

Some notes about BB- and GB-RBM implementation

Tomas Ukkonen, 2015 `tomas.ukkonen@iki.fi`

Restricted Boltzmann machines are tricky to implement as there doesn't seem to be very good reference material around describing all the details of the implementation. Here are some notes about RBM implementation in `..src/neuralnetwork/RBM.h` and `..src/neuralnetwork/RBM.cpp`. The learning method to be used is CD-10.

I couldn't find update rules for bias-terms used to calculate approximate gradient descent for (contrastive divergence) so I used the following approach:

- there are no separate bias-terms
- instead visible \mathbf{v} and hidden \mathbf{h} states are always *extended* to have one extra 1 after the data terms, so the actual terms given to the RBM are $[\mathbf{v}, 1]$ and $[\mathbf{h}, 1]$

This then makes the calculations relatively simple as you will only have to care about the weight matrix \mathbf{W} and can ignore other terms. How to handle these equal to 1 states during the stimulation phase of RBM in CD? The correct approach seems to handle them similarly to all other variables, in other words we just simulate and simulate their states normally. This then leads to situation where the bias terms automatically find rules that always set those values to 1.

TOY PROBLEM ANALYSIS

The function `rbm_test()` in `src/neuralnetwork/tst/test.cpp` generates a toy problem for which RBM CD-10 algorithm is tested. The data is generated as follows:

- hidden states (2) are selected randomly between 0 and 1
- after this visible states (4) are generated from hidden states by directly setting the k :th state exactly same as the $k/2$:th state in a hidden state
- in practice $[0, 1] \Rightarrow [0, 0, 1, 1], [1, 0] \Rightarrow [1, 1, 0, 0], [1, 1] \Rightarrow [1, 1, 1, 1]$ and $[0, 0] \Rightarrow [0, 0, 0, 0]$
- so the relationship between hidden states and input states is really simple, the easiest way is to simply to calculate mean value of every 2 items in visible layer and used it as a hidden layer

The results of for running CD-10 algorithm which discretizes or samples all values to 0 and 1 are always the same, the weight matrix for hidden to visible layers have the form:

$$\mathbf{W} = \begin{array}{|c|c|c|c|c|} \hline 4.6 & 4.4 & 4.3 & 4.5 & 5.6 \\ \hline 0.0 & -0.1 & -9.7 & -9.8 & 8.9 \\ \hline -9.7 & -9.7 & 0.09 & 0.2 & 8.7 \\ \hline \end{array}$$

This weight matrix clearly generates three hidden states where the first one is equal to the state 1 which is always on. Because binary values only take values 0 and 1 the sigmoid function only gets positive values (the last term is 5.6 meaning that the all 1 term in input layer sets value to be extremely likely to be 1 even when other terms are 0). After this, the second term sets the term to be 0 when the second part of the inputs are 1 and 0 when the inputs are 1 because in this case the sigmoidal value will be negative. And similarly for the last term.

Given this interpretation, the reversed weight matrix \mathbf{W}^T makes now also sense:

$$\mathbf{W}^T = \begin{array}{|c|c|c|} \hline 4.6 & +0.1 & -9.7 \\ \hline 4.4 & -0.1 & -9.7 \\ \hline 4.3 & -9.7 & 0.0 \\ \hline 4.5 & -9.7 & 0.2 \\ \hline 5.6 & 8.9 & 8.7 \\ \hline \end{array}$$

If we keep in mind that the first term is always 1 in our reversed. Then the first two terms are 1 only when the third term is zero and zero only when the third term is one (negative sigmoid value). Similar is true for the second set of variables and finally the 5th variable is always one due to fact that state values are always non-negative meaning that sigmoidal value will be always be positive.

*This then proofs and shows that the toy problem is properly solved by **the binary RBM with the CD-10 learning rule**. It can find proper two states and the one artificial always one state using the simple contrastive divergence algorithm.* This shows RBM is capable of data reduction.

Energy based model of RBM

The continuous RBM is more complicated to implement and understand. It seems that Gaussian-Bernoulli RBM could be the right choice (or maybe Beta-Bernoulli RBM) to model continuous input values. Because of this, I try to derive the whole BB-RBM and then continuous GB-RBM theory from energy based models (EBMs). It is important(?) to remember that Bernoulli-Bernoulli RBM model is related to Ising models and statistical physics but hidden variable models in general are not.

The energy of the restricted boltzman machine is:

$$E_{BB}(v, h) = -a^T v - v^T W h - b^T h$$

$$E_{GB}(v, h) = \frac{1}{2} \|v - a\|^2 - v^T W h - b^T h$$

Note that we could have added additional constant c term to these equations but it is not needed because probability distributions are normalized to have unit probability mass.

Now the probability of the (v, h) and only observed variables v is:

$$P(v, h) = \frac{1}{Z} e^{-E(v, h)}, \quad Z = \sum_{v, h} e^{-E(v, h)}$$

$$P(v) = \frac{1}{Z} \sum_h e^{-E(v, h)}$$

Now the hidden variables are Bernoulli distributed, that is, they take only values 0 and 1 which is a strong regularizer for the system.

We want to calculate probabilities of hidden h and visible v neurons:

$$P_{BB}(h|v) = \frac{P(v, h)}{P(v)} = \frac{\frac{1}{Z} e^{-E(v, h)}}{\frac{1}{Z} \sum_h e^{-E(v, h)}} = \frac{e^{v^T W h + a^T v + b^T h}}{\sum_h e^{v^T W h + a^T v + b^T h}}$$

$$P_{BB}(h_i = 1|v) = \frac{e^{+v^T w_i 1 + a^T v + b_i 1}}{e^{+v^T w_i 0 + a^T v + b_i 0} + e^{+v^T w_i 1 + a^T v + b_i 1}} = \frac{e^{+v^T w_i + b_i}}{1 + e^{+v^T w_i + b_i}} = \frac{1}{1 + e^{-v^T w_i - b_i}} = \text{sigmoid}(w_i^T v + b_i)$$

$$P_{BB}(\mathbf{h} = \mathbf{1}|\mathbf{v}) = \text{sigmoid}(\mathbf{W}\mathbf{v} + \mathbf{b})$$

$$P_{BB}(v_i = 1|h) = \frac{e^{+1 w_i^{(T)} h + a_i 1 + b^T h}}{e^{+0 w_i^{(T)} h + a_i 0 + b^T h} + e^{+1 w_i^{(T)} h + a_i 1 + b^T h}} = \frac{1}{1 + e^{-w_i^{(T)} h - a_i}} = \text{sigmoid}(w_i^{(T)} h + a_i)$$

$$P_{BB}(\mathbf{v} = \mathbf{1}|\mathbf{h}) = \text{sigmoid}(\mathbf{W}^T \mathbf{h} + \mathbf{a})$$

Free Energy

Next we use the definition of free energy and calculate its derivatives on parameters W , a and b . Free energy is defined to be:

$$F(v) = -\log \sum_h e^{-E(v, h)}$$

And the related probability distribution of data (visible states) is

$$P(v) = \frac{1}{Z} e^{-F(v)} = \frac{1}{Z} \sum_h e^{-E(v, h)}, \quad Z = \int e^{-F(v)} dv$$

