



Restricted Boltzmann Machines

(Apr 5, 2017)

YANG Jiancheng



Outline

- **I. Boltzmann Machines**
- **II. Restricted Boltzmann Machines**
- **III. Learning: CD and PCD**



• I. Boltzmann Machines

• Neural Network as Graph Models

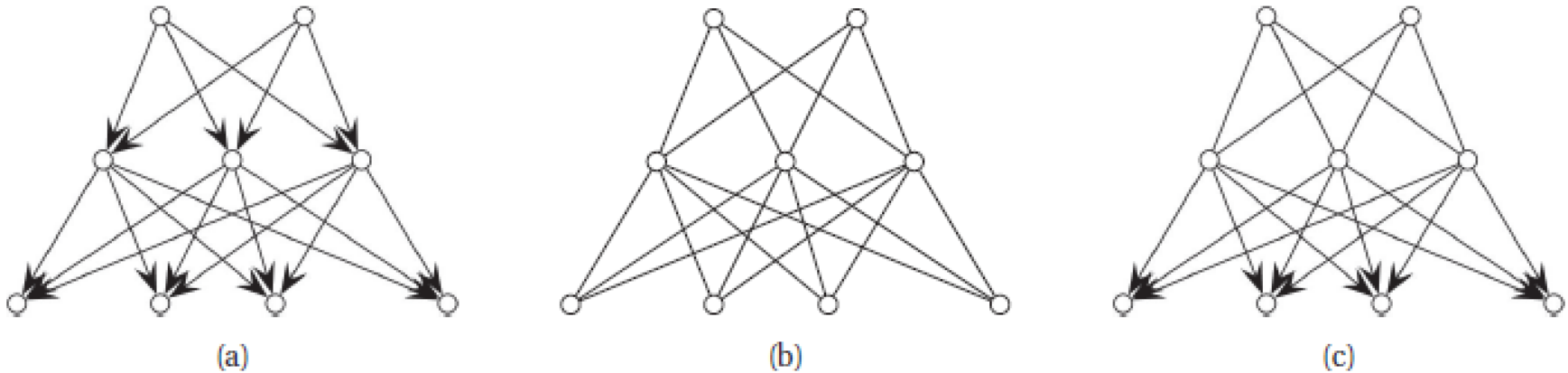


Figure 28.1 Some deep multi-layer graphical models. Observed variables are at the bottom. (a) A directed model. (b) An undirected model (deep Boltzmann machine). (c) A mixed directed-undirected model (deep belief net).

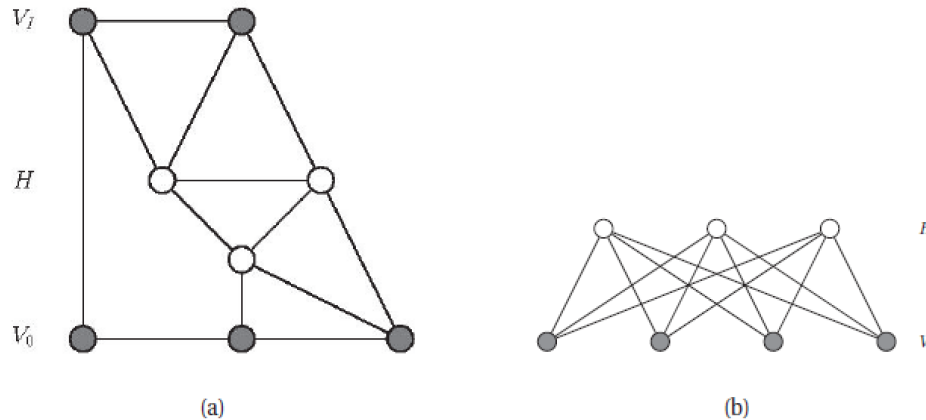
A common practice:

- Directed models are easy to train, but not easy to inference in parallel
- Undirected models are hard to train (need to inference when training), but easy to parallel



• I. Boltzmann Machines

- Boltzmann Machines to Restricted Boltzmann Machines



- a) Boltzmann Machines are arbitrary graphs with hidden and visible units, so-called “energy based” model
- b) Restricted versions are bipartite, i.e.
 - $h_i \perp h_j | \mathbf{v}$
 - $v_i \perp v_j | \mathbf{h}$



