

# An Infinite Restricted Boltzmann Machine

Marc-Alexandre Côté<sup>1</sup>, Hugo Larochelle<sup>1</sup>

<sup>1</sup>Department of Computer Science, Université de Sherbrooke, Sherbrooke, QC, Canada, J1K 2R1

**Keywords:** RBM, machine learning, unsupervised learning, neural network, generative model

## Abstract

We present a mathematical construction for the restricted Boltzmann machine (RBM) that doesn't require specifying the number of hidden units. In fact, the hidden layer size is adaptive and can grow during training. This is obtained by first extending the RBM to be sensitive to the ordering of its hidden units. Then, thanks to a carefully chosen definition of the energy function, we show that the limit of infinitely many hidden units is well defined. As with RBM, approximate maximum likelihood training can be performed, resulting in an algorithm that naturally and adaptively adds trained hidden units during learning. We empirically study the behaviour of this infinite RBM, showing that its performance is competitive to that of the RBM, while not requiring the tuning of a hidden layer size.

## 1 Introduction

Over the years, machine learning research has produced a large variety of latent variable probabilistic models. These include mixture models, factor analysis models, latent dynamical models, and many others. Such models usually require that the dimensionality of the latent representation be specified and fixed during learning. Adapting this quantity is then considered as a separate process, that takes the form of model selection and is normally treated as an additional hyper-parameter to tune.

For this reason, more recently, there has been a lot of work on extending these models such that the size of the representation can be treated as an adaptive quantity during training. These extensions, often referred to as "infinite" models, are non-parametric in nature where the latent space is infinite with probability 1 and can arbitrarily adapt their capacity to the training data (see Orbanz and Teh (2010) for a brief overview).

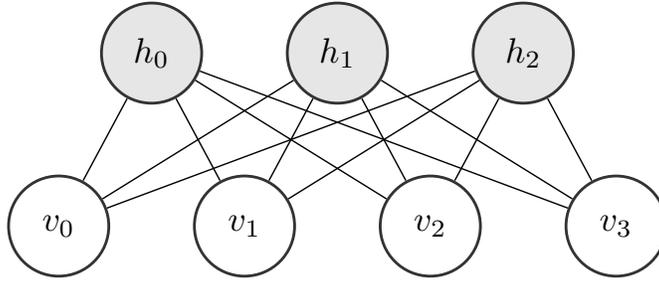


Figure 1: Graphical model of the restricted Boltzmann Machine. Interconnections between visible units and hidden units using symmetric weights.

While most latent variable models have been extended to one or more infinite variants, a notable exception is the restricted Boltzmann machine (RBM). The RBM is an undirected graphical model for binary vector observations, where the latent representation is itself a binary vector (i.e. hidden layer). The RBM (and its extensions to non-binary vectors) have been successfully applied to a variety of problems and data, such as images (Ranzato et al., 2010), movie user preferences (Salakhutdinov et al., 2007), motion capture (Taylor et al., 2011), text (Dahl et al., 2012) and many others. One explanation for the lack of literature on RBMs with an adaptive hidden layer size comes from its undirected nature. Indeed, undirected models tend to be less amenable to a Bayesian treatment of learning, on which relies the majority of the literature on infinite models.

Our main contribution in this paper is thus a proposal for an infinite RBM which can adapt the effective number of hidden units during training. While our proposal is not based on a Bayesian formulation, it does correspond to the infinite limit of a finite-sized model and behaves in such a way that it effectively adapts its capacity as training progresses.

First, we propose a finite extension of the RBM that is sensitive to the position of each unit in its hidden layer. This is achieved by introducing a random variable that represents the number of hidden units intervening in the RBM’s energy function. Then, thanks to the introduction of an energy cost for using each additional unit, we show that taking the infinite limit of the total number of hidden units is well defined. We describe an approximate maximum likelihood training algorithm for this infinite RBM, based on (Persistent) Contrastive Divergence, which results in a procedure where hidden units are implicitly added as training progresses. Finally, we empirically report how this model behaves in practice and show that it can achieve performance that is competitive to a traditional RBM on the binarized MNIST and Caltech101 Silhouettes datasets, while not requiring the tuning of a hyper-parameter for its hidden layer size.

## 2 Restricted Boltzmann Machine

We describe the basic RBM model, which we’ll build on to derive its ordered and infinite versions.

An RBM is a generative stochastic neural network composed of two layers:

visible  $\mathbf{v}$  and hidden  $\mathbf{h}$ . These layers are fully connected to each other, while connections within a layer are not allowed. This means each unit  $v_i$  is connected to all  $h_j$  units via undirected weighted connections (Figure 1).

Given a binary RBM with  $D$  visible units and  $K$  hidden units, the set of visible vectors is  $\mathcal{V} = \{0, 1\}^D$ , whereas the set of hidden vectors is  $\mathcal{H} = \{0, 1\}^K$ . In an RBM model, each configuration  $(\mathbf{v}, \mathbf{h}) \in \mathcal{V} \times \mathcal{H}$  has an associated energy value defined by the following function:

$$E(\mathbf{v}, \mathbf{h}) = -\mathbf{h}^T \mathbf{W} \mathbf{v} - \mathbf{v}^T \mathbf{b}^v - \mathbf{h}^T \mathbf{b}^h \quad (1)$$

The parameters  $\Theta = \{\mathbf{W}, \mathbf{b}^v, \mathbf{b}^h\}$  of this model are the weights  $\mathbf{W}$  ( $K \times D$  matrix), the visible unit biases  $\mathbf{b}^v$  ( $D \times 1$  vector) and the hidden unit biases  $\mathbf{b}^h$  ( $K \times 1$  vector).

A probability distribution over visible and hidden vectors is defined in terms of this energy function:

$$P(\mathbf{v}, \mathbf{h}) = \frac{1}{Z} e^{-E(\mathbf{v}, \mathbf{h})} \quad (2)$$

with

$$Z = \sum_{\mathbf{v}' \in \mathcal{V}} \sum_{\mathbf{h}' \in \mathcal{H}} e^{-E(\mathbf{v}', \mathbf{h}')}. \quad (3)$$

We see from Equation (3) that the partition function  $Z$  (normalizing constant) is intractable, as it requires summing over all possible  $2^{(D+K)}$  configurations.

The probability distribution of a visible vector is obtained by marginalizing over all configurations of hidden vectors. One property of the RBM is that the numerator of the marginal  $P(\mathbf{v})$  is tractable:

$$P(\mathbf{v}) = \frac{1}{Z} \sum_{\mathbf{h}' \in \mathcal{H}} e^{-E(\mathbf{v}, \mathbf{h}')} = \frac{1}{Z} e^{-F(\mathbf{v})} \quad (4)$$

with

$$F(\mathbf{v}) = -\mathbf{v}^T \mathbf{b}^v - \sum_{i=1}^K \text{soft}_+(\mathbf{W}_{i \cdot} \mathbf{v} + b_i^h) \quad (5)$$

where  $\text{soft}_+(x) = \ln(1 + e^x)$  and the notation  $\mathbf{W}_{i \cdot}$  designates the  $i^{\text{th}}$  row of  $\mathbf{W}$ , likewise for columns  $\mathbf{W}_{\cdot j}$ . This allows for an equivalent definition of the RBM model in terms of what is known as the free energy  $F(\mathbf{v})$ . However, the partition function still requires summing over all configurations of visible vectors, which is intractable even for moderate values of  $D$ .

