

# Restricted Boltzmann Machine Derivations

Jan Schlüter

last change: March 26, 2014

## Abstract

This document gives detailed derivations for the central quantities of Restricted Boltzmann Machines (RBMs): The conditional distributions of visible and hidden units, and the log likelihood gradient with respect to the model parameters. It handles the standard Bernoulli-Bernoulli RBM (with binary visible and hidden units) as well as different formulations of the Gaussian-Bernoulli RBM (with real-valued visible units). It is not meant as a general introduction to RBMs, but as a supplement helping to follow the mathematics. If you are not familiar with RBMs, introductions can be found in [3] (mathematical), [2] (practical) or [7, Section 4.4] (intuitive).

## 1 General equations

An RBM is defined by its energy function

$$E(\mathbf{v}, \mathbf{h}, \boldsymbol{\theta}) \tag{1}$$

defining the joint probability distribution

$$p(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta}) = \frac{\exp(-E(\mathbf{v}, \mathbf{h}, \boldsymbol{\theta}))}{Z(\boldsymbol{\theta})} = \frac{\exp(-E(\mathbf{v}, \mathbf{h}, \boldsymbol{\theta}))}{\sum_{\mathbf{u}, \mathbf{g}} \exp(-E(\mathbf{u}, \mathbf{g}, \boldsymbol{\theta}))}, \tag{2}$$

where  $\mathbf{v}$  is the vector of visible units,  $\mathbf{h}$  is the vector of hidden units, and  $\boldsymbol{\theta}$  is the set of model parameters (i.e., the connection weights and biases). In sums over configurations, we use  $\mathbf{u}$  for visible and  $\mathbf{g}$  for hidden unit states. For continuous visible units, all sums over visible configurations  $\mathbf{u}$  in this section need to be replaced by integrals.

### 1.1 Marginal distribution of visible and hidden units

From (2), we can derive the marginal distribution of the visible units

$$p(\mathbf{v}; \boldsymbol{\theta}) = \sum_{\mathbf{g}} p(\mathbf{v}, \mathbf{g}; \boldsymbol{\theta}) = \frac{\sum_{\mathbf{g}} \exp(-E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta}))}{Z(\boldsymbol{\theta})} \tag{3}$$

and the marginal distribution of the hidden units

$$p(\mathbf{h}; \boldsymbol{\theta}) = \sum_{\mathbf{u}} p(\mathbf{u}, \mathbf{h}; \boldsymbol{\theta}) = \frac{\sum_{\mathbf{u}} \exp(-E(\mathbf{u}, \mathbf{h}, \boldsymbol{\theta}))}{Z(\boldsymbol{\theta})}. \tag{4}$$

### 1.2 Conditional distribution of visible and hidden units

From (2) and (4), we can derive the conditional distribution of the visible units

$$p(\mathbf{v}|\mathbf{h}; \boldsymbol{\theta}) = \frac{p(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta})}{p(\mathbf{h}; \boldsymbol{\theta})} = \frac{\frac{1}{Z(\boldsymbol{\theta})} \exp(-E(\mathbf{v}, \mathbf{h}, \boldsymbol{\theta}))}{\frac{1}{Z(\boldsymbol{\theta})} \sum_{\mathbf{u}} \exp(-E(\mathbf{u}, \mathbf{h}, \boldsymbol{\theta}))} = \frac{\exp(-E(\mathbf{v}, \mathbf{h}, \boldsymbol{\theta}))}{\sum_{\mathbf{u}} \exp(-E(\mathbf{u}, \mathbf{h}, \boldsymbol{\theta}))}, \tag{5}$$

and from (2) and (3), we can derive the conditional distribution of the hidden units

$$p(\mathbf{h}|\mathbf{v}; \boldsymbol{\theta}) = \frac{p(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta})}{p(\mathbf{v}; \boldsymbol{\theta})} = \frac{\frac{1}{Z(\boldsymbol{\theta})} \exp(-E(\mathbf{v}, \mathbf{h}, \boldsymbol{\theta}))}{\frac{1}{Z(\boldsymbol{\theta})} \sum_{\mathbf{g}} \exp(-E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta}))} = \frac{\exp(-E(\mathbf{v}, \mathbf{h}, \boldsymbol{\theta}))}{\sum_{\mathbf{g}} \exp(-E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta}))}. \tag{6}$$

### 1.3 Log likelihood gradient

We will derive the general form of the gradient of the log likelihood w.r.t. an arbitrary model parameter  $\theta$  for a single training sample  $\mathbf{v}$  here and use it to obtain the gradients for specific types of RBMs and model parameters in the following sections. Some steps are marked with a number and explained below.

$$\begin{aligned}
& \frac{\partial}{\partial \theta} \log p(\mathbf{v}; \boldsymbol{\theta}) \\
\stackrel{1}{=} & \frac{\partial}{\partial \theta} \log \left( \frac{\sum_{\mathbf{g}} \exp(-E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta}))}{Z(\boldsymbol{\theta})} \right) \\
= & \frac{\partial}{\partial \theta} \log \left( \sum_{\mathbf{g}} \exp(-E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta})) - \log Z(\boldsymbol{\theta}) \right) \\
= & \frac{\partial}{\partial \theta} \log \sum_{\mathbf{g}} \exp(-E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta})) - \frac{\partial}{\partial \theta} \log Z(\boldsymbol{\theta}) \\
\stackrel{2}{=} & \frac{\frac{\partial}{\partial \theta} \sum_{\mathbf{g}} \exp(-E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta}))}{\sum_{\mathbf{g}} \exp(-E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta}))} - \frac{\frac{\partial}{\partial \theta} Z(\boldsymbol{\theta})}{Z(\boldsymbol{\theta})} \\
\stackrel{3}{=} & \frac{\sum_{\mathbf{g}} \exp(-E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta})) \cdot \frac{\partial}{\partial \theta} (-E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta}))}{\sum_{\mathbf{g}} \exp(-E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta}))} - \frac{\sum_{\mathbf{u}, \mathbf{g}} \exp(-E(\mathbf{u}, \mathbf{g}, \boldsymbol{\theta})) \cdot \frac{\partial}{\partial \theta} (-E(\mathbf{u}, \mathbf{g}, \boldsymbol{\theta}))}{Z(\boldsymbol{\theta})} \\
= & \sum_{\mathbf{g}} \left( \frac{\exp(-E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta}))}{\sum_{\mathbf{g}'} \exp(-E(\mathbf{v}, \mathbf{g}', \boldsymbol{\theta}))} \cdot \frac{\partial}{\partial \theta} (-E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta})) \right) - \sum_{\mathbf{u}, \mathbf{g}} \left( \frac{\exp(-E(\mathbf{u}, \mathbf{g}, \boldsymbol{\theta}))}{Z(\boldsymbol{\theta})} \cdot \frac{\partial}{\partial \theta} (-E(\mathbf{u}, \mathbf{g}, \boldsymbol{\theta})) \right) \\
\stackrel{4}{=} & \sum_{\mathbf{g}} p(\mathbf{g}|\mathbf{v}; \boldsymbol{\theta}) \cdot \frac{\partial}{\partial \theta} (-E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta})) - \sum_{\mathbf{u}, \mathbf{g}} p(\mathbf{u}, \mathbf{g}; \boldsymbol{\theta}) \cdot \frac{\partial}{\partial \theta} (-E(\mathbf{u}, \mathbf{g}, \boldsymbol{\theta})) \\
= & \sum_{\mathbf{g}} p(\mathbf{g}|\mathbf{v}; \boldsymbol{\theta}) \cdot \frac{\partial}{\partial \theta} (-E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta})) - \sum_{\mathbf{u}} p(\mathbf{u}; \boldsymbol{\theta}) \sum_{\mathbf{g}} p(\mathbf{g}|\mathbf{u}; \boldsymbol{\theta}) \cdot \frac{\partial}{\partial \theta} (-E(\mathbf{u}, \mathbf{g}, \boldsymbol{\theta})) \\
\stackrel{5}{=} & G_{\theta}(\mathbf{v}, \boldsymbol{\theta}) - \langle G_{\theta}(\mathbf{u}, \boldsymbol{\theta}) \rangle_{p(\mathbf{u}; \boldsymbol{\theta})} \quad \text{with } G_{\theta}(\mathbf{x}, \boldsymbol{\theta}) := \sum_{\mathbf{g}} p(\mathbf{g}|\mathbf{x}; \boldsymbol{\theta}) \cdot \frac{\partial}{\partial \theta} (-E(\mathbf{x}, \mathbf{g}, \boldsymbol{\theta})) \tag{7}
\end{aligned}$$

Explanation of steps marked with a number:

1. We substitute the definition of  $p(\mathbf{v}; \boldsymbol{\theta})$  (3).
2. We derive the log functions, applying the chain rule of differentiation.
3. We substitute the definition of the partition function  $Z(\boldsymbol{\theta})$  (see denominators of (2)), then derive the exponential functions, again applying the chain rule.
4. We back-substitute equations (6) and (2), respectively.
5. We note that the second term computes the expectation of the first term under the model's marginal distribution of visible units  $p(\mathbf{v}; \boldsymbol{\theta})$ .

## 2 Bernoulli-Bernoulli RBM

Energy function:

$$E(\mathbf{v}, \mathbf{h}, \boldsymbol{\theta}) = - \sum_i v_i a_i - \sum_{i,j} v_i W_{ij} h_j - \sum_j h_j b_j \quad (8)$$

### 2.1 Conditional distribution of the visible units

In the binary RBM, the conditional probability of a visible unit being active given the hidden unit states is given as:

$$p(v_k = 1 | \mathbf{h}; \boldsymbol{\theta}) = \sigma \left( a_k + \sum_j W_{kj} h_j \right), \quad (9)$$

where  $\sigma(x) = (1 + \exp(-x))^{-1}$  is the logistic sigmoid function. The derivation is symmetric to the one for the hidden units discussed in the next section.

