

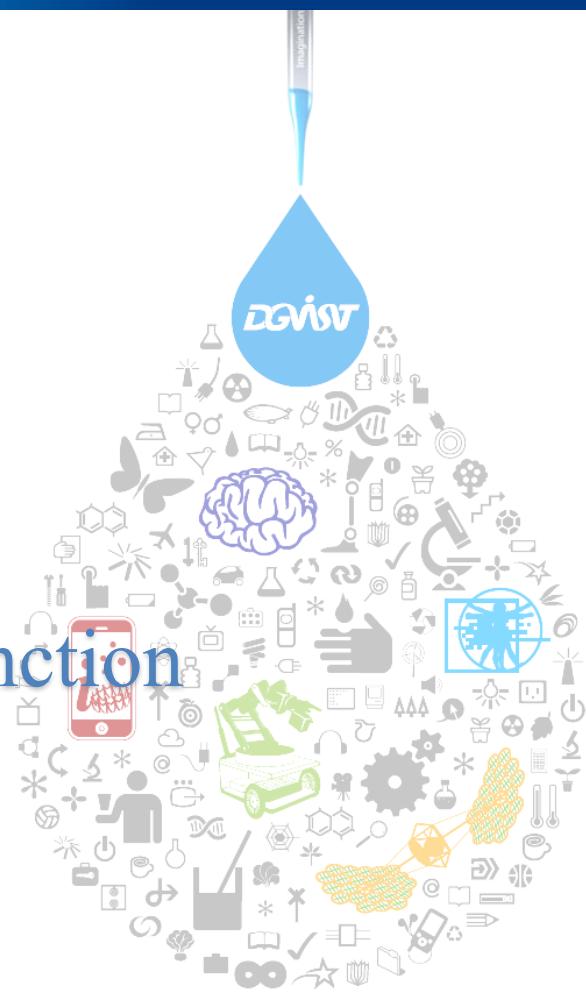


Deep Learning Seminar

CH 18. Confronting the Partition Function

2018-01-31

Eunjeong Yi



Chapter 18. Confronting the Partition Function

- 18.1 The Log-Likelihood Gradient**
- 18.2 Stochastic Maximum Likelihood and Contrastive Divergence**
- 18.3 Pseudolikelihood**
- 18.4 Score Matching and Ratio Matching**
- 18.5 Denoising Score Matching**
- 18.6 Noise-Contrastive Estimation**
- 18.7 Estimating the Partition Function**

Content

■ Undirected Models

■ Partition Function

■ Log-likelihood Gradient

■ Restricted Boltzmann Machine

■ Stochastic Maximum Likelihood

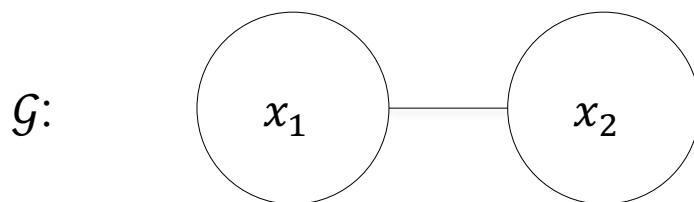
- Naïve MCMC for Maximizing Log-likelihood
- Contrastive Divergence

■ Estimating the Partition Function

- Importance Sampling

Undirected Models

- **Describe probability distribution of random variables which have no intrinsic direction or operate both directions**
 - there is no conditional probability distribution between random variables



- **Defined by an unnormalized probability distribution**

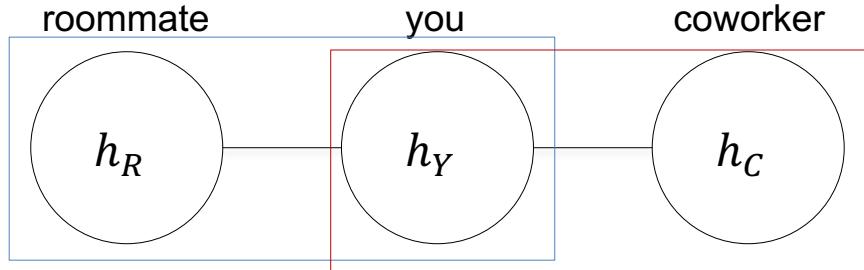
$$\tilde{p}(x) = \prod_{\mathcal{C} \in \mathcal{G}} \phi(\mathcal{C})$$

- a factor $\phi(\mathcal{C})$ measures the affinity of the variables in \mathcal{C} (non-negative)
 - \mathcal{G} : an undirected graph
 - \mathcal{C} : a clique in \mathcal{G}
 - ϕ : a factor (or clique potential)

Health State Example

■ Modeling the states of health of you, roommate, and coworker

➤ h_i : a state of health (1: good health, 0: with a cold)



