

1 Introduction

Restricted Boltzmann Machines (RBMs) are neural networks that can reproduce a probability distribution, given a set of training samples. They have the architectural peculiarity of being composed of one visible (input) layer and one hidden layer. As opposed to general Boltzmann Machines, they are restricted in the sense that the units, or neurons, are not connected to other units in the same layer, only ones in the other layer. The values themselves of the units are binary (i.e, 0 or 1).

visible layer

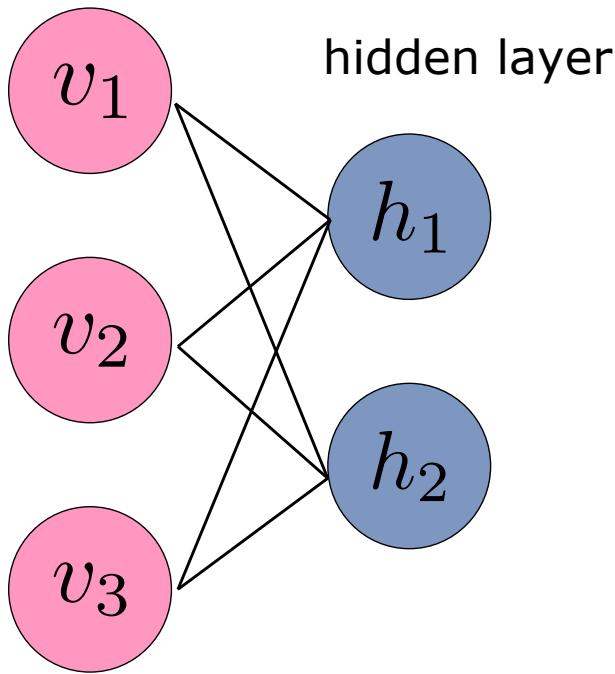


Figure 1: Example RBM with visible layer \mathbf{v} consisting of $D = 3$ units and hidden layer consisting of $n_h = 2$ units.

2 Theory

Consider an RBM with a visible layer \mathbf{v} of length D and hidden layer \mathbf{h} of length n_h , where each component of \mathbf{v} and \mathbf{h} can be either 0 or 1 and let the parameters of this network be $\lambda = \{\mathbf{a}, \mathbf{b}, \mathbf{w}\}$. The probability of a configuration consisting of visible layer \mathbf{v} and hidden

layer \mathbf{h} , with parameters λ , is given by the joint Boltzmann distribution:

$$P(\mathbf{v}, \mathbf{h}) = \frac{e^{-E(\mathbf{v}, \mathbf{h})}}{\mathcal{Z}_\lambda} \quad (1)$$

where the partition function is:

$$\mathcal{Z}_\lambda = \sum_{\mathbf{v}, \mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})} \quad (2)$$

and the "energy" is:

$$E(\mathbf{v}, \mathbf{h}) = - \sum_{i=1}^D a_i v_i - \sum_{j=1}^{n_h} b_j h_j - \sum_{i,j} v_i w_{ij} h_j \quad (3)$$

The marginalized probability of a visible vector can be obtained by summing out the hidden vector dependence:

$$P(\mathbf{v}) = \frac{1}{\mathcal{Z}_\lambda} \sum_{\mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})} \quad (4)$$

and likewise for the marginal probability of a hidden vector:

$$P(\mathbf{h}) = \frac{1}{\mathcal{Z}_\lambda} \sum_{\mathbf{v}} e^{-E(\mathbf{v}, \mathbf{h})} \quad (5)$$

The conditional probability of a visible vector given a hidden vector is:

$$P_\lambda(\mathbf{v}|\mathbf{h}) = \prod_{i=1}^D P_\lambda(v_i|\mathbf{h}) \quad (6)$$

where the factorization of probabilities is due to the fact the input unit activations v_i are mutually independent, since there are no connections between units in the same layer. Similarly, the probability of a hidden vector given an input is:

$$P_\lambda(\mathbf{h}|\mathbf{v}) = \prod_{j=1}^{n_h} P_\lambda(h_j|\mathbf{v}) \quad (7)$$

The activation probabilities of each unit are given by:

$$P_\lambda(h_j|\mathbf{v}) = \begin{cases} \sigma(z_j) & h_j = 1 \\ 1 - \sigma(z_j) & h_j = 0 \end{cases} \quad (8)$$

where $\sigma(z)$ is the logistic sigmoid function:

$$\sigma(z_j) = \frac{1}{1 + e^{-z_j}} \quad (9)$$

and z_j is the weighted sum of inputs:

$$z_j = b_j + \sum_{i=1}^D v_i w_{ij} \quad (10)$$

The activation probabilities of input units are obtained similarly, replacing the weighted sum of inputs with:

$$z_i = a_i + \sum_{j=1}^{n_h} h_j w_{ij} \quad (11)$$

Using these probabilities, a Markov chain can be produced that generates hidden vectors from visible ones, and from the generated hidden vector, a new visible one:

$$\mathbf{v} \rightarrow \mathbf{h}' \rightarrow \mathbf{v}' \rightarrow \mathbf{h}'' \rightarrow \mathbf{v}'' \rightarrow \mathbf{h}''' \rightarrow \dots$$

At the end of the random walk, a model for the probability distribution of input vectors is obtained.