Now we want to calculate gradient with respect to parameters of the distribution in order to maximize likelihood of the data $P(v)$:

$$\frac{\partial -\log p(v)}{\partial \theta} = -\frac{\partial p(v)/\partial \theta}{p(v)} = -\frac{p(v)(-\partial F(v)/\partial \theta) - p(v)(\partial Z/\partial \theta)/Z}{p(v)} = \frac{\partial F(v)}{\partial \theta} + \frac{1}{Z} \frac{\partial Z}{\partial \theta}$$

We calculate the term $\frac{1}{Z} \frac{\partial Z}{\partial \theta}$ separately:

$$\frac{1}{Z} \frac{\partial Z}{\partial \theta} = \int -\frac{\partial F(v)}{\partial \theta} \frac{1}{Z} e^{-F(v)} dv = -\int p(v) \frac{\partial F(v)}{\partial \theta} dv$$

The general form of the derivate is then:

$$\frac{\partial \log p(v)}{\partial \theta} = -\left(\frac{\partial F(v)}{\partial \theta} - \int p(v) \frac{\partial F(v)}{\partial \theta} dv \right) = -\left(\frac{\partial F(v)}{\partial \theta} - E_v\left[\frac{\partial F(v)}{\partial \theta}\right] \right)$$

And the latter term can be approximated using contrastive divergence algorithm. We use the training data to produce $p(h|v)$ and then sample $p(v|h)$ and repeat the procedure to get sample (v_i, h_i) and only keep v_i to get the approximate sample from distribution $p(v)$. This *maybe* special case of Gibbs sampling. Alternatively, it is possible to use MCMC, AIS, PT-MCMC or some other sampling technique to sample directly from $P(v)$ instead using CD-k Gibbs sampling.

Gradient descent

Parameters of the distribution are optimized using gradient descent algorithm (maximum likelihood method) so it is important to calculate actual derivatives of $p(v)$ for Bernoulli-Bernoulli RBM. (Also note that gradients are needed by approximative second order methods like L-BFGS).

First we further simplify the $F(v)$ terms

$$F(v) = -\log \sum_h e^{-E_{BB}(v,h)} = -a^T v - \log \sum_h e^{(W^T v + b)^T h} = -a^T v - \log \sum_h e^{\sum_{i,j} h_i (v_j w_{ij} + b_i)}$$

$$F(v) = -a^T v - \log \sum_h \prod_i e^{h_i (\sum_j v_j w_{ij} + b_i)} = -a^T v - \sum_i \log \sum_{h_i} e^{h_i (\sum_j v_j w_{ij} + b_i)}$$

If we decide $h = \{0, 1\}$ then the equation simplifies further into

$$F(v) = -a^T v - \sum_i \log (1 + e^{\sum_j v_j w_{ij} + b_i})$$

Calculating gradients leads into eqs:

$$\frac{\partial F(v)}{\partial a} = -v$$

$$\frac{\partial F(v)}{\partial b_i} = -\frac{e^{\sum_j v_j w_{ij} + b_i}}{1 + e^{\sum_j v_j w_{ij} + b_i}} = -\text{sigmoid}(\sum_j v_j w_{ij} + b_i)$$

$$\frac{\partial F(v)}{\partial w_{ij}} = -\frac{e^{\sum_j v_j w_{ij} + b_i} v_j}{1 + e^{\sum_j v_j w_{ij} + b_i}} = -v_j \text{sigmoid}(\sum_j v_j w_{ij} + b_i)$$

Gaussian distribution (Gaussian-Bernoulli RBM) - Continuous RBM

So far we have only discussed about Bernoulli-Bernoulli RBM. But what we really want is to process continuous valued input data (and then maybe use BB-RBM to further process its hidden variables). The possible models to use are Gaussian-Bernoulli and Beta-Bernoulli (*note that Beta and Bernoulli distributions are conjugate distributions, this might be useful..*). It seems that gaussian distribution is far more popular so I try to calculate Gaussian-Bernoulli RBM instead.

If we ignore (or fix it to unit value) the variance term of normal distribution, the logical energy function seem to be:

$$E_{GB}(v, h) = \frac{1}{2} \|v - a\|^2 - v^T W h - b^T h$$

To further justify this model, let's calculate marginalized distributions $p(v|h)$ and $p(h|v)$ for this model.

$$P(v|h) = \frac{P(v, h)}{P(h)} = \frac{e^{-E_{GB}(v, h)}}{\int_v e^{-E_{GB}(v, h)} dv} = \frac{e^{v^T W h - \frac{1}{2} \|v - a\|^2}}{\int_v e^{v^T W h - \frac{1}{2} \|v - a\|^2} dv}$$

$$\begin{aligned} \int_v e^{v^T W h - \frac{1}{2} \|v - a\|^2 + b^T h} dv &= e^{b^T h} \int e^{v^T W h - \frac{1}{2} \|v - a\|^2} dv = \int e^{-\frac{1}{2} \|v\|^2 + v^T (W h + a) - \frac{1}{2} \|a\|^2} dv \\ \int e^{-\frac{1}{2} \|v\|^2 + v^T (W h + a) - \frac{1}{2} \|a\|^2} dv &= e^{\frac{1}{2} \|W h + a\|^2 - \frac{1}{2} \|a\|^2} \int e^{-\frac{1}{2} \|v - (W h + a)\|^2} dv \\ P(v|h) &= \frac{e^{\frac{1}{2} \|W h + a\|^2 - \frac{1}{2} \|a\|^2} e^{-\frac{1}{2} \|v - (W h + a)\|^2}}{e^{\frac{1}{2} \|W h + a\|^2 - \frac{1}{2} \|a\|^2} \int e^{-\frac{1}{2} \|v - (W h + a)\|^2} dv} = \frac{1}{Z} e^{-\frac{1}{2} \|v - (W h + a)\|^2} \sim \text{Normal}(W h + a, I) \end{aligned}$$

And similar calculations can be done to calculate: $P(h|v) = \text{sigmoid}(v^T W + b)$.

The related free energy model is:

$$F(v) = -\log \sum_h e^{-E_{GB}(v, h)} = -\log \sum_h e^{-\frac{1}{2} \|v - a\|^2 + v^T W h + b^T h} = \frac{1}{2} \|v - a\|^2 - \log \sum_h e^{(W^T v + b)^T h}$$

And the related gradient of normal distribution parameter a is:

$$\frac{\partial F}{\partial a} = a - v$$

And the other gradients should be the same as in the Bernoulli-Bernoulli model.

Non-unit variance

Unit variance assumption of Gaussian-Bernoulli RBM can be a problem when the input data is not properly preprocessed (in practice we assume $\text{Var}[f(x) + e] = \text{Var}[f(x)] + N(0, 1)$). For example, maximum variance of Bernoulli distributed hidden variables is 0.25 which means that the noise in input data should be modelled to be much smaller.

NOTE: We should add a inverse-Wishart prior for Σ term assuming it to be diagonal I initially (a regularizer). After this we can normalize variance of each dimension to have unit variance meaning that prior is hopefully close to noise prior and integrate over the $p(\Sigma)$ prior.