### 2.2 Conditional distribution of the hidden units

This derivation is copied from [7, p. 143]. Some steps are marked with a number and explained below.

$$\begin{aligned} p(h_k = 1 | \mathbf{v}; \boldsymbol{\theta}) &\stackrel{1}{=} \frac{p(\mathbf{v}, h_k = 1; \boldsymbol{\theta})}{p(\mathbf{v}; \boldsymbol{\theta})} \\ &= \frac{\sum_{\{\mathbf{g} | g_k = 1\}} p(\mathbf{v}, \mathbf{g}; \boldsymbol{\theta})}{\sum_{\mathbf{g}} p(\mathbf{v}, \mathbf{g}; \boldsymbol{\theta})} \\ &\stackrel{2}{=} \frac{\sum_{\{\mathbf{g} | g_k = 1\}} \exp(-E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta}))}{\sum_{\mathbf{g}} \exp(-E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta}))} \\ &\stackrel{3}{=} \frac{\sum_{\{\mathbf{g} | g_k = 1\}} \exp(-E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta}))}{\sum_{\{\mathbf{g} | g_k = 1\}} \exp(-E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta})) + \sum_{\{\mathbf{g} | g_k = 0\}} \exp(-E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta}))} \\ &= \frac{1}{1 + \frac{\sum_{\{\mathbf{g} | g_k = 0\}} \exp(-E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta}))}{\sum_{\{\mathbf{g} | g_k = 1\}} \exp(-E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta}))}} \\ &\stackrel{4}{=} \frac{1}{1 + \frac{\sum_{\{\mathbf{g} | g_k = 0\}} \exp(\sum_{i,j} v_i W_{ij} g_j + \sum_i v_i a_i + \sum_{j \neq k} g_j b_j)}{\sum_{\{\mathbf{g} | g_k = 1\}} \exp(\sum_{i,j} v_i W_{ij} g_j + \sum_i v_i a_i + \sum_{j \neq k} g_j b_j)}} \\ &\stackrel{5}{=} \frac{1}{1 + \frac{\sum_{\{\mathbf{g} | g_k = 0\}} \exp(\sum_i \sum_{j \neq k} v_i W_{ij} g_j + \sum_i v_i a_i + \sum_{j \neq k} g_j b_j)}{\sum_{\{\mathbf{g} | g_k = 1\}} \exp(\sum_i \sum_{j \neq k} v_i W_{ij} g_j + \sum_i v_i a_i + \sum_{j \neq k} g_j b_j)} \cdot \frac{\exp(\sum_i v_i W_{ik} \cdot 0 + 0 \cdot b_k)}{\exp(\sum_i v_i W_{ik} \cdot 1 + 1 \cdot b_k)}} \\ &\stackrel{6}{=} \frac{1}{1 + 1 \cdot \frac{1}{\exp(\sum_i v_i W_{ik} + b_k)}} \\ &= \frac{1}{1 + \exp(-(b_k + \sum_i v_i W_{ik}))} \\ &= \sigma \left( b_k + \sum_i v_i W_{ik} \right) \end{aligned} \quad (10)$$

Explanation of steps marked with a number:

1. We start from the definition of conditional probability.
2. We substitute the definition of the joint probability density function (2).
3. We split the sum over all configurations  $\mathbf{g}$  into two sums over configurations with an active and inactive hidden unit  $g_k$ , respectively.
4. We substitute the energy function (8).
5. From the sums over  $j$ , we move out the terms for  $j = k$ , for which we substitute the known values of  $g_j = g_k$ .

6. As the terms depending on  $g_k$  have been moved out of the two long exponentials, the sums over  $\{\mathbf{g}|g_k = 0\}$  and  $\{\mathbf{g}|g_k = 1\}$  are equal so their quotient evaluates to 1.

### 2.3 Log likelihood gradient

Following the derivation of Section 1.3, we just need to compute  $G_{\theta}(\mathbf{v}, \boldsymbol{\theta})$  for the model parameters  $\boldsymbol{\theta}$  we are interested in: a connection weight  $W_{ij}$ , visible bias  $a_i$  and hidden bias  $b_j$ . The derivations are based on [3, p. 26]. As always, some steps are marked with a number and explained in the end of the section. We start with the connection weight:

$$\begin{aligned}
G_{W_{ij}}(\mathbf{v}, \boldsymbol{\theta}) &= \sum_{\mathbf{g}} p(\mathbf{g}|\mathbf{v}; \boldsymbol{\theta}) \cdot \frac{\partial}{\partial W_{ij}} \left( -E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta}) \right) \\
&\stackrel{1}{=} \sum_{\mathbf{g}} p(\mathbf{g}|\mathbf{v}; \boldsymbol{\theta}) \cdot v_i g_j \\
&\stackrel{2}{=} v_i \cdot \sum_{\mathbf{g}} p(\mathbf{g}|\mathbf{v}; \boldsymbol{\theta}) g_j \\
&\stackrel{3}{=} v_i \cdot \sum_{g_j} \sum_{\mathbf{g}_{-j}} p(g_j|\mathbf{v}; \boldsymbol{\theta}) p(\mathbf{g}_{-j}|\mathbf{v}; \boldsymbol{\theta}) g_j \\
&\stackrel{4}{=} v_i \cdot \sum_{g_j} p(g_j|\mathbf{v}; \boldsymbol{\theta}) g_j \sum_{\mathbf{g}_{-j}} p(\mathbf{g}_{-j}|\mathbf{v}; \boldsymbol{\theta}) \\
&\stackrel{5}{=} v_i \cdot \sum_{g_j} p(g_j|\mathbf{v}; \boldsymbol{\theta}) g_j \\
&\stackrel{6}{=} v_i \cdot \left( p(h_j = 1|\mathbf{v}; \boldsymbol{\theta}) \cdot 1 + p(h_j = 0|\mathbf{v}; \boldsymbol{\theta}) \cdot 0 \right) \\
&= v_i \cdot p(h_j = 1|\mathbf{v}; \boldsymbol{\theta}) \\
&\stackrel{7}{=} v_i \cdot \sigma \left( b_j + \sum_k v_k W_{kj} \right)
\end{aligned}$$

Revisiting the derivation, note that we have established the following between the third and the previous to last line:

$$\sum_{\mathbf{g}} p(\mathbf{g}|\mathbf{v}; \boldsymbol{\theta}) g_j = p(h_j = 1|\mathbf{v}; \boldsymbol{\theta}) \tag{11}$$

This holds whenever  $p(\mathbf{h}|\mathbf{v}; \boldsymbol{\theta})$  is factorial and will become useful in other derivations. We continue with the gradient for the visible bias:

$$\begin{aligned}
G_{a_i}(\mathbf{v}, \boldsymbol{\theta}) &= \sum_{\mathbf{g}} p(\mathbf{g}|\mathbf{v}; \boldsymbol{\theta}) \cdot \frac{\partial}{\partial a_i} \left( -E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta}) \right) \\
&\stackrel{1}{=} \sum_{\mathbf{g}} p(\mathbf{g}|\mathbf{v}; \boldsymbol{\theta}) \cdot v_i \\
&\stackrel{2}{=} v_i \cdot \sum_{\mathbf{g}} p(\mathbf{g}|\mathbf{v}; \boldsymbol{\theta}) \\
&\stackrel{5}{=} v_i
\end{aligned}$$

Finally, we compute the gradient for the hidden bias:

$$\begin{aligned}
G_{b_j}(\mathbf{v}, \boldsymbol{\theta}) &= \sum_{\mathbf{g}} p(\mathbf{g}|\mathbf{v}; \boldsymbol{\theta}) \cdot \frac{\partial}{\partial b_j} \left( -E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta}) \right) \\
&\stackrel{1}{=} \sum_{\mathbf{g}} p(\mathbf{g}|\mathbf{v}; \boldsymbol{\theta}) \cdot g_j \\
&\stackrel{8}{=} \sigma \left( b_j + \sum_k v_k W_{kj} \right)
\end{aligned}$$

Explanation of steps marked with a number:

1. We substitute the energy function (8) and differentiate.
2. The factor  $v_i$  is the same for every  $\mathbf{g}$ , so we can pull it out of the sum.
3. We split the sum over all configurations  $\mathbf{g}$  into the sum over the two states of  $g_j$  and the sum over the states of the remaining hidden units  $\mathbf{g}_{-j}$ . Furthermore, noting that the distribution  $p(\mathbf{g}|\mathbf{v};\boldsymbol{\theta})$  is factorial (the hidden units are mutually independent given  $\mathbf{v}$ ), we pull out the factor  $p(g_j|\mathbf{v};\boldsymbol{\theta})$ .
4. The factor  $p(g_j|\mathbf{v};\boldsymbol{\theta})g_j$  is the same for every  $\mathbf{g}_{-j}$ , so we can pull it out of the inner sum.
5. The sum over a probability distribution evaluates to 1 and can be dropped.
6. We write out the sum over the two possible values of  $g_j$  (and rename it to  $h_j$ ).
7. We substitute the conditional hidden unit activation probability (10).
8. We substitute (11), then (10).

Inserting our results for  $G_{W_{ij}}(\mathbf{v},\boldsymbol{\theta})$ ,  $G_{a_i}(\mathbf{v},\boldsymbol{\theta})$ ,  $G_{b_j}(\mathbf{v},\boldsymbol{\theta})$  into (7), we obtain:

$$\frac{\partial}{\partial W_{ij}} \log p(\mathbf{v};\boldsymbol{\theta}) = v_i \cdot \sigma\left(b_j + \sum_k v_k W_{kj}\right) - \left\langle u_i \cdot \sigma\left(b_j + \sum_k u_k W_{kj}\right) \right\rangle_{p(\mathbf{u};\boldsymbol{\theta})} \quad (12)$$

$$\frac{\partial}{\partial a_i} \log p(\mathbf{v};\boldsymbol{\theta}) = v_i - \langle u_i \rangle_{p(\mathbf{u};\boldsymbol{\theta})} \quad (13)$$

$$\frac{\partial}{\partial b_j} \log p(\mathbf{v};\boldsymbol{\theta}) = \sigma\left(b_j + \sum_k v_k W_{kj}\right) - \left\langle \sigma\left(b_j + \sum_k u_k W_{kj}\right) \right\rangle_{p(\mathbf{u};\boldsymbol{\theta})} \quad (14)$$

The first terms can be easily computed from the training example  $\mathbf{v}$ , the second terms require evaluating an expectation over  $p(\mathbf{v};\boldsymbol{\theta})$ , which is computationally infeasible due to the partition function in (3). In practice, the second terms are usually approximated by averaging over samples from  $p(\mathbf{v};\boldsymbol{\theta})$ , which in turn are drawn approximately using Markov-Chain Monte-Carlo (MCMC) methods such as Gibbs sampling. Please see one of the RBM tutorials mentioned in the abstract for details.