$\phi(h_R, h_Y)$

c_1	$h_Y = 0$	$h_Y = 1$
$h_R = 0$	4	1
$h_R = 1$	1	10

$\phi(h_Y, h_C)$

c_2	$h_Y = 0$	$h_Y = 1$
$h_C = 0$	2	1
$h_C = 1$	1	10

■ Unnormalized probability distributions of health state

$$\tilde{p}(x) = \phi(h_R, h_Y)\phi(h_Y, h_C)$$

Partition Function

- Normalize probability distribution which is not guaranteed to sum or integrate to 1
- Partition function Z

$$Z = \int \tilde{p}(x) dx$$

or

$$= \Sigma_x \tilde{p}(x)$$

- Normalized probability distribution $p(\mathbf{x}; \theta)$

$$p(\mathbf{x}; \theta) = \frac{1}{Z(\theta)} \tilde{p}(\mathbf{x}; \theta)$$

\mathbf{x} : sampled data
 θ : parameter

Probabilities for Health State

■ $Z = \sum_{x \in \{0,1\}^3} \tilde{p}(x)$

h_R	h_Y	h_C	$\phi(\mathcal{C}_1)\phi(\mathcal{C}_2)$
0	0	0	8
0	0	1	4
0	1	0	1
0	1	1	10
1	0	0	2
1	0	1	1
1	1	0	10
1	1	1	100
Z			136

$$\phi(h_R, h_Y)$$

\mathcal{C}_1	$h_Y = 0$	$h_Y = 1$
$h_R = 0$	4	1
$h_R = 1$	1	10

$$\phi(h_Y, h_C)$$

\mathcal{C}_2	$h_Y = 0$	$h_Y = 1$
$h_C = 0$	2	1
$h_C = 1$	1	10

- All three people are in good health

$$p(1,1,1) = \frac{1}{136} \times 10 \times 10 = 0.7353$$

- Roommate and you are coughing but coworker is fine

$$p(0,0,1) = \frac{1}{136} \times 4 \times 1 = 0.0294$$

Log-likelihood Gradient

■ Normalized probability distribution $p(\mathbf{x}; \boldsymbol{\theta})$

$$p(\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \tilde{p}(\mathbf{x}; \boldsymbol{\theta})$$

■ Log-likelihood

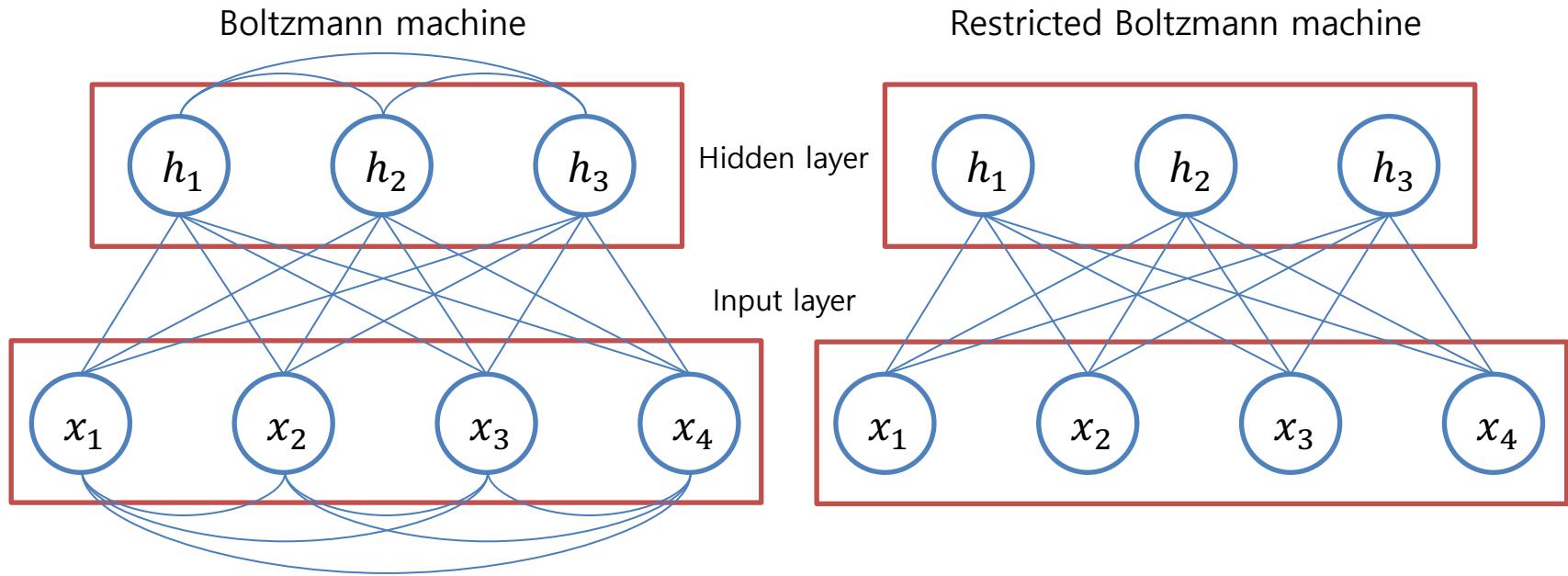
$$\nabla_{\boldsymbol{\theta}} \log p(\mathbf{x}; \boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x}; \boldsymbol{\theta}) - \nabla_{\boldsymbol{\theta}} \log Z(\boldsymbol{\theta})$$