3 Training

Typically, RBMs are trained using the method of contrastive divergence (CD), which performs Gibbs sampling, and inside a gradient descent (GD) procedure, which updates the

network weights.

Unbiased samples are obtained via alternating Gibbs sampling using the unit activation probabilities derived in the previous section.

The cost function to be used is:

$$C = - \sum_{\mathbf{v}} P(\mathbf{v}) \ln P_{\lambda}(\mathbf{v}) \quad (12)$$

which is known as the Kullback-Leibler divergence, without the constant term. (Alternatively, it is the difference between two Kullback-Leibler divergence functions)

For network parameters $\lambda = \{\mathbf{a}, \mathbf{b}, \mathbf{w}\}$, the gradients for GD are:

$$\delta_a = \frac{\partial C}{\partial a_i} = -\langle v_i \rangle + \langle v'_i \rangle \quad (13)$$

$$\delta_b = \frac{\partial C}{\partial b_j} = -\langle h_j \rangle + \langle h'_j \rangle \quad (14)$$

$$\delta_w = \frac{\partial C}{\partial w_{ij}} = -\langle v_i h_j \rangle + \langle v'_i h'_j \rangle \quad (15)$$

The averages that show up in the gradients are computed by taking the average of the unit over sample in the batch, or mini-batch:

$$\langle v_i \rangle = \frac{1}{N_{\text{mini}}} \sum_{n=1}^{N_{\text{mini}}} v_i^{(n)} \quad (16)$$

$$\langle h_j \rangle = \frac{1}{N_{\text{mini}}} \sum_{n=1}^{N_{\text{mini}}} h_j^{(n)} \quad (17)$$

$$\langle v_i h_j \rangle = \frac{1}{N_{\text{mini}}} \sum_{n=1}^{N_{\text{mini}}} v_i^{(n)} h_j^{(n)} \quad (18)$$

To track the training progress, the cost function needs to be computed for the various

iterations of network parameters. Computing it though, will not be tractable due to the exponentially large size of the configuration space of hidden-visible vector pairs:

$$C = - \sum_{\mathbf{v}} P_\lambda(\mathbf{v}) \ln P_\lambda(\mathbf{v}) \quad (19)$$

$$\approx -\frac{1}{N} \sum_{n=1}^N \ln P_\lambda(\mathbf{v}^{(n)}) \quad (20)$$

$$= -\frac{1}{N} \sum_{n=1}^N \left[\sum_{\mathbf{h}} \ln P_\lambda(\mathbf{v}^{(n)}, \mathbf{h}) \right] \quad (21)$$

$$C = -\frac{1}{N} \sum_{n=1}^N F(\mathbf{v}^{(n)}) + \ln \mathcal{Z}_\lambda \quad (22)$$

where the "free energy" $F(\mathbf{v}^{(n)})$ is defined such that:

$$e^{-F(\mathbf{v})} = \sum_{\mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})} \quad (23)$$

which is exponentially hard to compute.

There are three alternative approaches to keep track of the training progress:

Approach 1: Track squared error between data and reconstructions

$$C_1 = \frac{1}{N} \sum_{n=1}^N \|\mathbf{v}^{(n)} - \mathbf{v}^{(n)\prime}\|^2 \quad (24)$$

Approach 2: Track categorical cross-entropy of visible layer's binary units.

$$C_2 = -\frac{1}{N} \sum_{n=1}^N \left\{ \sum_{i=1}^D v_i \log \sigma(\tilde{z}_i) + (1 - v_i) \log [1 - \sigma(\tilde{z}_i)] \right\} \quad (25)$$

where $\tilde{z}_i = a_i + \sum_{j=1}^{n_h} w_{ij} h_j$.

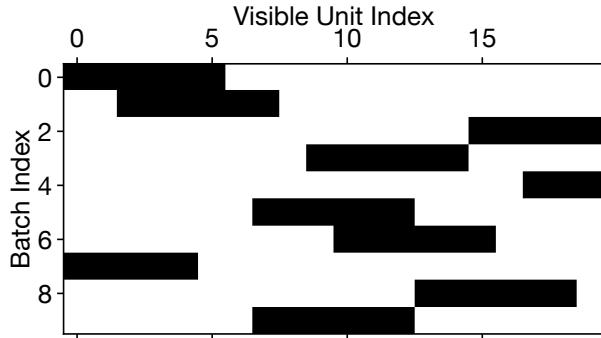
Approach 3: Compare current free-energy to free-energy of validation data.

