

Variational Autoencoder

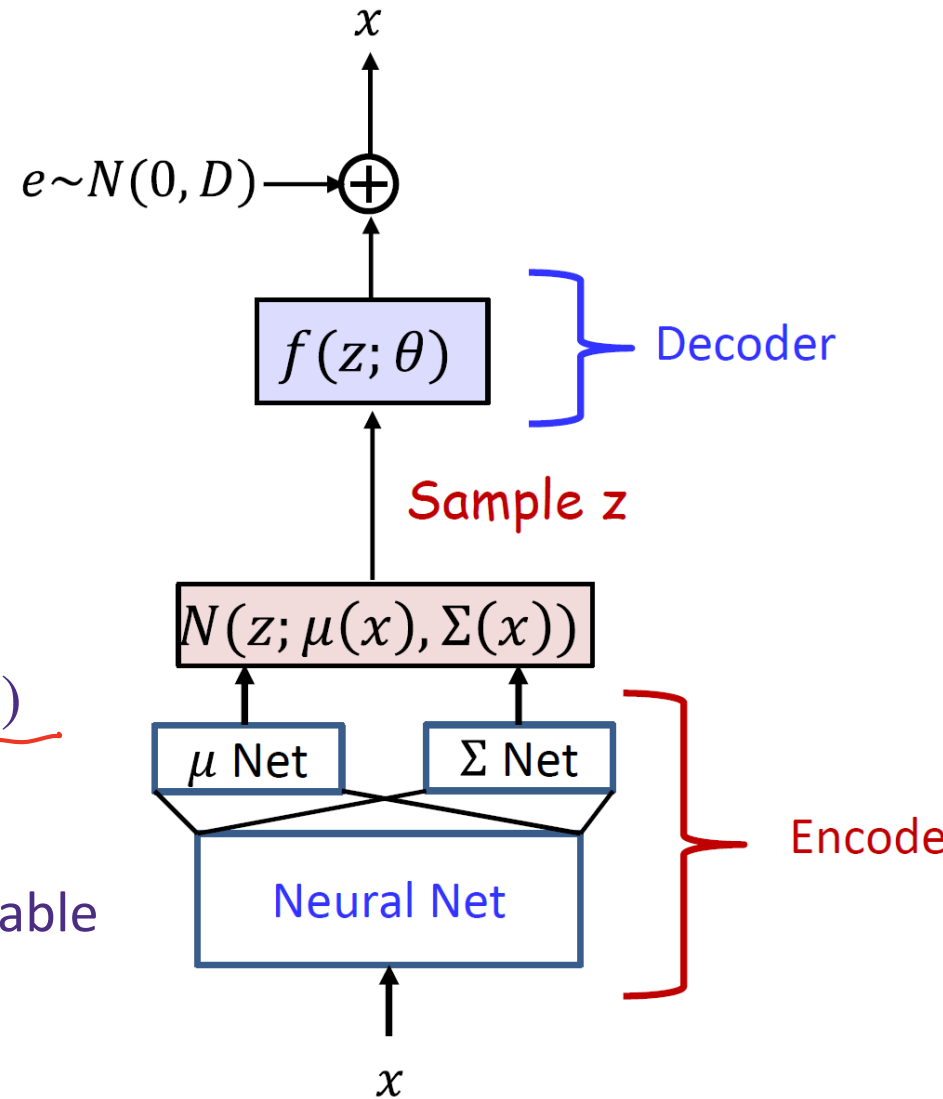


Architecture

- Auto-encoder: $x \rightarrow z \rightarrow x$
- Encoder: $q(z | x; \phi) : x \rightarrow z$
- Decoder: $p(x | z; \theta) : z \rightarrow x$

- Isomorphic Gaussian:
 $q(z | x; \phi) = N(\mu(x; \phi), \text{diag}(\exp(\sigma(x; \phi))))$
- Gaussian prior: $p(z) = N(0, I)$
- Gaussian likelihood: $p(x | z; \theta) \sim N(f(z; \theta), I)$

- Probabilistic model interpretation: latent variable model.



VAE Training

standard MLE: $\mathbb{E}_{x \sim p} \ell(f, y)$

$KL(p || q) = \mathbb{E}_p \log \frac{p}{q}$

- Training via optimizing ELBO

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

- Likelihood term + KL penalty

- KL penalty for Gaussians has closed form.