■ Gradient of $\log Z(\boldsymbol{\theta})$

$$\nabla_{\boldsymbol{\theta}} \log Z(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x}; \boldsymbol{\theta})$$

Restricted Boltzmann Machine (RBM)

- Quintessential example of how graphical models are used for deep learning
- Connection between visible units and hidden units



Stochastic Maximum Likelihood

■ Markov chain Monte Carlo (MCMC) method

- The method for sampling from target probability distributions using Markov chains

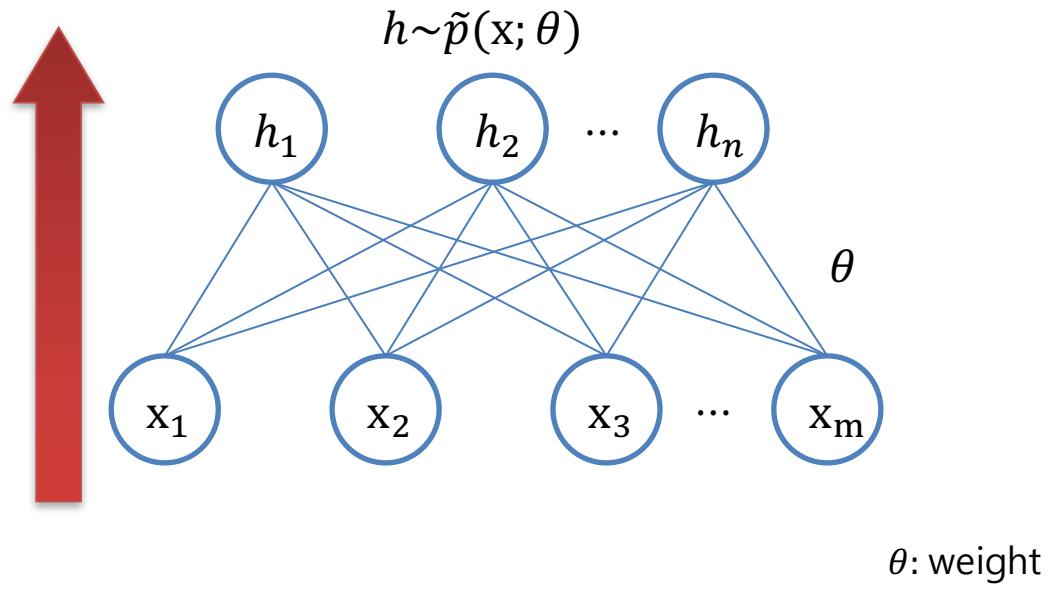
■ Algorithm

- 1) Sample $\{x^{(1)}, \dots, x^{(m)}\}$ from training set
- 2) Propagate sample to the hidden unit
- 3) Initialize $\{\tilde{x}^{(1)}, \dots, \tilde{x}^{(m)}\}$ to random value
- 4) Sample $\{\tilde{x}^{(1)}, \dots, \tilde{x}^{(m)}\}$ from model distribution
- 5) Propagate $\{\tilde{x}^{(1)}, \dots, \tilde{x}^{(m)}\}$

Naïve MCMC for Maximizing Log-likelihood

■ Forward step

- Sample $X = \{x^{(1)}, \dots, x^{(m)}\}$ from training set
- Model distribution $\tilde{p}(x; \theta)$ from hidden layer
- Calculate $\frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log \tilde{p}(x^{(i)}; \theta)$



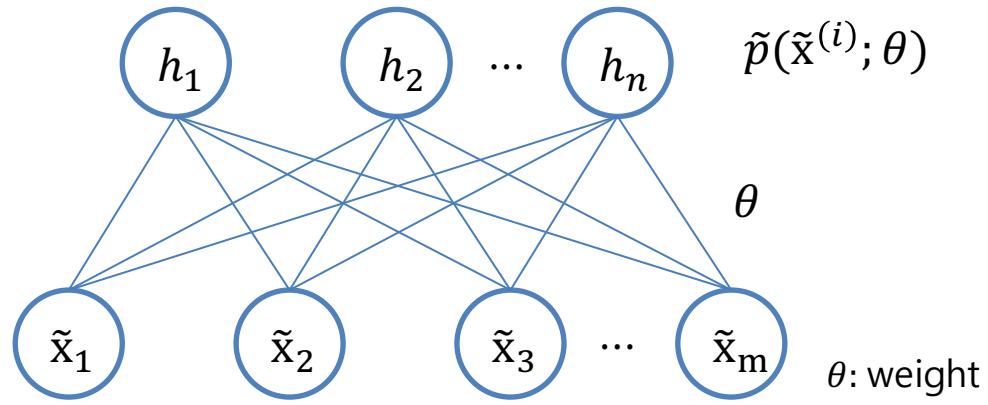
Naïve MCMC for Maximizing Log-likelihood