For notational brevity, let us introduce the following abbreviation:

$$\bar{h}_j(\mathbf{v},\boldsymbol{\theta}) := p(h_j = 1|\mathbf{v};\boldsymbol{\theta}) \quad (15)$$

This way, the gradients can be written as:

$$\frac{\partial}{\partial W_{ij}} \log p(\mathbf{v};\boldsymbol{\theta}) = v_i \cdot \bar{h}_j(\mathbf{v},\boldsymbol{\theta}) - \langle u_i \cdot \bar{h}_j(\mathbf{u},\boldsymbol{\theta}) \rangle_{p(\mathbf{u};\boldsymbol{\theta})} \quad (16)$$

$$\frac{\partial}{\partial a_i} \log p(\mathbf{v};\boldsymbol{\theta}) = v_i - \langle u_i \rangle_{p(\mathbf{u};\boldsymbol{\theta})} \quad (17)$$

$$\frac{\partial}{\partial b_j} \log p(\mathbf{v};\boldsymbol{\theta}) = \bar{h}_j(\mathbf{v},\boldsymbol{\theta}) - \langle \bar{h}_j(\mathbf{u},\boldsymbol{\theta}) \rangle_{p(\mathbf{u};\boldsymbol{\theta})} \quad (18)$$

In a practical implementation, one will usually draw samples  $\mathbf{u}$ , precompute all  $\bar{h}_j(\mathbf{v},\boldsymbol{\theta})$  and  $\bar{h}_j(\mathbf{u},\boldsymbol{\theta})$ , and then use those to compute the gradients. Note that  $\bar{h}_j(\mathbf{v},\boldsymbol{\theta})$  and  $\bar{h}_j(\mathbf{u},\boldsymbol{\theta})$  denote the real-valued activation probabilities of hidden units, not binary samples.

### 3 Gaussian-Bernoulli RBM (Krizhevsky's variant)

Energy function [5, p. 12]:

$$E(\mathbf{v}, \mathbf{h}, \boldsymbol{\theta}) = \sum_i \frac{(v_i - a_i)^2}{2\sigma_i^2} - \sum_{i,j} \frac{v_i}{\sigma_i} W_{ij} h_j - \sum_j h_j b_j \quad (19)$$

#### 3.1 Conditional distribution of the visible units

We insert (19) into (5). We follow [5, p. 13], fixing one mistake and proceeding a lot slower. Again, some steps are marked with a number and explained below.

$$\begin{aligned}
p(\mathbf{v}|\mathbf{h}; \boldsymbol{\theta}) &= \frac{\exp(-E(\mathbf{v}, \mathbf{h}, \boldsymbol{\theta}))}{\int_{\mathbf{u}} \exp(-E(\mathbf{u}, \mathbf{h}, \boldsymbol{\theta})) d\mathbf{u}} \\
&= \frac{\exp(-\sum_i \frac{(v_i - a_i)^2}{2\sigma_i^2} + \sum_{i,j} \frac{v_i}{\sigma_i} W_{ij} h_j + \sum_j h_j b_j)}{\int_{\mathbf{u}} \exp(-\sum_i \frac{(u_i - a_i)^2}{2\sigma_i^2} + \sum_{i,j} \frac{u_i}{\sigma_i} W_{ij} h_j + \sum_j h_j b_j) d\mathbf{u}} \\
&= \frac{\exp(-\sum_i \frac{(v_i - a_i)^2}{2\sigma_i^2} + \sum_{i,j} \frac{v_i}{\sigma_i} W_{ij} h_j) \exp(\sum_j h_j b_j)}{\int_{\mathbf{u}} \exp(-\sum_i \frac{(u_i - a_i)^2}{2\sigma_i^2} + \sum_{i,j} \frac{u_i}{\sigma_i} W_{ij} h_j) \exp(\sum_j h_j b_j) d\mathbf{u}} \\
&\stackrel{1}{=} \frac{\exp(-\sum_i \frac{(v_i - a_i)^2}{2\sigma_i^2} + \sum_i \frac{v_i}{\sigma_i} \sum_j W_{ij} h_j)}{\int_{\mathbf{u}} \exp(-\sum_i \frac{(u_i - a_i)^2}{2\sigma_i^2} + \sum_i \frac{u_i}{\sigma_i} \sum_j W_{ij} h_j) d\mathbf{u}} \\
&\stackrel{2}{=} \frac{\exp(-\sum_i \frac{(v_i - a_i)^2}{2\sigma_i^2} + \sum_i \frac{v_i}{\sigma_i} B_i)}{\int_{\mathbf{u}} \exp(-\sum_i \frac{(u_i - a_i)^2}{2\sigma_i^2} + \sum_i \frac{u_i}{\sigma_i} B_i) d\mathbf{u}} \quad \text{with } B_i := \sum_j W_{ij} h_j \\
&= \frac{\prod_i \exp(-\frac{(v_i - a_i)^2}{2\sigma_i^2} + \frac{v_i}{\sigma_i} B_i)}{\int_{\mathbf{u}} \prod_i \exp(-\frac{(u_i - a_i)^2}{2\sigma_i^2} + \frac{u_i}{\sigma_i} B_i) d\mathbf{u}} \\
&\stackrel{3}{=} \frac{\dots}{\prod_i \int_{u_i} \exp(-\frac{(u_i - a_i)^2}{2\sigma_i^2} + \frac{u_i}{\sigma_i} B_i) du_i} \\
&= \frac{\dots}{\prod_i \int_{u_i} \exp(-\frac{u_i^2 - 2u_i a_i + a_i^2}{2\sigma_i^2} + \frac{u_i}{\sigma_i} B_i) du_i} \\
&= \frac{\dots}{\prod_i \int_{u_i} \exp(-u_i^2 \frac{1}{2\sigma_i^2} + u_i(\frac{a_i}{\sigma_i^2} + \frac{B_i}{\sigma_i}) - \frac{a_i^2}{2\sigma_i^2}) du_i} \\
&\stackrel{4}{=} \frac{\dots}{\prod_i \sqrt{\pi 2\sigma_i^2} \exp(\frac{2\sigma_i^2}{4} (\frac{a_i}{\sigma_i^2} + \frac{B_i}{\sigma_i})^2 - \frac{a_i^2}{2\sigma_i^2})} \\
&= \frac{\dots}{\prod_i \sigma_i \sqrt{2\pi} \exp(\frac{1}{2} (\frac{a_i}{\sigma_i} + B_i)^2 - \frac{a_i^2}{2\sigma_i^2})} \\
&= \frac{\dots}{\prod_i \sigma_i \sqrt{2\pi} \exp(\frac{1}{2} (\frac{a_i^2}{\sigma_i^2} + 2\frac{a_i}{\sigma_i} B_i + B_i^2) - \frac{a_i^2}{2\sigma_i^2})} \\
&= \frac{\prod_i \exp(-\frac{(v_i - a_i)^2}{2\sigma_i^2} + \frac{v_i}{\sigma_i} B_i)}{\prod_i \sigma_i \sqrt{2\pi} \exp(\frac{a_i}{\sigma_i} B_i + \frac{1}{2} B_i^2)} \\
&= \prod_i \frac{1}{\sigma_i \sqrt{2\pi}} \exp(-\frac{(v_i - a_i)^2}{2\sigma_i^2} + \frac{v_i}{\sigma_i} B_i - \frac{a_i}{\sigma_i} B_i - \frac{1}{2} B_i^2) \\
&= \prod_i \frac{1}{\sigma_i \sqrt{2\pi}} \exp\left(-\frac{1}{2\sigma_i^2} \left((v_i - a_i)^2 - 2\sigma_i(v_i - a_i)B_i + \sigma_i^2 B_i^2\right)\right) \\
&= \prod_i \frac{1}{\sigma_i \sqrt{2\pi}} \exp\left(-\frac{1}{2\sigma_i^2} \left((v_i - a_i) - \sigma_i B_i\right)^2\right) \\
&\stackrel{5}{=} \prod_i \mathcal{N}\left(v_i; a_i + \sigma_i \sum_j W_{ij} h_j, \sigma_i^2\right) \quad (20)
\end{aligned}$$

Explanation of steps marked with a number:

1. The  $\exp(\sum_j h_j b_j)$  term in the denominator is constant w.r.t.  $\mathbf{u}$ , thus can be dragged out of the integral as a factor and cancels with the same factor in the nominator.
2. Just to make notation easier, we substitute  $B_i := \sum_j W_{ij} h_j$ .
3. Note that  $\int_{\mathbf{u}} \prod_i f(u_i) d\mathbf{u} = \int_{u_1} \int_{u_2} \dots \int_{u_n} f(u_1) f(u_2) \dots f(u_n) du_1 du_2 \dots du_n$ , which can be freely rearranged to  $\int_{u_1} f(u_1) du_1 \int_{u_2} f(u_2) du_2 \dots \int_{u_n} f(u_n) du_n = \prod_i \int_{u_i} f(u_i) du_i$  because any of the integrals is just a constant factor w.r.t. the other integration variables.
4. We solve the Gaussian function integral  $\int_x \exp(-x^2 a + x b + c) dx = \sqrt{\frac{\pi}{a}} \exp\left(\frac{b^2}{4a} + c\right)$ .
5. For your convenience, the Normal distribution is  $\mathcal{N}(x; \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$ .

### 3.2 Conditional distribution of the hidden units

We start from an intermediate step of the derivation in Section 2.2 (the one between the steps marked with 3 and 4). Again, some steps are marked with a number and explained below.

$$\begin{aligned}
p(h_k = 1 | \mathbf{v}) &= \frac{1}{1 + \frac{\sum_{\{g|g_k=0\}} \exp(-E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta}))}{\sum_{\{g|g_k=1\}} \exp(-E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta}))}} \\
&\stackrel{1}{=} \frac{1}{1 + \frac{\sum_{\{g|g_k=0\}} \exp(\sum_{i,j} \frac{v_i}{\sigma_i} W_{ij} g_j - \sum_i \frac{(v_i - a_i)^2}{\sigma_i^2} + \sum_j g_j b_j)}{\sum_{\{g|g_k=1\}} \exp(\sum_{i,j} \frac{v_i}{\sigma_i} W_{ij} g_j - \sum_i \frac{(v_i - a_i)^2}{\sigma_i^2} + \sum_j g_j b_j)}} \\
&\stackrel{2}{=} \frac{1}{1 + \frac{\exp(\sum_i \frac{v_i}{\sigma_i} W_{ik} \cdot 0 + 0 \cdot b_k)}{\exp(\sum_i \frac{v_i}{\sigma_i} W_{ik} \cdot 1 + 1 \cdot b_k)}} \\
&= \frac{1}{1 + \frac{1}{\exp(\sum_i \frac{v_i}{\sigma_i} W_{ik} + b_k)}} \\
&= \frac{1}{1 + \exp(-(b_k + \sum_i \frac{v_i}{\sigma_i} W_{ik}))} \\
&= \sigma\left(b_k + \sum_i \frac{v_i}{\sigma_i} W_{ik}\right) \tag{21}
\end{aligned}$$

Explanation of steps marked with a number:

1. We substitute the energy function (19).
2. We substitute the known values for  $g_j = g_k$  and cancel all terms not depending on  $g_k$  (cf. steps 5 and 6 in Section 2.2).