The energy model for the non-unit variance that can make sense is:

$$E_{GB}(v, h) = \frac{1}{2} (v - a)^T \Sigma^{-1} (v - a) - (\Sigma^{-0.5} v)^T W h - b^T h - \log (|\Sigma|^{-(v+p+1)/2}) + \frac{1}{2} \text{tr}(\Sigma^{-1})$$

which can be further justified by calculating $p(v|h)$ and $p(h|v)$ distributions:

$$\begin{aligned} P(v|h) &= \frac{P(v, h)}{P(h)} = \frac{e^{-E_{GB}(v, h)}}{\int_v e^{-E_{GB}(v, h)} dv} = \frac{e^{(\Sigma^{-0.5} v)^T W h - \frac{1}{2} (v - a)^T \Sigma^{-1} (v - a)}}{\int_v e^{(\Sigma^{-0.5} v)^T W h - \frac{1}{2} (v - a)^T \Sigma^{-1} (v - a)} dv} \\ \int_v e^{(\Sigma^{-0.5} v)^T W h - \frac{1}{2} (v - a)^T \Sigma^{-1} (v - a)} dv &= \int e^{-\frac{1}{2} v^T \Sigma^{-1} v + v^T \Sigma^{-1} (\Sigma^{0.5} W h + a) - \frac{1}{2} \|a\|^2} dv \\ \int e^{-\frac{1}{2} v^T \Sigma^{-1} v + v^T \Sigma^{-1} (\Sigma^{0.5} W h + a) - \frac{1}{2} \|a\|^2} dv &= e^{\frac{1}{2} \|\Sigma^{-0.5} W h + a\|_{\Sigma}^2 - \frac{1}{2} \|a\|^2} \int e^{-\frac{1}{2} \|v - (\Sigma^{0.5} W h + a)\|_{\Sigma}^2} dv \\ P(v|h) &= \frac{1}{Z(\Sigma)} e^{-\frac{1}{2} \|v - (\Sigma^{0.5} W h + a)\|_{\Sigma}^2} \sim \text{Normal}(v | \Sigma^{1/2} W h + a, \Sigma) \\ P(v|h, \Sigma) &\sim \text{Normal}(v | \Sigma^{1/2} W h + a, \Sigma) \end{aligned}$$

To generate normally distributed variables with wanted covariance matrix we notice that for $x \sim N(0, I)$ data we have:

$$\Sigma_x = \text{Cov}[A x] = E[A x x^T A^T] = A I A^T = A A^T = X \Lambda X^T \Rightarrow A = X \Lambda^{1/2}$$

And similarly we calculate

$$P(h|v) = \frac{P(v, h)}{P(v)} = \frac{e^{-E_{GB}(v, h)}}{\sum_h e^{-E_{GB}(v, h)}} = \frac{e^{(\Sigma^{-0.5}v)^T W h + b^T h}}{\sum_h e^{(\Sigma^{-0.5}v)^T W h + b^T h}} = \frac{e^{(\Sigma^{-0.5}v)^T W h + b^T h}}{\sum_h e^{(\Sigma^{-0.5}v)^T W h + b^T h}}$$

The term can be further rewritten as:

$$e^{(\Sigma^{-0.5}v)^T W h + b^T h} = e^{\sum_d ((\Sigma^{-0.5}v)^T w_d + b_d) h_d} = \prod_d e^{((\Sigma^{-0.5}v)^T w_d + b_d) h_d}$$

And if we only restrict to a single variable we have

$$P(h_d|v) = \frac{e^{((\Sigma^{-0.5}v)^T w_d + b_d) h_d}}{\sum_{h_d} e^{((\Sigma^{-0.5}v)^T w_d + b_d) h_d}} = \text{sigmoid}((\Sigma^{-0.5}v)^T w_d + b_d).$$

$$P(\mathbf{h}|\mathbf{v}, \Sigma) = \text{sigmoid}((\Sigma^{-0.5} \mathbf{v})^T \mathbf{W} + \mathbf{b})$$

Free Energy gradient of Σ

Next we want to calculate gradient to Σ which can be very tricky. In the continuous RBM paper I managed to find, the authors didn't use matrix form and only calculated gradient of σ and then *approximated* it. This result further shows that calculation of gradient of full covariance matrix Σ in this case is probably analytically extremely difficult or impossible.

Initially, we notice that Σ is a symmetric matrix. This mean that if we do a change of basis we can always find a space where Σ is a diagonal matrix, $\Sigma = \text{diag}(\sigma_1 \dots \sigma_D)$. Additionally, the RBM model implicitly assumes that variables are independent when given hidden or visible variables. Therefore, to continue to make this independence assumption (note: RBM can be seen as somewhat similar to ICA where we estimate "independent" or "semi-independent" components from the data), we assume there is no correlations between elements of the visible units given hidden units.

$$\begin{aligned} F(v) &= -\log \sum_h e^{-E_{GB}(v, h)} = -\log \sum_h e^{-\frac{1}{2}\|v-a\|_{\Sigma}^2 + (\Sigma^{-0.5}v)^T W h + b^T h + \log(|\Sigma|^{-(v+p+1)/2}) - \frac{1}{2}\text{tr}(\Sigma^{-1})} \\ &= \frac{1}{2}(v-a)^T \Sigma^{-1}(v-a) - \log(|\Sigma|^{-(v+p+1)/2}) + \frac{1}{2}\text{tr}(\Sigma^{-1}) - \log \sum_h e^{(W^T \Sigma^{-0.5}v + b)^T h} \\ &= \frac{1}{2}(v-a)^T \Sigma^{-1}(v-a) - \log(|\Sigma|^{-(v+p+1)/2}) + \frac{1}{2}\text{tr}(\Sigma^{-1}) - \log \prod_i \sum_{h_i} e^{(W^T \Sigma^{-0.5}v + b)_i h_i} \\ &= \frac{1}{2}(v-a)^T \Sigma^{-1}(v-a) + ((v+p+1)/2) \log(|\Sigma|) + \frac{1}{2}\text{tr}(\Sigma^{-1}) - \sum_i \log(1 + e^{(W^T \Sigma^{-0.5}v + b)_i}) \end{aligned}$$

Its derivate, assuming the covariance matrix is diagonal, is:

$$\begin{aligned} \frac{\partial F}{\partial \sigma_k} &= -\frac{(v_k - a_k)^2}{\sigma_k^3} + (v+p+1) \sum_i \frac{\partial}{\partial \sigma_k} \log(\sigma_i) - \frac{1}{\sigma_k^3} - \sum_i \frac{e^{(W^T \Sigma^{-0.5}v + b)_i}}{1 + e^{(W^T \Sigma^{-0.5}v + b)_i}} \frac{\partial}{\partial \sigma_k} (\sum_k w_{ki} \frac{v_k}{\sigma_k}) \\ &= -\frac{(v_k - a_k)^2}{\sigma_k^3} + (v+p+1) \frac{1}{\sigma_k} - \frac{1}{\sigma_k^3} + v_k / \sigma_k^2 \sum_i w_{ki} \text{sigmoid}(W^T \Sigma^{-0.5}v + b)_i \\ \text{diag}[\frac{\partial F}{\partial \Sigma}] &= \\ \text{diag}[-\Sigma^{-3/2}((v-a)(v-a)^T + I) + (v+p+1)\Sigma^{-1/2} + (\Sigma^{-1}v)(W \text{sigmoid}(W^T \Sigma^{-0.5}v + b))^T] \end{aligned}$$