The hope is that these should be approximately equal.

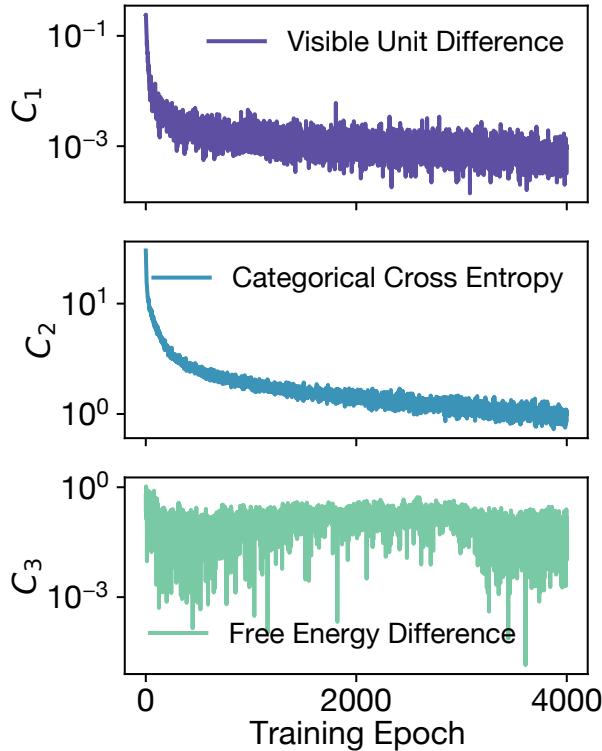
$$C_3 = \left| \frac{1}{N_{\text{val}}} \sum_{n=1}^{N_{\text{val}}} F(\mathbf{v}_{\text{val}}^{(n)}) - \frac{1}{N} \sum_{n=1}^N F(\mathbf{v}^{(n)}) \right| \quad (26)$$

4 Results

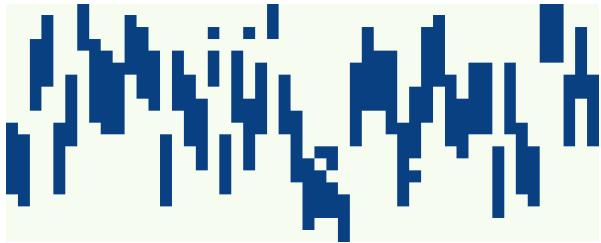
Example visible vectors for training



Comparison of cost function alternative approaches



Sampled visible vectors



5 Ising model

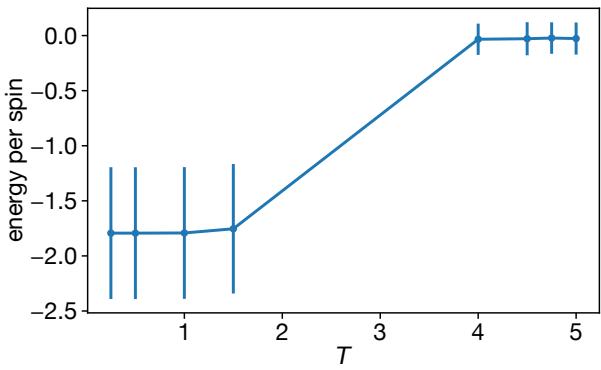


Figure 2: Energy per spin on a 10×10 lattice.

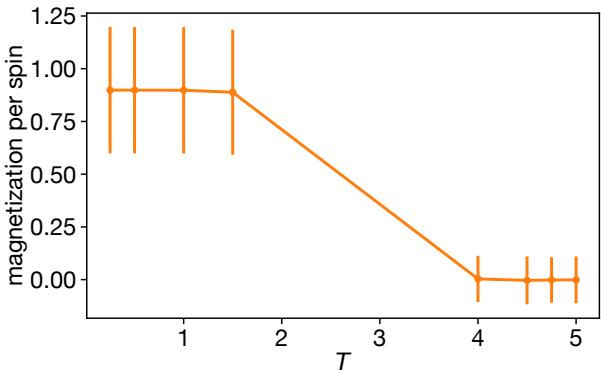


Figure 3: Magnetization per spin on a 10×10 lattice.

6 References

- Adrian's Machine Learning Lecture notes

- https://en.wikipedia.org/wiki/Restricted_Boltzmann_machine
- <http://www.cs.toronto.edu/~hinton/absps/guideTR.pdf>
- <https://www.nature.com/articles/s41567-019-0545-1>