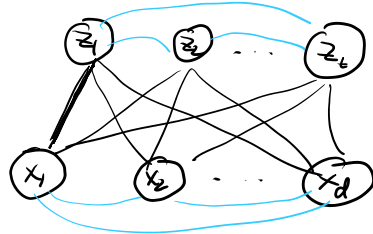


Recap: BM / RBM

parameterization: $P_w(s_1, s_2, \dots, s_m) \propto \exp(\vec{s}^T W \vec{s} - A(w))$

Gibbs sampling: $P_w(s_j | s_{-j}) = \text{sgm}(4s_j \langle W_{j,j}, s_j \rangle)$

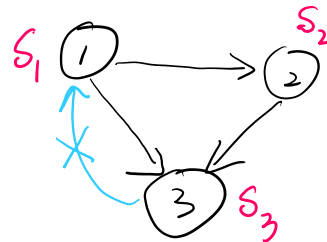
Graph: RBM



undirected graph
W is symmetric

PAG: directed acyclic graph

$pa(v) = \{u : \vec{uv} \in E\}$

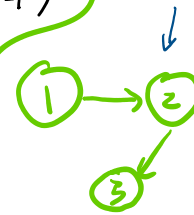


$pa(1) = \emptyset$
 $pa(2) = 1$
 $pa(3) = \{1, 2\}$

Belief Net (BN): $\vec{S} = (s_1, s_2, \dots, s_m)$

$P(s_1, s_2, s_3) = P(s_1) P(s_2 | s_1) P(s_3 | s_1, s_2)$

$P(s_1, s_2, \dots, s_m) = \prod_{j=1}^m P(s_j | pa(s_j))$



$P(s_1, s_2, s_3) = P(s_1) P(s_2 | s_1) P(s_3 | s_2)$

Thm: Fix a DAG. Given $\{P(s_j | pa(s_j))\}_{j=1}^m$, can construct

$P(s_1, s_2, \dots, s_m) = \prod_{j=1}^m P(s_j | pa(s_j))$

$\times \begin{cases} P(s_1) \\ P(s_2 | s_1) \\ P(s_3 | s_2) \end{cases}$

Ex: if the graph is not a DAG, then thm is not true.

Rem: $S \in \{\pm 1\}^m$, need $2^m - 1$ parameters to specify $P(s_1, s_2, \dots, s_m)$

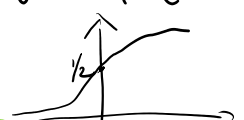
DAG with in-degree max K: $(2^K - 1)m$

Sigmoid BN (Neal '92)



$\text{Sgm}(t) = \frac{1}{1 + e^{-t}}$

specify $P_w(s_j | s_{<j}) = \text{sgm}(s_j \cdot (\sum_{k=1}^{j-1} W_{jk} s_k + b))$



specify $P_w(s_j | s_{<j}) = \text{sgm}(s_j \cdot (\sum_{k=1}^{j-1} W_{jk} s_k + b))$



$W \in \mathbb{R}^{m \times (m+1)}$ is lower-triangular

$W_{jk} = 0$ if $j \leq k \leq m$

$$\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & b_1 \\ x & 0 & 0 & 0 & 0 & b_2 \\ x & x & 0 & 0 & 0 & \vdots \\ x & x & x & 0 & 0 & \vdots \\ x & x & x & 0 & 0 & b_m \end{bmatrix}$$

$$= \text{sgm}(s_j \cdot W_{j:} \vec{s})$$

$$P_w(s_1, s_2, \dots, s_m) = \prod_{j=1}^m P_w(s_j | s_{<j})$$

$$= \prod_{j=1}^m \text{sgm}(s_j \cdot W_{j:} \vec{s})$$

Universality: $S = (X, Z)$, $P_w(x)$ can approx any discrete distr. for BN.