RBM's can be learned as generative models, to assign high probability (i.e. low energy) to training observations and low probability otherwise. One approach is to minimize the average negative log-likelihood (NLL) for a set of examples  $\mathcal{D} = \{\mathbf{v}_n\}_{n=1}^N$ :

$$f(\Theta, \mathcal{D}) = \frac{1}{N} \sum_{n=1}^N -\ln P(\mathbf{v}_n). \quad (6)$$

The gradient of this objective has a simple form, which is often referred to as the combination of positive and negative phases:

$$\nabla_{\theta} f(\Theta, \mathcal{D}) = \underbrace{\frac{1}{N} \sum_{n=1}^N \nabla_{\theta} F(\mathbf{v}_n)}_{\text{Positive phase}} - \underbrace{\sum_{\mathbf{v}' \in \mathcal{V}} P(\mathbf{v}') \nabla_{\theta} F(\mathbf{v}')}_{\text{Negative phase}} \quad (7)$$

where

$$\nabla_{\mathbf{W}} F(\mathbf{v}) = -\mathbb{E}[\mathbf{h}|\mathbf{v}] \mathbf{v}^T = -\widehat{\mathbf{h}}(\mathbf{v}) \mathbf{v}^T \quad (8)$$

$$\nabla_{\mathbf{b}^h} F(\mathbf{v}) = -\mathbb{E}[\mathbf{h}|\mathbf{v}] = -\widehat{\mathbf{h}}(\mathbf{v}) \quad (9)$$

$$\nabla_{\mathbf{b}^v} F(\mathbf{v}) = -\mathbf{v} \quad (10)$$

and where  $\widehat{\mathbf{h}}(\mathbf{v}) = \boldsymbol{\sigma}(\mathbf{W}\mathbf{v} + \mathbf{b}^h)$  with  $\boldsymbol{\sigma}(\cdot)$  being the sigmoid function  $\sigma(x) = \frac{1}{1+e^{-x}}$  applied element-wise. Derivation for the partial derivatives can be found in Appendix A.

Intuitively, the positive phase pushes up the probability of examples coming from our training set, whereas the negative phase lowers the probability of examples generated by the model. Much like the partition function, the negative phase is intractable. To overcome this we approximate the expectation under  $P(\mathbf{v})$  with an average of  $S$  samples  $\mathcal{S} = \{\hat{\mathbf{v}}_s\}_{s=1}^S$  drawn from  $P(\mathbf{v})$  i.e. the model.

$$\nabla_{\theta} f(\Theta, \mathcal{D}) \approx \underbrace{\frac{1}{N} \sum_{n=1}^N \nabla_{\theta} F(\mathbf{v}_n)}_{\text{Positive phase}} - \underbrace{\frac{1}{S} \sum_{s=1}^S \nabla F(\hat{\mathbf{v}}_s)}_{\text{Negative phase}} \quad (11)$$

Moreover, mini-batch training is usually employed and consists in replacing the positive phase average by one over a small subset of the training set, different for every training update.

Sampling from  $P(\mathbf{v})$  can be achieved using block Gibbs sampling, by alternating between sampling  $\mathbf{v} \sim P(\mathbf{v}|\mathbf{h})$  and  $\mathbf{h} \sim P(\mathbf{h}|\mathbf{v})$ . It can be done efficiently because RBMs have no connections within a layer, meaning that hidden units are conditionally independent given the visible units and vice versa. The conditional distributions of a binary RBM are Bernoulli distributions with parameters

$$P(h_i = 1|\mathbf{v}) = \sigma(\mathbf{W}_{i \cdot} \mathbf{v} + b_i^h) \quad (12)$$

$$P(v_j = 1|\mathbf{h}) = \sigma(\mathbf{h}^T \mathbf{W}_{\cdot j} + b_j^v) \quad (13)$$

In theory, the Markov chain should be run until equilibrium before drawing a sample for every training update, which is highly inefficient. Thus, Contrastive Divergence (CD) learning is often employed, where we initialize the update's Gibbs chains to the training examples and only perform  $T$  steps of Gibbs sampling (Hinton, 2002). Another approach, referred to as stochastic approximation or Persistent CD (PCD) (Tieleman, 2008), is to not reinitialize the Gibbs chains between updates.

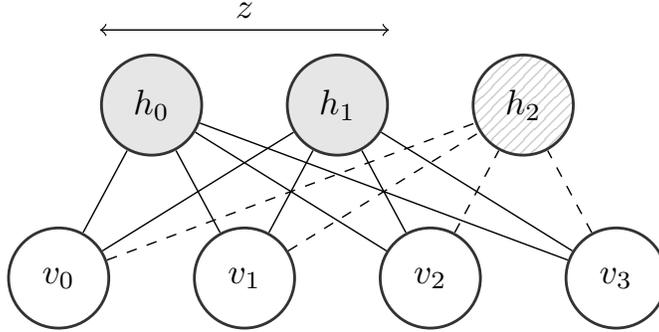


Figure 2: Illustration of the ordered RBM. Since  $z = 2$  only the first two hidden units are selected.

### 3 Ordered Restricted Boltzmann Machine

The model we propose is a variant of the RBM where the hidden units  $\mathbf{h}$  are ordered from left to right, with this order being taken into account by the energy function. We refer to this model as an ordered RBM (oRBM). As shown in Figure 2, the oRBM takes hidden unit order into account by introducing a random variable  $z$  that can be understood as the effective number of hidden units participating to the energy. Hidden units are selected starting from the left and the selection of each hidden unit is associated with an incremental cost in energy.

Concretely, we define the energy function of the oRBM as

$$E(\mathbf{v}, \mathbf{h}, z) = -\mathbf{v}^T \mathbf{b}^v - \sum_{i=1}^z (h_i(\mathbf{W}_i \cdot \mathbf{v} + b_i^h) - \beta_i) \quad (14)$$

where  $z$  represents the number of selected hidden units that are active and  $\beta_i$  is a energy penalty for selecting each  $i^{\text{th}}$  hidden unit. As we will see, carefully parametrizing the per unit energy penalty will allow us to consider the case of an infinite pool of hidden units.

In our experiments, as we wanted the filters of each unit to be the dominating factor in a unit being selected, we parametrized it as  $\beta_i = \beta \text{soft}_+(b_i^h)$ , where  $\beta$  is a global hyper-parameter (critically, as we’ll discuss later, this hyper-parameter doesn’t actually require tuning and a generic value for it works fine). Intuitively, the penalty term acts as a form of regularization since it forces the model to avoid using more hidden units than needed, prioritizing smaller networks.

