditional Models

- ◆ We're more interested in conditional distributions in many cases

- Predictive modeling:  $P \left[ y = \text{cat} \mid x = \text{ }$  

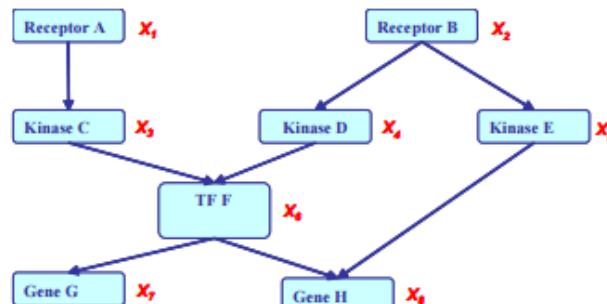
- Better performance with fewer training samples

- Contextual generation:

$$P \left[ y = \text{ } \left| \begin{array}{l} \text{cat wearing glasses and tie} \\ \text{ }$$

$$x = \text{cat with glass and tie} \right. \right]$$

- Building large networks:



# Conditional Moment-Matching Networks

- ◆ Two sets of data samples (one from training set, one from the model)

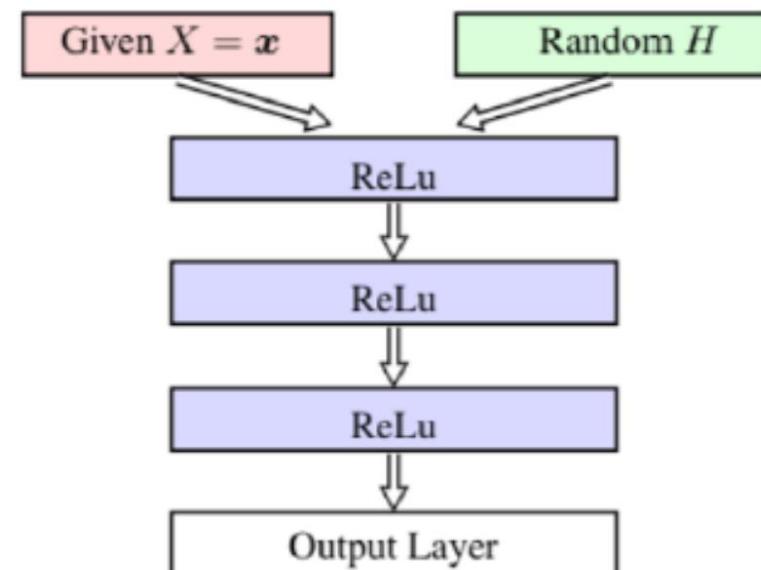
$$\mathcal{D}_{XY}^s = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$$

$$\mathcal{D}_{XY}^d = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^M$$

- ◆ A DGM:

$$p(\mathbf{h}) = \prod_{i=1}^d U(h_i)$$

$$\mathbf{y} = f(\mathbf{x}, \mathbf{h} | \mathbf{w})$$



# Conditional Moment-Matching Networks

- ◆ We define conditional max-min discrepancy (CMMMD):

$$\begin{aligned}\widehat{\mathcal{L}}_{\text{CMMMD}}^2 &= \left\| \Phi_d(K_d + \lambda I)^{-1} \Upsilon_d^\top - \Phi_s(K_s + \lambda I)^{-1} \Upsilon_s^\top \right\|_{\mathcal{F} \otimes \mathcal{G}}^2 \\ &= \text{Tr} \left( K_d \widetilde{K}_d^{-1} L_d \widetilde{K}_d^{-1} \right) + \text{Tr} \left( K_s \widetilde{K}_s^{-1} L_s \widetilde{K}_s^{-1} \right) \\ &\quad - 2 \cdot \text{Tr} \left( K_{sd} \widetilde{K}_d^{-1} L_{ds} \widetilde{K}_s^{-1} \right),\end{aligned}$$

- the superscripts  $s$  and  $d$  denote the two sets of samples
- $\widetilde{K} = K + \lambda I$
- $\Phi_d := (\phi(\mathbf{y}_1^d), \dots, \phi(\mathbf{y}_N^d))$ ,  $\Upsilon_d := (\phi(\mathbf{x}_1^d), \dots, \phi(\mathbf{x}_N^d))$ ;  $\Phi_s, \Upsilon_s$ : similarly for  $\mathcal{D}_{XY}^s$ .
- $K_d = \Upsilon_d^\top \Upsilon_d, K_s = \Upsilon_s^\top \Upsilon_s$ : gram matrices for input variables;  
 $L_d = \Phi_d^\top \Phi_d, L_s = \Phi_s^\top \Phi_s$ : gram matrices for output variables
- $K_{sd} = \Upsilon_s^\top \Upsilon_d, L_{ds} = \Phi_d^\top \Phi_s$ : gram matrices between the two datasets on input and output variables.