SBN: Max Likelihood

Given $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n \in \mathcal{P}(X)$, $\min_w \text{KL}(\hat{P}(x) \| P_w(x))$

$$\equiv \min_w -\frac{1}{n} \sum_{i=1}^n \log P_w(\vec{x}_i)$$

$$\frac{\partial}{\partial W_{jk}} = -\frac{1}{n} \sum_{i=1}^n \frac{\partial_{W_{jk}} P_w(\vec{x}_i)}{P_w(\vec{x}_i)} = -\frac{1}{n} \sum_{i=1}^n \frac{\partial_{W_{jk}} \sum_{\vec{z}_i} P_w(\vec{x}_i, \vec{z}_i)}{P_w(\vec{x}_i)}$$

$$= -\frac{1}{n} \sum_{i=1}^n \sum_{\vec{z}_i} \frac{\partial_{W_{jk}} P_w(\vec{x}_i, \vec{z}_i)}{P_w(\vec{x}_i)}$$

$$= -\frac{1}{n} \sum_{i=1}^n \sum_{\vec{z}_i} \frac{P_w(\vec{x}_i, \vec{z}_i)}{P_w(\vec{x}_i)} \cdot \frac{\partial_{W_{jk}} \log P_w(\vec{x}_i, \vec{z}_i)}{\log P_w(\vec{x}_i, \vec{z}_i)}$$

$$\prod_j \text{sgm}(s_j \cdot W_{j:} \vec{s})$$

$$P_w(\vec{z}_i | \vec{x}_i)$$

$$= -E_{\hat{P}(s)} \frac{\partial_{W_{jk}} \log \prod_j \text{sgm}(s_j \cdot W_{j:} \vec{s})}{\log \prod_j \text{sgm}(s_j \cdot W_{j:} \vec{s})}, \quad \hat{P}(s) = \hat{P}(\vec{x}, \vec{z}) = \hat{P}(\vec{x}) \cdot P_w(\vec{x} | \vec{z})$$

$$= -E_{\hat{P}_w} \frac{\partial_{W_{jk}} \log \text{sgm}(s_j \cdot W_{j:} \vec{s})}{\log \text{sgm}(s_j \cdot W_{j:} \vec{s})}$$

$$= -E_{\hat{P}_w} \frac{\text{sgm}'(s_j \cdot W_{j:} \vec{s})}{\text{sgm}(s_j \cdot W_{j:} \vec{s})} \cdot s_j \cdot s_k$$

$$\text{sgm}'(t) = \text{sgm}(t) \cdot \text{sgm}(-t)$$

$$= -E_{\hat{P}_w} \text{sgm}(-s_j W_j \vec{s}) s_j s_k$$

Gibbs sampling: sample x from training data

$$P(z_j = z | z_{-j}, x) = \text{sgm}(x W_z \vec{s}) \cdot \prod_{k > j} \text{sgm}(s_k [W_{kz} s + W_{kz}(z - s_j)])$$

$$\propto P(z_j = z, z_{-j}, x) = \prod_k P(s_k | \text{pa}(s_k))$$

Ex: Can we use EM?

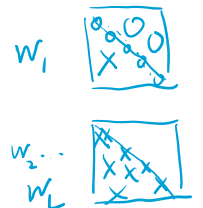
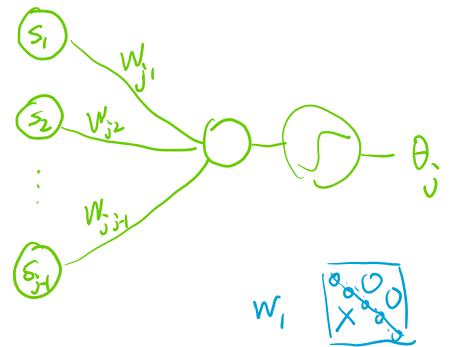
DBN (Bengio & Bengio '99)

specify m univariate densities $P(s_j | s_{<j}) = P_j(s_j; \theta_j(s_{<j}))$

for SBN: Bernoulli: for each $P_j(s_j = 1) = \theta_j(s_{<j})$ Gaussian

$$\theta_j(s_{<j}) = \text{sgm}(\sum_{k < j} W_{jk} s_k)$$

$$\text{DBN} \begin{cases} \vec{h}_0 = \vec{s} \\ \vec{h}_l = \sigma(W_l \vec{h}_{l-1}), l = 1, 2, \dots, L-1 \\ \vec{\theta} = \sigma(W_L \vec{h}_{L-1}) \end{cases}$$



θ_j can only depend on s_1, s_2, \dots, s_{j-1}