Moreover, having the penalty depending on the hidden biases also implies that the selection of a hidden units (i.e. influencing the outcome of the random variable  $z$ ) will be mostly controlled by the values taken by the connections  $\mathbf{W}$ . Higher values of the bias of a hidden unit will not increase its probability of being selected. In other words, for the model to increase its capacity and better fit the training data, it will have to learn better filters. Note that alternative parametrizations could certainly be considered.

As with the RBM,  $P(\mathbf{v})$  is defined in terms of its energy function. For this, we have to specify the set of legal values for  $\mathbf{v}$ ,  $\mathbf{h}$  and  $z$ . Since, for a given  $z$ , the value

of the energy is irrelevant for the dimensions of  $\mathbf{h}$  from  $z$  to  $K$ , we will assume they are set to 0. There is thus a coupling between the value of  $z$  and the legal values of  $\mathbf{h}$ . We will note  $\mathcal{H}_z = \{\mathbf{h} \in \mathcal{H} | h_k = 0 \ \forall k > z\}$  the legal values of  $\mathbf{h}$  for a given  $z$ . As for  $z$ , it can vary in  $\{1, \dots, K\}$ , and  $\mathbf{v} \in \mathcal{V}$  as usual.

The joint probability over  $\mathbf{v}$ ,  $\mathbf{h}$  and  $z$  is thus:

$$P(\mathbf{v}, \mathbf{h}, z) = \frac{1}{Z} e^{-E(\mathbf{v}, \mathbf{h}, z)} \quad (15)$$

where

$$Z = \sum_{z'=1}^K \sum_{\mathbf{v}' \in \mathcal{V}} \sum_{\mathbf{h}' \in \mathcal{H}_{z'}} e^{-E(\mathbf{v}', \mathbf{h}', z')}. \quad (16)$$

As for the marginal distribution  $P(\mathbf{v})$  of the oRBM model, it can also be written in terms of a free energy. Indeed, in a derivation similar to the case of the RBM, we can show:

$$P(\mathbf{v}) = \frac{1}{Z} \sum_{z=1}^K \sum_{\mathbf{h} \in \mathcal{H}_z} e^{-E(\mathbf{v}, \mathbf{h}, z)} = \frac{1}{Z} \sum_{z=1}^K e^{-F(\mathbf{v}, z)} \quad (17)$$

$$F(\mathbf{v}, z) = -\mathbf{v}^T \mathbf{b}^v - \sum_{i=1}^z (\text{soft}_+(\mathbf{W}_i \cdot \mathbf{v} + b_i^h) - \beta_i) \quad (18)$$

This gives us a free energy where only the hidden units have been marginalized. We can also derive a formulation where the free energy depends only on  $\mathbf{v}$ :

$$P(\mathbf{v}) = \frac{1}{Z} \sum_{z=1}^K e^{-F(\mathbf{v}, z)} = \frac{1}{Z} e^{-F(\mathbf{v})} \quad \text{with} \quad F(\mathbf{v}) = -\ln \left( \sum_{z=1}^K e^{-F(\mathbf{v}, z)} \right) \quad (19)$$

It should be noticed that, in the oRBM,  $z$  does not correspond to the number of hidden units assumed to have generated *all* observations. Instead, the model allows for different observations having been generated by a different number of hidden units. Specifically, for a given  $\mathbf{v}$ , the conditional distribution over the corresponding value of  $z$  is

$$P(z|\mathbf{v}) = \frac{\exp(-F(\mathbf{v}, z))}{\sum_{z'=1}^K \exp(-F(\mathbf{v}, z'))}. \quad (20)$$

As for the conditional distribution over the hidden units, given a value of  $z$  it takes the same form as for the regular RBM, except for unselected hidden units which are forced to zero. Similarly, the distribution of  $\mathbf{v}$  given a value of the hidden layer and  $z$  reflects that of the RBM:

$$P(h_i = 1 | \mathbf{v}, z) = \begin{cases} \sigma(\mathbf{W}_i \cdot \mathbf{v} + b_i^h) & \text{if } i \leq z \\ 0 & \text{otherwise} \end{cases} \quad (21)$$

$$P(v_j = 1 | \mathbf{h}, z) = \sigma \left( \sum_{i=1}^z W_{ij} h_i + b_j^v \right) \quad (22)$$

To train the oRBM, we can also rely on CD or PCD for estimating the gradients based on Equation 11 but using  $F(\mathbf{v})$  as defined in equation 19. Defining

$\mathbf{1}_z = \overbrace{[1, \dots, 1, 0, \dots, 0]}^z$  and  $\mathbf{cdf}(z|\mathbf{v}) = [P(z < 1|\mathbf{v}), \dots, P(z < K|\mathbf{v})]^T$  with  $\odot$  denoting the element-wise product, the free energy gradients are then slightly modified as follows:

$$\nabla_{\mathbf{w}} F(\mathbf{v}) = - \mathbb{E}_{\mathbf{h}, z} [\mathbf{h} \odot \mathbf{1}_z | \mathbf{v}] \mathbf{v}^T = -(\widehat{\mathbf{h}}(\mathbf{v}) \odot (1 - \mathbf{cdf}(z|\mathbf{v}))) \mathbf{v}^T \quad (23)$$

$$\nabla_{\mathbf{b}^h} F(\mathbf{v}) = - \mathbb{E}_{\mathbf{h}, z} [(\mathbf{h} - \beta \boldsymbol{\sigma}(\mathbf{b}^h)) \odot \mathbf{1}_z | \mathbf{v}] = -(\widehat{\mathbf{h}}(\mathbf{v}) - \beta \boldsymbol{\sigma}(\mathbf{b}^h)) \odot (1 - \mathbf{cdf}(z|\mathbf{v})) \quad (24)$$

$$\nabla_{\mathbf{b}^v} F(\mathbf{v}) = -\mathbf{v} \quad (25)$$

with  $\widehat{\mathbf{h}}(\mathbf{v}) = \boldsymbol{\sigma}(\mathbf{W}\mathbf{v} + \mathbf{b}^h)$ . Derivation for the partial derivatives can be found in Appendix A.

Compared to the RBM, computing these gradients requires one additional quantity: the vector of cumulative probabilities  $\mathbf{cdf}(z|\mathbf{v})$ . Fortunately, this quantity can be efficiently computed, in  $O(K)$ , by first computing the required probabilities vector  $P(z|\mathbf{v})$  and performing a cumulative sum.

Sampling from  $P(\mathbf{v})$  slightly differs from the RBM as we need to consider  $z$  in the Markov chain. With the oRBM, Gibbs steps alternate between sampling  $(\mathbf{h}, z) \sim P(\mathbf{h}, z|\mathbf{v})$  and  $\mathbf{v} \sim P(\mathbf{v}|\mathbf{h}, z)$ . Sampling from  $P(\mathbf{h}, z|\mathbf{v})$  is done in two steps:  $z \sim P(z|\mathbf{v})$  followed by  $\mathbf{h} \sim P(\mathbf{h}|\mathbf{v}, z)$ .

During training, what we observe is that the hidden units are each trained gradually, in sequence, from left to right. This effect is mainly due to the multiplicative term  $(1 - \mathbf{cdf}(z|\mathbf{v}))$  in the hidden unit parameter updates of Equations 23 and 24, which is monotonically decreasing. Effectively, the model is thus growing in capacity during training, until its maximum capacity of  $K$  hidden units.

