

November 24, 2017



Inter-University Research Institute Corporation /
Research Organization of Information and Systems

National Institute of Informatics

Boltzmann Machines

Data Mining 05 (データマイニング)

Mahito Sugiyama (杉山磨人)

Today's Outline

- Boltzmann Machines (Ising models):
A fundamental probabilistic model of deep learning
- Probabilistic models on Posets
 - Relationship with pattern mining
- Relationship to the deep architecture
 - DBM (Deep Boltzmann Machines)

Boltzmann Machines

- A **Boltzmann machine** (BM) is represented as an undirected graph $G = (V, E)$ with $V = \{1, 2, \dots, n\}$ and $E \subseteq \{\{i, j\} \mid i, j \in V\}$
- The **energy function** $\Phi: \{0, 1\}^n \rightarrow \mathbb{R}$ of a BM G is defined as

$$\Phi(\mathbf{x}; \boldsymbol{\theta}) = - \sum_{i \in V} \theta_i x_i - \sum_{\{i, j\} \in E} \theta_{ij} x_i x_j$$

- $\mathbf{x} = (x_1, x_2, \dots, x_n) \in \{0, 1\}^n$
- $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_n, \theta_{12}, \theta_{13}, \dots, \theta_{n-1n})$ is a **parameter vector** for vertices (bias) $\theta_1, \dots, \theta_n$ and edges (weight) $\theta_{12}, \dots, \theta_{n-1n}$
- $\theta_{ij} = 0$ if $\{i, j\} \notin E$

Boltzmann Distribution

- The probability $p(\mathbf{x}; \boldsymbol{\theta})$ is obtained for each $\mathbf{x} \in \{0, 1\}^n$ as

$$p(\mathbf{x}; \boldsymbol{\theta}) = \frac{\exp(-\Phi(\mathbf{x}, \boldsymbol{\theta}))}{Z(\boldsymbol{\theta})}$$

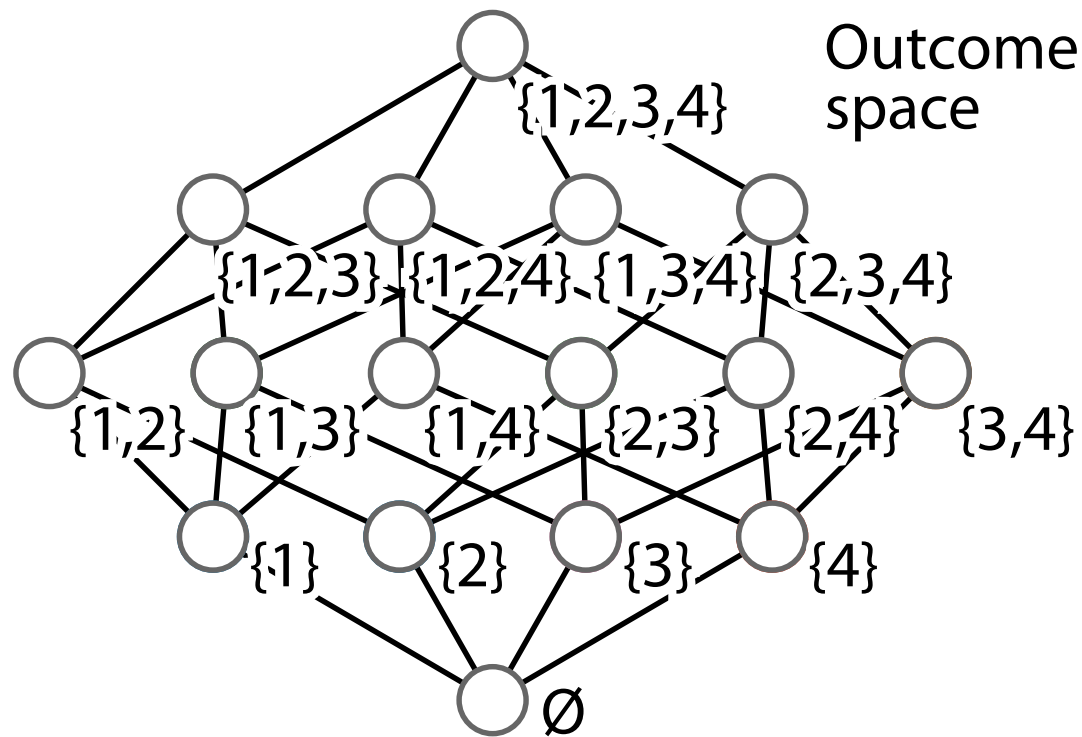
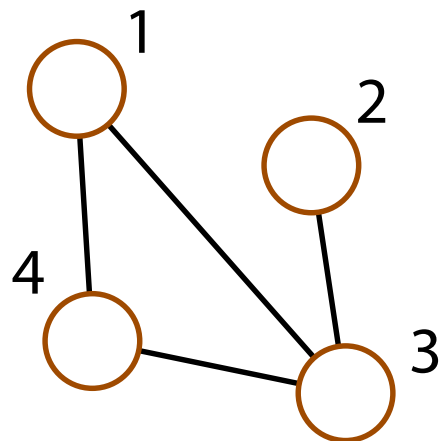
- $Z(\boldsymbol{\theta})$ is a **partition function** such that