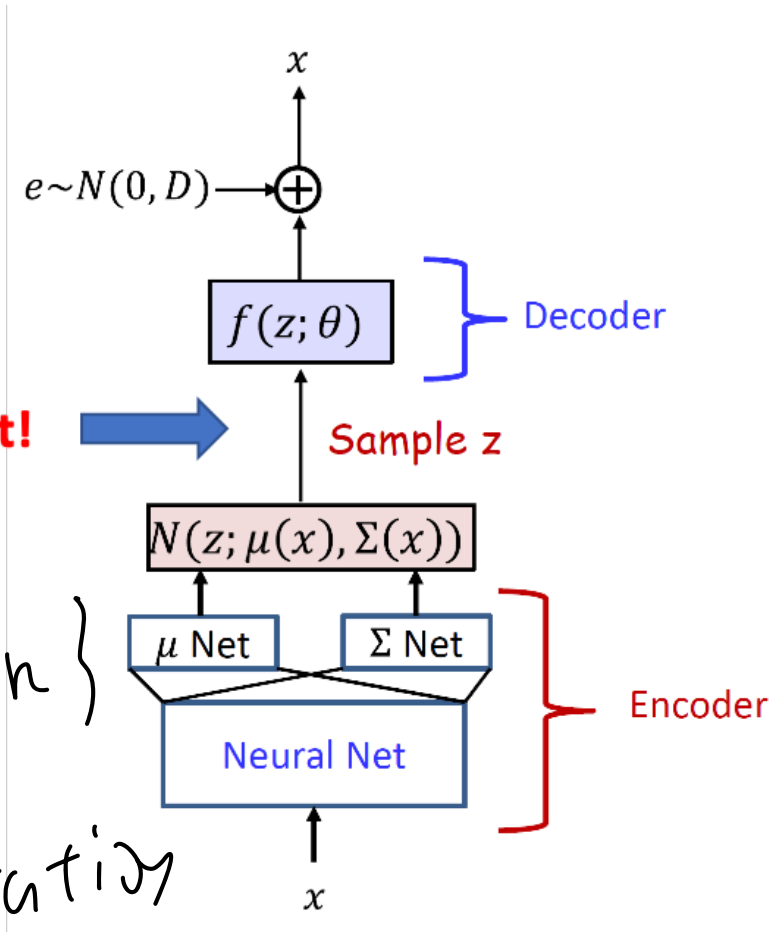
- Likelihood term (reconstruction loss):

- Monte-Carlo estimation
- Draw samples from $q(z|x; \phi)$
- Compute gradient of θ :

- $x \sim N(f(z; \theta); I)$

- $p(x) = \frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2} \|x - f(z; \theta)\|_2^2)$

$x \mapsto$ sample $\{z_1, \dots, z_n\}$
 \rightarrow approximate expectation



VAE Training

- Likelihood term (reconstruction loss):

- Gradient for ϕ . Loss: $L(\phi) = \mathbb{E}_{z \sim q(z; \phi)} [\log p(x | z)]$

- Reparameterization trick:

- $z \sim N(\mu, \Sigma) \Leftrightarrow z = \mu + \epsilon, \epsilon \sim N(0, \Sigma)$

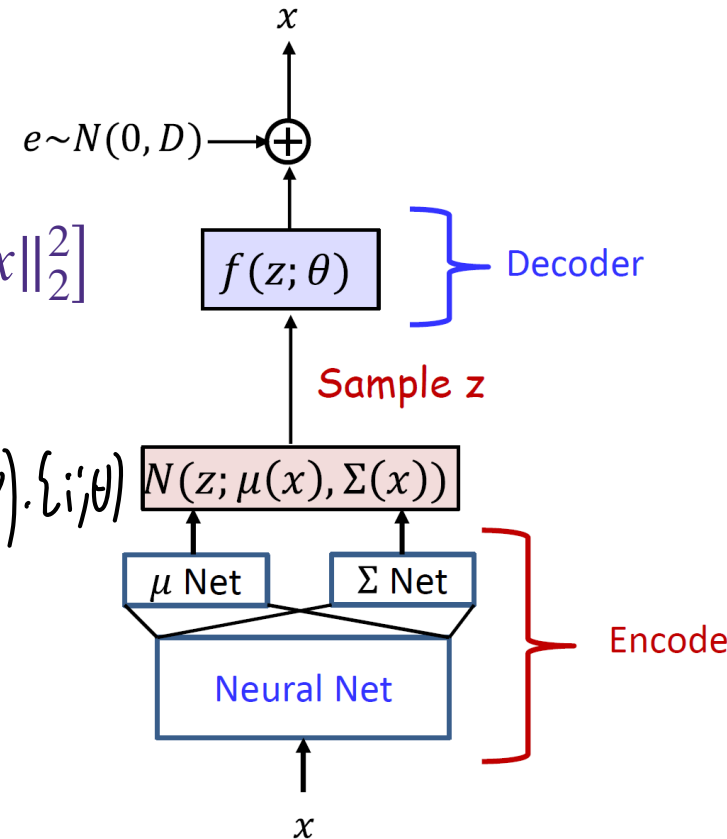
- $L(\phi) \propto \mathbb{E}_{z \sim q(z|\phi)} [\|f(z; \theta) - x\|_2^2]$
 - $\propto \mathbb{E}_{\epsilon \sim N(0, I)} [\|f(\mu(x; \phi) + \sigma(x; \phi) \cdot \epsilon; \theta) - x\|_2^2]$

