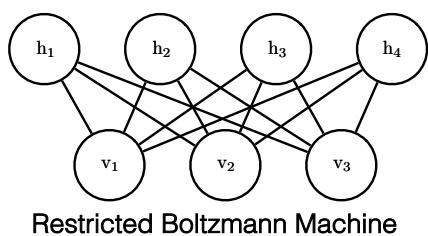


# Learning in Probabilistic Models: RBMs, DBMs, DBNs (chapters 18, 20)

University of Padova, A.A. 2022/23

## Restricted Boltzmann Machines: recap



### Energy-based Undirected Model

$$P(\mathbf{v} = \mathbf{v}, \mathbf{h} = \mathbf{h}) = \frac{1}{Z} \exp(-E(\mathbf{v}, \mathbf{h}))$$

$$E(\mathbf{v}, \mathbf{h}) = -\mathbf{b}^\top \mathbf{v} - \mathbf{c}^\top \mathbf{h} - \mathbf{v}^\top \mathbf{W} \mathbf{h}$$

hidden units:  
not observable variables

visible units:  
observable variables

partition function

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

Binary visible and hidden variables

$p(\mathbf{v})$  is intractable, however both  
 $p(\mathbf{h}|\mathbf{v})$  and  $p(\mathbf{v}|\mathbf{h})$   
are *factorial* and easy to compute!



$$p(\mathbf{h}|\mathbf{v}) = \prod_{j=1}^{n_h} p(h_j|\mathbf{v})$$

$$p(h_j = 1|\mathbf{v}) = \sigma(c_j + \mathbf{v}^\top \mathbf{W}_{:,j})$$

# Training Restricted Boltzmann Machines

Now that we know how to perform Gibbs sampling, we are ready to understand how a RBM can be trained!

- Given training data  $\mathcal{T} = \{\mathbf{x}^{(i)}\}_{i=1}^n$  (to be clamped to the visible units) we want to find parameters  $\theta = \{\mathbf{W}, \mathbf{b}, \mathbf{c}\}$  that maximize  $p(\mathcal{T}|\theta)$
- We know it is easier to maximize the Log-Likelihood

$$\log p(\mathcal{T}|\theta) = \log \frac{1}{Z(\theta)} \tilde{p}(\mathcal{T}|\theta) = \log \tilde{p}(\mathcal{T}|\theta) - \log Z(\theta),$$

where

$$\log \tilde{p}(\mathcal{T}|\theta) = \log \prod_{i=1}^n \tilde{p}(\mathbf{x}^{(i)}|\theta) = \sum_{i=1}^n \underbrace{\log \tilde{p}(\mathbf{x}^{(i)}|\theta)}_{i.i.d. \text{ samples}}$$

- We proceed by gradient ascent, so we need to compute

$$\begin{aligned} \nabla_{\theta} \log p(\mathcal{T}|\theta) &= \nabla_{\theta} \log \tilde{p}(\mathcal{T}|\theta) - \nabla_{\theta} \log Z(\theta) \\ &= \sum_{i=1}^n \nabla_{\theta} \log \tilde{p}(\mathbf{x}^{(i)}|\theta) - \nabla_{\theta} \log Z(\theta) \end{aligned}$$

# Training Restricted Boltzmann Machines

- The term  $\sum_{i=1}^n \nabla_{\theta} \log \tilde{p}(\mathbf{x}^{(i)}|\theta)$  is not difficult to deal with