• II. Restricted Boltzmann Machines

- Binary RBMs

a) Most common

b) Binary hidden and binary visible units

c) Formula

$$p(\mathbf{v}, \mathbf{h} | \theta) = \frac{1}{Z(\theta)} \exp(-E(\mathbf{v}, \mathbf{h}; \theta))$$

$$\begin{aligned} E(\mathbf{v}, \mathbf{h}; \theta) &\triangleq -\sum_{r=1}^R \sum_{k=1}^K v_r h_k W_{rk} - \sum_{r=1}^R v_r b_r - \sum_{k=1}^K h_k c_k \\ &= -(\mathbf{v}^T \mathbf{W} \mathbf{h} + \mathbf{v}^T \mathbf{b} + \mathbf{h}^T \mathbf{c}) \end{aligned}$$

$$Z(\theta) = \sum_{\mathbf{v}} \sum_{\mathbf{h}} \exp(-E(\mathbf{v}, \mathbf{h}; \theta))$$

d) Inference

$$\mathbb{E}[\mathbf{h} | \mathbf{v}, \theta] = \text{sigm}(\mathbf{W}^T \mathbf{v})$$

$$\mathbb{E}[\mathbf{v} | \mathbf{h}, \theta] = \text{sigm}(\mathbf{W} \mathbf{h})$$

$$p(h|v, \theta) = \frac{p(h, v | \theta)}{p(v | \theta)} = \frac{e^{-v^T w h}}{\sum_h e^{-v^T w h}} = \frac{1}{1 + e^{-w^T v}}$$



• II. Restricted Boltzmann Machines

- Variants of RBMs

a) Categorical RBMs

$$E(\mathbf{v}, \mathbf{h}; \theta) \triangleq - \sum_{r=1}^R \sum_{k=1}^K \sum_{c=1}^C v_r^c h_k W_{rk}^c - \sum_{r=1}^R \sum_{c=1}^C v_r^c b_r^c - \sum_{k=1}^K h_k c_k$$

$$p(v_r | \mathbf{h}, \theta) = \text{Cat}(\mathcal{S}(\{b_r^c + \sum_k h_k W_{rk}^c\}_{c=1}^C)))$$

$$p(h_k = 1 | \mathbf{c}, \theta) = \text{sigm}(c_k + \sum_r \sum_c v_r^c W_{rk}^c)$$

$$\mathcal{S} = \text{softmax}$$

b) Gaussian RBMs

$$E(\mathbf{v}, \mathbf{h} | \theta) = - \sum_{r=1}^R \sum_{k=1}^K W_{rk} h_k v_r - \frac{1}{2} \sum_{r=1}^R (v_r - b_r)^2 - \sum_{k=1}^K a_k h_k$$

$$p(v_r | \mathbf{h}, \theta) = \mathcal{N}(v_r | b_r + \sum_k w_{rk} h_k, 1)$$

$$p(h_k = 1 | \mathbf{v}, \theta) = \text{sigm}\left(c_k + \sum_r w_{rk} v_r\right)$$



• III. Learning: CD and PCD

- Objective: Maxent models (1)

a) If all units are visible, then it's a Markov Random Field (MRF), in log-linear form:

$$p(v|\theta) = \frac{1}{Z(\theta)} \exp\left(\sum_c \theta_c^T \phi_c(v)\right)$$

$$\ell(\theta) = \frac{1}{N} \sum \log(p|\theta) = \frac{1}{N} \sum \left(\sum_c \theta_c^T \phi_c(v) - \log Z(\theta)\right)$$

$$\frac{\partial \ell}{\partial \theta_c} = \frac{1}{N} \sum_i \left[\phi_c(v_i) - \frac{\partial}{\partial \theta_c} \log Z(\theta)\right]$$

$$= \frac{1}{N} \sum_i \phi_c(v_i) - \mathbb{E}_v(\phi_c(v)|\theta)$$

$$= E_{emp} \phi_c(v) - E_{model} \phi_c(v)$$



- **III. Learning: CD and PCD**

- **Objective: Maxent models (2)**

Proof:

$$\begin{aligned}\frac{\partial}{\partial \theta_c} \log Z(\theta) &= \frac{1}{Z(\theta)} \sum_v \frac{\partial}{\partial \theta_c} \exp \left(\sum_c \theta_c^T \phi_c(v) \right) \\ &= \sum_v \phi_c(v) \frac{\exp(\sum_c \theta_c^T \phi_c(v))}{Z(\theta)} = \sum_v \phi_c(v) p(v|\theta) \\ &= \mathbb{E}_v(\phi_c(v)|\theta)\end{aligned}$$



• III. Learning: CD and PCD

- Objective: Partially observed maxent models

b) If some units are hidden:

$$p(v, h|\theta) = \frac{1}{Z(\theta)} \exp\left(\sum_c \theta_c^T \phi_c(v, h)\right)$$

$$\ell(\theta) = \frac{1}{N} \sum \log(p|\theta)$$

Similarly (but with some efforts),

$$\frac{\partial \ell}{\partial \theta_c} = \frac{1}{N} \sum_i \mathbb{E}_h \phi_c(v_i, h|\theta) - \mathbb{E}_{h,v} \phi_c(v, h|\theta)$$



• III. Learning: CD and PCD

- Objective: RBMs

c) RBMs are bipartite with visible and hidden units

$$\begin{aligned}\frac{\partial \ell}{\partial \theta_c} &= \frac{1}{N} \sum_i \mathbb{E}_h \phi_c(v_i, h | \theta) - \mathbb{E}_{h,v} \phi_c(v, h | \theta) \\ &= \frac{1}{N} \sum_i v_{c,i} \text{sigm}(v_i, \theta) - \mathbb{E}_{h,v} \phi_c(v, h | \theta)\end{aligned}$$

The first part is easy, but what about the second part?