## 4 Infinite Restricted Boltzmann Machine

The growing behaviour of the oRBM begs for the question: could we achieve a similar effect without having to specify a maximum capacity to the model? Indeed, while Montufar and Ay (2010) have shown that with  $2^{V-1} - 1$  hidden units an RBM is a universal approximator, a variant of the RBM that could automatically increase its capacity until it is sufficiently high is likely to yield much smaller models in practice.

It turns out that this is possible, by taking the limit of  $K \rightarrow \infty$ . For this reason, we refer to this model as the infinite RBM (iRBM).

This limit is made possible thanks to two modeling choices. The first is the assumption that a finite (but variable!) number of hidden units have non-zero weights and biases. This is trivial to ensure, for any optimization procedure, using any amount of any type of weight decay (e.g. L2 or L1 regularization) on all the weights and hidden biases. An infinite number of non-zero weights and biases

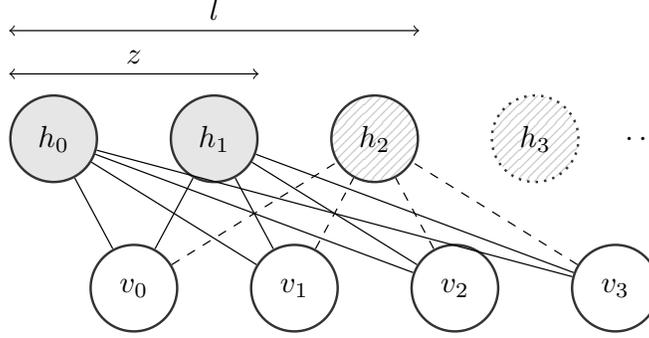


Figure 3: Illustration of the infinite RBM. With  $z = 2$ , only the first two hidden units are currently selected. The dashed lines illustrate that there are connections that are trained (non-zero) with the third hidden unit. All (infinitely many) hidden units after the third have zero-valued weights, which correspond to  $l$  being equal to 3.

could then correspond to an infinite penalty, so no proper optimization would ever diverge to this solution, no matter the initialization. This is guaranteed when using L1 regularization, thanks to its sparsity inducing property. As for L2 regularization, while it could theoretically lead to an infinite number of hidden units (e.g. if the L2 norm of the parameters associated with each hidden unit decreases exponentially with respect to the position of the hidden unit), in practice the floating precision would clip very small parameters to zero, thus having a finite number of hidden units.

The second key choice is our parametrization of the per-unit energy penalty  $\beta_i$ , which will ensure that the infinite sums required in computing probabilities will be convergent. For instance, consider the conditional  $P(z|\mathbf{v})$ :

$$P(z|\mathbf{v}) = \frac{\exp(-F(\mathbf{v}, z))}{Z(\mathbf{v})} = \frac{\exp(-F(\mathbf{v}, z))}{\sum_{z'=1}^{\infty} \exp(-F(\mathbf{v}, z'))} \quad (26)$$

Let's note  $l$  the number of effectively trained hidden units, i.e. where all hidden units  $> l$  have zero weights and biases. This is guaranteed to happen thanks to the growing behaviour that ensures hidden units are ordered from left to right. Then, we can split the normalization constant  $Z(\mathbf{v})$  of Equation 26 into two parts, split at  $z = l$ , as follows:

$$\begin{aligned} Z(\mathbf{v}) &= \sum_{z=1}^l \exp(-F(\mathbf{v}, z)) + \sum_{z=l+1}^{\infty} \exp(-F(\mathbf{v}, z)) \\ &= \sum_{z=1}^l \exp(-F(\mathbf{v}, z)) + \sum_{z=l+1}^{\infty} \exp\left(-F(\mathbf{v}, l) + \sum_{i=l+1}^z \text{soft}_+(\mathbf{W}_i \cdot \mathbf{v} + b_i^h) - \beta_i\right) \\ &= \sum_{z=1}^l \exp(-F(\mathbf{v}, z)) + \exp(-F(\mathbf{v}, l)) \underbrace{\sum_{z=1}^{\infty} \exp((1 - \beta)\text{soft}_+(0))^z}_{\text{Geometric series}} \end{aligned} \quad (27)$$

where Equation 27 is obtained by exploiting the fact that all weights and biases of hidden units at position  $l + 1$  and higher are zero. By ensuring that  $\beta > 1$ , the geometric series of Equation 27 is finite and can be analytically computed. This in turn implies that  $P(z|\mathbf{v})$  is tractable and can be sampled from. Following a similar reasoning, the global partition function  $Z$  can be shown to be finite (see Appendix B), thus yielding a properly defined joint distribution for any configurations with a finite number of non-zero weights and hidden biases.

One could think that, compared to a regular RBM, we have merely traded the hyper-parameter of the hidden layer size with the hyper-parameter  $\beta$ . However, crucially,  $\beta$ 's role is only to ensure that the iRBM is properly defined, and the penalty it imposes in the energy function can be compensated by the learned parameters. The extent to which the parameters can grow enough to compensate for that penalty is then controlled by the strength of weight decay, a hyper-parameter the iRBM shares with the RBM. We've thus effectively removed one hyper-parameter. Moreover, we've indeed observed that results are robust to the choice of  $\beta$ , that is finely tuning beta was not necessary to ultimately achieving good performance. While the choice of  $\beta$  can impact the number of epochs it would take for the weights to compensate for the penalty, this (the number of epochs) is a quantity that must be tuned anyways, even in regular RBMs.

The question of the identifiability of the binary RBM is a complex one, which has been studied (Cueto et al., 2010). Unlike the RBM, the iRBM is sensitive to the ordering of its hidden units, thanks to the penalty term. This means permutations of iRBM's hidden units do not correspond to the same distribution, making its parametrization more identifiable.

As for learning, it can be done mostly by following the procedure of the oRBM, i.e. minimizing the NLL with stochastic gradient descent using (Persistent) CD to approximate the gradients. One slight modification is required however. Indeed, since the free energy gradient for the hidden weights and biases can be non-zero for all (infinite) hidden units, we cannot use the gradient of Equations 23 and 24 for all hidden units.

To avoid this issue, we consider the following observation. Instead of using the derivative of  $F(\mathbf{v})$ , we could instead use the derivative of  $F(\mathbf{v}, z)$ , where  $z$  is obtained by sampling from  $P(z|\mathbf{v})$ :

$$\nabla_{\mathbf{w}} F(\mathbf{v}, z) = -\mathbb{E}_{\mathbf{h}}[\mathbf{h} \odot \mathbf{1}_z | z, \mathbf{v}] \mathbf{v}^T = -(\widehat{\mathbf{h}}(\mathbf{v}) \odot \mathbf{1}_z) \mathbf{v}^T \quad (28)$$

$$\nabla_{\mathbf{b}^h} F(\mathbf{v}, z) = -\mathbb{E}_{\mathbf{h}}[(\mathbf{h} - \beta \boldsymbol{\sigma}(\mathbf{b}^h)) \odot \mathbf{1}_z | z, \mathbf{v}] = -(\widehat{\mathbf{h}}(\mathbf{v}) - \beta \boldsymbol{\sigma}(\mathbf{b}^h)) \odot \mathbf{1}_z . \quad (29)$$

In this case, all weights and biases with an index greater than the sampled  $z$  have a gradient of zero, i.e. do not require any update. Moreover, the expectation of these gradients with respect to  $z$  (conditioned on  $\mathbf{v}$ ) are the gradients of  $F(\mathbf{v})$ , making them unbiased in this respect. This comes at the cost of higher variance in the updates. But thanks to this observation, we are justified to use a hybrid approach, where we use the  $F(\mathbf{v})$  gradients only for the units with index less or equal than  $l$ , and "use" the gradient of  $F(\mathbf{v}, z)$  for the other units, i.e. leave them set to zero.