$$\begin{aligned} \sum_{i=1}^n \nabla_{\theta} \log \tilde{p}(\mathbf{x}^{(i)}|\theta) &= \sum_{i=1}^n \frac{1}{\tilde{p}(\mathbf{x}^{(i)}|\theta)} \nabla_{\theta} \tilde{p}(\mathbf{x}^{(i)}|\theta) \\ &= (\mathbf{x}^{(i)} \text{ is clamped to } \mathbf{v}^{(i)}, \text{ i.e. } \mathbf{x}^{(i)} = \mathbf{v}^{(i)}) \\ &= \sum_{i=1}^n \frac{1}{\sum_{\mathbf{h}} \exp(-E(\mathbf{v}^{(i)}, \mathbf{h}))} \nabla_{\theta} \sum_{\mathbf{h}} \exp(-E(\mathbf{v}^{(i)}, \mathbf{h})) \\ &= - \sum_{i=1}^n \sum_{\mathbf{h}} \underbrace{\frac{\exp(-E(\mathbf{v}^{(i)}, \mathbf{h}))}{\sum_{\mathbf{h}} \exp(-E(\mathbf{v}^{(i)}, \mathbf{h}))}}_{= p(\mathbf{h}|\mathbf{v}^{(i)})} \nabla_{\theta} E(\mathbf{v}^{(i)}, \mathbf{h}) \\ &= \boxed{\frac{\exp(-E(\mathbf{v}^{(i)}, \mathbf{h}))}{\sum_{\mathbf{h}} \exp(-E(\mathbf{v}^{(i)}, \mathbf{h}))} = \frac{\tilde{p}(\mathbf{h}, \mathbf{v}^{(i)})}{\tilde{p}(\mathbf{v}^{(i)})} \frac{Z}{Z}} \\ &= \boxed{\frac{p(\mathbf{h}, \mathbf{v}^{(i)})}{p(\mathbf{v}^{(i)})} = p(\mathbf{h}|\mathbf{v}^{(i)})} \\ &= - \sum_{i=1}^n \sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{v}^{(i)}) \nabla_{\theta} E(\mathbf{v}^{(i)}, \mathbf{h}) \end{aligned}$$

# Training Restricted Boltzmann Machines

Knowing that  $p(\mathbf{h}|\mathbf{v}^{(i)}) = \prod_{j=1}^{n_h} p(h_j|\mathbf{v}^{(i)})$ , we have (recall that  $\mathbf{x}^{(i)} = \mathbf{v}^{(i)}$ )

$$\sum_{i=1}^n \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x}^{(i)}|\boldsymbol{\theta}) = - \sum_{i=1}^n \sum_{\mathbf{h}} \prod_{j=1}^{n_h} p(h_j|\mathbf{v}^{(i)}) \nabla_{\boldsymbol{\theta}} E(\mathbf{v}^{(i)}, \mathbf{h})$$

Recalling that  $\forall j, h_j \in \{0, 1\}$ ,  $\nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x}^{(i)}|\boldsymbol{\theta}) = \nabla_{\mathbf{W}, \mathbf{b}, \mathbf{c}} \log \tilde{p}(\mathbf{x}^{(i)}|\boldsymbol{\theta})$ , if we now focus on  $\nabla_{\mathbf{W}}$ , we have for any  $q$  and  $z$  (notice that  $\frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial w_{qz}} = \frac{\partial(-\mathbf{b}^\top \mathbf{v} - \mathbf{c}^\top \mathbf{h} - \mathbf{v}^\top \mathbf{W} \mathbf{h})}{\partial w_{qz}} = -v_q h_z$ )

$$\begin{aligned} \sum_{i=1}^n \frac{\partial \log \tilde{p}(\mathbf{x}^{(i)}|\boldsymbol{\theta})}{\partial w_{qz}} &= - \sum_{i=1}^n \sum_{\mathbf{h}} \prod_{j=1}^{n_h} p(h_j|\mathbf{v}^{(i)}) \frac{\partial E(\mathbf{v}^{(i)}, \mathbf{h})}{\partial w_{qz}} \\ &= \sum_{i=1}^n \sum_{h_1, \dots, h_{n_h}} \prod_{j=1}^{n_h} p(h_j|\mathbf{v}^{(i)}) v_q^{(i)} h_z \\ &= \sum_{i=1}^n \sum_{h_z \in \{0, 1\}} p(h_z|\mathbf{v}^{(i)}) v_q^{(i)} h_z \underbrace{\sum_{\substack{h_1, \dots, \\ h_{z-1}, h_{z+1}, \\ \dots, h_{n_h}}} \prod_{\substack{j=1 \\ j \neq z}}^{n_h} p(h_j|\mathbf{v}^{(i)})}_{= 1} \end{aligned}$$