■ Backward step

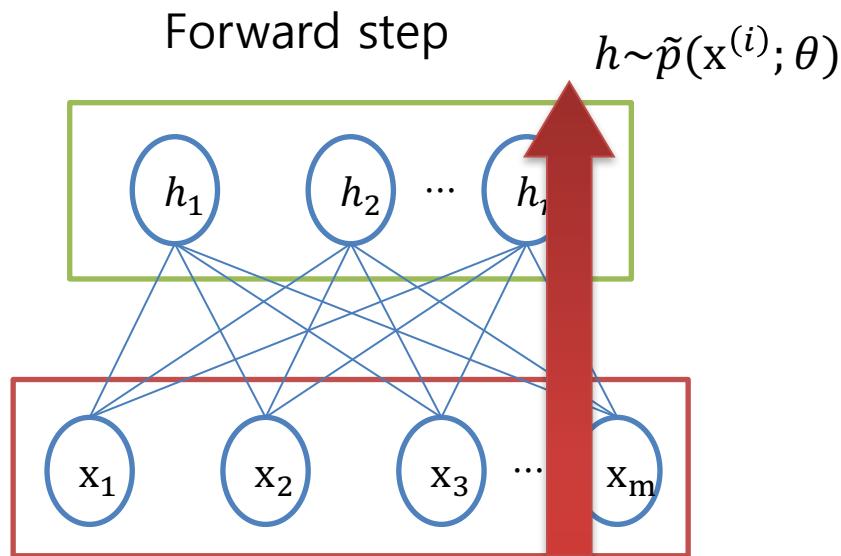
- $\tilde{X} = \{\tilde{x}^{(1)}, \dots, \tilde{x}^{(m)}\}$ initialized to random values
- \tilde{X} from $\tilde{p}(x^{(i)}; \theta)$ with Gibbs sampling
- Calculate $\frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log \tilde{p}(\tilde{x}^{(i)}; \theta)$
- Weight update

$$\theta \leftarrow \theta + \epsilon \left(\frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log \tilde{p}(x^{(i)}; \theta) - \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log \tilde{p}(\tilde{x}^{(i)}; \theta) \right)$$

■ iteration until converge $\tilde{p}(\tilde{x}^{(i)}; \theta)$ with $\tilde{p}(x^{(i)}; \theta)$

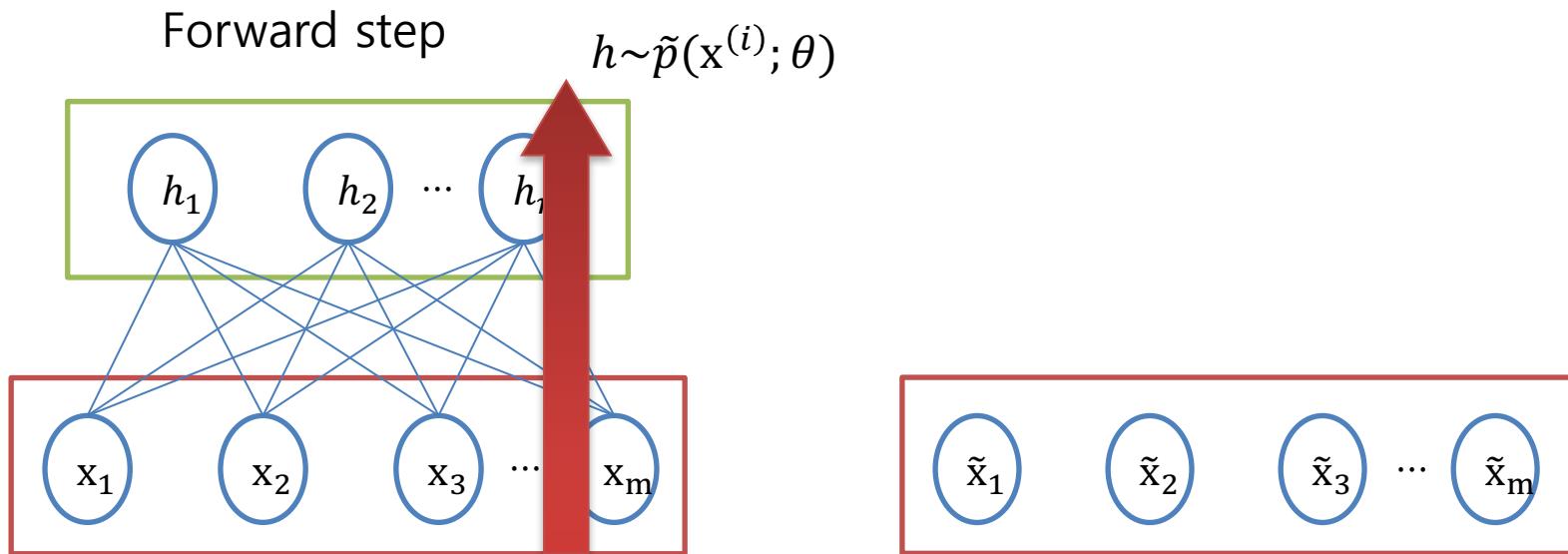


Naïve MCMC for Maximizing Log-likelihood



θ : weight

Naïve MCMC for Maximizing Log-likelihood



\tilde{x}_1 initialized to random values

$\tilde{X} = \{\tilde{x}^{(1)}, \dots, \tilde{x}^{(m)}\}$ from $\tilde{p}(x^{(i)}; \theta)$
Gibbs sampling