- Monte-Carlo estimate for $\nabla L(\phi)$

for approx $\epsilon_1, \dots, \epsilon_k \sim N(0, I)$

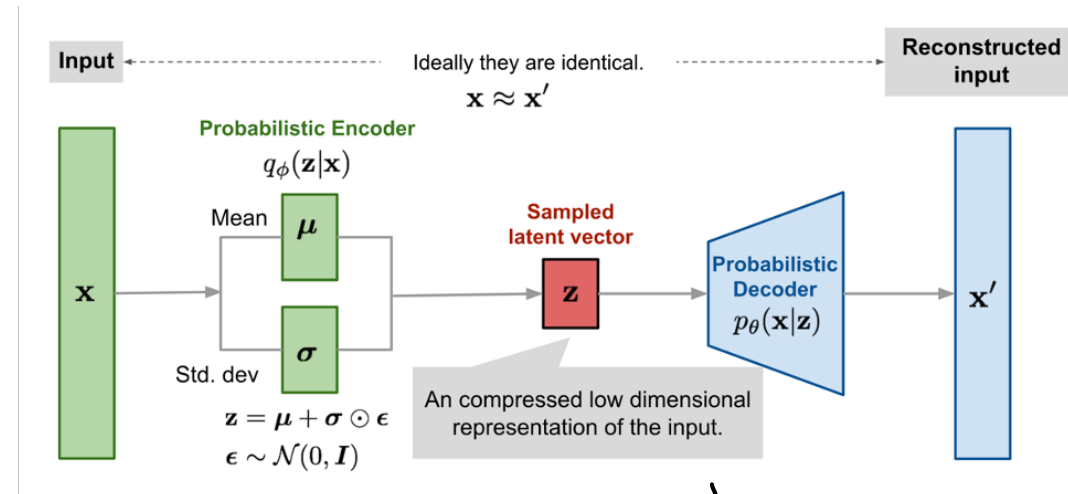
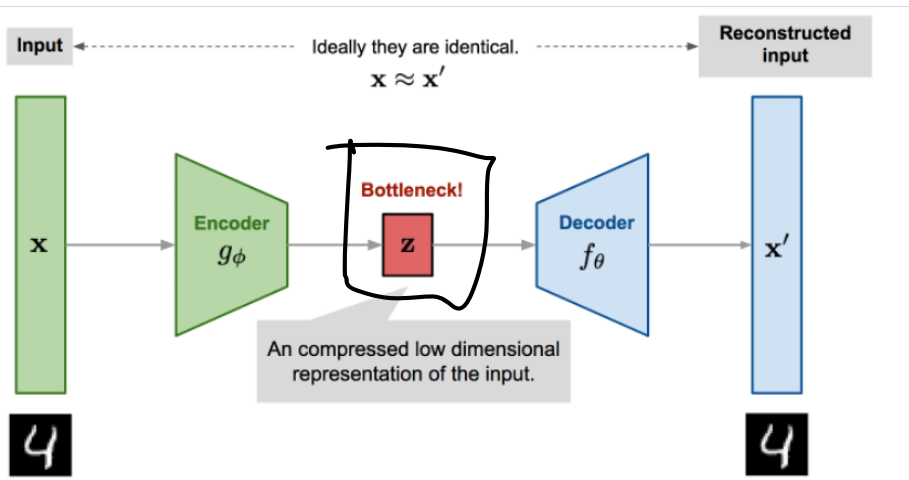
- End-to-end training

$\sum_i \|f(\mu(x, \phi) + \sigma(x, \phi) \cdot \epsilon_i; \theta) - x\|_2^2$



VAE vs. AE

- AE: classical unsupervised representation learning method.
- VAE: a probabilistic model of AE
 - AE + Gaussian noise on z
 - KL penalty: L_2 constraint on the latent vector z



Handwritten notes:

$$z \sim \mathcal{N}(0, I)$$

apply ϵ $P(x|z; \theta)$
 $x \sim f(z; \theta) + \epsilon$

Conditioned VAE

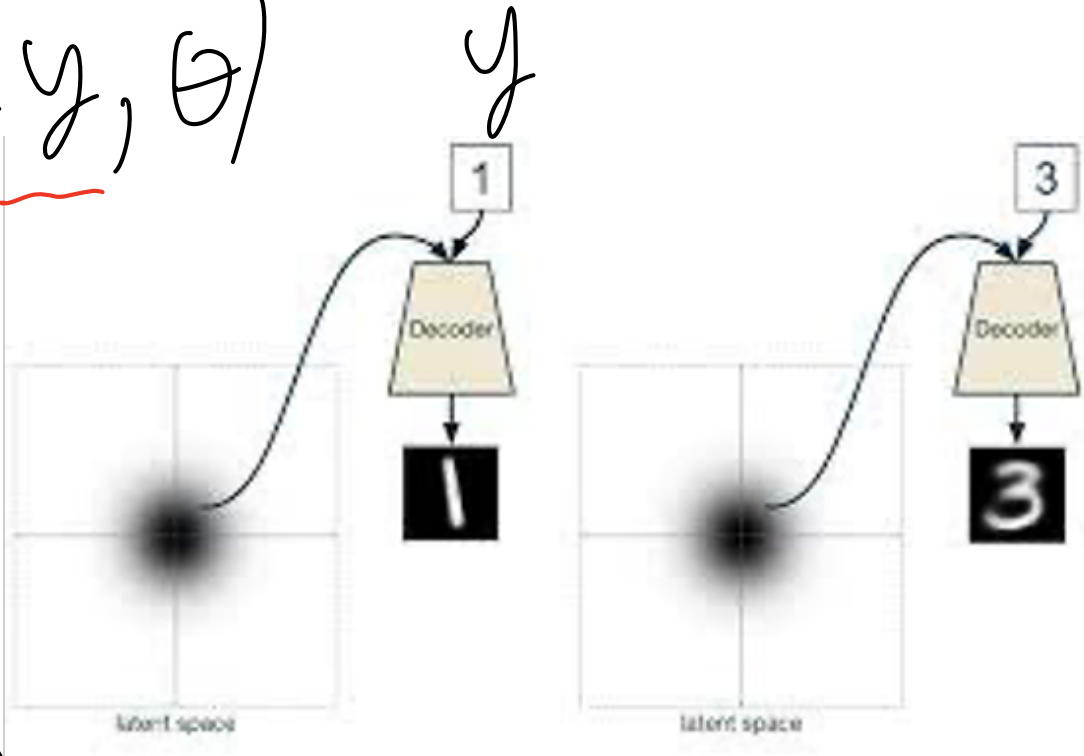
- Semi-supervised learning: some labels are also available