Other derivates of the free energy are:

$$\begin{aligned} \frac{\partial F}{\partial a} &= -\Sigma^{-1}(v-a) \\ \frac{\partial F(v)}{\partial b_i} &= -\text{sigmoid}(\sum_j w_{ij} v_j / \sigma_j + b_i), \quad \frac{\partial F}{\partial b} = -\text{sigmoid}(W^T \Sigma^{-1/2}v + b) \\ \frac{\partial F(v)}{\partial w_{ij}} &= -\frac{v_j}{\sigma_j} \text{sigmoid}(\sum_j w_{ij} v_j / \sigma_j + b_i), \quad \frac{\partial F}{\partial W} = -\text{sigmoid}(W^T \Sigma^{-1/2}v + b) (\Sigma^{-1/2}v)^T \end{aligned}$$

Now the open question is whether these derivatives are valid for **any** covariance matrix Σ as it would mean that we could then estimate correlations between visible layer inputs and not expect them to be statistically independent given hidden variables.

RESULTS: This does not seem to work for any Σ matrix, but it is uncertain as the primary problem here is that variances can become negative meaning that the method is not robust or a complex method is needed to dynamically adjust learning rates so that Σ matrix (eigenvalues) can never become negative.

It is possible to get this working but it seems that the convergence is really slow. The primary idea is to check if updated covariance matrix has positive diagonal and is otherwise good looking (non complex numbers, NaNs, Infs etc) and always drop the learning rate by 50% if it seems that the update would make covariance matrix “unhealthy”. But this is not robust or a good solution as it “filters out” updates that would turn matrix into bad one. Additionally, the update rule seem to be really SLOW here. (gbrbm-variance-model-v1)

The direct variance update don't work, we instead try to fix the situation with the change of variable $e^{-z_i} = 1/\sigma_i^2$ (as recommended below), leading into the equation:

$$\begin{aligned} F(v) &= \frac{1}{2}(v-a)^T \Sigma^{-1}(v-a) + ((v+p+1)/2) \log(|\Sigma|) + \frac{1}{2} \text{tr}(\Sigma^{-1}) - \sum_i \log(1 + e^{(W^T \Sigma^{-0.5} v + b)_i}) \\ &= \frac{1}{2} \sum_i (1 + (v_i - a_i)^2) e^{-z_i} + ((v+p+1)/2) \sum_i z_i - \sum_j \log(1 + \exp(\sum_k w_{kj} v_k e^{-z_k/2} + b_j)) \\ \frac{\partial F}{\partial z_i} &= -\frac{1}{2} (1 + (v_i - a_i)^2) e^{-z_i} + ((v+p+1)/2) - \sum_j \text{sigmoid}(W^T \Sigma^{-0.5} v + b)_j \left(-\frac{1}{2} w_{ij} v_i e^{-z_i/2} \right) \\ \frac{\partial F}{\partial z_i} &= -\frac{1}{2} e^{-z_i} (1 + (v_i - a_i)^2) + ((2 \dim(v) + 1)/2) + \frac{1}{2} e^{-z_i/2} v_i \sum_j w_{ij} \text{sigmoid}(W^T \Sigma^{-0.5} v + b)_j \end{aligned}$$

Wishart prior's degrees of freedom is chosen to be low so we choose the minimum $v = p$ so that mean of inverse wishart prior/regularizer so degrees of freedom is minimized and variance is maximized (maximum uncertainty).

This works somewhat better but **still** leads into NaNs and other problems when the larger number of hidden nodes are used. (This may be because of small covariance matrixes leading into eqs near Infinities.)

ALTERNATIVE GAUSSIAN MODEL

Because their basic approach to introduce variances into model seem to cause problems I now try another slightly modified model introduced by Aalto university researchers (Cho et al. 2011).

Now the energy of the system is defined to be

$$\begin{aligned} E_{GB}(v, h) &= \frac{1}{2} \alpha \sum_i (v_i - a_i)^2 / \sigma_i^2 - \sum_{i,j} \frac{v_i}{\sigma_i^2} w_{ij} h_j - \sum_i b_i h_i \\ &= \frac{1}{2} (v-a)(\alpha D^{-2})(v-a) - v^T D^{-2} W h - b^T h \\ &= \sum_i \frac{1}{2} v_i^2 \alpha / \sigma_i^2 + \frac{1}{2} a_i^2 \alpha / \sigma_i^2 - \sum_i (a_i + \sum_j w_{ij} h_j / \alpha) v_i \alpha / \sigma_i^2 - \sum_i b_i h_i \\ &= \frac{1}{2} \sum_i (v_i - (a_i + \sum_j w_{ij} h_j / \alpha))^2 \alpha / \sigma_i^2 + \frac{1}{2} \sum_i a_i^2 \alpha / \sigma_i^2 - \frac{1}{2} \sum_i (a_i + \sum_j w_{ij} h_j / \alpha)^2 \alpha / \sigma_i^2 - \sum_i b_i h_i \\ &= \frac{1}{2} (v - (a + W h / \alpha))^T (\alpha D^{-2}) (v - (a + W h / \alpha)) + \frac{1}{2} a^T (\alpha D^{-2}) a - \frac{1}{2} \|a + W h / \alpha\|_{D^2/\alpha}^2 - b^T h \end{aligned}$$

After this it is again quite straightforward to calculate the conditional probabilities given hidden or visible neurons.

$$P(v|h) = \frac{e^{-\frac{1}{2}\sum_i (v_i - (a_i + \sum_j w_{ij}h_j/\alpha))^2 \alpha / \sigma_i^2}}{\int e^{-\frac{1}{2}\sum_i (v_i - (a_i + \sum_j w_{ij}h_j/\alpha))^2 \alpha / \sigma_i^2 + \dots} dv} \propto e^{-\frac{1}{2}\sum_i (v_i - (a_i + \sum_j w_{ij}h_j/\alpha))^2 \alpha / \sigma_i^2}$$

$$\sim \prod_i N(a_i + \sum_j w_{ij}h_j/\alpha, \sigma_i^2/\alpha)$$

$$P(v|h) \sim \text{Normal}(a + Wh/\alpha, D^2/\alpha), \quad D^2 = \text{diag}[\sigma_1^2 \dots \sigma_D^2]$$

$$P(h_j|v) = \frac{e^{\sum_{i,j} \frac{v_i}{\sigma_i^2} w_{ij}h_j + \sum_j b_j h_j}}{\sum_{h_j} e^{\sum_{i,j} \frac{v_i}{\sigma_i^2} w_{ij}h_j + \sum_j b_j h_j}} = \frac{e^{(\sum_i \frac{v_i}{\sigma_i^2} w_{ij} + b_j)h_j}}{\sum_{h_j} e^{(\sum_i \frac{v_i}{\sigma_i^2} w_{ij} + b_j)h_j}} = \text{sigmoid}(\sum_i \frac{v_i}{\sigma_i^2} w_{ij} + b_j)$$

$$P(h|v) = \text{sigmoid}(b + v^T D^{-2} W)$$

Notice here that this time the mean is not related to variances which can be useful.