$$Z(\boldsymbol{\theta}) = \sum_{\mathbf{x} \in \{0, 1\}^n} \exp(-\Phi(\mathbf{x}; \boldsymbol{\theta}))$$

to ensure the condition $\sum_{\mathbf{x} \in \{0, 1\}^n} p(\mathbf{x}) = 1$

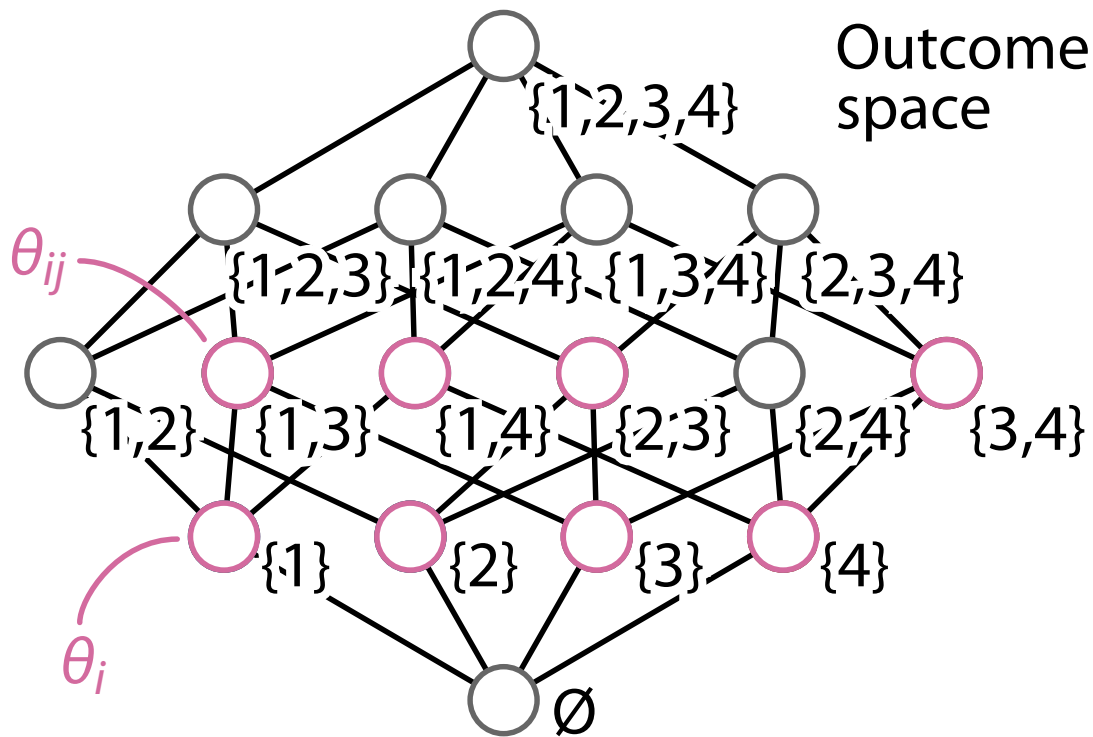
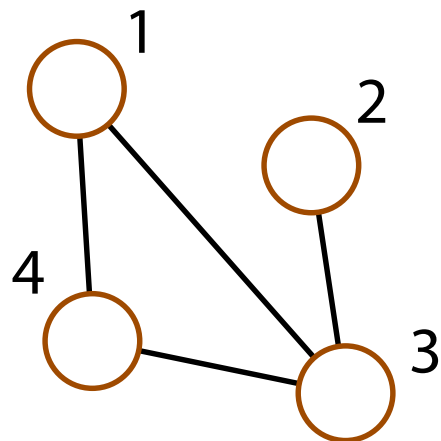
Outcome Space of BM