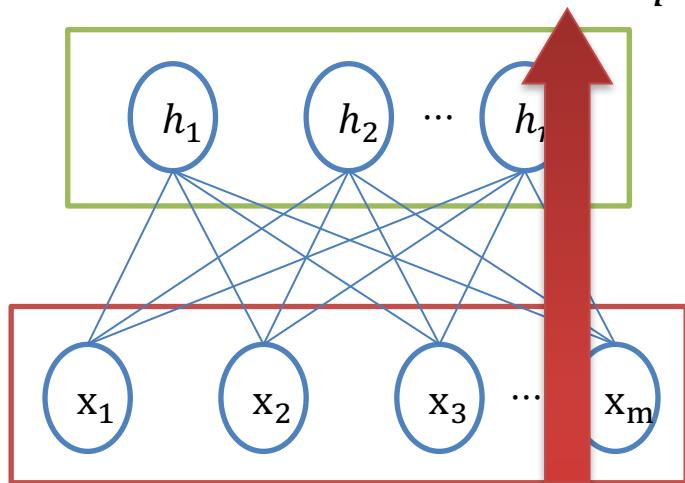
θ : weight

Naïve MCMC for Maximizing Log-likelihood

$$\theta \leftarrow \theta + \epsilon \left(\frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log \tilde{p}(x^{(i)}; \theta) - \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log \tilde{p}(\tilde{x}^{(i)}; \theta) \right)$$

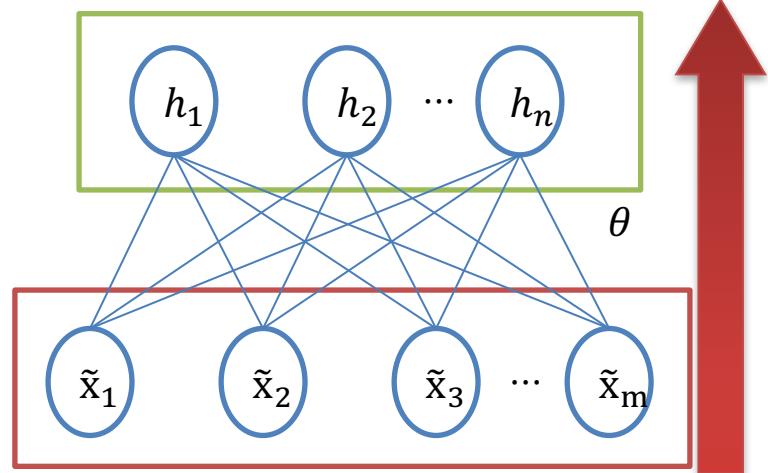
Forward step

$$h \sim \tilde{p}(x^{(i)}; \theta)$$



Backward step

$$\tilde{p}(\tilde{x}^{(i)}; \theta)$$



\tilde{x}_1 initialized to random values

$$\tilde{X} = \{\tilde{x}^{(1)}, \dots, \tilde{x}^{(m)}\} \text{ from } \tilde{p}(x^{(i)}; \theta)$$

Gibbs sampling

θ : weight

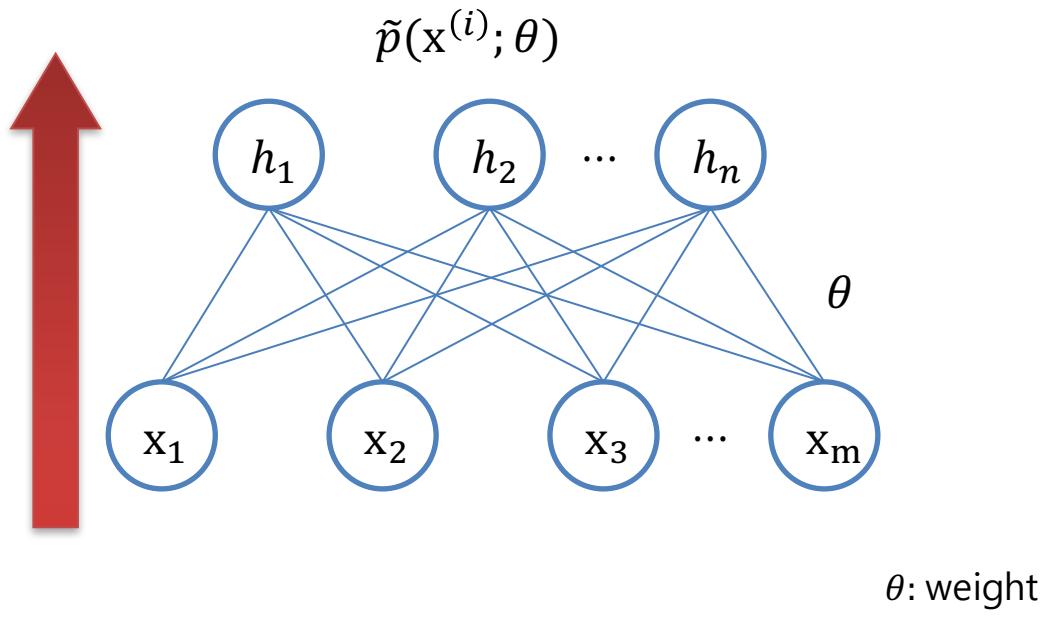
Contrastive Divergence

- Initialization of sample from the training set instead of the Markov chain
- Algorithm