decoder
 $P(x|z, \underline{y}, \theta)$

$$f(z, y; \theta)$$

$$z \sim N(\mu, \Sigma)$$

$$f(z, y, \theta)$$



conditioned generation

Comments on VAE

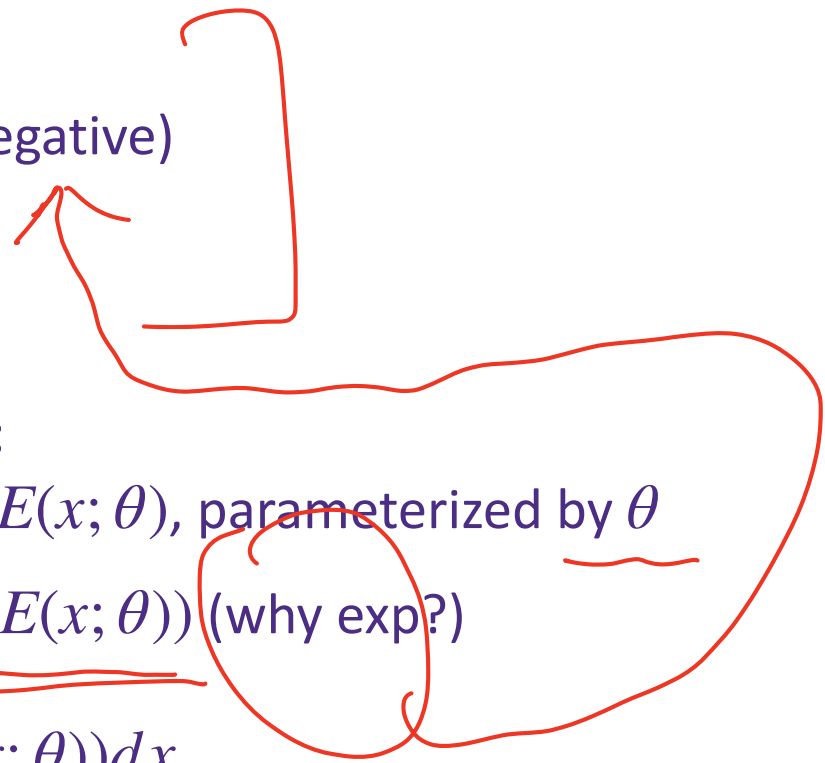
- Pros:
 - Flexible architecture
 - Stable training
- Cons:
 - Inaccurate probability evaluation (approximate inference)

$$p(q_b(y))$$

Energy-Based Models

W

Energy-based Models

- Goal of generative models:
 - a probability distribution of data: $P(x)$
 - Requirements
 - $P(x) \geq 0$ (non-negative)
 - $\int_x P(x)dx = 1$
 - Energy-based model:
 - Energy function: $E(x; \theta)$, parameterized by θ
 - $P(x) = \frac{1}{z} \exp(-E(x; \theta))$ (why exp?)
 - $z = \int_x \exp(-E(x; \theta))dx$
- 

Boltzmann Machine

$$Z = \sum_S E(S)$$

- Generative model

- $E(y) = \frac{1}{2} y^T W y$

- $P(y) = \frac{1}{Z} \exp\left(-\frac{E(y)}{T}\right)$, T : temperature hyper-parameter

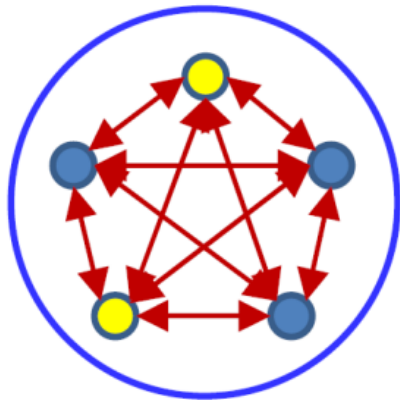
- W : parameter to learn

- When y_i is binary, patterns are affecting each other through W

$y_i \in \{0, 1\}$

$$S = y$$

$S \in \mathbb{R}^d$



$$z_i = \frac{1}{T} \sum_j w_{ji} s_j$$

$$P(s_i = 1 | s_{j \neq i}) = \frac{1}{1 + e^{-z_i}}$$

Boltzmann Machine: Training

- Objective: maximum likelihood learning (assume $T=1$):
 - Probability of one sample:

$$P(y) = \frac{\exp(\frac{1}{2}y^T W y)}{\sum_{y'} \exp(\frac{1}{2}y'^T W y')}$$

$$y' \in \{0, 1\}$$

- Maximum log-likelihood:

$$L(W) = \frac{1}{N} \sum_{y \in D} \frac{1}{2} y^T W y - \log \sum_{y'} \exp(\frac{1}{2} y'^T W y')$$

$$D = \{y_1, \dots, y_N\}$$

$$2^d: \# \text{ of } y'$$

Boltzmann Machine: Training

$$L(W) = \frac{1}{N} \sum_{y \in Y} \frac{1}{2} y^T W y - \log \sum_{y'} \exp\left(\frac{1}{2} y'^T W y'\right)$$