As previously mentioned, we use weight decay to ensure that the number of non-zero parameters cannot diverge to infinity. For practical reasons, our implementation also used a capacity-limiting heuristic. If the Gibbs sampling chain ever sampled a value for  $z$  that is greater than  $l$ , then we clamped it to  $l + 1$ . Intuitively, this corresponds to "adding" a single hidden unit. This avoids filling all the memory in the (unlikely) event where we'd draw a large value for  $z$ . When adding a hidden unit, its associated weights and biases are initialized to zero.

We emphasize that these were not required to avoid divergence (weight decay is sufficient): it merely ensured a practical and efficient implementation of the model on the GPU. Note also that when using L1 regularisation,  $l$  can decrease in value, thanks to the sparsity promoting property of the L1 norm. Again, we highlight that while a finite number of weights and biases is maintained, that number of such weights does vary and is learned, while the implicit number of hidden units is indeed infinite (infinitely many contribute to the partition function).

## 5 Related Work

This work falls within the research literature on discovering extensions of the original RBM model to different contexts and objectives. Of note here is the implicit mixture of RBMs (Nair and Hinton, 2008). Indeed, the oRBM can be interpreted as a special case of an implicit mixture of RBMs. Writing  $P(\mathbf{v})$  as  $\sum_{z=1}^K P(z)P(\mathbf{v}|z)$  we see that the oRBM is an implicit mixture of  $K$  RBMs, where each RBM has a different number of hidden units (from 1 to  $K$ ) and the weights are tied between RBMs. The prior  $P(z)$  represents the probability of using the  $z^{\text{th}}$  RBM and is also derived from the energy function. However, as in the implicit mixture of RBMs,  $P(z)$  is intractable as it would require the value of the partition function. That said, the work of Nair and Hinton (2008) is otherwise very different and did not address the question of having an RBM with adaptive capacity.

Another related work is that of the Cardinality RBMs proposed by Swersky et al. (2012). They used a cardinality potential to control the sparsity of the RBM, i.e. limiting the number of hidden units that can be active. In the oRBM and the iRBM,  $z$  effectively acts as an upper bound on the number of hidden units  $h_i$  that can be equal to 1, since we are limiting  $\mathbf{h}$  to be in  $\mathcal{H}_z$ , a subset of  $\mathcal{H}$ . In their work, Swersky et al. (2012) use cardinality potentials that allow only configurations having at most  $k$  active hidden units. One difference with our work however is that their cardinality potential is order agnostic, meaning that the active hidden units can be positioned anywhere within the hidden layer while still satisfying the cardinality potential. On the other hand, in the oRBM, all units with index higher than  $z$  must be set to zero, with only the previous hidden units being allowed to be active. In addition, their parameter  $k$  is fixed during training whereas our number of active hidden units  $z$  changes depending on the input.

The oRBM also bears some similarity with autoencoders trained by a nested version of dropout (Rippel et al., 2014). Nested dropout works by stochastically selecting the number of hidden units used to reconstruct an input example at

training time, and so independently for each update and example. Rippel et al. (2014) showed that this defines a learning objective that makes the solution identifiable and no longer invariant to hidden unit permutation. In addition to being concerned with a different type of model, this work doesn't discuss the case of an unbounded and adaptive hidden layer size.

Welling et al. (2003) proposed a self supervised boosting approach, which is applicable to the RBM and in which hidden units are sequentially added and trained. However, like boosting in general and unlike the iRBM, this procedure trains each hidden unit greedily instead of jointly, which could lead to much larger networks than necessary. Moreover, it is not easily generalizable to online learning.

While the work on unsupervised neural networks with adaptive hidden layer size is otherwise relatively scarce, there's been much more work in the context of supervised learning. There is the well known work of Fahlman and Lebiere (1990) on Cascade-Correlation networks. More recently, Zhou et al. (2012) proposed a procedure for learning discriminative features with a denoising autoencoder (a model related to the RBM). The procedure is also applicable to the online setting. It relies on invoking two heuristics that either add or merge hidden units during training. We note that the iRBM framework could easily be generalized to discriminative and hybrid training as in Zhou et al. (2012). The corresponding mechanisms for adding and merging units would then be implicitly derived from gradient descent on the corresponding supervised training objective.

Finally, we highlight that our model is not based on a Bayesian formulation, as most of the literature on infinite models. On the other hand, it does correspond to the infinite limit of a finite-sized model and yields a model that can learn its size with training.

## 6 Experiments

We compare the performance of the oRBM and the iRBM with the classic RBM on two datasets: binarized MNIST (Salakhutdinov and Murray, 2008) and Cal-Tech101 Silhouettes (Marlin et al.). We aim to demonstrate that the iRBM effectively removes the need of tuning an hyper-parameter for the hidden layer size while still achieving comparable performance to the standard RBM. The code to reproduce the experiments of the paper is available on GitHub<sup>1</sup>. Our implementation is done using Theano (Bastien et al., 2012; Bergstra et al., 2010).

For completeness, we wish to mention that more sophisticated or deep models have reported results on one or both of these datasets (e.g. EoNADE (Uria et al., 2014), DBNs (Murray and Salakhutdinov, 2009), Deep autoregressive networks (Gregor et al., 2014), Iterative Neural Autoregressive Distribution Estimator (Raiko and Bengio, 2014)) that improve on the standard RBM. However, since our objective with the iRBM is to effectively remove a hyper-parameter of the RBM, instead of achieving improved performances, we focus our comparison on this baseline.

---

<sup>1</sup><http://github.com/MarcCote/iRBM>

All NLL results of this section were obtained by estimating the log-partition function  $\ln \hat{Z}$  using Annealed Importance Sampling (AIS) (Salakhutdinov and Murray, 2008) with 100,000 intermediate distributions and 5000 chains. As an additional validation step, samples were generated from best models and visually inspected.

Each model was trained with mini-batch stochastic gradient descent using batch size of 64 examples and using PCD with 10 Gibbs steps between parameter updates. We used the ADAGRAD stochastic gradient update (Duchi et al., 2011), a per-dimension learning rate method, to train the oRBMs and the iRBMs. We found that having different learning rates for different hidden units was very beneficial, since units positioned earlier in the hidden layer will approach convergence faster than units to their right, and thus will benefit from a learning rate decaying more rapidly. We tried several learning rates  $lr \in \{5 \times 10^{-1}, 10^{-1}, 5 \times 10^{-2}, 10^{-2}\}$  and always set ADAGRAD’s epsilon parameter to  $10^{-6}$ .