### 3.3 Log likelihood gradient

As in Section 2.3, we compute  $G_{\theta}(\mathbf{v}, \boldsymbol{\theta})$  for the model parameters  $\theta$  we are interested in: a connection weight  $W_{ij}$ , visible bias  $a_i$ , hidden bias  $b_j$  and visible standard deviation  $\sigma_i$ . Again, some steps are marked with a number and explained in the end of the section.

We start with the connection weight:

$$\begin{aligned}
G_{W_{ij}}(\mathbf{v}, \boldsymbol{\theta}) &= \sum_{\mathbf{g}} p(\mathbf{g} | \mathbf{v}; \boldsymbol{\theta}) \cdot \frac{\partial}{\partial W_{ij}} \left( -E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta}) \right) \\
&\stackrel{1}{=} \sum_{\mathbf{g}} p(\mathbf{g} | \mathbf{v}; \boldsymbol{\theta}) \cdot \frac{v_i}{\sigma_i} g_j \\
&\stackrel{2}{=} \frac{v_i}{\sigma_i} \cdot \sum_{\mathbf{g}} p(\mathbf{g} | \mathbf{v}; \boldsymbol{\theta}) g_j
\end{aligned}$$

$$\begin{aligned}
&\stackrel{3}{=} \frac{v_i}{\sigma_i} \cdot p(h_j = 1 | \mathbf{v}; \boldsymbol{\theta}) \\
&\stackrel{4}{=} \frac{v_i}{\sigma_i} \cdot \bar{h}_j(\mathbf{v}, \boldsymbol{\theta})
\end{aligned}$$

We continue with the gradient for the visible bias:

$$\begin{aligned}
G_{a_i}(\mathbf{v}, \boldsymbol{\theta}) &= \sum_{\mathbf{g}} p(\mathbf{g} | \mathbf{v}; \boldsymbol{\theta}) \cdot \frac{\partial}{\partial a_i} \left( -E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta}) \right) \\
&\stackrel{1}{=} \sum_{\mathbf{g}} p(\mathbf{g} | \mathbf{v}; \boldsymbol{\theta}) \cdot \frac{v_i - a_i}{\sigma_i^2} \\
&\stackrel{5}{=} \frac{v_i - a_i}{\sigma_i^2} \cdot \sum_{\mathbf{g}} p(\mathbf{g} | \mathbf{v}; \boldsymbol{\theta}) \\
&\stackrel{6}{=} \frac{v_i - a_i}{\sigma_i^2}
\end{aligned}$$

We proceed with the gradient for the hidden bias:

$$\begin{aligned}
G_{b_j}(\mathbf{v}, \boldsymbol{\theta}) &= \sum_{\mathbf{g}} p(\mathbf{g} | \mathbf{v}; \boldsymbol{\theta}) \cdot \frac{\partial}{\partial b_j} \left( -E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta}) \right) \\
&\stackrel{1}{=} \sum_{\mathbf{g}} p(\mathbf{g} | \mathbf{v}; \boldsymbol{\theta}) \cdot g_j \\
&\stackrel{3,4}{=} \bar{h}_j(\mathbf{v}, \boldsymbol{\theta})
\end{aligned}$$

Finally, we compute the gradient for the visible standard derivation:

$$\begin{aligned}
G_{\sigma_i}(\mathbf{v}, \boldsymbol{\theta}) &= \sum_{\mathbf{g}} p(\mathbf{g} | \mathbf{v}; \boldsymbol{\theta}) \cdot \frac{\partial}{\partial \sigma_i} \left( -E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta}) \right) \\
&\stackrel{1}{=} \sum_{\mathbf{g}} p(\mathbf{g} | \mathbf{v}; \boldsymbol{\theta}) \cdot \left( \frac{(v_i - a_i)^2}{\sigma_i^3} - \frac{v_i}{\sigma_i^2} \sum_j W_{ij} g_j \right) \\
&= \sum_{\mathbf{g}} p(\mathbf{g} | \mathbf{v}; \boldsymbol{\theta}) \frac{(v_i - a_i)^2}{\sigma_i^3} - \sum_{\mathbf{g}} p(\mathbf{g} | \mathbf{v}; \boldsymbol{\theta}) \frac{v_i}{\sigma_i^2} \sum_j W_{ij} g_j \\
&\stackrel{7}{=} \frac{(v_i - a_i)^2}{\sigma_i^3} \cdot \sum_{\mathbf{g}} p(\mathbf{g} | \mathbf{v}; \boldsymbol{\theta}) - \frac{v_i}{\sigma_i^2} \sum_j W_{ij} \cdot \sum_{\mathbf{g}} p(\mathbf{g} | \mathbf{v}; \boldsymbol{\theta}) g_j \\
&\stackrel{6,3}{=} \frac{(v_i - a_i)^2}{\sigma_i^3} - \frac{v_i}{\sigma_i^2} \sum_j W_{ij} \cdot p(h_j = 1 | \mathbf{v}; \boldsymbol{\theta}) \\
&\stackrel{4}{=} \frac{(v_i - a_i)^2}{\sigma_i^3} - \frac{v_i}{\sigma_i^2} \sum_j W_{ij} \cdot \bar{h}_j(\mathbf{v}, \boldsymbol{\theta})
\end{aligned}$$

Explanation of steps marked with a number:

1. We substitute the energy function (19) and differentiate.
2. The factor  $\frac{v_i}{\sigma_i}$  is the same for every  $\mathbf{g}$ , so we can pull it out of the sum.
3. We substitute (11) from Section 2.3.
4. We substitute the abbreviation (15).
5. The factor  $\frac{v_i - a_i}{\sigma_i^2}$  is the same for every  $\mathbf{g}$ , so we can pull it out of the sum.
6. The sum over a probability distribution evaluates to 1 and can be dropped.
7. We pull out of the sums over  $\mathbf{g}$  all factors that are independent of  $\mathbf{g}$ .



Inserting our results for  $G_{W_{ij}}(\mathbf{v}, \boldsymbol{\theta})$ ,  $G_{a_i}(\mathbf{v}, \boldsymbol{\theta})$ ,  $G_{b_j}(\mathbf{v}, \boldsymbol{\theta})$ ,  $G_{\sigma_i}(\mathbf{v}, \boldsymbol{\theta})$  into (7), we obtain:

$$\frac{\partial}{\partial W_{ij}} \log p(\mathbf{v}; \boldsymbol{\theta}) = \frac{v_i}{\sigma_i} \cdot \bar{h}_j(\mathbf{v}, \boldsymbol{\theta}) - \left\langle \frac{u_i}{\sigma_i} \cdot \bar{h}_j(\mathbf{u}, \boldsymbol{\theta}) \right\rangle_{p(\mathbf{u}; \boldsymbol{\theta})} \quad (22)$$

$$\begin{aligned} \frac{\partial}{\partial a_i} \log p(\mathbf{v}; \boldsymbol{\theta}) &= \frac{v_i - a_i}{\sigma_i^2} - \left\langle \frac{u_i - a_i}{\sigma_i^2} \right\rangle_{p(\mathbf{u}; \boldsymbol{\theta})} \\ &= \frac{v_i}{\sigma_i^2} - \frac{a_i}{\sigma_i^2} - \left\langle \frac{u_i}{\sigma_i^2} - \frac{a_i}{\sigma_i^2} \right\rangle_{p(\mathbf{u}; \boldsymbol{\theta})} \\ &= \frac{v_i}{\sigma_i^2} - \frac{a_i}{\sigma_i^2} + \frac{a_i}{\sigma_i^2} - \left\langle \frac{u_i}{\sigma_i^2} \right\rangle_{p(\mathbf{u}; \boldsymbol{\theta})} \\ &= \frac{v_i}{\sigma_i^2} - \left\langle \frac{u_i}{\sigma_i^2} \right\rangle_{p(\mathbf{u}; \boldsymbol{\theta})} \end{aligned} \quad (23)$$

$$\frac{\partial}{\partial b_j} \log p(\mathbf{v}; \boldsymbol{\theta}) = \bar{h}_j(\mathbf{v}, \boldsymbol{\theta}) - \langle \bar{h}_j(\mathbf{u}, \boldsymbol{\theta}) \rangle_{p(\mathbf{u}; \boldsymbol{\theta})} \quad (24)$$

$$\begin{aligned} \frac{\partial}{\partial \sigma_i} \log p(\mathbf{v}; \boldsymbol{\theta}) &= \frac{(v_i - a_i)^2}{\sigma_i^3} - \frac{v_i}{\sigma_i^2} \sum_j W_{ij} \cdot \bar{h}_j(\mathbf{v}, \boldsymbol{\theta}) \\ &\quad - \left\langle \frac{(u_i - a_i)^2}{\sigma_i^3} - \frac{u_i}{\sigma_i^2} \sum_j W_{ij} \cdot \bar{h}_j(\mathbf{u}, \boldsymbol{\theta}) \right\rangle_{p(\mathbf{u}; \boldsymbol{\theta})} \end{aligned} \quad (25)$$

## 4 Gaussian-Bernoulli RBM (Cho's variant)

Energy function [1, p. 2]:

$$E(\mathbf{v}, \mathbf{h}, \boldsymbol{\theta}) = \sum_i \frac{(v_i - a_i)^2}{2\sigma_i^2} - \sum_{i,j} \frac{v_i}{\sigma_i^2} W_{ij} h_j - \sum_j h_j b_j \quad (26)$$

### 4.1 Conditional distribution of the visible units

Note that we can obtain the energy function (26) from (19) by simply replacing  $W_{ij}$  with  $\frac{1}{\sigma_i} W_{ij}$ . This substitution carries through the derivation (Section 3.1) and results in:

$$p(\mathbf{v}|\mathbf{h}; \boldsymbol{\theta}) = \prod_i \mathcal{N}\left(v_i; a_i + \sum_j W_{ij} h_j, \sigma_i^2\right) \quad (27)$$

### 4.2 Conditional distribution of the hidden units

Similarly, we can substitute  $W_{ij}$  with  $\frac{1}{\sigma_i} W_{ij}$  in the derivation of Section 3.2 to arrive at:

$$p(h_k = 1|\mathbf{v}) = \sigma\left(b_k + \sum_i \frac{v_i}{\sigma_i^2} W_{ik}\right) \quad (28)$$

### 4.3 Log likelihood gradient

The gradient derivations are very similar to Section 3.3. For the biases  $a_i$  and  $b_j$ , nothing changes. The connection weight  $W_{ij}$  and visible standard deviation  $\sigma_i$  gradients are affected by the slightly different  $\frac{\partial}{\partial W_{ij}}(-E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta}))$  and  $\frac{\partial}{\partial \sigma_i}(-E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta}))$ , respectively.