$W \in \mathbb{R}^{d \times d}$, focus on W_{ij}

$$\nabla_{W_{ij}} L = \frac{1}{N} \sum_y y_i \cdot y_j - \sum_{y'} \frac{\exp\left(\frac{1}{2} y'^T W y'\right)}{Z} \cdot y_i \cdot y_j$$

$$\mathbb{E}[y_i \cdot y_j]$$

\Leftrightarrow Monte-Carlo
sample $S = \{y^1, \dots, y^K\}$

$$\Rightarrow \nabla_{W_{ij}} L \approx \frac{1}{N} \sum_y y_i \cdot y_j - \frac{1}{K} \sum_{y' \in S} y_i \cdot y_j$$

Boltzmann Machine: ~~Training~~ Sampling

Sampling: $M \times M$ d

Initialize

$$y(0) \in \mathcal{R}^d$$

for $t=1, \dots, T$

iterate

$j \in \{1, \dots, d\}$, conditioned sampling

$$y_j(t) \sim P(\cdot | y_{j \neq i}(t-1))$$

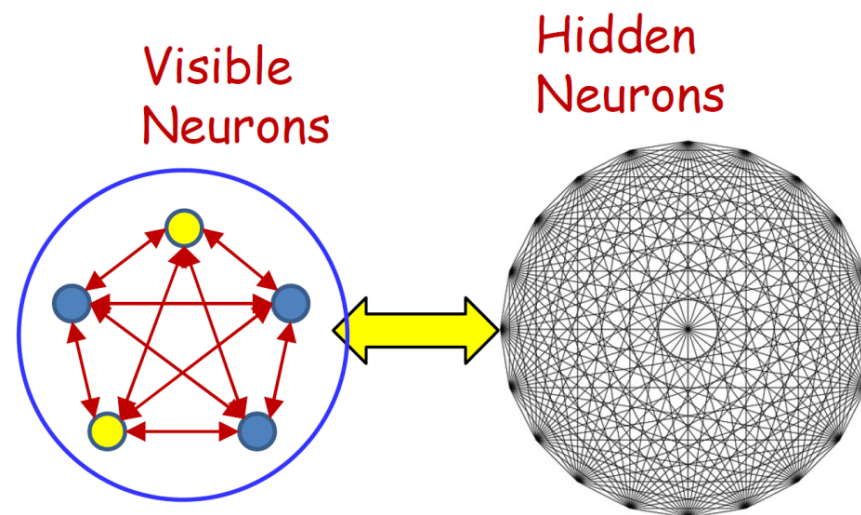
$$\Rightarrow \{y(0), \dots, y(t)\}$$

Boltzmann Machine with Hidden Neurons

- Visible and hidden neurons:

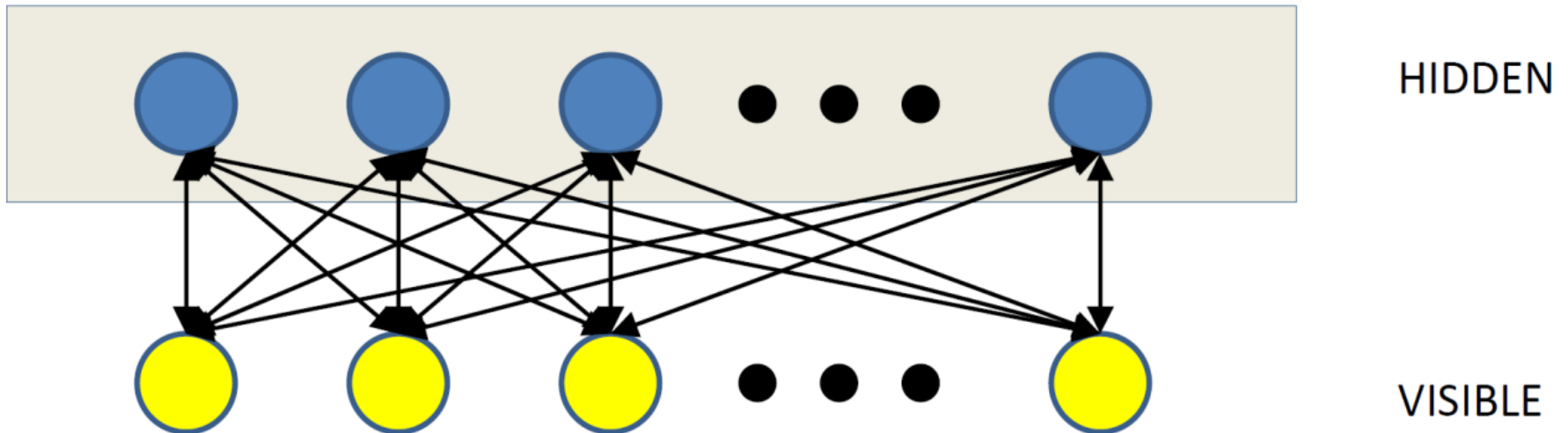
- y : visible, h : hidden

- $$P(y) = \sum_v P(y, v)$$



Restricted Boltzmann Machine

- A structured Boltzmann Machine
 - Hidden neurons are only connected to visible neurons
 - No intra-layer connections
 - Invented by Paul Smolensky in '89
 - Became more practical after Hinton invested fast learning algorithms in mid 2000



Restricted Boltzmann Machine

- Computation Rules

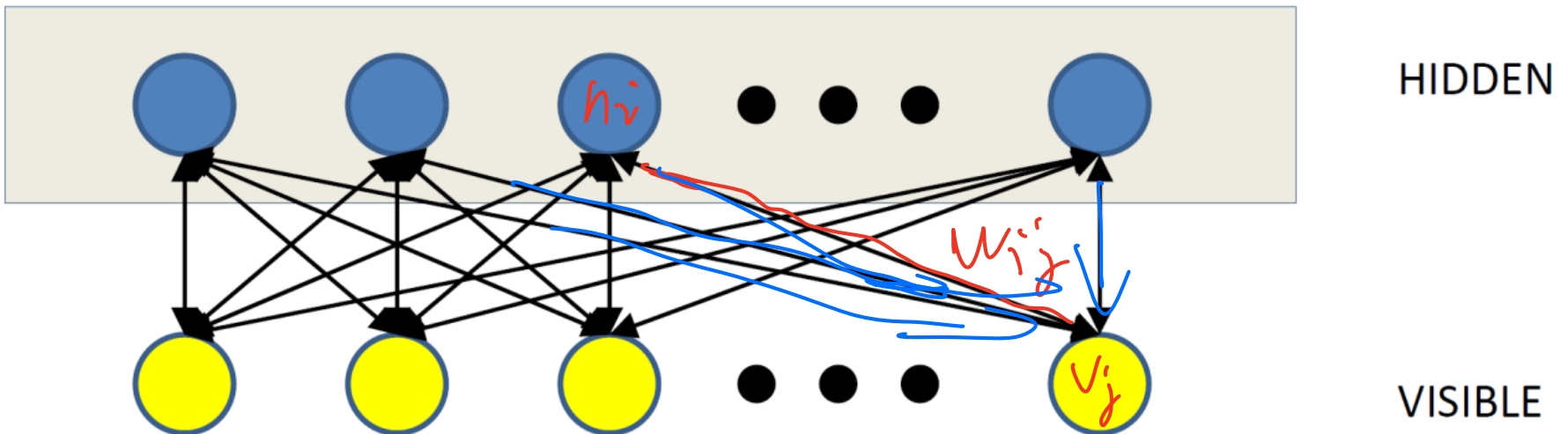
- Iterative sampling

$M(M)$

init $u(0), v(0)$
 generate $h(t)$ condition on $v(t-1)$
 $v(p)$ condition on $h(t)$

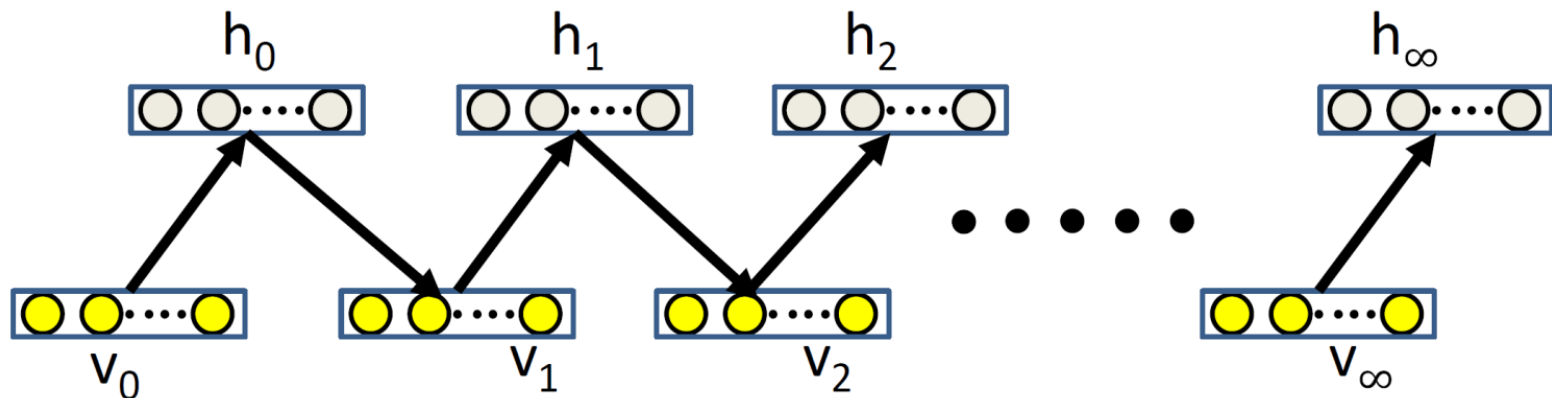
- Hidden neurons $h_i: z_i = \sum_j w_{ij} v_j, P(h_i | v) = \frac{1}{1 + \exp(-z_i)}$