Boltzmann
Machine (BM)



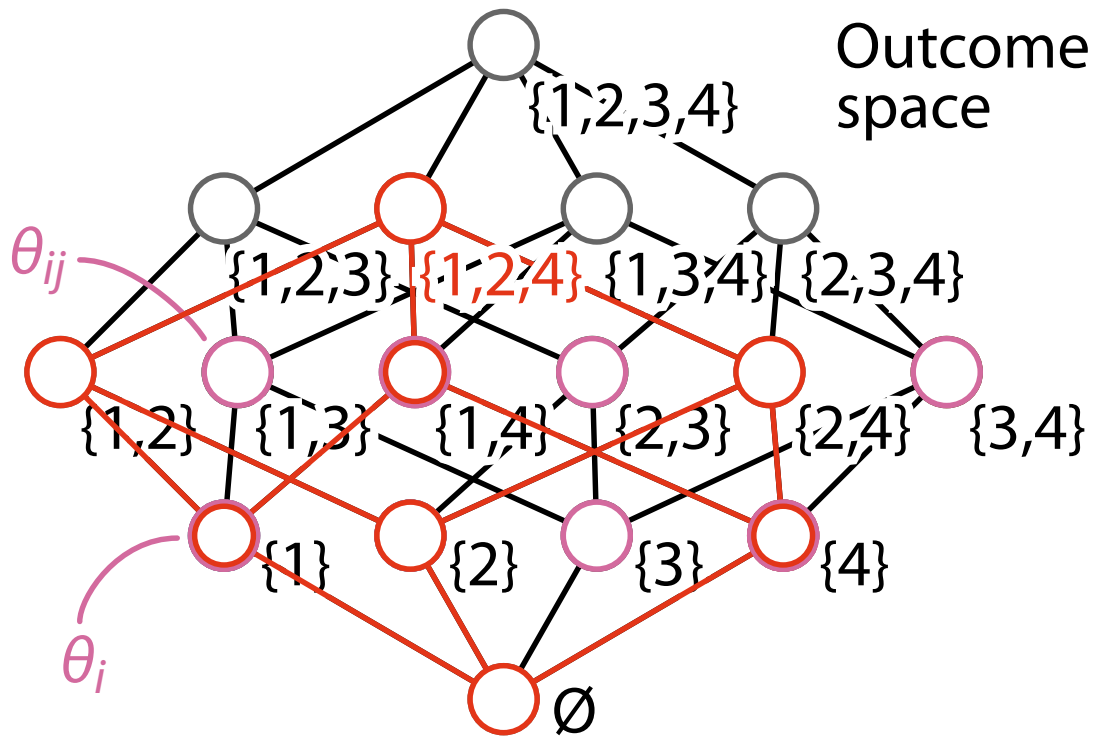
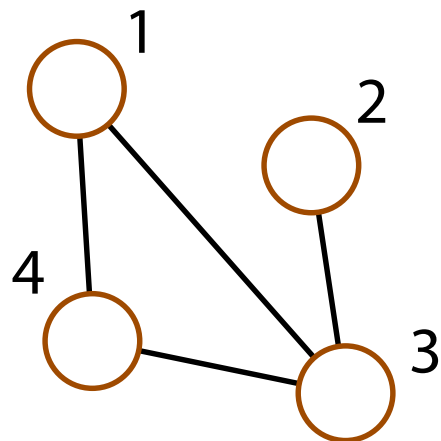
Parameters θ