After calculating the conditional probabilities, we again calculate free energy of the energy function (and then calculate its partial derivatives).

$$\begin{aligned} F(v) &= -\log \sum_h e^{-E_{GB}(v,h)} = -\log \sum_h e^{-\frac{1}{2}\alpha(v-a)^T D^{-2}(v-a) + v^T D^{-2} W h + b^T h} \\ &= \frac{1}{2}\alpha(v-a)^T D^{-2}(v-a) - \sum_i \log \sum_{h_i} e^{(W^T D^{-2} v + b)_i h_i} \\ &= \frac{1}{2}\alpha(v-a)^T D^{-2}(v-a) - \sum_i \log(1 + e^{(W^T D^{-2} v + b)_i}) \\ &= \frac{1}{2}\sum_i \alpha \frac{(v_i - a_i)^2}{\sigma_i^2} - \sum_i \log(1 + \exp(\sum_j w_{ji} v_j / \sigma_j^2 + b_i)) \end{aligned}$$

From this formulation we can calculate

$$\frac{\partial F}{\partial a} = -\alpha D^{-2}(a - v)$$

$$\frac{\partial F}{\partial b_i} = -\frac{e^{(W^T D^{-2} v + b)_i}}{1 + e^{(W^T D^{-2} v + b)_i}} = -\text{sigmoid}(W^T D^{-2} v + b)_i$$

$$\frac{\partial F}{\partial b} = -\text{sigmoid}(W^T D^{-2} v + b) = -h$$

$$\frac{\partial F}{\partial w_{ij}} = -\frac{e^{(W^T D^{-2} v + b)_i}}{1 + e^{(W^T D^{-2} v + b)_i}} (v_j / \sigma_j^2) = -\frac{v_j}{\sigma_j^2} \text{sigmoid}(W^T D^{-2} v + b)_i = -D^{-2} v h^T$$

But the crucial derivate is the derivate of σ_i terms, we alter energy function by the change of terms: $1/\sigma_i^2 = e^{-z_i}$ leading into formula:

$$F(v) = \frac{1}{2}\sum_i \alpha (v_i - a_i)^2 e^{-z_i} - \sum_j \log(1 + \exp(\sum_i w_{ij} v_i e^{-z_i} + b_j))$$

And then derivate with respect to z_i :

$$\begin{aligned} \frac{\partial F}{\partial z_i} &= -\frac{1}{2}\alpha (v_i - a_i)^2 e^{-z_i} + \sum_j \frac{\exp(\sum_i w_{ij} v_i e^{-z_i} + b_j)}{1 + \exp(\sum_i w_{ij} v_i e^{-z_i} + b_j)} w_{ij} v_i e^{-z_i} \\ &= -e^{-z_i} [\frac{1}{2}\alpha (v_i - a_i)^2 - v_i \sum_j w_{ij} \text{sigmoid}(W^T D^{-2} v + b)_j] = e^{-z_i} S(a, b, D) \end{aligned}$$

But because variance is difficult parameter to optimize for, we want to calculate its Hessian matrix and use second order derivation

This is the exactly same formula as the one given in the paper of Aalto uni people (*Improved Learning of Gaussian-Bernoulli Restricted Boltzmann Machines. ICANN 2011.*)

Important: The analysis of results seem to imply to following result. Aalto university people model performance with data having multiple modes (circle data etc) have **bad performance** if gradient descent method is used and variance should be learned - cannot learn the variance (Parallel Tempering MCMC may be possible to learn the model variance correctly). However, it has **BETTER** performance if correct variance is known and gradient descent method is used.

1. This seem to imply to following approach: first optimize “original” GB-RBM so that you can learn the variances (correctly)
2. After learning the variances switch to “Aalto-university” GB-RBM model which can give better reconstruction error. With correct variances known it is now possible to learn multiple modes of the distribution easily.
3. Alternatively try to use more sophisticated methods.

Parallel Tempering Annihilated Importance Sampling (AIS)

Parallel tempering is a powerful approach that was used by Aalto university people in 2011 (also see *On the Quantitative Analysis of Deep Belief Networks. Salakhutdinov.* about AIS). I will try to use same method to generate intermediate distributions from which to estimate partition function Z which is needed to estimate exact probabilities.

$$P(v) = \frac{1}{Z} \sum_h e^{-E(v,h)}$$

We generate series of temperatures $\beta = [0...1]$ and define intermediate models as (*Cho et al. 2011*):

$$\mathbf{W}^{(\beta)} = \beta \mathbf{W}, \quad a_i = \beta a_i + (1 - \beta) m_i$$

$$\mathbf{b}^{(\beta)} = \beta \mathbf{b}, \quad \sigma_i^{(\beta)} = \sqrt{\beta \sigma_i^2 + (1 - \beta) s_i^2}$$

So with $\beta = 0$ we simply sample from unscaled normal probability distribution $v_i \sim \mathcal{N}(m_i, s_i^2)$ and then define **stochastic** transition operation as

$$T_\beta(v', v) = p_\beta(v'|h) p_\beta(h|v)$$

Because we know (**proof?**) that iterative Gibbs sampling $T(v', v)$ from distribution will lead to distribution $p(v)$ given any distribution $p(v)$ for any v (assuming all $p(v) > 0$) then the transition from probability distribution $p(v)$ of values will lead to limiting distribution of $p(v)$ of values?? (Convergence point of markov chain).

After this we use AIS to estimate Z -ratios (which mean is used as an estimate) using un-scaled distributions by using samples:

$$\frac{Z_{\beta_K}}{Z_{\beta_0}} = \frac{p'_{\beta_1}(\mathbf{v}_1)}{p'_{\beta_0}(\mathbf{v}_1)} \frac{p'_{\beta_2}(\mathbf{v}_2)}{p'_{\beta_1}(\mathbf{v}_2)} \dots \frac{p'_{\beta_K}(\mathbf{v}_K)}{p'_{\beta_{K-1}}(\mathbf{v}_K)}.$$

(NOTE: is \mathbf{v}_K generated this way also distributed as $p_K(\mathbf{v})$?)