- Visible neurons $v_j: z_j = \sum_i w_{ij} h_i, P(v_j | h) = \frac{1}{1 + \exp(-z_j)}$



Restricted Boltzmann Machine

- Sampling:
 - Randomly initialize visible neurons v_0
 - Iterative sampling between hidden neurons and visible neurons
 - Get final sample (v_∞, h_∞)
- Training:
 - MLE
 - Sampling to approximate gradient

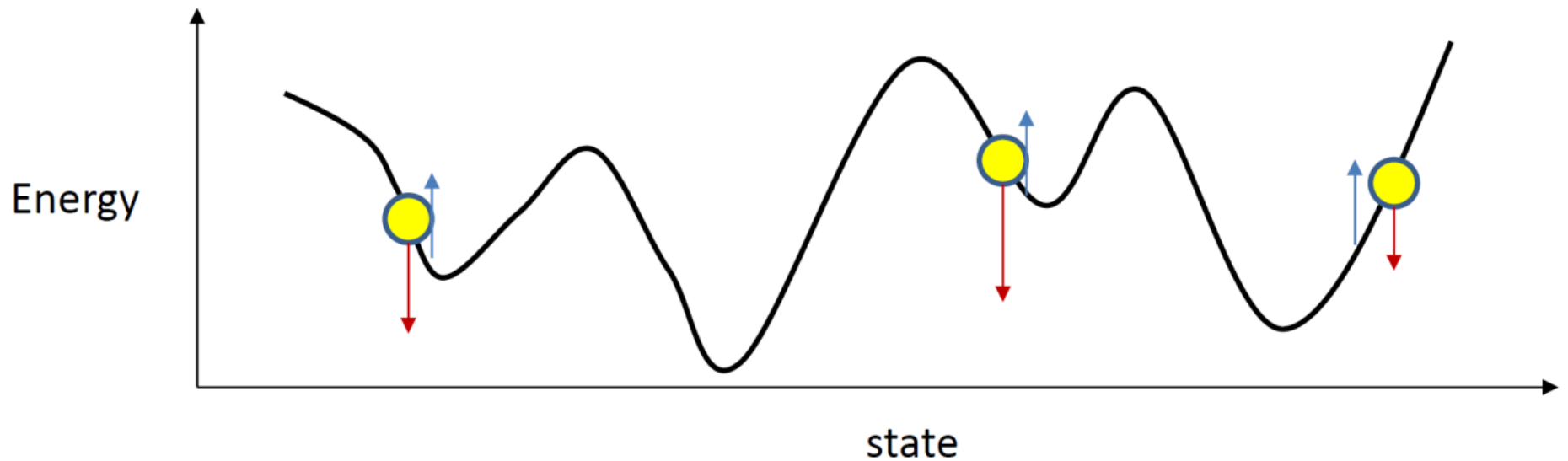


Restricted Boltzmann Machine

- Maximum likelihood estimated:

- $$\nabla_{w_{ij}} L(W) = \frac{1}{N_P K} \sum_{v \in P} v_{0i} h_{0j} - \frac{1}{M} \sum v_{\infty i} h_{\infty j}$$

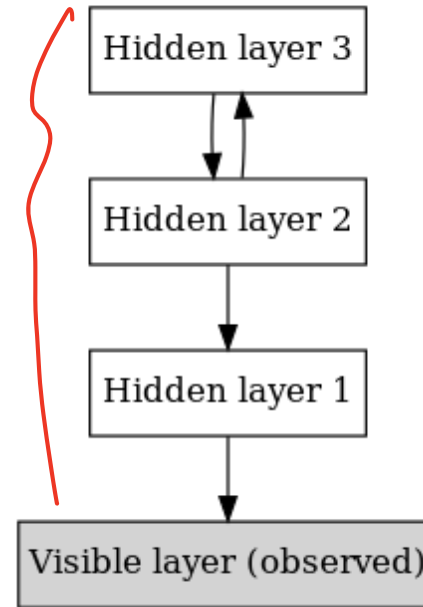
- No need to lift up the entire energy landscape!
 - Raising the neighborhood of desired patterns is sufficient