In summary, we obtain:

$$\frac{\partial}{\partial W_{ij}} \log p(\mathbf{v}; \boldsymbol{\theta}) = \frac{v_i}{\sigma_i^2} \cdot \bar{h}_j(\mathbf{v}, \boldsymbol{\theta}) - \left\langle \frac{u_i}{\sigma_i^2} \cdot \bar{h}_j(\mathbf{u}, \boldsymbol{\theta}) \right\rangle_{p(\mathbf{u}; \boldsymbol{\theta})} \quad (29)$$

$$\frac{\partial}{\partial a_i} \log p(\mathbf{v}; \boldsymbol{\theta}) = \frac{v_i}{\sigma_i^2} - \left\langle \frac{u_i}{\sigma_i^2} \right\rangle_{p(\mathbf{u}; \boldsymbol{\theta})} \quad (30)$$

$$\frac{\partial}{\partial b_j} \log p(\mathbf{v}; \boldsymbol{\theta}) = \bar{h}_j(\mathbf{v}, \boldsymbol{\theta}) - \langle \bar{h}_j(\mathbf{u}, \boldsymbol{\theta}) \rangle_{p(\mathbf{u}; \boldsymbol{\theta})} \quad (31)$$

$$\begin{aligned} \frac{\partial}{\partial \sigma_i} \log p(\mathbf{v}; \boldsymbol{\theta}) &= \frac{(v_i - a_i)^2}{\sigma_i^3} - \frac{v_i}{2\sigma_i^3} \sum_j W_{ij} \cdot \bar{h}_j(\mathbf{v}, \boldsymbol{\theta}) \\ &\quad - \left\langle \frac{(u_i - a_i)^2}{\sigma_i^3} - \frac{u_i}{2\sigma_i^3} \sum_j W_{ij} \cdot \bar{h}_j(\mathbf{u}, \boldsymbol{\theta}) \right\rangle_{p(\mathbf{u}; \boldsymbol{\theta})} \end{aligned} \quad (32)$$

To make learning of the visible standard deviations (or variances) more stable, Cho et al. [1, pp. 2–3] propose to reparameterize the energy function by substituting  $\sigma_i^2$  with  $\exp(z_i)$  and learn  $z_i$ , the logarithms of the variances. The reparameterized energy function becomes:

$$E(\mathbf{v}, \mathbf{h}, \boldsymbol{\theta}) = \sum_i \frac{(v_i - a_i)^2}{2 \exp(z_i)} - \sum_{i,j} \frac{v_i}{\exp(z_i)} W_{ij} h_j - \sum_j h_j b_j \quad (33)$$

The corresponding gradient becomes:

$$\begin{aligned} G_{z_i}(\mathbf{v}, \boldsymbol{\theta}) &= \sum_{\mathbf{g}} p(\mathbf{g}|\mathbf{v}; \boldsymbol{\theta}) \cdot \frac{\partial}{\partial z_i} (-E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta})) \\ &= \sum_{\mathbf{g}} p(\mathbf{g}|\mathbf{v}; \boldsymbol{\theta}) \cdot \left( \frac{(v_i - a_i)^2}{2 \exp(z_i)} - \frac{v_i}{\exp(z_i)} \sum_j W_{ij} g_j \right) \\ &= \frac{(v_i - a_i)^2}{2 \exp(z_i)} - \frac{v_i}{\exp(z_i)} \sum_j W_{ij} \cdot \bar{h}_j(\mathbf{v}, \boldsymbol{\theta}) \end{aligned}$$

$$\begin{aligned}
\frac{\partial}{\partial z_i} \log p(\mathbf{v}; \boldsymbol{\theta}) &= \frac{(v_i - a_i)^2}{2 \exp(z_i)} - \frac{v_i}{\exp(z_i)} \sum_j W_{ij} \cdot \bar{h}_j(\mathbf{v}, \boldsymbol{\theta}) \\
&\quad - \left\langle \frac{(u_i - a_i)^2}{2 \exp(z_i)} - \frac{u_i}{\exp(z_i)} \sum_j W_{ij} \cdot \bar{h}_j(\mathbf{u}, \boldsymbol{\theta}) \right\rangle_{p(\mathbf{u}; \boldsymbol{\theta})}
\end{aligned} \tag{34}$$

## 5 Gaussian-Bernoulli RBM (Goodfellow's variant)

Energy function [4, p. 8]:

$$E(\mathbf{v}, \mathbf{h}, \boldsymbol{\theta}) = \sum_i \frac{(v_i - a_i)^2}{2\sigma_i^2} - \sum_{i,j} \frac{v_i - a_i}{\sigma_i^2} W_{ij} h_j - \sum_j h_j b_j \quad (35)$$

### 5.1 Conditional distribution of the visible units

Note that the energy function can be rewritten as:

$$E(\mathbf{v}, \mathbf{h}, \boldsymbol{\theta}) = \sum_i \frac{(v_i - a_i)^2}{2\sigma_i^2} - \sum_{i,j} \frac{v_i}{\sigma_i^2} W_{ij} h_j + \sum_{i,j} \frac{a_i}{\sigma_i^2} W_{ij} h_j - \sum_j h_j b_j \quad (36)$$

Compared to (26) this just adds a term that is independent of  $\mathbf{v}$ . In the derivation of  $p(\mathbf{v}|\mathbf{h}; \boldsymbol{\theta})$  this term cancels just like the  $\sum_j h_j b_j$  term (cf. Section 3.1), so the resulting conditional distribution is identical to Cho's variant:

$$p(\mathbf{v}|\mathbf{h}; \boldsymbol{\theta}) = \prod_i \mathcal{N}\left(v_i; a_i + \sum_j W_{ij} h_j, \sigma_i^2\right) \quad (37)$$

### 5.2 Conditional distribution of the hidden units

The derivation is very similar to the one in Section 3.2:

$$\begin{aligned} p(h_k = 1|\mathbf{v}) &= \frac{1}{1 + \frac{\sum_{\{\mathbf{g}|g_k=0\}} \exp(-E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta}))}{\sum_{\{\mathbf{g}|g_k=1\}} \exp(-E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta}))}} \\ &\stackrel{1}{=} \frac{1}{1 + \frac{\sum_{\{\mathbf{g}|g_k=0\}} \exp(\sum_{i,j} \frac{v_i - a_i}{\sigma_i^2} W_{ij} g_j - \sum_i \frac{(v_i - a_i)^2}{\sigma_i^2} + \sum_j g_j b_j)}{\sum_{\{\mathbf{g}|g_k=1\}} \exp(\sum_{i,j} \frac{v_i - a_i}{\sigma_i^2} W_{ij} g_j - \sum_i \frac{(v_i - a_i)^2}{\sigma_i^2} + \sum_j g_j b_j)}} \\ &\stackrel{2}{=} \frac{1}{1 + \frac{\exp(\sum_i \frac{v_i - a_i}{\sigma_i^2} W_{ik} \cdot 0 + 0 \cdot b_k)}{\exp(\sum_i \frac{v_i - a_i}{\sigma_i^2} W_{ik} \cdot 1 + 1 \cdot b_k)}} \\ &= \frac{1}{1 + \frac{1}{\exp(\sum_i \frac{v_i - a_i}{\sigma_i^2} W_{ik} + b_k)}} \\ &= \sigma\left(b_k + \sum_i \frac{v_i - a_i}{\sigma_i^2} W_{ik}\right) \end{aligned} \quad (38)$$

Explanation of steps marked with a number:

1. We substitute the energy function (35).
2. We substitute the known values for  $g_j = g_k$  and cancel all terms not depending on  $g_k$ .

### 5.3 Log likelihood gradient

Again, the gradient derivations are similar to Section 3.3. For the connection weight  $W_{ij}$ , hidden bias  $b_j$  and visible standard deviation  $\sigma_i$ , the gradients only change a little. For the visible bias, we will take a more detailed look at the derivation:

$$\begin{aligned} G_{a_i}(\mathbf{v}, \boldsymbol{\theta}) &= \sum_{\mathbf{g}} p(\mathbf{g}|\mathbf{v}; \boldsymbol{\theta}) \cdot \frac{\partial}{\partial a_i} (-E(\mathbf{v}, \mathbf{g}, \boldsymbol{\theta})) \\ &\stackrel{1}{=} \sum_{\mathbf{g}} p(\mathbf{g}|\mathbf{v}; \boldsymbol{\theta}) \cdot \left( \frac{v_i - a_i}{\sigma_i^2} - \frac{1}{\sigma_i^2} \sum_j W_{ij} g_j \right) \end{aligned}$$

$$\begin{aligned}
&\stackrel{2}{=} \frac{v_i - a_i}{\sigma_i^2} \cdot \sum_{\mathbf{g}} p(\mathbf{g}|\mathbf{v}; \boldsymbol{\theta}) - \frac{1}{\sigma_i^2} \sum_j W_{ij} \cdot \sum_{\mathbf{g}} p(\mathbf{g}|\mathbf{v}; \boldsymbol{\theta}) g_j \\
&\stackrel{3,4}{=} \frac{v_i - a_i}{\sigma_i^2} - \frac{1}{\sigma_i^2} \sum_j W_{ij} \cdot p(h_j = 1|\mathbf{v}; \boldsymbol{\theta}) \\
&\stackrel{5}{=} \frac{v_i - a_i}{\sigma_i^2} - \frac{1}{\sigma_i^2} \sum_j W_{ij} \cdot \bar{h}_j(\mathbf{v}, \boldsymbol{\theta})
\end{aligned}$$

Explanation of steps marked with a number:

1. We substitute the energy function (35) and differentiate.
2. We pull out of the sums over  $\mathbf{g}$  all factors that are independent of  $\mathbf{g}$ .
3. The sum over a probability distribution evaluates to 1 and can be dropped.
4. We substitute (11) from Section 2.3.
5. We substitute the abbreviation (15).