- 1) Sample $\{x^{(1)}, \dots, x^{(m)}\}$ from training set
- 2) Propagate sample to the hidden unit
- 3) Sample $\{\tilde{x}^{(1)}, \dots, \tilde{x}^{(m)}\}$ from model distribution
- 4) Propagate $\{\tilde{x}^{(1)}, \dots, \tilde{x}^{(m)}\}$

■ Forward step

- Sample $X = \{x^{(1)}, \dots, x^{(m)}\}$ from training set
- Model distribution $\tilde{p}(x^{(i)}; \theta)$ from hidden layer
- Calculate $\frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log \tilde{p}(x^{(i)}; \theta)$

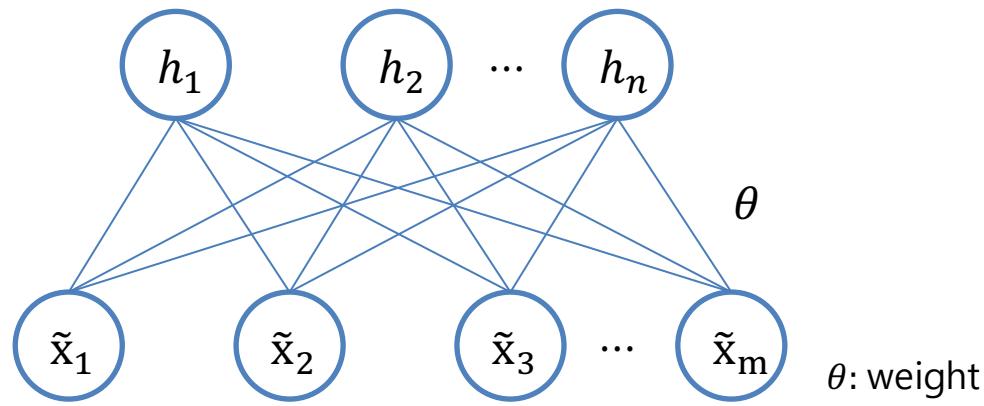


■ Backward step

- Sample $\tilde{X} = \{\tilde{x}^{(1)}, \dots, \tilde{x}^{(m)}\} = \{x^{(1)}, \dots, x^{(m)}\}$ with Gibbs sampling
- Calculate $\frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log \tilde{p}(\tilde{x}^{(i)}; \theta)$
- Weight update

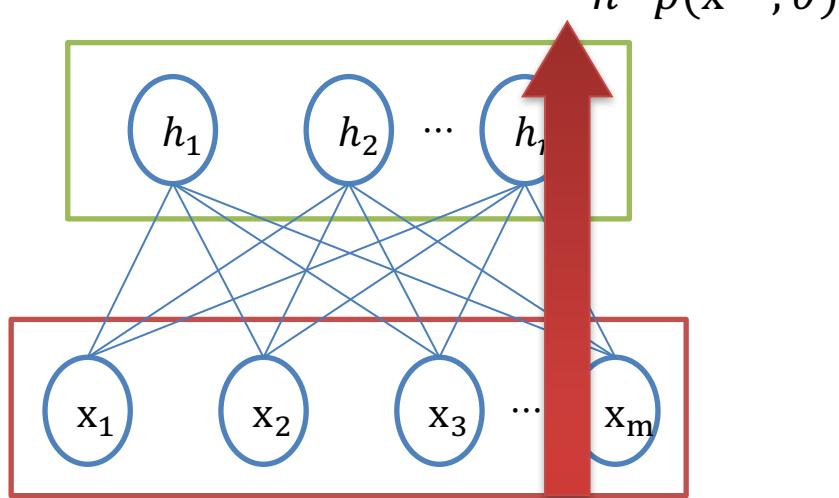
$$\theta \leftarrow \theta + \epsilon \left(\frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log \tilde{p}(x^{(i)}; \theta) - \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log \tilde{p}(\tilde{x}^{(i)}; \theta) \right)$$

■ iteration until converge $\tilde{p}(\tilde{x}^{(i)}; \theta)$ with $\tilde{p}(x^{(i)}; \theta)$