# Training Restricted Boltzmann Machines

$$\begin{aligned} \sum_{i=1}^n \frac{\partial \log \tilde{p}(\mathbf{x}^{(i)}|\boldsymbol{\theta})}{\partial w_{qz}} &= \sum_{i=1}^n \sum_{h_z \in \{0, 1\}} p(h_z|\mathbf{v}^{(i)}) v_q^{(i)} h_z \\ &= \sum_{i=1}^n p(h_z = 1|\mathbf{v}^{(i)}) v_q^{(i)} \end{aligned}$$

and recalling that  $p(h_z = 1|\mathbf{v}) = \sigma(c_z + \mathbf{v}^\top \mathbf{W}_{:,z})$

$$\sum_{i=1}^n \frac{\partial \log \tilde{p}(\mathbf{x}^{(i)}|\boldsymbol{\theta})}{\partial w_{qz}} = \sum_{i=1}^n \sigma(c_z + \mathbf{v}^{(i)\top} \mathbf{W}_{:,z}) v_q^{(i)}$$

# Training Restricted Boltzmann Machines

- As we already discussed, the term  $\nabla_{\theta} \log Z(\theta)$  is difficult to deal with. In fact, recalling that  $Z(\theta) = \sum_{v,h} \exp(-E(v, h))$ , we have

$$\begin{aligned}
 \nabla_{\theta} \log Z(\theta) &= \frac{1}{\sum_{v,h} \exp(-E(v, h))} \nabla_{\theta} \sum_{v,h} \exp(-E(v, h)) \\
 &= - \sum_{v,h} \underbrace{\frac{\exp(-E(v, h))}{\sum_{v,h} \exp(-E(v, h))}}_{= p(v, h)} \nabla_{\theta} E(v, h) \\
 &= - \sum_{v,h} p(v, h) \nabla_{\theta} E(v, h) \\
 &= - \sum_{v,h} \underbrace{p(v)p(h|v)}_{\text{product rule}} \nabla_{\theta} E(v, h) \quad (\text{alternative: } p(h)p(v|h)) \\
 &= - \sum_v p(v) \sum_h p(h|v) \nabla_{\theta} E(v, h)
 \end{aligned}$$

A.A. 2022/23: Deep Learning

# Training Restricted Boltzmann Machines

If we now focus on  $\nabla_{W_q}$ , we have for any  $q$  and  $z$

$$\begin{aligned}
 \frac{\partial \log Z(\theta)}{\partial w_{qz}} &= - \sum_v p(v) \underbrace{\sum_h p(h|v) \frac{\partial E(v, h)}{\partial w_{qz}}}_{= -\sigma(c_z + v^T W_{:,z}) v_q} \\
 &= \sum_v p(v) \sigma(c_z + v^T W_{:,z}) v_q
 \end{aligned}$$

which is intractable for regular sized RBMs because its complexity is still exponential in the size of  $v$  (or  $h$ )

Recalling that we are assuming binary input/hidden units:

$\text{size}(v) \rightarrow$  with  $k$  input units the size is  $2^k$

A.A. 2022/23: Deep Learning

# Training Restricted Boltzmann Machines

Summing up

$$\frac{\partial \log p(\mathbf{x}|\boldsymbol{\theta})}{\partial w_{qz}} = \underbrace{\sum_{i=1}^n \sigma(c_z + \mathbf{v}^{(i)\top} \mathbf{W}_{:,z}) v_q^{(i)}}_{\text{data}} - \underbrace{\sum_{\mathbf{v}} p(\mathbf{v}) \sigma(c_z + \mathbf{v}^\top \mathbf{W}_{:,z}) v_q}_{\text{model}}$$

$$\propto \langle v_q h_z \rangle_{\text{data}} - \langle v_q h_z \rangle_{\text{model}}$$

due to partition function  
needs Gibbs sampling

and in a similar way it is possible to get the  $\nabla_b$  and  $\nabla_c$  components of the gradient