• III. Learning: CD and PCD

• Block Gibbs Sampling

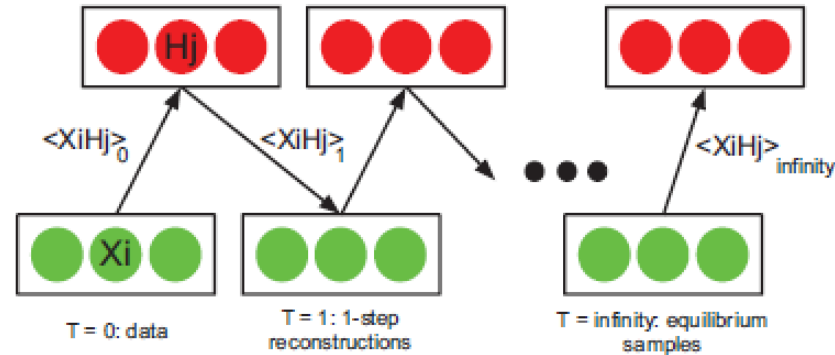


Figure 27.31 Illustration of Gibbs sampling in an RBM. The visible nodes are initialized at a datavector, then we sample a hidden vector, then another visible vector, etc. Eventually (at “infinity”) we will be producing samples from the joint distribution $p(\mathbf{v}, \mathbf{h}|\theta)$.

Key ideas of the second part: MCMC

• Sampling: fantasy data



• III. Learning: CD and PCD

- Contrastive divergence (CD)

Algorithm 27.3: CD-1 training for an RBM with binary hidden and visible units

```
1 Initialize weights  $\mathbf{W} \in \mathbb{R}^{R \times K}$  randomly;
2  $t := 0$ ;
3 for each epoch do
4      $t := t + 1$  ;
5     for each minibatch of size  $B$  do
6         Set minibatch gradient to zero,  $\mathbf{g} := \mathbf{0}$  ;
7         for each case  $\mathbf{v}_i$  in the minibatch do
8             Compute  $\mu_i = \mathbb{E}[\mathbf{h}|\mathbf{v}_i, \mathbf{W}]$ ;
9             Sample  $\mathbf{h}_i \sim p(\mathbf{h}|\mathbf{v}_i, \mathbf{W})$ ;
10            Sample  $\mathbf{v}'_i \sim p(\mathbf{v}|\mathbf{h}_i, \mathbf{W})$ ;
11            Compute  $\mu'_i = \mathbb{E}[\mathbf{h}|\mathbf{v}'_i, \mathbf{W}]$ ;
12            Compute gradient  $\nabla_{\mathbf{W}} = (\mathbf{v}_i)(\mu_i)^T - (\mathbf{v}'_i)(\mu'_i)^T$  ;
13            Accumulate  $\mathbf{g} := \mathbf{g} + \nabla_{\mathbf{W}}$ ;
14         Update parameters  $\mathbf{W} := \mathbf{W} + (\alpha_t/B)\mathbf{g}$ 
```



• III. Learning: CD and PCD

- Persistent CD (PCD)

Algorithm 27.4: Persistent CD for training an RBM with binary hidden and visible units

```
1 Initialize weights  $\mathbf{W} \in \mathbb{R}^{D \times L}$  randomly;
2 Initialize chains  $(\mathbf{v}_s, \mathbf{h}_s)_{s=1}^S$  randomly ;
3 for  $t = 1, 2, \dots$  do
4     // Mean field updates ;
5     for each case  $i = 1 : N$  do
6          $\mu_{ik} = \text{sigm}(\mathbf{v}_i^T \mathbf{w}_{:,k})$ 
7     // MCMC updates ;
8     for each sample  $s = 1 : S$  do
9         Generate  $(\mathbf{v}_s, \mathbf{h}_s)$  by brief Gibbs sampling from old  $(\mathbf{v}_s, \mathbf{h}_s)$ 
10    // Parameter updates ;
11     $\mathbf{g} = \frac{1}{N} \sum_{i=1}^N \mathbf{v}_i (\mu_i)^T - \frac{1}{S} \sum_{s=1}^S \mathbf{v}_s (\mathbf{h}_s)^T$  ;
12     $\mathbf{W} := \mathbf{W} + \alpha_t \mathbf{g}$ ;
13    Decrease  $\alpha_t$ 
```



Bibliography

- **MLaPP: Chapter 27.7 Restricted Boltzman machines (RBMs)**
- **Neural Networks for Machine Learning by Geoffrey Hinton (Coursera): Week 12,13.14**



Thanks for listening!