We also tested different values for both L1 and L2 regularization’s factor  $\lambda \in \{0, 10^{-2}, 10^{-3}, 10^{-4}, 10^{-5}\}$ . Note that we allow the iRBM to shrink only if L1 regularization is used.

We did try varying the  $\beta$  found in the penalty term and as expected we’ve found results to be robust to its value. Since  $\beta$  must be greater than 1, we explored positive constants to add to 1, on a log scale (1, 0.25, 0.1, 0.01, 0.001, etc.). We settled on using  $\beta = 1.01$  for all experiments as it provides a penalty high enough to have a growing behavior and requires around five hundred epochs for the weights to compensate for the penalty.

Finally, we note that improved performances could certainly have been achieved using an improved sampler (e.g. parallel tempering (Desjardins et al., 2010)) or parametrization (e.g. enhanced gradient parametrization (Cho et al., 2013)). However, these changes would equally improve the baseline RBM, so we decided to concentrate on this more common learning setup.

## 6.1 Binarized MNIST

The MNIST dataset<sup>2</sup> is composed of 70,000 images of size 28x28 pixels representing handwritten digits (0-9). Images have been stochastically binarized according to their pixel intensity as in Salakhutdinov and Murray (2008). We use the same split as in Larochelle and Murray (2011), corresponding to 50,000 examples for training, 10,000 for validation and 10,000 for testing.

Each model was trained up to 5000 epochs but we performed AIS evaluation every 1000 epochs and kept the model having the best NLL approximation on the valid set. We report the associated NLL approximations obtained on the test set. Taking after past studies assessing RBM results on binarized MNIST, we fixed the number of hidden units to 500 for the RBM and the oRBM. Best results for the RBM, oRBM and iRBM are reported in Table 1. The oRBM and the iRBM models reach competitive performance compared to the RBM. Samples from all three models are illustrated in Figure 5.

---

<sup>2</sup><http://yann.lecun.com/exdb/mnist>

Table 1: Average NLL on binarized MNIST test set for best RBMs, oRBM and iRBM. Partition functions were estimated using AIS with 100,000 intermediate distributions and 5000 chains. The confidence interval on the average NLL assumes  $\ln \hat{Z}$  has no variance and reflects the confidence of a finite sample average. By taking the uncertainty about the partition function into account, the interval would be larger.

MODEL	SIZE	BINARIZED MNIST		
		$\ln \hat{Z}$	$\ln(\hat{Z} \pm 3\sigma)$	AVG. NLL
RBM	100	600.92	[600.88, 600.95]	$98.17 \pm 0.52$
RBM	500	613.28	[613.24, 613.31]	$86.50 \pm 0.44$
RBM	2000	1099.07	[1098.94, 1099.17]	$85.03 \pm 0.42$
oRBM	500	40.06	[39.90, 40.19]	$88.15 \pm 0.46$
iRBM	1208	40.32	[40.03, 40.54]	$85.65 \pm 0.44$

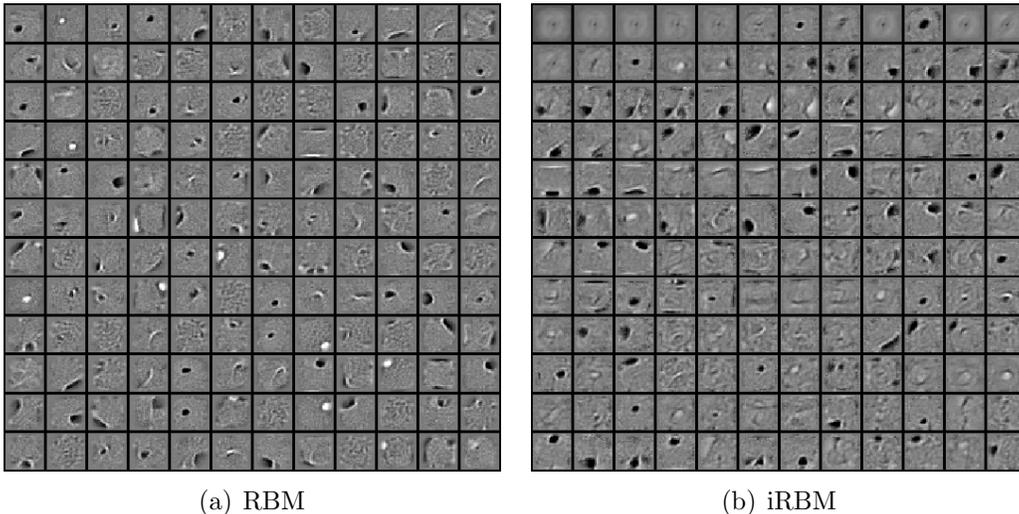


Figure 4: Comparing the filters of an RBM and an iRBM both trained on binarized MNIST. The first 96 filters are shown starting from the top-left corner and incrementing across columns first.

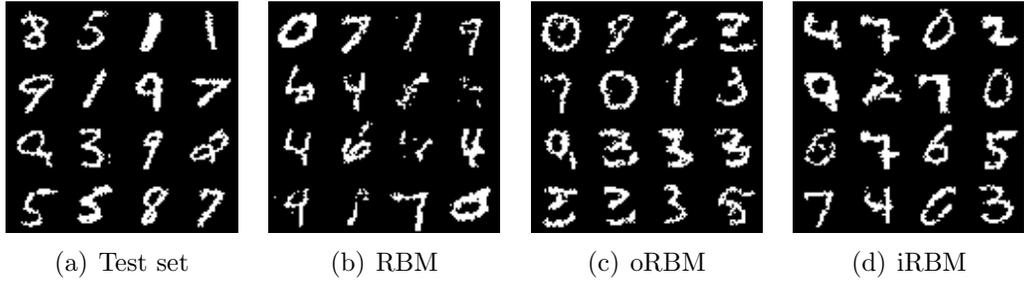


Figure 5: Comparison between data from binarized MNIST and random samples generated from the three models by randomly initializing visible units and running 10,000 Gibbs steps. The RBM and oRBM both have 500 hidden units, whereas the iRBM final size is 1208 hidden units.