$$\frac{\partial \log p(\mathbf{x}|\boldsymbol{\theta})}{\partial b_q} = \sum_{i=1}^n v_q^{(i)} - \underbrace{\sum_{\mathbf{v}} p(\mathbf{v}) v_q}_{\text{model}}$$

$$\frac{\partial \log p(\mathbf{x}|\boldsymbol{\theta})}{\partial c_z} = \sum_{i=1}^n \sigma(c_z + \mathbf{v}^\top \mathbf{W}_{:,z}) - \underbrace{\sum_{\mathbf{v}} p(\mathbf{v}) \sigma(c_z + \mathbf{v}^\top \mathbf{W}_{:,z})}_{\text{model}}$$

due to partition function  
needs Gibbs sampling

A.A. 2022/23: Deep Learning

# Training Restricted Boltzmann Machines

**Algorithm 18.1** A naive MCMC algorithm for maximizing the log-likelihood with an intractable partition function using gradient ascent

Set  $\epsilon$ , the step size, to a small positive number.

Set  $k$ , the number of Gibbs steps, high enough to allow burn in. Perhaps 100 to train an RBM on a small image patch.

**while** not converged **do**

    Sample a minibatch of  $m$  examples  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  from the training set  
 $\mathbf{g} \leftarrow \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x}^{(i)}; \boldsymbol{\theta}).$

    Initialize a set of  $m$  samples  $\{\tilde{\mathbf{x}}^{(1)}, \dots, \tilde{\mathbf{x}}^{(m)}\}$  to random values (e.g., from a uniform or normal distribution, or possibly a distribution with marginals matched to the model's marginals)

**for**  $i = 1$  to  $k$  **do**

**for**  $j = 1$  to  $m$  **do**

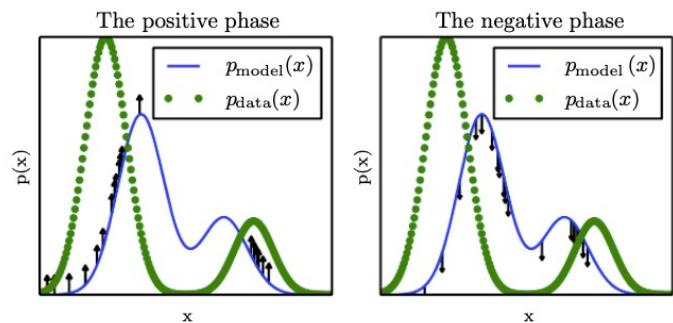
$\tilde{\mathbf{x}}^{(j)} \leftarrow \text{gibbs\_update}(\tilde{\mathbf{x}}^{(j)}).$

**end for**

**end for**

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

**end while**



A.A. 2022/23: Deep Learning

# Training Restricted Boltzmann Machines

---

**Algorithm 18.2** The contrastive divergence algorithm, using gradient ascent as the optimization procedure

---

Set  $\epsilon$ , the step size, to a small positive number.

Set  $k$ , the number of Gibbs steps, high enough to allow a Markov chain sampling from  $p(\mathbf{x}; \boldsymbol{\theta})$  to mix when initialized from  $p_{\text{data}}$ . Perhaps 1–20 to train an RBM on a small image patch.

**while** not converged **do**

    Sample a minibatch of  $m$  examples  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  from the training set

$$\mathbf{g} \leftarrow \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x}^{(i)}; \boldsymbol{\theta}).$$

**for**  $i = 1$  to  $m$  **do**

$$\tilde{\mathbf{x}}^{(i)} \leftarrow \mathbf{x}^{(i)}.$$

**end for**

**for**  $i = 1$  to  $k$  **do**

**for**  $j = 1$  to  $m$  **do**

$$\tilde{\mathbf{x}}^{(j)} \leftarrow \text{gibbs\_update}(\tilde{\mathbf{x}}^{(j)}).$$