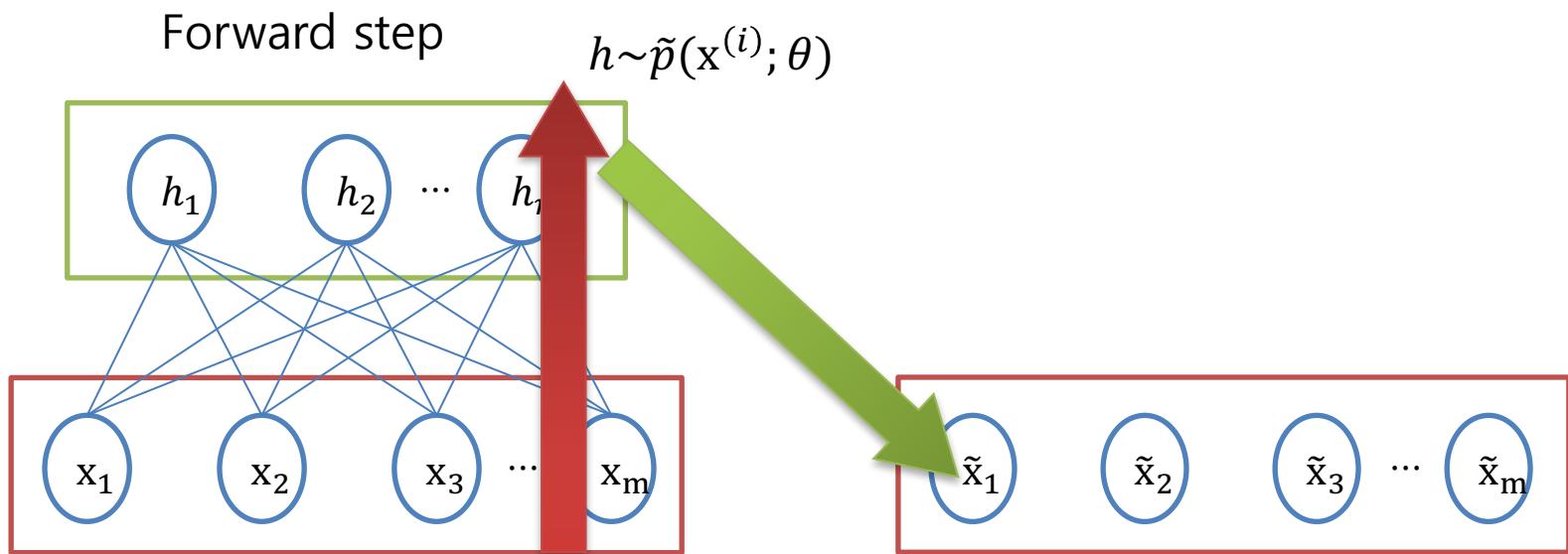
CD

Forward step



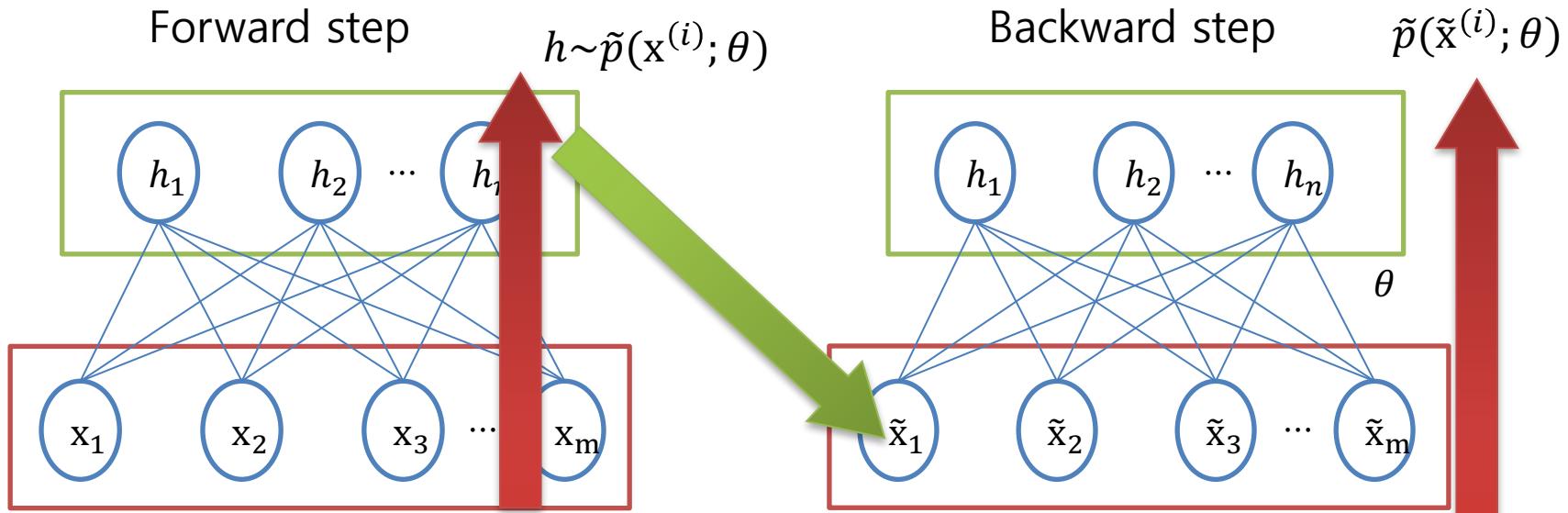
θ : weight

CD



CD

$$\theta \leftarrow \theta + \epsilon \left(\frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log \tilde{p}(x^{(i)}; \theta) - \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log \tilde{p}(\tilde{x}^{(i)}; \theta) \right)$$