But because we know $p_{\beta_0}(\mathbf{v})$ to be exactly normal it is trivial to sample from and its normalizing constant Z can be computed directly.

$$Z_{\beta_0} = (2\pi)^{D/2} \sigma_1 \dots \sigma_D$$

What is still needed to estimate Z is the calculation of un-scaled probability distribution values $p'(\mathbf{v})$ given parameters of GB-RBM. This needs summation over all values of \mathbf{h} :

$$P'(\mathbf{v}) = e^{-\frac{1}{2}(\mathbf{v}-\mathbf{a})^T \mathbf{\Sigma}^{-1}(\mathbf{v}-\mathbf{a})} \sum_{\mathbf{h}} e^{\mathbf{v}^T \mathbf{\Sigma}^{-0.5} \mathbf{W} \mathbf{h} + \mathbf{b}^T \mathbf{h}}$$

For \mathbf{h} this simplifies further into summation:

$$S = \sum_{\mathbf{h}} e^{\alpha^T \mathbf{h}} = \sum_{\mathbf{h}} \prod_i e^{\alpha_i h_i}$$

The next comes the tricky/smart step, because \mathbf{h} takes only values 0 and 1 this summation can be simplified without going through all the $2^{\dim(\mathbf{h})}$ possible values by using property of multiplication $(\alpha + \beta)(\zeta + \delta) \dots$ “generates”/“goes through” all possible 2^D states **in linear time!**

$$S = \sum_{\mathbf{h}} \prod_i e^{\alpha_i h_i} = \prod_i (e^{\alpha_i(h_i=0)} + e^{\alpha_i(h_i=1)}) = \prod_i (1 + e^{\alpha_i})$$

Therefore, the unscaled log probability is:

$$\begin{aligned} \log(P'(\mathbf{v})) &= -\frac{1}{2}(\mathbf{v}-\mathbf{a})^T \mathbf{\Sigma}^{-1}(\mathbf{v}-\mathbf{a}) + \sum_i \log(1 + e^{\alpha_i(\mathbf{v}, \mathbf{W}, \mathbf{\Sigma}, \mathbf{b})}), \\ \alpha &= \mathbf{W}^T \mathbf{\Sigma}^{-0.5} \mathbf{v} + \mathbf{b}. \end{aligned}$$

About numerical stability: in practice, the exponent of alpha can become too large. If e^{α_i} becomes floating point infinity and destroys numerical accuracy, we use approximation

$$\log(1 + e^{\alpha}) \approx \log(e^{\alpha}) = \alpha$$

which is good approximation if e^{α} is near infinity (even at floating point accuracy).

Building Neural Network from Stacked RBM

After the problem of learning GB-RBM and BB-RBM has been solved it is naturally interesting to try to use it to construct *feedforward neural network*. This is rather straight forward as the probability functions always has sigmoidal probability function (even at the top layer) and we can directly insert weights and biases from the RBM. The only problem is that our RBM is a stochastic/probabilistic meaning that we will now make significant error when approximating hidden neurons to have continuous values.

TODO

Hamiltonian Monte Carlo sampler based approach

Because the direct optimization method doesn't seem to work very well. It seems that an interesting approach could be try to use Monte Carlo sampling as seen in many papers as the given probability model fits into HMC approach. For hamiltonian we will use function

$$H(\mathbf{q}, \mathbf{p}) = E_T(\mathbf{q}) + K(\mathbf{p}), \quad K(\mathbf{p}) = \mathbf{p}^T \mathbf{M}^{-1} \mathbf{p} / 2$$

And we want to use HMC to sample from distribution:

$$p(\boldsymbol{\theta}|\mathbf{v}) = \frac{1}{Z(\mathbf{v})} \exp(-E_T(\boldsymbol{\theta}|\mathbf{v})) \propto p(\mathbf{v}|\boldsymbol{\theta}) p(\boldsymbol{\theta})$$

But in practice we have only know the data likelihood $p(\mathbf{v}|\boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} e^{-F(\mathbf{v}|\boldsymbol{\theta})}$ term from the previous sections. This will lead to the following eqs (for the gradient):

$$E_T(\boldsymbol{\theta}|\mathbf{v}) = -\log\left(\frac{Z(\mathbf{v})}{Z(\boldsymbol{\theta})}\right) + F_T(\mathbf{v}|\boldsymbol{\theta}) - \log\left(\frac{p(\boldsymbol{\theta})}{p(\mathbf{v})}\right)$$

$$\nabla_{\boldsymbol{\theta}} E_T(\boldsymbol{\theta}|\mathbf{v}) = \nabla_{\boldsymbol{\theta}} F_T(\mathbf{v}|\boldsymbol{\theta}) - E_{\mathbf{v}}[\nabla_{\boldsymbol{\theta}} F_T(\mathbf{v}|\boldsymbol{\theta})] - \nabla_{\boldsymbol{\theta}} \log(p(\boldsymbol{\theta}))$$

We can decide for the flat priors meaning that $p(\boldsymbol{\theta}) \propto 1$ and the last term will be zero. This means that we can use directly gradient computed for the $-\log(P(\mathbf{v}))$. But HMC also calculates difference between energy functions

$$E_T(\boldsymbol{\theta}_{n+1}|\mathbf{v}) - E_T(\boldsymbol{\theta}_n|\mathbf{v}) = F_T(\mathbf{v}|\boldsymbol{\theta}_{n+1}) - F_T(\mathbf{v}|\boldsymbol{\theta}_n) + \log\left(\frac{Z_T(\boldsymbol{\theta}_{n+1})}{Z_T(\boldsymbol{\theta}_n)}\right) - \log\left(\frac{p(\boldsymbol{\theta}_{n+1})}{p(\boldsymbol{\theta}_n)}\right)$$

The last term can be assumed to be zero but the problem is ratio of partition functions. We know that ratio of partition functions will be close to 1 if the change between parameters is not too large so it can be initially approximated to be zero but in practice this seem to lead to too small variances causing problems. But we can approximate by using the previous results (by using first order derivate approximation from $\boldsymbol{\theta}_{n+1}$ to $\boldsymbol{\theta}_n$ and from $\boldsymbol{\theta}_n$ to $\boldsymbol{\theta}_{n+1}$ and calculating the mean):

$$\log\left(\frac{Z_T(\boldsymbol{\theta}_{n+1})}{Z_T(\boldsymbol{\theta}_n)}\right) \approx -\frac{1}{2}(\boldsymbol{\theta}_{n+1} - \boldsymbol{\theta}_n)^T (E_{\mathbf{v}}[\nabla_{\boldsymbol{\theta}} F_T(\mathbf{v}|\boldsymbol{\theta}_{n+1})] + E_{\mathbf{v}}[\nabla_{\boldsymbol{\theta}} F_T(\mathbf{v}|\boldsymbol{\theta}_n)])$$

This approximation does not seem to work well in practice and cannot be used.