Boltzmann Machine (BM)



Probability Computation

Boltzmann
Machine (BM)



Log-Linear Model

- BM is a special case of the **log-linear model** in statistics
- Let each set $x \in 2^V$ be the set of indices of “1” of $\mathbf{x} \in \{0, 1\}^n$
- The **parameter set** $B = \{x \in 2^V \setminus \{\emptyset\} \mid |x| = 1 \text{ or } x \in E\}$
- The Boltzmann distribution is described as

$$\log p(x) = \sum_{s \subseteq x} \theta(s) - \psi(\theta), \quad \psi(\theta) = -\theta(\emptyset) = \log Z(\boldsymbol{\theta})$$

- **Log** of the probability is obtained by the **linear** combination of coefficients θ
- $\theta(s) = 0$ if $s \notin B$

Learning of BM by MLE

- Given a dataset $D = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$, the objective of learning Boltzmann machines is to maximize the (log-)likelihood (Maximum Likelihood Estimation; MLE)

Find θ that maximizes $\prod_{i=1}^N p(\mathbf{x}_i; \theta) = p(\mathbf{x}_1; \theta) \cdot p(\mathbf{x}_2; \theta) \cdot \dots \cdot p(\mathbf{x}_N; \theta)$

- The probability of generating the given dataset by a BM

- The log-likelihood is usually treated

$$L_D(\theta) = \log \prod_{i=1}^N p(\mathbf{x}_i; \theta) = \sum_{i=1}^N \log p(\mathbf{x}_i; \theta)$$

Gradient of θ

- The log-likelihood is

$$L_D(\boldsymbol{\theta}) = \sum_{i=1}^N \log p(\mathbf{x}_i; \boldsymbol{\theta}) = \sum_{i=1}^N \sum_{s \subseteq \mathbf{x}_i} \theta(s) - N\psi(\boldsymbol{\theta})$$

- The **gradient** of $\theta(s)$ for each $s \in B$ is obtained as

$$\frac{\partial L_D(\boldsymbol{\theta})}{\partial \theta(s)} = |\{x_i \in D \mid x_i \supseteq s\}| - N \eta(s),$$

$$\eta(s) = \sum_{x \supseteq s} p(x) = \begin{cases} \mathbf{E}_{\boldsymbol{\theta}}[s^i] = \sum_{\mathbf{s}} s^i p(\mathbf{s}; \boldsymbol{\theta}) = \Pr(s^i = 1) \\ \mathbf{E}_{\boldsymbol{\theta}}[s^i s^j] = \sum_{\mathbf{s}} s^i s^j p(\mathbf{s}; \boldsymbol{\theta}) = \Pr(s^i = 1 \text{ and } s^j = 1) \end{cases}$$