$\tilde{X} = \{\tilde{x}^{(1)}, \dots, \tilde{x}^{(m)}\}$ from $\tilde{p}(x^{(i)}; \theta)$
Gibbs sampling

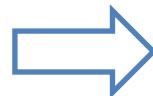
θ : weight

Importance sampling

- A variance reduction technique that can be used in the Monte Carlo method
- Estimating properties of a particular distribution using samples generated from a different distribution
- Input $x \sim p(x)$

$$\mathbb{E}_{x \sim p(x)} f(x) = \int p(x) f(x) dx$$

$$p(x)f(x) = q(x) \frac{p(x)f(x)}{q(x)}$$



$$\mathbb{E}_{x \sim p(x)} f(x) = \mathbb{E}_{x \sim q(x)} \frac{p(x)f(x)}{q(x)}$$

Importance sampling

- **Sample** $X = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}\}, (\mathbf{x}^{(i)} \sim q(\mathbf{x}))$
- **Importance sampling estimator** \hat{s}_q

$$\hat{s}_q = \frac{1}{n} \sum_{i=1}^n \frac{p(\mathbf{x}^{(i)})f(\mathbf{x}^{(i)})}{q(\mathbf{x}^{(i)})}$$

- **Minimum variance of** \hat{s}_q

$$Var(\hat{s}_q) = \frac{1}{n} Var\left(\frac{p(\mathbf{x})f(\mathbf{x})}{q(\mathbf{x})}\right)$$

Estimating the Partition Function

■ Estimate partition function $Z_B(\theta)$ of model M_B

- Use different probability distribution $p_A(x; \theta_A) = \frac{1}{Z_A(\theta)} \tilde{p}_A(x; \theta_A)$
- Estimate the ratio $r = \frac{Z(\theta_B)}{Z(\theta_A)}$ using importance sampling

■ Sample $X = \{x^{(1)}, \dots, x^{(K)}\}$ from $p_A(x; \theta_A)$

Model M_A

$$p_A(x; \theta_A) = \frac{1}{Z_A} \tilde{p}_A(x; \theta_A)$$

Model M_B

$$p_B(x; \theta_B) = \frac{1}{Z_B} \tilde{p}_B(x; \theta_B)$$

■ Comparing model M_A with M_B

$$\sum_i \log p_A(x^{(i)}; \theta_A) - \sum_i \log p_B(x^{(i)}; \theta_B) = \sum_i \log \frac{\tilde{p}_A(x^{(i)}; \theta_A)}{\tilde{p}_B(x^{(i)}; \theta_B)} - m \log \frac{Z(\theta_A)}{Z(\theta_B)}$$

- Estimate $\log \frac{Z(\theta_A)}{Z(\theta_B)}$ using importance sampling

Estimating the Partition Function

■ Partition function $Z_B(\theta)$

$$Z_B(\theta) = \int \tilde{p}_B(x; \theta) dx = \int \frac{p_A(x; \theta)}{\tilde{p}_A(x; \theta)} \tilde{p}_B(x; \theta) dx = Z_A(\theta) \int p_A(x; \theta) \frac{\tilde{p}_B(x; \theta)}{\tilde{p}_A(x; \theta)} dx$$

■ Ratio r

$$r = \frac{Z_B(\theta)}{Z_A(\theta)} = \int p_A(x; \theta) \frac{\tilde{p}_B(x; \theta)}{\tilde{p}_A(x; \theta)} dx$$

■ Ratio Estimator \hat{r} using sample $X = \{x^{(1)}, \dots, x^{(K)}\}$

$$\hat{r} = \frac{1}{K} \sum_{k=1}^K \frac{\tilde{p}_B(x^{(k)}; \theta)}{\tilde{p}_A(x^{(k)}; \theta)}, \quad (x^{(k)} \sim p_A(x; \theta))$$

Estimating the Partition Function

■ Estimate $Z_B(\theta)$

$$\hat{Z}_B(\theta) = \frac{Z_A(\theta)}{K} \sum_{k=1}^K \frac{\tilde{p}_B(x^{(k)}; \theta)}{\tilde{p}_A(x^{(k)}; \theta)}, \quad (x^{(k)} \sim p_A(x; \theta))$$

■ Minimize variance of $\hat{Z}_B(\theta)$

$$Var(\hat{Z}_B(\theta)) = \frac{Z_A(\theta)}{K^2} \sum_{k=1}^K \left(\frac{\tilde{p}_B(x^{(k)}; \theta)}{\tilde{p}_A(x^{(k)}; \theta)} - \hat{Z}_B \right)^2$$

Thank you