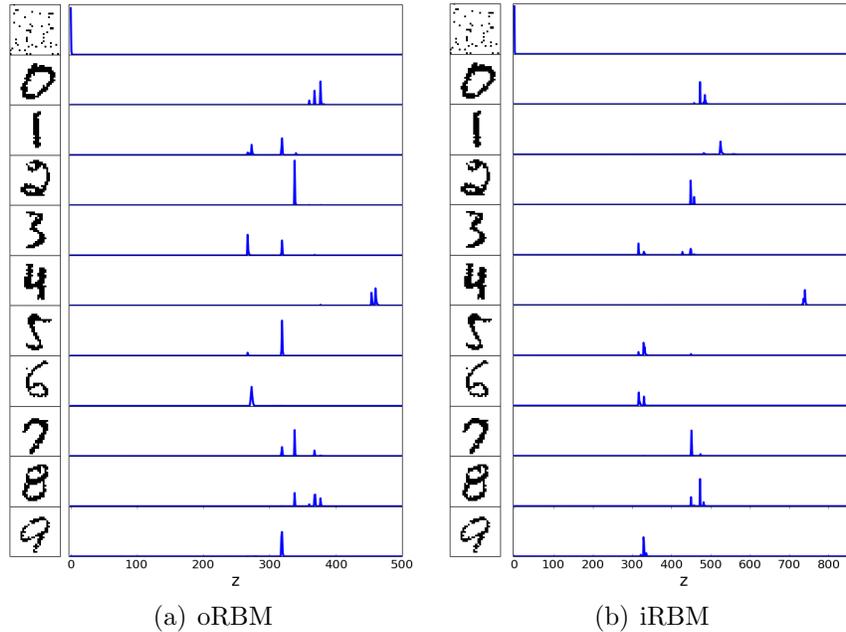


Figure 6: Each row shows a plot of  $P(z|\mathbf{v})$  where  $\mathbf{v}$  is a given example from MNIST test set and is displayed to the left. The first row illustrates the impact of a noisy image on sampling  $z$ . As explained in Section 3 of the paper, we see that different input images are related to different values for the number  $z$  of used units.

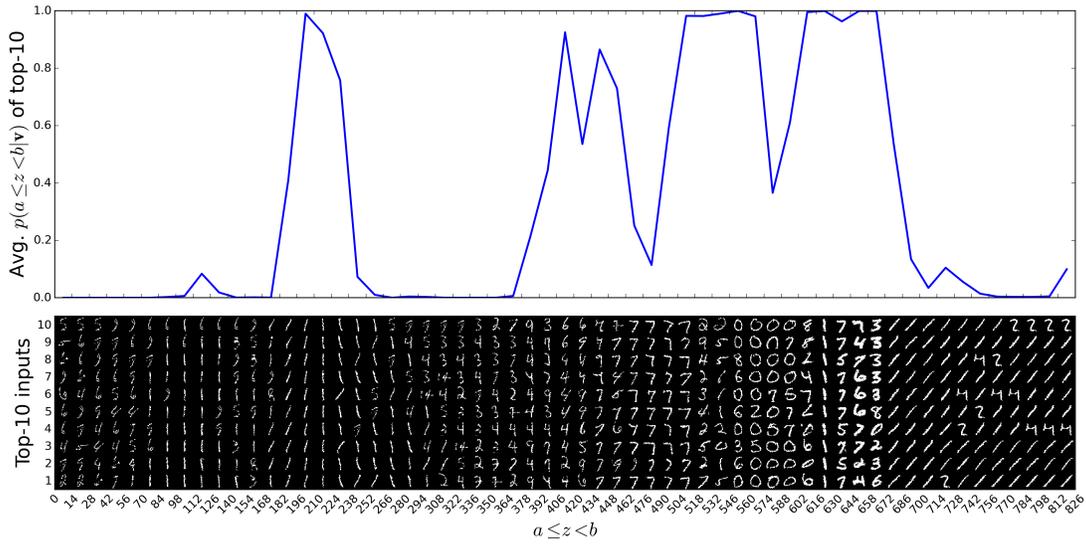


Figure 7: **(Bottom)** Top 10 inputs from the test set with highest value of  $P(z|\mathbf{v})$  within different intervals for  $z$ , i.e.  $\text{argmax}_{\mathbf{v}} P(a \leq z < b|\mathbf{v})$  for different intervals  $[a, b)$ . Interestingly, bolder inputs seem to be related to bigger values for the number  $z$  of used units. Also, simpler characters (e.g. "ones") tend to favor smaller values of  $z$  compared to more complex characters. **(Top)** Average of  $P(a \leq z < b|\mathbf{v})$  over the top 10 inputs. Low values highlights regions in the hidden layer where the hidden units are only useful when taken together with hidden units further right in the layer.

The best RBM (500 hidden units) was trained without any regularization and  $lr = 10^{-2}$  for 5000 epochs. We used our own implementation to train the RBM, which is why our result slightly differs from what is reported by Salakhutdinov and Murray (2008). The difference can be justified by the fact that they used the full 60,000 training set images, instead of a 50,000 subset. Also, they use a custom schedule to gradually increase the number of CD steps during training. That said, the oRBM and the iRBM would probably also benefit from having more training data and an improved sampling strategy.

The best oRBM (500 hidden units) was trained without any regularization and  $lr = 3 \times 10^{-2}$  for 500 epochs. After 3000 epochs, the best iRBM had 1208 hidden units with non-zero weights. It was trained with L1 regularization using a regularization factor of  $\lambda = 10^{-4}$  and  $lr = 5 \times 10^{-2}$ .

To show that our best iRBM does find an appropriate number of hidden units, we compared it with two other RBMs having respectively 100 and 2000 hidden units. Both were trained for 5000 epochs without any regularization and respectively with  $lr = 10^{-1}$  and  $lr = 10^{-2}$ . Results are reported in Table 1 where we can see the oRBM and the iRBM still achieve competitive results compared to the RBM with 2000 hidden units.

Figure 4 shows the ordering effect on the filters obtained with an iRBM. The ordering is even more apparent when observing the hidden unit filters during training. We generated a video of this visualization, illustrating the filter values and the generated negative samples at epochs 1, 10, 50 and 100. See link: <http://goo.gl/LGQDaI>.

Interestingly, we’ve observed that Gibbs sampling can mix much more slowly with the oRBM. The reason is the addition of variable  $z$  increases the dependence between states and thus hurts the convergence of Gibbs sampling. In particular, we observed that when the Gibbs chain is in a state corresponding to a noisy image without any structure, it can require many steps before stepping out of this region of the input space. Yet, comparing the free energy of such random images and images that resemble digits confirmed that these random images have significantly higher free energy (and thus are unlikely samples of the model). Figure 6 also confirms the high dependence between  $z$  and  $\mathbf{v}$ : the distribution of the unstructured image is peaked at  $z = 1$ , while all digits prefer values of  $z$  greater than 250. To fix this issue, we’ve found that simply initializing the Gibbs chain to  $z = K$  was sufficient. We used this when sampling from a trained oRBM model.

The iRBM doesn’t seem to suffer as much from a low mixing rate and thus doesn’t require the  $z = K$  initialization heuristic for sampling. In fact, using the heuristic when sampling from an iRBM has almost no impact on the final samples when running 10,000 Gibbs steps. This could be an artefact of the model being trained progressively, i.e. we only add one hidden unit when sampling a large value for  $z$  bigger than  $l$ . Understanding how the lower mixing rate affects the proposed models and if a heuristic such as the one we mentioned earlier could be used to improve training is a topic left for future work.