Learning Equation of BM

- $L_D(\boldsymbol{\theta})$ is maximized when the gradient is zero \iff

$$\frac{1}{N} |\{x_i \in D \mid x_i \supseteq s\}| = \eta(s)$$

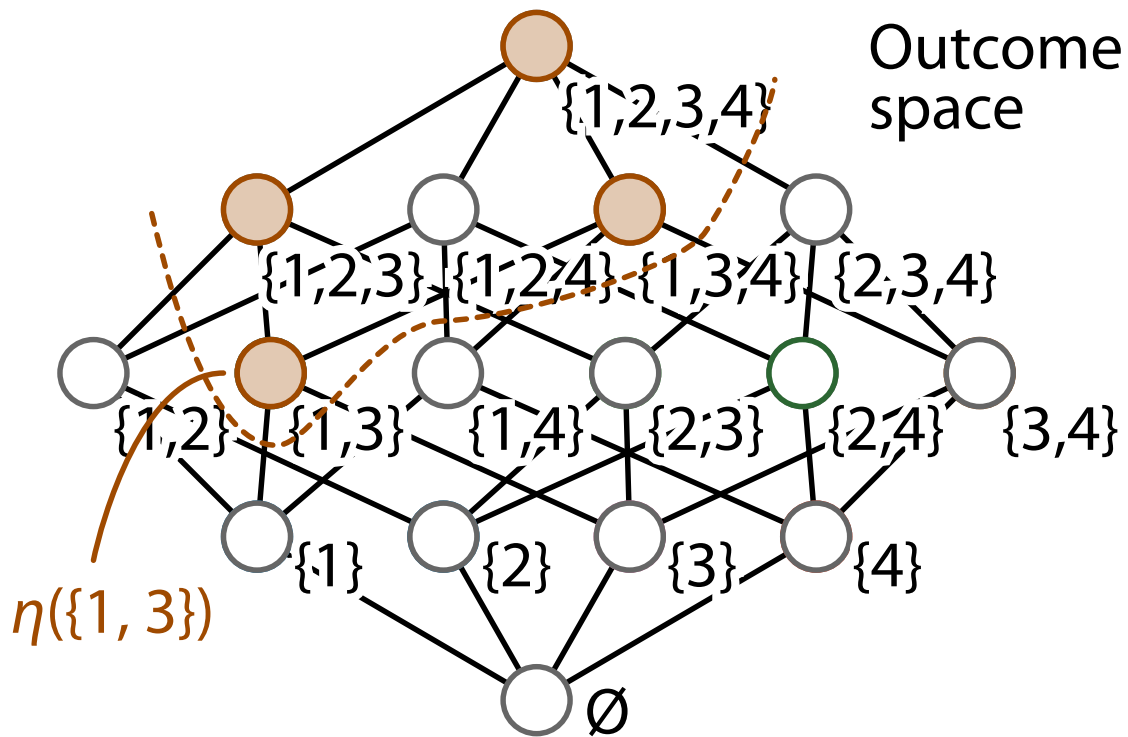
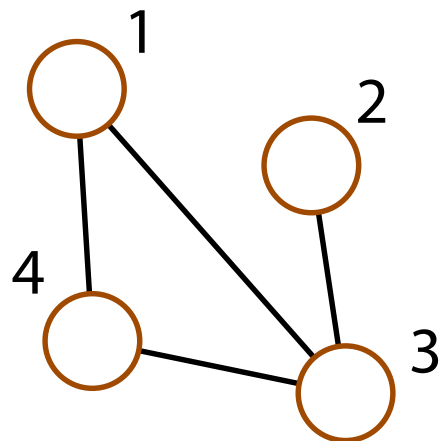
$$\hat{\eta}(s) = \eta(s)$$

for all $s \in B$

- This is known as **learning equation of BM**
- $\hat{\eta}(s)$ coincides with the **frequency** of a pattern s used in itemset mining

Frequency η

Boltzmann
Machine (BM)