In summary, we obtain:

$$\frac{\partial}{\partial W_{ij}} \log p(\mathbf{v}; \boldsymbol{\theta}) = \frac{v_i - a_i}{\sigma_i^2} \cdot \bar{h}_j(\mathbf{v}, \boldsymbol{\theta}) - \left\langle \frac{u_i - a_i}{\sigma_i^2} \cdot \bar{h}_j(\mathbf{u}, \boldsymbol{\theta}) \right\rangle_{p(\mathbf{u}; \boldsymbol{\theta})} \quad (39)$$

$$\begin{aligned}
\frac{\partial}{\partial a_i} \log p(\mathbf{v}; \boldsymbol{\theta}) &= \frac{v_i - a_i}{\sigma_i^2} - \frac{1}{\sigma_i^2} \sum_j W_{ij} \cdot \bar{h}_j(\mathbf{v}, \boldsymbol{\theta}) \\
&\quad - \left\langle \frac{u_i - a_i}{\sigma_i^2} - \frac{1}{\sigma_i^2} \sum_j W_{ij} \cdot \bar{h}_j(\mathbf{u}, \boldsymbol{\theta}) \right\rangle_{p(\mathbf{u}; \boldsymbol{\theta})} \\
&= \frac{v_i}{\sigma_i^2} - \frac{1}{\sigma_i^2} \sum_j W_{ij} \cdot \bar{h}_j(\mathbf{v}, \boldsymbol{\theta}) \\
&\quad - \left\langle \frac{u_i}{\sigma_i^2} - \frac{1}{\sigma_i^2} \sum_j W_{ij} \cdot \bar{h}_j(\mathbf{u}, \boldsymbol{\theta}) \right\rangle_{p(\mathbf{u}; \boldsymbol{\theta})} \quad (40)
\end{aligned}$$

$$\frac{\partial}{\partial b_j} \log p(\mathbf{v}; \boldsymbol{\theta}) = \bar{h}_j(\mathbf{v}, \boldsymbol{\theta}) - \langle \bar{h}_j(\mathbf{u}, \boldsymbol{\theta}) \rangle_{p(\mathbf{u}; \boldsymbol{\theta})} \quad (41)$$

$$\begin{aligned}
\frac{\partial}{\partial \sigma_i} \log p(\mathbf{v}; \boldsymbol{\theta}) &= \frac{(v_i - a_i)^2}{\sigma_i^3} - \frac{v_i - a_i}{2\sigma_i^3} \sum_j W_{ij} \cdot \bar{h}_j(\mathbf{v}, \boldsymbol{\theta}) \\
&\quad - \left\langle \frac{(u_i - a_i)^2}{\sigma_i^3} - \frac{u_i - a_i}{2\sigma_i^3} \sum_j W_{ij} \cdot \bar{h}_j(\mathbf{u}, \boldsymbol{\theta}) \right\rangle_{p(\mathbf{u}; \boldsymbol{\theta})} \quad (42)
\end{aligned}$$

As in Section 4.3, we can reparameterize  $\sigma_i^2$  as  $\exp(z_i)$  and learn  $z_i$  instead:

$$\begin{aligned}
\frac{\partial}{\partial z_i} \log p(\mathbf{v}; \boldsymbol{\theta}) &= \frac{(v_i - a_i)^2}{2 \exp(z_i)} - \frac{v_i - a_i}{\exp(z_i)} \sum_j W_{ij} \cdot \bar{h}_j(\mathbf{v}, \boldsymbol{\theta}) \\
&\quad - \left\langle \frac{(u_i - a_i)^2}{2 \exp(z_i)} - \frac{u_i - a_i}{\exp(z_i)} \sum_j W_{ij} \cdot \bar{h}_j(\mathbf{u}, \boldsymbol{\theta}) \right\rangle_{p(\mathbf{u}; \boldsymbol{\theta})} \quad (43)
\end{aligned}$$

An interesting aspect of this GRBM variant is that the visible bias is subtracted in all occurrences of the visible unit states, similar to the *centering trick* [6] discussed in the next section. We will see that this formulation is ill-suited for centering the hidden units, though.

## 6 Centering Trick

The *centering trick* [6] is a temporary reparameterization of an RBM's energy function that helps learning without changing the modelled probability distribution. It introduces a scalar offset for each visible and hidden unit that is updated to subtract its mean activation during training, and (optionally) reset to zero after training to obtain a standard uncentered RBM. In order to keep the probability distribution unchanged and not create a moving target, any changes in the offsets have to be countered by updates to other model parameters. In this section, we will show how to derive these reparameterization rules for the RBM variants considered before.

### 6.1 Bernoulli-Bernoulli RBM

We begin by introducing offset vectors  $\boldsymbol{\alpha}$  and  $\boldsymbol{\beta}$  for the visible and hidden units, respectively, resulting in the following energy function:

$$E(\mathbf{v}, \mathbf{h}, \boldsymbol{\theta}) = -\sum_i (v_i - \alpha_i) a_i - \sum_{i,j} (v_i - \alpha_i) W_{ij} (h_j - \beta_j) - \sum_j (h_j - \beta_j) b_j \quad (44)$$

$$= -(\mathbf{v} - \boldsymbol{\alpha})^T \mathbf{a} - (\mathbf{v} - \boldsymbol{\alpha})^T \mathbf{W} (\mathbf{h} - \boldsymbol{\beta}) - (\mathbf{h} - \boldsymbol{\beta})^T \mathbf{b} \quad (45)$$

The substitution of  $v_i$  by  $(v_i - \alpha_i)$  and of  $h_j$  by  $(h_j - \beta_j)$  simply carries through the derivations of the conditional probabilities (9) and (10) and the log likelihood gradients in Section 2.3.

When recentering the hidden units, we set  $\boldsymbol{\beta}$  to a new value  $\boldsymbol{\beta}'$  (during training, this will be the mean activation of hidden units, usually computed as a running average as per line 20 of Fig. 3 in [6], and after training this may be zero to eliminate the offset terms). We want to ensure that the modelled probability distribution does not change, so we need to choose the remaining model parameters in  $\boldsymbol{\theta}'$  to fulfill

$$p(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta}) = p(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta}'). \quad (46)$$

One way to satisfy this constraint is by requiring the energy function to only change by a constant term (which will drop out as a common factor in the nominator and denominator of (2)), as in [6, Sect. 3.1]:

$$E(\mathbf{v}, \mathbf{h}, \boldsymbol{\theta}) = E(\mathbf{v}, \mathbf{h}, \boldsymbol{\theta}') + \text{const} \quad (47)$$

Thus, we insert (45) into (47) and solve for  $\mathbf{a}'$ ,  $\mathbf{b}'$ . We set  $\boldsymbol{\alpha}' = \boldsymbol{\alpha}$  because we need to be able to choose the visible unit offsets independently of the hidden unit offsets. We set  $\mathbf{W}' = \mathbf{W}$  because we cannot use it to counter the offsets.

$$\begin{aligned} E(\mathbf{v}, \mathbf{h}, \boldsymbol{\theta}) &= E(\mathbf{v}, \mathbf{h}, \boldsymbol{\theta}') + \text{const} \\ -(\mathbf{v} - \boldsymbol{\alpha})^T \mathbf{a} - (\mathbf{v} - \boldsymbol{\alpha})^T \mathbf{W} (\mathbf{h} - \boldsymbol{\beta}) &= -(\mathbf{v} - \boldsymbol{\alpha})^T \mathbf{a}' - (\mathbf{v} - \boldsymbol{\alpha})^T \mathbf{W} (\mathbf{h} - \boldsymbol{\beta}') \\ -(\mathbf{h} - \boldsymbol{\beta})^T \mathbf{b} &= -(\mathbf{h} - \boldsymbol{\beta}')^T \mathbf{b}' + \text{const} \\ \frac{(\mathbf{v} - \boldsymbol{\alpha})^T (\mathbf{a}' - \mathbf{a})}{-(\mathbf{h} - \boldsymbol{\beta})^T \mathbf{b}} &\stackrel{1}{=} \frac{(\mathbf{v} - \boldsymbol{\alpha})^T \mathbf{W} (\mathbf{h} - \boldsymbol{\beta} - \mathbf{h} + \boldsymbol{\beta}')}{-(\mathbf{h} - \boldsymbol{\beta}')^T \mathbf{b}' + \text{const}} \\ \frac{\mathbf{v}^T (\mathbf{a}' - \mathbf{a}) - \boldsymbol{\alpha}^T (\mathbf{a}' - \mathbf{a})}{-\mathbf{h}^T \mathbf{b} + \boldsymbol{\beta}^T \mathbf{b}} &= \frac{\mathbf{v}^T \mathbf{W} (\boldsymbol{\beta}' - \boldsymbol{\beta}) - \boldsymbol{\alpha}^T \mathbf{W} (\boldsymbol{\beta}' - \boldsymbol{\beta})}{-\mathbf{h}^T \mathbf{b}' + \boldsymbol{\beta}'^T \mathbf{b}' + \text{const}} \\ \frac{\mathbf{v}^T (\mathbf{a}' - \mathbf{a})}{-\mathbf{h}^T \mathbf{b}} &\stackrel{2}{=} \frac{\mathbf{v}^T \mathbf{W} (\boldsymbol{\beta}' - \boldsymbol{\beta})}{-\mathbf{h}^T \mathbf{b}' + \text{const}} \\ \mathbf{a}' - \mathbf{a} &= \mathbf{W} (\boldsymbol{\beta}' - \boldsymbol{\beta}) \\ \mathbf{b}' &= \mathbf{b} \end{aligned}$$

Explanation of steps marked with a number:

1. We add  $(\mathbf{v} - \boldsymbol{\alpha})^T \mathbf{a}'$  and  $(\mathbf{v} - \boldsymbol{\alpha})^T \mathbf{W} (\mathbf{h} - \boldsymbol{\beta})$  to both sides of the equation.
2. We absorb any terms not involving  $\mathbf{v}$  or  $\mathbf{h}$  in the constant term.