You can calculate gradient quite easily but you cannot calculate difference between energy functions as it requires estimating the ratio of partition functions. For this case I will try to use *importance sampling* (as the distributions should be close to each other and we want only the ratio and not the actual Z_T -values. Theory:

$$\frac{Z_T(\boldsymbol{\theta}_{n+1})}{Z_T(\boldsymbol{\theta}_n)} = \int \frac{p'_T(\mathbf{v}|\boldsymbol{\theta}_{n+1})}{Z_T(\boldsymbol{\theta}_n)} d\mathbf{v} = \int \frac{p'_T(\mathbf{v}|\boldsymbol{\theta}_{n+1})}{p'_n(\mathbf{v}|\boldsymbol{\theta}_n)} p(\mathbf{v}|\boldsymbol{\theta}_n) d\mathbf{v} = E_{\boldsymbol{\theta}_n} \left[\frac{p'_T(\mathbf{v}|\boldsymbol{\theta}_{n+1})}{p'_n(\mathbf{v}|\boldsymbol{\theta}_n)} \right]$$

So we need to sample \mathbf{v} :s from the distribution using AIS and “parallel tempering” approach described by Cho et. al. in 2011. This seem to work rather well after which we calculate ratios of unscaled distributions $p'(\mathbf{v}|\boldsymbol{\theta}) = e^{-F(\mathbf{v}|\boldsymbol{\theta})}$ and calculate the mean value:

$$E_{\mathbf{v}|\boldsymbol{\theta}_n} \left[\frac{p'_T(\mathbf{v}|\boldsymbol{\theta}_{n+1})}{p'_n(\mathbf{v}|\boldsymbol{\theta}_n)} \right] = E_{\mathbf{v}|\boldsymbol{\theta}_n} [e^{F(\mathbf{v}|\boldsymbol{\theta}_n) - F(\mathbf{v}|\boldsymbol{\theta}_{n+1})}].$$

This is easy to parallelize and seem to work with enough accuracy, **in low dimensions** (<256). TODO: Probably a good way to estimate $Z_T(\boldsymbol{\theta}_n)$ would be to use variational methods (variational bayes but currently it is too complicated for me).

Latent variable model

In general, the RBM method works by first defining probability function $p(\mathbf{v}, \mathbf{h}|\boldsymbol{\theta})$ and then maximizing visible states (visible data) probability

$$\max_{\boldsymbol{\theta}} p(\mathbf{v}|\boldsymbol{\theta}) = \int p(\mathbf{v}, \mathbf{h}|\boldsymbol{\theta}) d\mathbf{h}$$

But instead of calculating maximum likelihood solution, one could also calculate bayesian samples and estimate distribution or sample from

$$p(\boldsymbol{\theta}|\mathbf{v}) \propto p(\mathbf{v}|\boldsymbol{\theta}) p(\boldsymbol{\theta})$$

But sampling from this distribution can be difficult but one can maybe use MCMC sampling which is only interested in relative probabilities

$$r = \frac{p(\boldsymbol{\theta}_{n+1}|\mathbf{v}) p(\boldsymbol{\theta})}{p(\boldsymbol{\theta}_n|\mathbf{v}) p(\boldsymbol{\theta})} = \frac{\sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h}|\boldsymbol{\theta}_{n+1})}{\sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h}|\boldsymbol{\theta}_n)}$$

If we furthermore assume that $Z(\boldsymbol{\theta}_{n+1}) \approx Z(\boldsymbol{\theta}_n)$ then we get a formula:

$$r(\mathbf{v}) \approx \frac{\sum_{\mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h}|\boldsymbol{\theta}_{n+1})}}{\sum_{\mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h}|\boldsymbol{\theta}_n)}} \approx$$

After which hidden states are sampled or calculated using a formula

$$p(\mathbf{h}|\mathbf{v}) = \int p(\mathbf{h}|\boldsymbol{\theta}, \mathbf{v}) p(\boldsymbol{\theta}|\mathbf{v}) d\boldsymbol{\theta} \approx \frac{1}{N} \sum_i p(\mathbf{h}|\boldsymbol{\theta}_i, \mathbf{v})$$

Now for multivariate probability distribution $\mathbf{x} = [\mathbf{v}, \mathbf{h}]^T$, $p(\mathbf{x}) \sim N(\boldsymbol{\mu}_{\mathbf{x}}, \boldsymbol{\Sigma}_{\mathbf{x}})$ it is trivially easy to compute $p(\mathbf{v}|\boldsymbol{\mu}_{\mathbf{v}}, \boldsymbol{\Sigma}_{\mathbf{v}})$ as the solution is yet another normal distribution. But now because marginal distribution do not have dependency to $\boldsymbol{\mu}_{\mathbf{h}}$ or $\boldsymbol{\Sigma}_{\mathbf{h}}$ or other variables given known variable, it is not possible to calculate relationship between \mathbf{h} and \mathbf{v} .

If one does not want to use bayesian approach when maximizing probability of visible states given parameters, another approach could be still use gradient descent but now giving the gradient to L-BFGS optimizer method (conjugate gradient) method which should then be able to do higher quality maximization of visible states probability by using only gradient information.

Solving variance model with bayesian prior distribution

We add inverse scaled chi-squared distribution as prior/regularizer distribution for σ_i^2 terms in order to keep variance as low as possible

$$\begin{aligned} E_{GB}(v, h) + \text{ScaledInv}\chi^2(\boldsymbol{\sigma}^2) \\ = \frac{1}{2} \sum_i (v_i - a_i)^2 / \sigma_i^2 - \sum_{i,j} \frac{v_i}{\sigma_i^2} w_{ij} h_j - \sum_i b_i h_i + \frac{1}{2} (\alpha - 1) \sum_i (v_i - a_i)^2 / \sigma_i^2 \\ = \frac{1}{2} \alpha (v - (a + Wh))^T D^{-2} (v - (a + Wh)) + \frac{1}{2} a^T D^{-2} a - \frac{1}{2} \|a + Wh\|^2 - b^T h \end{aligned}$$

After this It is possible to calculate free energy derivate of σ_i^2 .

$$\begin{aligned} F(v) = -\log \sum_h e^{-E_{GB}(v, h)} = -\log \sum_h e^{-\frac{1}{2} (v - a)^T D^{-2} (v - a) - \alpha \frac{1}{2} s^T D s + v^T D^{-2} Wh + b^T h} \\ = \frac{1}{2} (v - a)^T D^{-2} (v - a) + \frac{1}{2} s^T D^2 s - \sum_i \log \sum_{h_i} e^{(W^T D^{-2} v + b)_i h_i} \\ = \frac{1}{2} (v - a)^T D^{-2} (v - a) + \frac{1}{2} s^T D^{-2} s - \sum_i \log (1 + e^{(W^T D^{-2} v + b)_i}) \\ = \frac{1}{2} \sum_i \frac{(v_i - a_i)^2}{\sigma_i^2} + \frac{1}{2} \sum_i s_i^2 \sigma_i^2 - \sum_i \log (1 + \exp(\sum_j w_{ji} v_j / \sigma_j^2 + b_i)) \end{aligned}$$

And the derivate with respect to z_i is:

$$\frac{\partial F}{\partial z_i} = -e^{-z_i} [\frac{1}{2} (v_i - a_i)^2 - v_i \sum_j w_{ij} \text{sigmoid}(W^T D^{-2} v + b)_j] + \frac{1}{2} s_i^2 e^{z_i}$$

And s_i is a sample square statistic $s_i^2 = \sum_{j=1}^N (v_i(j) - \bar{v}_i(j))^2$ where j :s are indexes for individual observations from the data. Now from the data it is easy to calculate s_i^2 but from the model this is a different thing. If we approximate $s_i^2 = \alpha (v_i - a_i)^2$ we get regularized gradients.