KL Divergence Minimization

- Given two distributions P, Q , the KL divergence from P to Q :

$$D_{\text{KL}}(P, Q) = \sum_x p(x) \log \frac{p(x)}{q(x)}$$

- Given a dataset $D = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$, the **empirical distribution** \hat{P} is

$$\hat{p}(x) = \frac{1}{N} |\{x_i \in D \mid x_i = x\}|$$

- Maximizing the (log)likelihood is equivalent to minimizing the **Kullback–Leibler (KL) divergence**: $\min_{P \in \mathcal{S}(B)} D_{\text{KL}}(\hat{P}, P)$
 - $\mathcal{S}(B)$: the set of Boltzmann distributions

Optimization: Gradient Ascent

Algorithm 1: Learning of BM by gradient ascent

```
1 Initialize  $\theta$  with some values;  
2  $t \leftarrow 0$ ;  
3 repeat  
4   |   foreach  $s \in B$  do  
5   |   |    $\theta^{(t+1)}(s) \leftarrow \theta^{(t)}(s) + \varepsilon(\hat{\eta}(s) - \eta(s))$ ;  
6   |   |    $t \leftarrow t + 1$   
7 until  $\theta^{(t)} = \theta^{(t+1)}$ ;
```

Combinatorial Explosion

- The serious problem of learning BMs: **combinatorial explosion!!**

- The time complexity of computation of $\eta(x)$:

$$\eta(x) = \sum_{s \supseteq x} p(s),$$

is $2^{O(n)}$ and it is impossible to evaluate

- This is required to get the gradient $\hat{\eta}(x) - \eta(x)$
- Solution: approximate it by **Gibbs sampling**

Gibbs Sampling (1/2)

- A Markov chain Monte Carlo (MCMC) algorithm
- We can generate samples from the current Boltzmann distribution
 - n variables are dependent with each other
 - The partition function is not needed
- After obtaining an enough sample $S = \{s_1, s_2, \dots, s_M\}$ by Gibbs sampling, $\eta(s)$ can be approximated as

$$\eta(s) \approx \frac{1}{M} |\{s_i \in S \mid s_i \ni s\}|$$

Gibbs Sampling (2/2)

- For $\mathbf{x} = (x^1, x^2, \dots, x^n)$, the conditional probability of the i th variable being x^i with fixing others is

$$p_i = \frac{p(x^1, \dots, x^{i-1}, x^i, x^{i+1}, \dots, x^n)}{p(x^1, \dots, x^{i-1}, 0, x^{i+1}, \dots, x^n) + p(x^1, \dots, x^{i-1}, 1, x^{i+1}, \dots, x^n)}$$
$$= \frac{\exp(\lambda_i x_i)}{1 + \exp(\lambda_i)}$$

$$\lambda_i = \theta_i + \sum_{j \neq i} \theta_{ij} x^j$$

Algorithm 2: Gibbs Sampling

```
1 Initialize  $\mathbf{x}$  with some values;  
2 repeat  
3   foreach  $i \in \{1, 2, \dots, n\}$  do  
4     if  $p_i \geq$  a random value  $u \in [0, 1]$  then  
5        $x^i \leftarrow 1$   
6     else  
7        $x^i \leftarrow 0$   
8     Output  $\mathbf{x}$  and use it for the next initial vector  
9 until getting enough sample;
```