So we find that when we change the hidden unit offsets from  $\boldsymbol{\beta}$  to  $\boldsymbol{\beta}'$ , we need to update the visible biases to  $\mathbf{a}' = \mathbf{a} + \mathbf{W} (\boldsymbol{\beta}' - \boldsymbol{\beta})$  and leave the hidden unit biases unchanged – i.e., we use the visible unit biases to exactly compensate the change in top-down input caused by changing the hidden unit offsets.

An alternative way to satisfy (46) is by requiring the following set of equations to hold:

$$p(\mathbf{v}|\mathbf{h};\boldsymbol{\theta}') = p(\mathbf{v}|\mathbf{h};\boldsymbol{\theta}) \quad (48a)$$

$$p(\mathbf{h}|\mathbf{v};\boldsymbol{\theta}') = p(\mathbf{h}|\mathbf{v};\boldsymbol{\theta}) \quad (48b)$$

We insert the centered version of (9) into (48a) and the centered version of (10) into (48b) and solve the system of equations for  $\mathbf{a}, \mathbf{b}$ , again setting  $\mathbf{W}' = \mathbf{W}$  and  $\boldsymbol{\alpha}' = \boldsymbol{\alpha}$  beforehand.

$$\begin{cases} \forall_k & \sigma\left(a'_k + \sum_j W_{kj}(h_j - \beta'_j)\right) = \sigma\left(a_k + \sum_j W_{kj}(h_j - \beta_j)\right) \\ \forall_k & \sigma\left(b'_k + \sum_i (v_i - \alpha_i)W_{ik}\right) = \sigma\left(b_k + \sum_i (v_i - \alpha_i)W_{ik}\right) \\ \forall_k & a'_k + \sum_j W_{kj}(h_j - \beta'_j) = a_k + \sum_j W_{kj}(h_j - \beta_j) \\ \forall_k & b'_k = b_k \\ \forall_k & \begin{cases} a'_k = a_k + \sum_j W_{kj}(h_j - \beta_j - h_j + \beta'_j) \\ = a_k + \sum_j W_{kj}(\beta'_j - \beta_j) \end{cases} \\ \forall_k & b'_k = b_k \end{cases}$$

Not surprisingly, we obtain the same solution  $\mathbf{a}' = \mathbf{a} + \mathbf{W}(\boldsymbol{\beta}' - \boldsymbol{\beta})$  and  $\mathbf{b}' = \mathbf{b}$ . However, for the Gaussian-Binary RBMs, this way of deriving the reparameterization rule results in far less writing, so we will use this for the remaining derivations.

When recentering the visible units, we need to set  $\mathbf{b}' = \mathbf{b} + (\boldsymbol{\alpha}' - \boldsymbol{\alpha})^T \mathbf{W}$  and  $\mathbf{a}' = \mathbf{a}$ . The derivation is symmetric to the ones above. Note that when the offsets are only learned from data-dependent states as in [6] and not from model samples,  $\boldsymbol{\alpha}$  can be initialized to the mean of training data and left untouched during training. For an RBM (as opposed to a Deep Boltzmann Machine), this also avoids the need of reparameterizing the hidden unit biases  $\mathbf{b}$ .

## 6.2 Gaussian-Bernoulli RBMs

We start with the standard Gaussian-Bernoulli RBM as used by Krizhevsky [5] and discussed in Section 3. As before, we introduce offset vectors  $\boldsymbol{\alpha}$  and  $\boldsymbol{\beta}$  for the visible and hidden units, respectively, obtaining the following energy function:

$$E(\mathbf{v}, \mathbf{h}, \boldsymbol{\theta}) = \sum_i \frac{(v_i - a_i - \alpha_i)^2}{2\sigma_i^2} - \sum_{i,j} \frac{v_i - \alpha_i}{\sigma_i} W_{ij}(h_j - \beta_j) - \sum_j (h_j - \beta_j)b_j \quad (49)$$

$$= \frac{1}{2}(\mathbf{v} - \mathbf{a} - \boldsymbol{\alpha})^T \boldsymbol{\Lambda}(\mathbf{v} - \mathbf{a} - \boldsymbol{\alpha}) - (\mathbf{v} - \boldsymbol{\alpha})^T \boldsymbol{\Lambda}^{\frac{1}{2}} \mathbf{W}(\mathbf{h} - \boldsymbol{\beta}) - (\mathbf{h} - \boldsymbol{\beta})^T \mathbf{b}, \quad (50)$$

where  $\boldsymbol{\Lambda}$  is a diagonal matrix of precisions, or inverse variances,  $\sigma_i^{-2}$ . The substitutions of  $v_i$  by  $(v_i - \alpha_i)$  and of  $h_j$  by  $(h_j - \beta_j)$  carry through the derivations in Sections 3–5 except for a minor change in  $p(\mathbf{v}|\mathbf{h})$ . To derive the reparameterization rule for recentering the hidden units, we insert the centered version of (20) into (48a) and the centered version of (21) into (48b) and solve the system of equations for  $\mathbf{a}, \mathbf{b}, \boldsymbol{\Lambda}$ , setting  $\mathbf{W}' = \mathbf{W}$  and  $\boldsymbol{\alpha}' = \boldsymbol{\alpha}$  beforehand:

$$\begin{cases} \prod_i \mathcal{N}\left(v_i; a'_i + \alpha_i + \sigma'_i \sum_j W_{ij}(h_j - \beta'_j), \sigma_i'^2\right) = \prod_i \mathcal{N}\left(v_i; a_i + \alpha_i + \sigma_i \sum_j W_{ij}(h_j - \beta_j), \sigma_i^2\right) \\ \forall_k & \sigma\left(b'_k + \sum_i \frac{v_i - \alpha_i}{\sigma'_i} W_{ik}\right) = \sigma\left(b_k + \sum_i \frac{v_i - \alpha_i}{\sigma_i} W_{ik}\right) \\ \forall_i & \sigma'_i = \sigma_i \\ \forall_i & a'_i + \alpha_i + \sigma_i \sum_j W_{ij}(h_j - \beta'_j) = a_i + \alpha_i + \sigma_i \sum_j W_{ij}(h_j - \beta_j) \\ \forall_k & b'_k + \sum_i \frac{v_i - \alpha_i}{\sigma_i} W_{ik} = b_k + \sum_i \frac{v_i - \alpha_i}{\sigma_i} W_{ik} \end{cases}$$

$$\left\{ \begin{array}{l} \forall_i \quad \sigma'_i = \sigma_i \\ \forall_i \quad a'_i = a_i + \sigma_i \sum_j W_{ij}(h_j - \beta_j - h_j + \beta'_j) \\ \qquad \qquad = a_i + \sigma_i \sum_j W_{ij}(\beta'_j - \beta_j) \\ \forall_k \quad b'_k = b_k \end{array} \right.$$

We see that the hidden biases and the visible standard deviations stay unmodified, and we need to set  $\mathbf{a}' = \mathbf{a} + \mathbf{\Lambda}^{-\frac{1}{2}} \mathbf{W}(\boldsymbol{\beta}' - \boldsymbol{\beta})$  to counter the change of hidden offsets.

For recentering the visible units, we solve the same system of equations for  $\mathbf{a}, \mathbf{b}, \mathbf{\Lambda}$ , this time setting  $\mathbf{W}' = \mathbf{W}$  and  $\boldsymbol{\beta}' = \boldsymbol{\beta}$  beforehand:

$$\left\{ \begin{array}{l} \prod_i \mathcal{N}(v_i; a'_i + \alpha'_i + \sigma'_i \sum_j W_{ij}(h_j - \beta_j), \sigma_i'^2) = \prod_i \mathcal{N}(v_i; a_i + \alpha_i + \sigma_i \sum_j W_{ij}(h_j - \beta_j), \sigma_i^2) \\ \forall_k \quad \sigma(b'_k + \sum_i \frac{v_i - \alpha'_i}{\sigma'_i} W_{ik}) = \sigma(b_k + \sum_i \frac{v_i - \alpha_i}{\sigma_i} W_{ik}) \end{array} \right.$$

$$\left\{ \begin{array}{l} \forall_i \quad \sigma'_i = \sigma_i \\ \forall_i \quad a'_i + \alpha'_i + \sigma_i \sum_j W_{ij}(h_j - \beta_j) = a_i + \alpha_i + \sigma_i \sum_j W_{ij}(h_j - \beta_j) \\ \forall_k \quad b'_k + \sum_i \frac{v_i - \alpha'_i}{\sigma_i} W_{ik} = b_k + \sum_i \frac{v_i - \alpha_i}{\sigma_i} W_{ik} \end{array} \right.$$

$$\left\{ \begin{array}{l} \forall_i \quad \sigma'_i = \sigma_i \\ \forall_i \quad a'_i = a_i + \alpha_i - \alpha'_i \\ \forall_k \quad b'_k = b_k + \sum_i \frac{v_i - \alpha_i - v_i + \alpha'_i}{\sigma_i} W_{ik} \\ \qquad \qquad = b_k + \sum_i \frac{v_i + \alpha'_i - \alpha_i}{\sigma_i} W_{ik} \end{array} \right.$$

Here, both the visible and hidden unit biases have to be adapted. When changing the visible unit offsets from  $\boldsymbol{\alpha}$  to  $\boldsymbol{\alpha}'$ , we have to counter that by setting  $\mathbf{a}' = \mathbf{a} + \boldsymbol{\alpha} - \boldsymbol{\alpha}'$  and  $\mathbf{b}' = \mathbf{b} + (\boldsymbol{\alpha}' - \boldsymbol{\alpha})^T \mathbf{\Lambda}^{\frac{1}{2}} \mathbf{W}$ .