**end for**

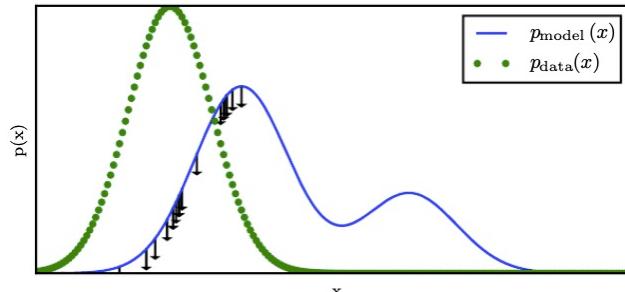
**end for**

$$\mathbf{g} \leftarrow \mathbf{g} - \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\tilde{\mathbf{x}}^{(i)}; \boldsymbol{\theta}).$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \epsilon \mathbf{g}.$$

**end while**

---



# Training Restricted Boltzmann Machines

---

**Algorithm 18.3** The stochastic maximum likelihood / persistent contrastive divergence algorithm using gradient ascent as the optimization procedure

---

Set  $\epsilon$ , the step size, to a small positive number.

Set  $k$ , the number of Gibbs steps, high enough to allow a Markov chain sampling from  $p(\mathbf{x}; \boldsymbol{\theta} + \epsilon \mathbf{g})$  to burn in, starting from samples from  $p(\mathbf{x}; \boldsymbol{\theta})$ . Perhaps 1 for RBM on a small image patch, or 5–50 for a more complicated model like a DBM.

Initialize a set of  $m$  samples  $\{\tilde{\mathbf{x}}^{(1)}, \dots, \tilde{\mathbf{x}}^{(m)}\}$  to random values (e.g., from a uniform or normal distribution, or possibly a distribution with marginals matched to the model's marginals).

**while** not converged **do**

    Sample a minibatch of  $m$  examples  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  from the training set

$$\mathbf{g} \leftarrow \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x}^{(i)}; \boldsymbol{\theta}).$$

**for**  $i = 1$  to  $k$  **do**

**for**  $j = 1$  to  $m$  **do**

$$\tilde{\mathbf{x}}^{(j)} \leftarrow \text{gibbs\_update}(\tilde{\mathbf{x}}^{(j)}).$$

**end for**

**end for**

$$\mathbf{g} \leftarrow \mathbf{g} - \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\tilde{\mathbf{x}}^{(i)}; \boldsymbol{\theta}).$$

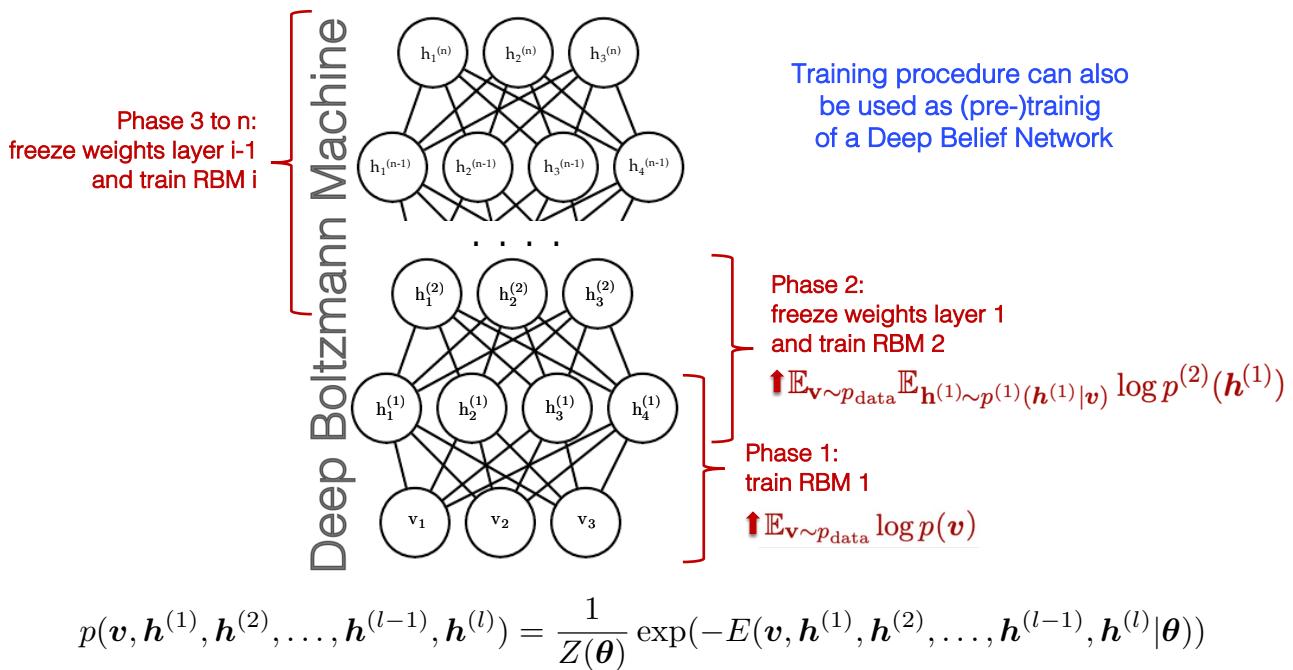
$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \epsilon \mathbf{g}.$$