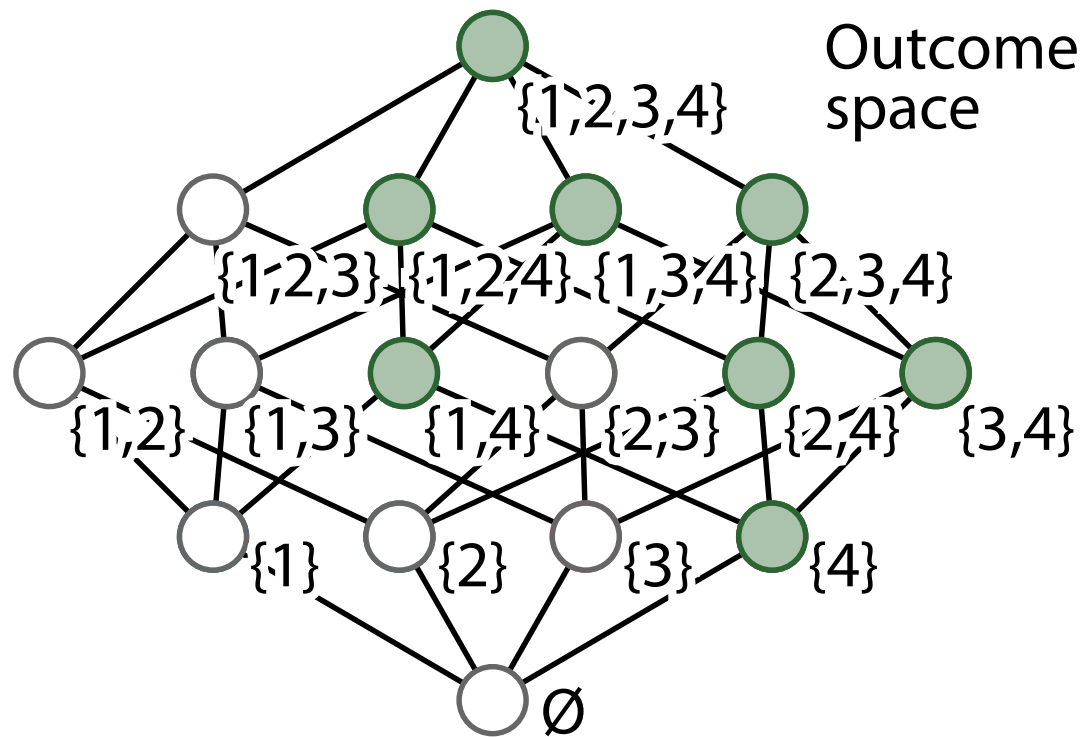
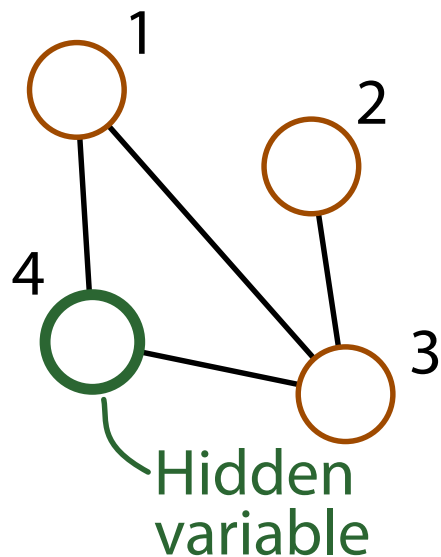
Introducing Hidden Variables

- To increase the representation power of the BM, we can introduce **hidden variables**
- When there are hidden nodes, the (log-)likelihood is maximized with respect to the distribution in which the hidden variables are marginalized out
- Let V and H be visible and hidden nodes
- The learning equation for $x = v \cup h$ with $v \subseteq V$ and $h \subseteq H$ is