For Cho's variant of the Gaussian-Bernoulli RBM discussed in Section 4, nothing much changes. As noted in Section 4.1, its energy function can be derived from Krizhevsky's variant by simply replacing  $W_{ij}$  with  $\frac{1}{\sigma_i} W_{ij}$  (or  $\mathbf{W}$  with  $\mathbf{\Lambda}^{\frac{1}{2}} \mathbf{W}$ ), resulting in:

$$E(\mathbf{v}, \mathbf{h}, \boldsymbol{\theta}) = \sum_i \frac{(v_i - a_i - \alpha_i)^2}{2\sigma_i^2} - \sum_{i,j} \frac{v_i - \alpha_i}{\sigma_i^2} W_{ij}(h_j - \beta_j) - \sum_j (h_j - \beta_j) b_j \quad (51)$$

$$= \frac{1}{2} (\mathbf{v} - \mathbf{a} - \boldsymbol{\alpha})^T \mathbf{\Lambda} (\mathbf{v} - \mathbf{a} - \boldsymbol{\alpha})^T - (\mathbf{v} - \boldsymbol{\alpha})^T \mathbf{\Lambda} \mathbf{W} (\mathbf{h} - \boldsymbol{\beta}) - (\mathbf{h} - \boldsymbol{\beta})^T \mathbf{b} \quad (52)$$

This substitution carries through the derivation of the reparameterization rule: When updating the visible and hidden offsets from  $\boldsymbol{\alpha}$  to  $\boldsymbol{\alpha}'$  and  $\boldsymbol{\beta}$  to  $\boldsymbol{\beta}'$ , the biases need to be set according to  $\mathbf{a}' = \mathbf{a} + \boldsymbol{\alpha} - \boldsymbol{\alpha}' + \mathbf{W}(\boldsymbol{\beta}' - \boldsymbol{\beta})$  and  $\mathbf{b}' = \mathbf{b} + (\boldsymbol{\alpha}' - \boldsymbol{\alpha})^T \mathbf{\Lambda} \mathbf{W}$ .

For Goodfellow's variant discussed in Section 5, the centered energy function would be:

$$E(\mathbf{v}, \mathbf{h}, \boldsymbol{\theta}) = \sum_i \frac{(v_i - a_i - \alpha_i)^2}{2\sigma_i^2} - \sum_{i,j} \frac{v_i - a_i - \alpha_i}{\sigma_i^2} W_{ij}(h_j - \beta_j) - \sum_j (h_j - \beta_j) b_j \quad (53)$$

$$= \frac{1}{2} (\mathbf{v} - \mathbf{a} - \boldsymbol{\alpha})^T \mathbf{\Lambda} (\mathbf{v} - \mathbf{a} - \boldsymbol{\alpha})^T - (\mathbf{v} - \mathbf{a} - \boldsymbol{\alpha})^T \mathbf{\Lambda} \mathbf{W} (\mathbf{h} - \boldsymbol{\beta}) - (\mathbf{h} - \boldsymbol{\beta})^T \mathbf{b} \quad (54)$$

Note that the visible bias  $\mathbf{a}$  and offset  $\boldsymbol{\alpha}$  assume the same role, so they could be merged. We will leave them separated for clarity. As before, we derive the reparameterization rule for updating the hidden



offsets by inserting the centered versions of (37) and (38) into (48), setting  $\mathbf{W}' = \mathbf{W}$ ,  $\boldsymbol{\alpha}' = \boldsymbol{\alpha}$  and solving for  $\mathbf{a}$ ,  $\mathbf{b}$ ,  $\boldsymbol{\Lambda}$ :

$$\left\{ \begin{array}{l} \prod_i \mathcal{N}\left(v_i; a'_i + \alpha_i + \sum_j W_{ij}(h_j - \beta'_j), \sigma_i'^2\right) = \prod_i \mathcal{N}\left(v_i; a_i + \alpha_i + \sum_j W_{ij}(h_j - \beta_j), \sigma_i^2\right) \\ \forall_k \quad \sigma\left(b'_k + \sum_i \frac{v_i - a'_i - \alpha_i}{\sigma_i'^2} W_{ik}\right) = \sigma\left(b_k + \sum_i \frac{v_i - a_i - \alpha_i}{\sigma_i^2} W_{ik}\right) \end{array} \right.$$

$$\left\{ \begin{array}{l} \forall_i \quad \sigma'_i = \sigma_i \\ \forall_i \quad a'_i + \alpha_i + \sum_j W_{ij}(h_j - \beta'_j) = a_i + \alpha_i + \sum_j W_{ij}(h_j - \beta_j) \\ \forall_k \quad b'_k + \sum_i \frac{v_i - a'_i - \alpha_i}{\sigma_i'^2} W_{ik} = b_k + \sum_i \frac{v_i - a_i - \alpha_i}{\sigma_i^2} W_{ik} \end{array} \right.$$

$$\left\{ \begin{array}{l} \forall_i \quad \sigma'_i = \sigma_i \\ \forall_i \quad a'_i = a_i + \sum_j W_{ij}(\beta'_j - \beta_j) \\ \forall_k \quad b'_k = b_k + \sum_i \frac{a'_i - a_i}{\sigma_i^2} W_{ik} \\ \quad \quad = b_k + \sum_i \frac{\sum_j W_{ij}(\beta'_j - \beta_j)}{\sigma_i^2} W_{ik} \end{array} \right.$$

That is, when changing the hidden offsets from  $\boldsymbol{\beta}$  to  $\boldsymbol{\beta}'$ , we need to set  $\mathbf{a}' = \mathbf{a} + \mathbf{W}(\boldsymbol{\beta}' - \boldsymbol{\beta})$  and  $\mathbf{b}' = \mathbf{b} + (\mathbf{a} - \mathbf{a}')^T \boldsymbol{\Lambda} \mathbf{W} = \mathbf{b} + (\mathbf{W}(\boldsymbol{\beta}' - \boldsymbol{\beta}))^T \boldsymbol{\Lambda} \mathbf{W} = \mathbf{b} + (\boldsymbol{\beta}' - \boldsymbol{\beta})^T \mathbf{W}^T \boldsymbol{\Lambda} \mathbf{W}$ .

For recentering the visible units, we solve the same system of equations for  $\mathbf{a}$ ,  $\mathbf{b}$ ,  $\boldsymbol{\Lambda}$ , setting  $\mathbf{W}' = \mathbf{W}$  and  $\boldsymbol{\beta}' = \boldsymbol{\beta}$  beforehand:

$$\left\{ \begin{array}{l} \prod_i \mathcal{N}\left(v_i; a'_i + \alpha'_i + \sum_j W_{ij}(h_j - \beta_j), \sigma_i'^2\right) = \prod_i \mathcal{N}\left(v_i; a_i + \alpha_i + \sum_j W_{ij}(h_j - \beta_j), \sigma_i^2\right) \\ \forall_k \quad \sigma\left(b'_k + \sum_i \frac{v_i - a'_i - \alpha'_i}{\sigma_i'^2} W_{ik}\right) = \sigma\left(b_k + \sum_i \frac{v_i - a_i - \alpha_i}{\sigma_i^2} W_{ik}\right) \end{array} \right.$$

$$\left\{ \begin{array}{l} \forall_i \quad \sigma'_i = \sigma_i \\ \forall_i \quad a'_i + \alpha'_i + \sum_j W_{ij}(h_j - \beta_j) = a_i + \alpha_i + \sum_j W_{ij}(h_j - \beta_j) \\ \forall_k \quad b'_k + \sum_i \frac{v_i - a'_i - \alpha'_i}{\sigma_i'^2} W_{ik} = b_k + \sum_i \frac{v_i - a_i - \alpha_i}{\sigma_i^2} W_{ik} \end{array} \right.$$

$$\left\{ \begin{array}{l} \forall_i \quad \sigma'_i = \sigma_i \\ \forall_i \quad a'_i = a_i + \alpha_i - \alpha'_i \\ \forall_k \quad b'_k = b_k + \sum_i \frac{a_i - a'_i + \alpha_i - \alpha'_i}{\sigma_i^2} W_{ik} \\ \quad \quad = b_k + \sum_i \frac{a_i - a_i - \alpha_i + \alpha'_i + \alpha_i - \alpha'_i}{\sigma_i^2} W_{ik} \\ \quad \quad = b_k \end{array} \right.$$

That is, when changing the visible offsets from  $\boldsymbol{\alpha}$  to  $\boldsymbol{\alpha}'$ , we need to set  $\mathbf{a}' = \mathbf{a} - \boldsymbol{\alpha}'$  and  $\mathbf{b}' = \mathbf{b}$ . This is of course a much simpler reparameterization rule than for the other two Gaussian-Bernoulli RBM variants, but given that the visible offsets  $\boldsymbol{\alpha}$  are typically set in advance and not changed during training, this is only a virtual advantage and does not outweigh the disadvantage of the much more involved reparameterization rule for the hidden offsets derived before.

## Acknowledgements

Thanks to Alex Krizhevsky and Asja Fischer et al. for publishing their respective derivations [5, 3], and thanks to Karen Ullrich for helping to get the derivation in 3.1 straight.

## References

- [1] KyungHyun Cho, Alexander Ilin, and Tapani Raiko. Improved learning of gaussian-bernoulli restricted boltzmann machines. In *Proceedings of the 21st International Conference on Artificial Neural Networks (ICANN)*, pages 10–17, Espoo, Finland, 2011.
- [2] Deep learning tutorials: Restricted Boltzmann Machines (RBM). <http://deeplearning.net/tutorial/rbm.html>, January 2014.
- [3] Asja Fischer and Christian Igel. An introduction to restricted boltzmann machines. In *Proceedings of the 17th Iberoamerican Congress on Pattern Recognition (CIARP)*, Buenos Aires, Argentina, 2012.
- [4] Ian Goodfellow, Mehdi Mirza, Aaron Courville, and Yoshua Bengio. Multi-prediction deep boltzmann machines. In *Advances in Neural Information Processing Systems (NIPS) 26*, pages 548–556. 2013.
- [5] Alex Krizhevsky. Learning multiple layers of features from tiny images. Master’s thesis, Dept. of Comp. Science, Univ. of Toronto, 2009.
- [6] Grégoire Montavon and Klaus-Robert Müller. Deep boltzmann machines and the centering trick. In Grégoire Montavon, Geneviève B. Orr, and Klaus-Robert Müller, editors, *Neural Networks: Tricks of the Trade*, volume 7700 of *Lecture Notes in Computer Science*, pages 621–637. Springer Berlin Heidelberg, 2012.
- [7] Jan Schlüter. Unsupervised Audio Feature Extraction for Music Similarity Estimation. Master’s thesis, Technische Universität München, Munich, Germany, October 2011.