**end while**

---

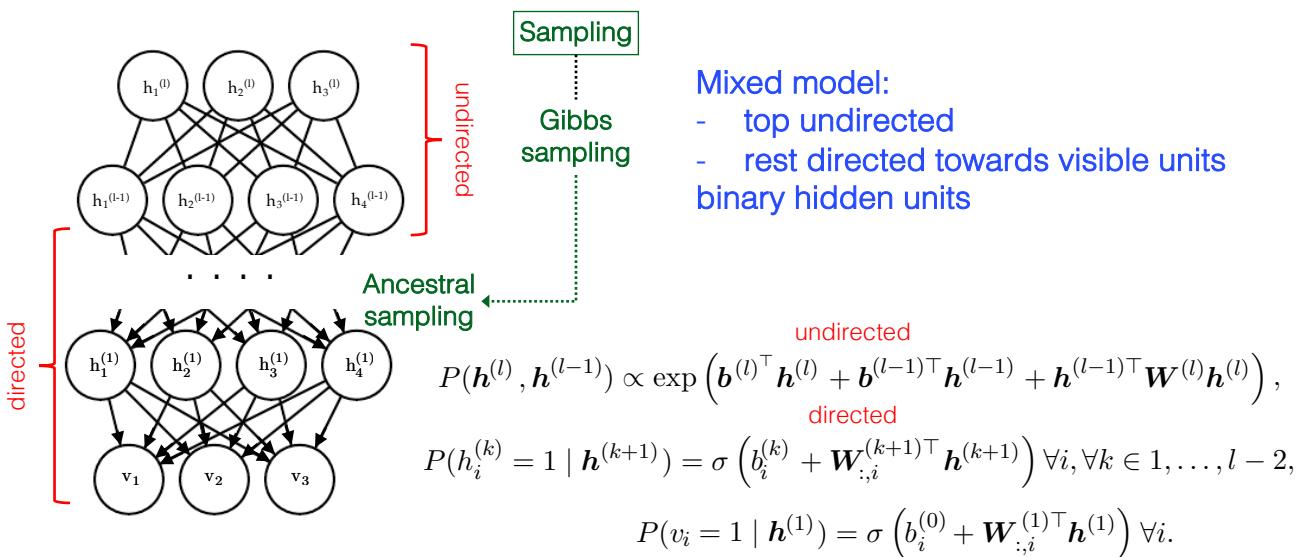
# RBM-based Deep Models

Layer-wise approach: greatly stack RBMs  
(suboptimal)



A.A. 2022/23: Deep Learning

## Deep Belief Networks



In the case of real-valued visible units, substitute

$$\mathbf{v} \sim \mathcal{N} \left( \mathbf{v}; \mathbf{b}^{(0)} + \mathbf{W}^{(1)\top} \mathbf{h}^{(1)}, \boldsymbol{\beta}^{-1} \right)$$

diagonal

A.A. 2022/23: Deep Learning