$$\sum_{s \ni v, h} p(h | s) \hat{p}(s) = \eta(x)$$

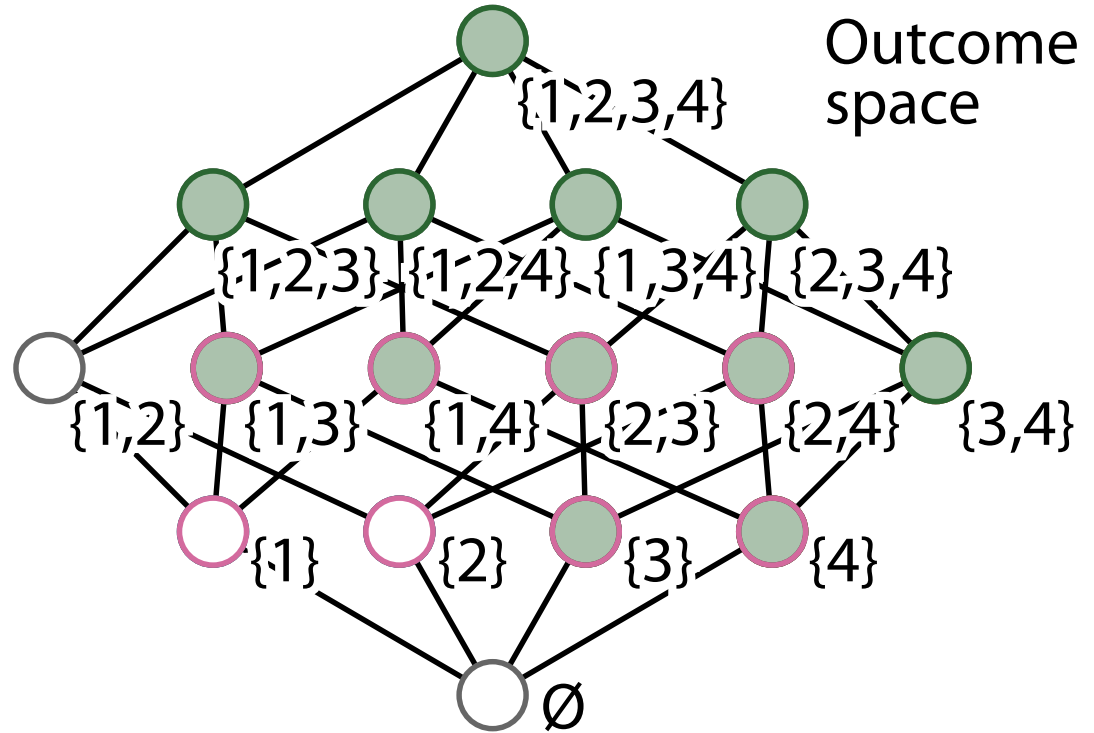
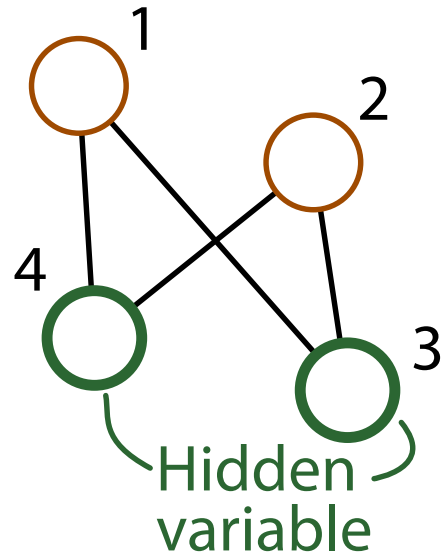
Outcome Space with Hidden Variable

Boltzmann Machine (BM)



Restricted Boltzmann Machines (RBMs)

Boltzmann Machine (BM)



Learning of RBMs

- Given a dataset $D = \{\mathbf{v}_{(1)}, \dots, \mathbf{v}_{(N)}\}$, learning equations in RBMs are

$$\hat{p}(v) = \eta(v)$$

$$\frac{1}{N} \sum_{\mu=1}^N \text{sig}(\lambda_{(\mu)}^j) = \eta(h)$$

$$\frac{1}{N} \sum_{\mu=1}^N v_{(\mu)}^i \text{sig}(\lambda_{(\mu)}^j) = \eta(v \cup h), \quad \text{where}$$

$$\text{sig}(\lambda_{(\mu)}^j) = \frac{\exp(\lambda_{(\mu)}^j)}{1 + \exp(\lambda_{(\mu)}^j)}, \quad \lambda_{(\mu)}^j = \theta_j + \sum_{i \in V} \theta_{ij} v_{(\mu)}^i$$

Deep Boltzmann Machines (DBMs)