Multiple visible elements

In real world applications, it is often interesting to learn distribution $p(\mathbf{y}|\mathbf{x})$ and in RBM context this means we have multiple different visible elements $\mathbf{v} = [\mathbf{v}_1, \mathbf{v}_2]$ and we need to be able to sample from $p(\mathbf{v}_2|\mathbf{v}_1)$ in order to predict value of \mathbf{v}_2 given elements \mathbf{v}_1 . In this model, the learning of RBM elements can be done normally. However, in reconstruction phase we need to sample from $p(\mathbf{h}|\mathbf{v}_1)$ and then from $p(\mathbf{v}_2|\mathbf{h})$ in order to approximate sample from integral

$$p(\mathbf{v}_2|\mathbf{v}_1) = \int p(\mathbf{v}_2|\mathbf{h}) p(\mathbf{h}|\mathbf{v}_1) d\mathbf{h}$$

In GB-RBM we need that $p(\mathbf{v}|\mathbf{h})$ has a normal distribution and calculating conditional normal distribution $p(\mathbf{v}_1|\mathbf{h})$ and $p(\mathbf{v}_2|\mathbf{h})$ is "easy":

$$\begin{aligned} p(\mathbf{v}|\mathbf{h}) \sim \text{Normal}(\mathbf{v}|\boldsymbol{\Sigma}^{1/2} \mathbf{W} \mathbf{h} + \mathbf{a}, \boldsymbol{\Sigma}), \boldsymbol{\Sigma} = \text{diag}([\sigma_1 \dots \sigma_D]) \\ p(\mathbf{v}_2|\mathbf{v}_1 = \mathbf{x}, \mathbf{h}) \sim \text{Normal}(\mathbf{v}_2|\boldsymbol{\mu}_2 + \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1} (\mathbf{x} - \boldsymbol{\mu}_1), \boldsymbol{\Sigma}_{22} - \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1} \boldsymbol{\Sigma}_{12}) \end{aligned}$$

But because Σ is diagonal conditional distribution does not directly depend on values of \mathbf{v}_1 .

$$p(\mathbf{v}_1|\mathbf{h}) \sim \text{Normal}(\mathbf{v}_1|\boldsymbol{\mu}_1(\mathbf{h}), \Sigma_{11})$$

$$p(\mathbf{v}_2|\mathbf{h}) \sim \text{Normal}(\mathbf{v}_2|\boldsymbol{\mu}_2(\mathbf{h}), \Sigma_{22})$$

What is left is to calculate $p(\mathbf{h}|\mathbf{v}_1)$ to generate hidden variables from input values $\mathbf{v}_1 = \mathbf{x}$, which is more difficult.

$$p(\mathbf{h}|\mathbf{v}) = \text{sigmoid}((\Sigma^{-0.5} \mathbf{v})^T \mathbf{W} + \mathbf{b})$$

By generating \mathbf{v}_2 's randomly we can generate candidate(s) \mathbf{h}_0 which can be then improved and we can use bayes rule to write

$$p(\mathbf{h}|\mathbf{v}_1) \propto p(\mathbf{v}_1|\mathbf{h}) p(\mathbf{h})$$

And if we assume each configuration is equally likely in our model (without any data hidden states should have equal probability of being either 0 or 1 \Rightarrow so the prior should maybe still have some effect - it should force values of \mathbf{h} to be either 0 or 1 though the values should be equally likely), then the prior $p(\mathbf{h})$ is flat constant. This makes sampling from $p(\mathbf{h}|\mathbf{v}_1)$ using HMC rather easy. Our $U(\mathbf{h})$ terms are (Γ_1 -matrix is simple used to select $\boldsymbol{\mu}_1$ from $\boldsymbol{\mu}$):

$$U(\mathbf{h}) = \frac{1}{2}(\mathbf{v}_1 - \boldsymbol{\mu}_1(\mathbf{h}))^T \Sigma_{11}^{-1}(\mathbf{v}_1 - \boldsymbol{\mu}_1(\mathbf{h})), \boldsymbol{\mu}_1(\mathbf{h}) = \Gamma_1(\Sigma^{1/2} \mathbf{W} \mathbf{h} + \mathbf{a})$$

$$\nabla_{\mathbf{h}} U(\mathbf{h}) = (\boldsymbol{\mu}_1(\mathbf{h}) - \mathbf{v}_1)^T \Sigma_{11}^{-1} \Gamma_1 \Sigma^{1/2} \mathbf{W}$$

Here we have assumed \mathbf{h} can take continuous values between 0 and 1 in practice, the sampling should maybe happen in continuous \mathbf{h} -space but the exact samples generated should be discretized. Alternatives:

- keep \mathbf{h} values discretized to 0, 1 values always. Now there is problem that with too small epsilon the sampler stays in a same state and cannot never try to change state
- keep values of \mathbf{h} continuous internally and only discretize to 0, 1 states when emitting a sample
- use continuous \mathbf{h} values, everywhere, this is not well-defined as everything in our theory assumes samples must be discretized and other sources recommend against completely continuous variables..

Another approach would generate candidate \mathbf{h}_0 states and then do gradient descent until convergence to the best state and discretize \mathbf{h} . After discretization calculate gradient descent again until convergence and discretize again \mathbf{h} . Stop until there is no changes in \mathbf{h} . Repeat N times. This will be maximum data likelihood estimate of $p(\mathbf{h}|\mathbf{v}_1)$. After this sample $p(\mathbf{v}_1|\mathbf{h})$.

Brute force (modular) approach

Assume we can divide \mathbf{h} into parts $[\mathbf{h}_1 \dots \mathbf{h}_L]$ and go through all the combinations separately. If we have perfect parallel machine this happens in $O(2^{\max_i \dim(\mathbf{h}_i)})$ time or in $O(2^{\dim(\mathbf{h})/K})$ on average if we can divide problem perfectly into K parts. We can choose, for example, $\dim(\mathbf{h}_i) = 10$ and the computations scale now linearly (although the results are not equally good).

General multiple elements (modular) approach

In practice, the memory requirements $O(\dim(\mathbf{v}) \times \dim(\mathbf{h}))$ and computation of huge RBM's grow too quickly if we use perfect model and use full matrix \mathbf{W} . Because of this we consider modular approach (also suggested by neuroscience research) where vectors are divided into modules/subvectors. This leads into following energy equations:

$$E_{BB}(\mathbf{v}_1 \dots \mathbf{v}_K, \mathbf{h}_1 \dots \mathbf{h}_L) = -\mathbf{a}^T \mathbf{v} - \sum_{k,l} \mathbf{v}_k^T \mathbf{W}_{kl} \mathbf{h}_l - \mathbf{b}^T \mathbf{h}$$

$$E_{GB}(\mathbf{v}_1 \dots \mathbf{v}_K, \mathbf{h}_1 \dots \mathbf{h}_L) = \frac{1}{2} \|\mathbf{v} - \mathbf{a}\|^2 - \sum_{k,l} \mathbf{v}_k^T \mathbf{W}_{kl} \mathbf{h}_l - \mathbf{b}^T \mathbf{h}$$

This reduces amount of memory required during computations but doesn't lead into other improvements. But we want that $p(\mathbf{h}) = p(\mathbf{h}_1)p(\mathbf{h}_2) \dots p(\mathbf{h}_L)$ in order to be able to brute force through \mathbf{h} when looking for optimum because now we can look for optimum for each \mathbf{h}_i vector separately.