

---

**Algorithm 20.1** The variational stochastic maximum likelihood algorithm for training a DBM with two hidden layers.

---

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

Set  $k$ , the number of Gibbs steps, high enough to allow a Markov chain of  $p(\mathbf{v}, \mathbf{h}^{(1)}, \mathbf{h}^{(2)}; \boldsymbol{\theta} + \epsilon \Delta_{\boldsymbol{\theta}})$  to burn in, starting from samples from  $p(\mathbf{v}, \mathbf{h}^{(1)}, \mathbf{h}^{(2)}; \boldsymbol{\theta})$ . Initialize three matrices,  $\tilde{\mathbf{V}}$ ,  $\tilde{\mathbf{H}}^{(1)}$  and  $\tilde{\mathbf{H}}^{(2)}$  each with  $m$  rows set to random values (e.g., from Bernoulli distributions, possibly with marginals matched to the model's marginals).

**while** not converged (learning loop) **do**

Sample a minibatch of  $m$  examples from the training data and arrange them as the rows of a design matrix  $\mathbf{V}$ .

Initialize matrices  $\hat{\mathbf{H}}^{(1)}$  and  $\hat{\mathbf{H}}^{(2)}$ , possibly to the model's marginals.

**while** not converged (mean field inference loop) **do**

$$\hat{\mathbf{H}}^{(1)} \leftarrow \sigma \left( \mathbf{V} \mathbf{W}^{(1)} + \hat{\mathbf{H}}^{(2)} \mathbf{W}^{(2)\top} \right).$$

$$\hat{\mathbf{H}}^{(2)} \leftarrow \sigma \left( \hat{\mathbf{H}}^{(1)} \mathbf{W}^{(2)} \right).$$

**end while**

$$\Delta \mathbf{W}^{(1)} \leftarrow \frac{1}{m} \mathbf{V}^\top \hat{\mathbf{H}}^{(1)}$$

$$\Delta \mathbf{W}^{(2)} \leftarrow \frac{1}{m} \hat{\mathbf{H}}^{(1)\top} \hat{\mathbf{H}}^{(2)}$$

**for**  $l = 1$  to  $k$  (Gibbs sampling) **do**

Gibbs block 1:

$$\forall i, j, \tilde{V}_{i,j} \text{ sampled from } P(\tilde{V}_{i,j} = 1) = \sigma \left( \mathbf{W}_{j,:}^{(1)} \left( \tilde{\mathbf{H}}_{i,:}^{(1)} \right)^\top \right).$$

$$\forall i, j, \tilde{H}_{i,j}^{(2)} \text{ sampled from } P(\tilde{H}_{i,j}^{(2)} = 1) = \sigma \left( \tilde{\mathbf{H}}_{i,:}^{(1)} \mathbf{W}_{:,j}^{(2)} \right).$$

Gibbs block 2:

$$\forall i, j, \tilde{H}_{i,j}^{(1)} \text{ sampled from } P(\tilde{H}_{i,j}^{(1)} = 1) = \sigma \left( \tilde{\mathbf{V}}_{i,:} \mathbf{W}_{:,j}^{(1)} + \tilde{\mathbf{H}}_{i,:}^{(2)} \mathbf{W}_{j,:}^{(2)\top} \right).$$

**end for**

$$\Delta \mathbf{W}^{(1)} \leftarrow \Delta \mathbf{W}^{(1)} - \frac{1}{m} \mathbf{V}^\top \tilde{\mathbf{H}}^{(1)}$$

$$\Delta \mathbf{W}^{(2)} \leftarrow \Delta \mathbf{W}^{(2)} - \frac{1}{m} \tilde{\mathbf{H}}^{(1)\top} \tilde{\mathbf{H}}^{(2)}$$

$\mathbf{W}^{(1)} \leftarrow \mathbf{W}^{(1)} + \epsilon \Delta \mathbf{W}^{(1)}$  (this is a cartoon illustration, in practice use a more effective algorithm, such as momentum with a decaying learning rate)

$$\mathbf{W}^{(2)} \leftarrow \mathbf{W}^{(2)} + \epsilon \Delta \mathbf{W}^{(2)}$$

**end while**

---