# Connections with MMD

◆ MMD:

$$\mathcal{L}_{MMD}^2 = \text{Tr}(L_d \cdot \mathbf{1}) + \text{Tr}(L_s \cdot \mathbf{1}) - 2 \cdot \text{Tr}(L_{ds} \cdot \mathbf{1}),$$

where  $\mathbf{1}$  is the matrix with all entities equal to 1.

◆ CMMMD:

$$\mathcal{L}_{CMMMD}^2 = \text{Tr}(L_d \cdot C_1) + \text{Tr}(L_s \cdot C_2) - 2 \cdot \text{Tr}(L_{ds} \cdot C_3),$$

where  $C_1, C_2, C_3$  are some matrix based on the training and generated data.

Instead of putting uniform weights on the gram matrix, CMMMD applies **non-uniform weight**, reflecting the influence of conditional variables.

# Mini-batch Training Algorithm

- ◆ BP to compute the gradients:

**Input:** Dataset  $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$

**Output:** Learned parameters  $\mathbf{w}$

Randomly divide training dataset  $\mathcal{D}$  into mini batches

**While** Stopping criterion not met **do**

    Draw a minibatch  $\mathcal{B}$  from  $\mathcal{D}$ ;

    For each  $\mathbf{x} \in \mathcal{B}$ , generate a  $\mathbf{y}$ ; and set  $\mathcal{B}'$  to contain all the generated  $(\mathbf{x}, \mathbf{y})$ ;

    Compute the gradient  $\frac{\partial \widehat{\mathcal{L}}_{\text{CMMMD}}^2}{\partial \mathbf{w}}$  on  $\mathcal{B}$  and  $\mathcal{B}'$ ;

    Update  $\mathbf{w}$  using the gradient with proper regularizer.

**EndWhile**

# Some Results – Classification

Classification accuracy on MNIST:

Model	Error Rate
VA+Pegasos	1.04
MMVA	0.90
CGMMN	0.97
CVA + Pegasos	1.35
CGMMN-CNN	0.47
Stochastic Pooling	0.47
Network in Network	0.47
Maxout Network	0.45
CMMVA	0.45
DSN	0.39

Classification accuracy on SVHN:

Model	Error Rate
CVA+Pegasos	25.3
CGMMN-CNN	3.13
CNN	4.9
CMMVA	3.09
Stochastic Pooling	2.80
Network in Network	2.47
Maxout Network	2.35
DSN	1.92

# Some Results – Generation

Generating performance on MNIST:



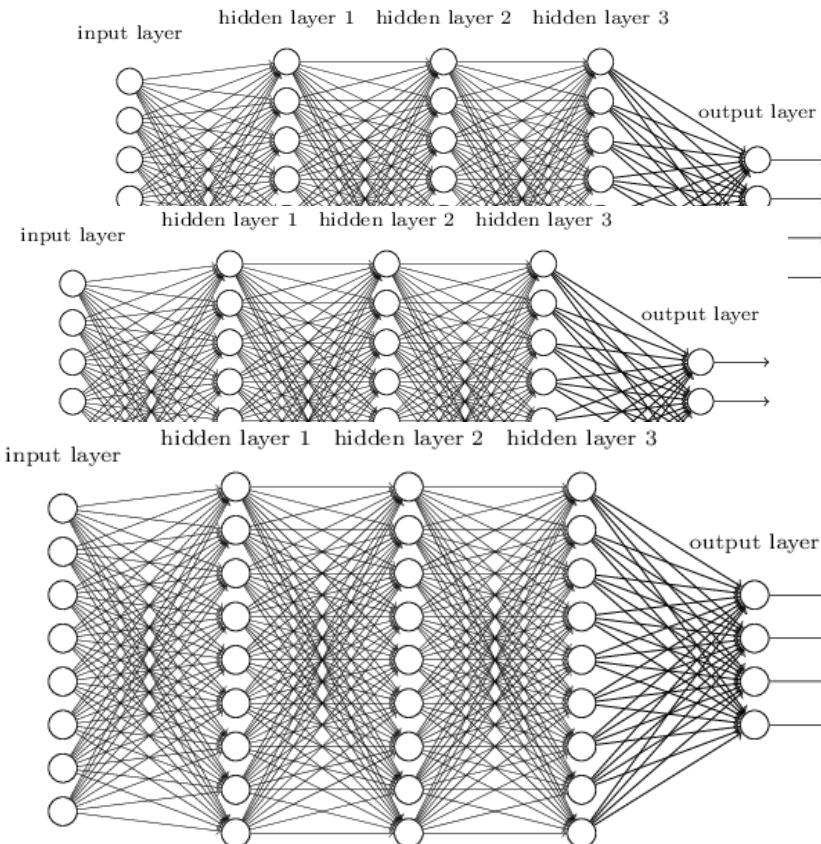
(a) MNIST samples

(b) Random samples

(c) Samples conditioned on label 0

# Some Results – Bayesian Dark Knowledge

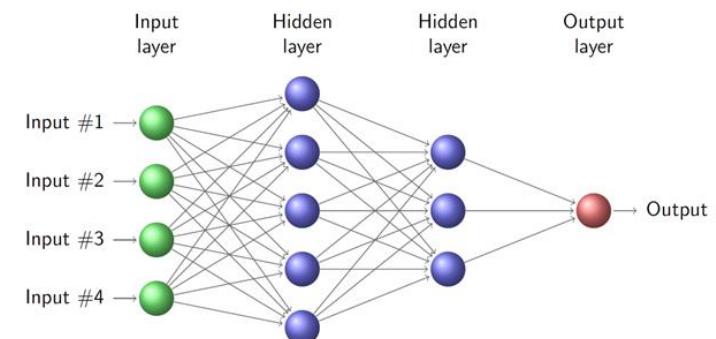
- ◆ Distilling the knowledge (Hinton et al., 2015) & Bayesian dark knowledge (Korattikara et al., 2015)



$$p(\mathbf{y}|\mathbf{x}, \mathcal{D}) = \int p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})p(\boldsymbol{\theta}|\mathcal{D})d\boldsymbol{\theta}$$

Bayes prediction is expensive!

Teach a student network for prediction



# Some Results – Bayesian Dark Knowledge

Distilling knowledge in Bayesian networks on Boston housing dataset:

- one-dimensional regression: 506 data points where each data is of dimension 13
- distill a PBP model
- test whether the distilled model will degrade the prediction performance.

PBP prediction	Distilled by CGMMN
$2.574 \pm 0.089$	$2.580 \pm 0.093$

- measured by MSE
- MLP network with three hidden layers and (100, 50, 50) hidden units for middle layers.
- $N = 3,000$  sample pairs  $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$
- $\mathbf{x}_i$  is generated by adding noise into training data to avoid fitting the testing data directly

# Summary

- ◆ Deep generative models are flexible in density estimation, sampling generation, etc.
  - VAE, GAN, etc.
- ◆ Memory and attention are helpful for DGMs
- ◆ Supervised/Semi-supervised learning for DGMs with state-of-the-art performance
- ◆ Moment-matching networks can learn conditional distributions well
- ◆ Code: <https://github.com/thu-ml>

# References

- ◆ Kingma, D. P. and Welling, M. *Auto-encoding variational Bayes*. In ICLR, 2013.
- ◆ Rezende, D. J., Mohamed, S., and Wierstra, D. *Stochastic backpropagation and approximate inference in deep generative models*. In ICML, 2014.
- ◆ Burda, Y., Grosse, R., and Salakhutdinov, R. *Importance weighted autoencoders*. In arXiv:1509.00519, 2015.
- ◆ Goodfellow, I. J., Abadie, J. P., Mirza, M., Xu, B., Farley, D. W., S.ozair, Courville, A., and Bengio, Y. *Generative adversarial nets*. In NIPS, 2014
- ◆ Bengio, Y., Thihodeau-Laufer, E., Alain, G., and Yosinski, J. *Deep generative stochastic networks trainable by backprop*. In ICML, 2014.
- ◆ Li, C., Zhu, J., and Zhang, B. *Learning to generate with memory*. In ICML, 2016
- ◆ Li, C., Zhu, J., and Zhang, B. *Max-margin deep generative models*. In NIPS, 2015
- ◆ Li, C., Zhu, J., and Zhang, B. *Max-margin deep generative models for (semi-)supervised learning*, arXiv, 2016
- ◆ Ren, Y., Li, J., Luo, Y., and Zhu, J. *Conditional generative moment-matching networks*, NIPS, 2016.

Thanks