

There are many ways of parametrizing Gaussian-Bernoulli RBMs. One choice is whether to use a covariance matrix or a precision matrix for the Gaussian distribution. Here we present the precision formulation. The modification to obtain the covariance formulation is straightforward. We wish to have the conditional distribution

$$p(\mathbf{v} \mid \mathbf{h}) = \mathcal{N}(\mathbf{v}; \mathbf{W}\mathbf{h}, \boldsymbol{\beta}^{-1}). \quad (20.38)$$

We can find the terms we need to add to the energy function by expanding the unnormalized log conditional distribution:

$$\log \mathcal{N}(\mathbf{v}; \mathbf{W}\mathbf{h}, \boldsymbol{\beta}^{-1}) = -\frac{1}{2} (\mathbf{v} - \mathbf{W}\mathbf{h})^\top \boldsymbol{\beta} (\mathbf{v} - \mathbf{W}\mathbf{h}) + f(\boldsymbol{\beta}). \quad (20.39)$$

Here  $f$  encapsulates all the terms that are a function only of the parameters and not the random variables in the model. We can discard  $f$  because its only role is to normalize the distribution, and the partition function of whatever energy function we choose will carry out that role.

If we include all of the terms (with their sign flipped) involving  $\mathbf{v}$  from equation 20.39 in our energy function and do not add any other terms involving  $\mathbf{v}$ , then our energy function will represent the desired conditional  $p(\mathbf{v} \mid \mathbf{h})$ .

We have some freedom regarding the other conditional distribution,  $p(\mathbf{h} \mid \mathbf{v})$ . Note that equation 20.39 contains a term

$$\frac{1}{2} \mathbf{h}^\top \mathbf{W}^\top \boldsymbol{\beta} \mathbf{W} \mathbf{h}. \quad (20.40)$$

This term cannot be included in its entirety because it includes  $h_i h_j$  terms. These correspond to edges between the hidden units. If we included these terms, we would have a linear factor model instead of a restricted Boltzmann machine. When designing our Boltzmann machine, we simply omit these  $h_i h_j$  cross terms. Omitting them does not change the conditional  $p(\mathbf{v} \mid \mathbf{h})$  so equation 20.39 is still respected. However, we still have a choice about whether to include the terms involving only a single  $h_i$ . If we assume a diagonal precision matrix, we find that for each hidden unit  $h_i$  we have a term

$$\frac{1}{2} h_i \sum_j \beta_j W_{j,i}^2. \quad (20.41)$$

In the above, we used the fact that  $h_i^2 = h_i$  because  $h_i \in \{0, 1\}$ . If we include this term (with its sign flipped) in the energy function, then it will naturally bias  $h_i$  to be turned off when the weights for that unit are large and connected to visible units with high precision. The choice of whether or not to include this bias term does not affect the family of distributions the model can represent (assuming that