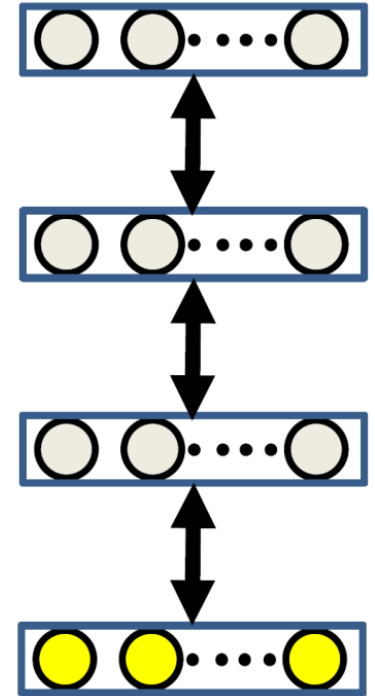
Deep Boltzmann Machine

- Can we have a **deep** version of RBM?
 - Deep Belief Net ('06)
 - Deep Boltzmann Machine ('09)
- Sampling?
 - Forward pass: bottom-up
 - Backward pass: top-down
- Deep Boltzmann Machine
 - The very first deep generative model
 - Salakhudinov & Hinton

fix h_1, h_2, h_3, V



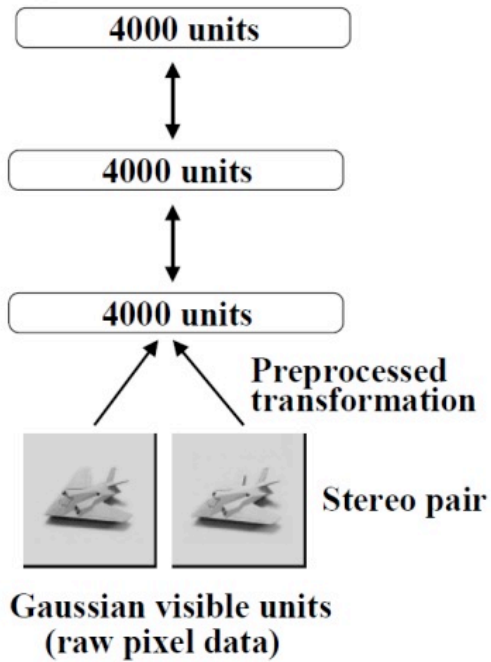
deep belief net



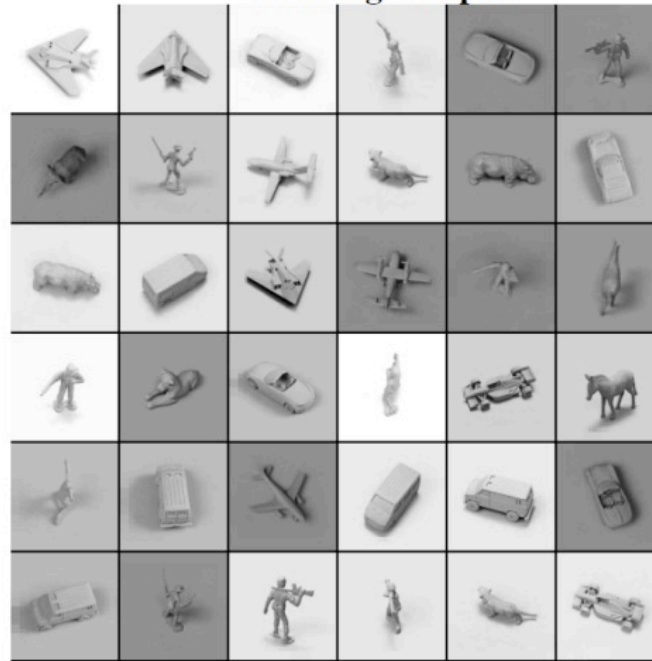
Deep Boltzmann Machine

Deep Boltzmann Machine

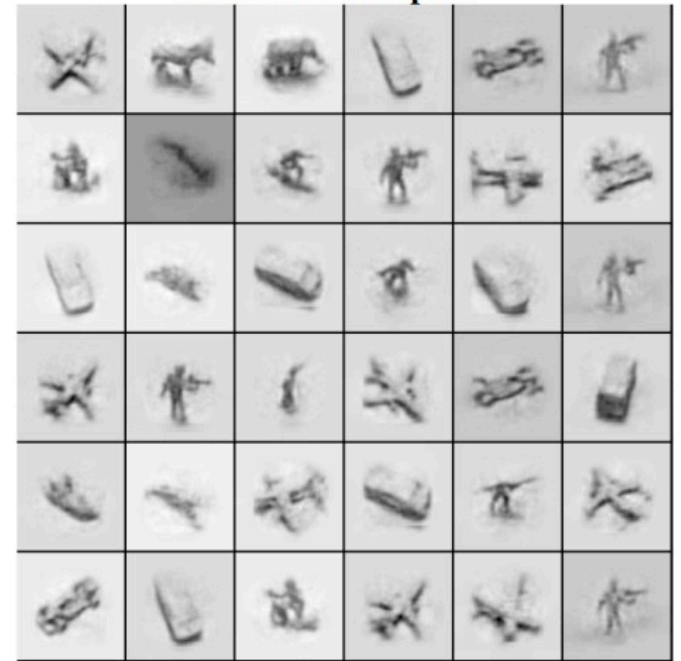
Deep Boltzmann Machine



Training Samples



Generated Samples



Summary

- Pros: powerful and flexible

- An arbitrarily complex density function $p(x) = \frac{1}{Z} \exp(-E(x))$

- Cons: hard to sample / train

- Hard to sample:

- MCMC sampling

- Partition function

- No closed-form calculation for likelihood
- Cannot optimize MLE loss exactly
- MCMC sampling



VAE or GAN

$z \sim N(\mu, \Sigma)$, $f(z)$

$\{x_1, x_2, \dots, x_T\} \approx p(x)$