We’ve also investigated what kind of inputs are maximizing  $P(z|\mathbf{v})$ , for different values of  $z$ . Using our best iRBM model trained with L1 regularization, we generated Figure 7. It highlights the fact that  $P(z|\mathbf{v})$  does capture some structure

Table 2: Average NLL on CalTech101 Silhouettes test set estimated using AIS with 100,000 intermediate distributions and 5000 chains. The confidence interval on the average NLL assumes  $\ln \hat{Z}$  has no variance and reflects the confidence of a finite sample average. By taking the uncertainty about the partition function into account, the interval would be larger.

CALTECH101 SILHOUETTES				
MODEL	SIZE	$\ln \hat{Z}$	$\ln(\hat{Z} \pm 3\sigma)$	AVG. NLL
RBM	100	2512.20	[2511.62, 2512.56]	$177.37 \pm 2.81$
RBM	500	2385.91	[2385.68, 2386.10]	$119.05 \pm 2.27$
RBM	2000	3353.47	[3349.85, 3354.15]	$118.29 \pm 2.25$
oRBM	500	1782.96	[1782.88, 1783.02]	$114.99 \pm 1.97$
iRBM	915	2000.08	[1999.93, 2000.22]	$121.47 \pm 2.07$

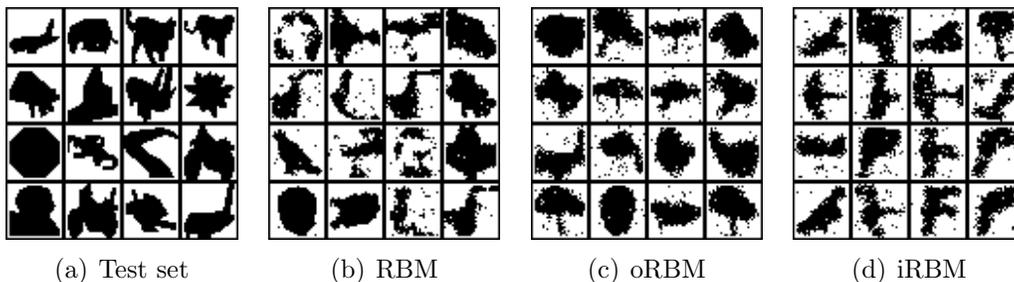


Figure 8: Comparison between data from CalTech101 Silhouettes and random samples generated from three models by randomly initializing visible units and running 10,000 Gibbs steps. The RBM and oRBM both have 500 hidden units, whereas the iRBM final size is 915 hidden units.

about the data, as the identity of the character with highest  $P(z|\mathbf{v})$  vary between different values of  $z$ .

## 6.2 CalTech101 Silhouettes

The CalTech101 Silhouettes dataset<sup>3</sup> (Marlin et al.) is composed of 8,671 images of size 28x28 binary pixels, representing object silhouettes (101 classes). The dataset is divided in three subsets: 4,100 examples for training, 2,264 for validation and 2,307 for testing.

Following a protocol similar to the one used for MNIST, each model was trained up to 5000 epochs, AIS evaluation was done every 1000 epochs. We report the NLL approximations obtained on the test set. Best results for the RBM, oRBM and iRBM are reported in Table 2. Again, the oRBM and the iRBM models reach

<sup>3</sup><http://people.cs.umass.edu/~marlin/data.shtml>

competitive performance compared to the RBM. Samples from all three models are illustrated in Figure 8.

The best RBM (500 hidden units) was trained without any regularization and  $lr = 10^{-2}$  for 3000 epochs. We used our own implementation to train the RBM. The best oRBM (500 hidden units) was trained with L1 regularization using a regularization factor of  $\lambda = 10^{-3}$  and  $lr = 10^{-2}$  for 5000 epochs. After 4000 epochs, the best iRBM had 915 hidden units with non-zero weights. It was trained with L1 regularization using a regularization factor of  $\lambda = 10^{-3}$  and  $lr = 5 \times 10^{-2}$ .

Again, to show that our best iRBM does find an appropriate number of hidden units, we compared it with two other RBMs having respectively 100 and 2000 hidden units. Both were trained without any regularization and respectively with  $lr = 10^{-1}$  for 5000 epochs and  $lr = 10^{-2}$  for 2000 epochs. Results are reported in Table 2 where we can see the oRBM and the iRBM still achieve competitive results compared to the RBM with 2000 hidden units.

## Conclusion

We proposed a novel extension of the RBM, the infinite RBM, which obviates the need to specify the hidden layer size. The iRBM is derived from the ordered RBM by taking the infinite limit of its hidden layer size. We presented a training procedure, derived from Contrastive Divergence, such that training the iRBM yields a learning procedure where the effective hidden layer size can grow.

In future work, we are interested in generalizing the idea of a growing latent representation to structures other than a flat vector representation. We are currently exploring extensions of the RBM allowing for a tree-structured latent representation. We believe a similar construction, involving a similar  $z$  random variable, should allow us to derive a training algorithm that also learns the latent representation's size.

## Acknowledgments

We thank NSERC for supporting this research, Nicolas Le Roux for discussions and comments, and Stanislas Lauly for making the iRBM's training video.

# A Partial derivatives

## Partial derivatives related to the RBM

Recall equation (5) representing the free energy of the RBM:

$$F(\mathbf{v}) = -\mathbf{v}^T \mathbf{b}^v - \sum_{i=1}^K \text{soft}_+(\mathbf{W}_i \cdot \mathbf{v} + b_i^h)$$

Taking the partial derivatives of  $F(\mathbf{v})$  w.r.t.  $W_{ij}$ ,  $b_i^h$  and  $b_j^v$  respectively, we obtain the following:

$$\frac{\partial F(\mathbf{v})}{\partial W_{ij}} = - \sum_{k=1}^K \frac{\partial \text{soft}_+(\mathbf{W}_k \cdot \mathbf{v} + b_k^h)}{\partial W_{ij}} = - \sum_{k=1}^K \sigma(\mathbf{W}_k \cdot \mathbf{v} + b_k^h) \frac{\partial \mathbf{W}_k \cdot \mathbf{v}}{\partial W_{ij}} = -\sigma(\mathbf{W}_i \cdot \mathbf{v} + b_i^h) v_j \quad (30)$$

$$\frac{\partial F(\mathbf{v})}{\partial b_i^h} = - \sum_{k=1}^K \frac{\partial \text{soft}_+(\mathbf{W}_k \cdot \mathbf{v} + b_k^h)}{\partial b_i^h} = - \sum_{k=1}^K \sigma(\mathbf{W}_k \cdot \mathbf{v} + b_k^h) \frac{\partial b_k^h}{\partial b_i^h} = -\sigma(\mathbf{W}_i \cdot \mathbf{v} + b_i^h) \quad (31)$$

$$\frac{\partial F(\mathbf{v})}{\partial b_j^v} = - \frac{\mathbf{v}^T \mathbf{b}^v}{\partial b_j^v} = -v_j \quad (32)$$

where  $\sigma(\mathbf{W}_i \cdot \mathbf{v} + b_i^h)$  can be expressed as a conditional expectation over  $h_i$  using equation (12)

$$\sigma(\mathbf{W}_i \cdot \mathbf{v} + b_i^h) = P(h_i = 1 | \mathbf{v}) = \sum_{h_i \in \{0,1\}} P(h_i = 1 | \mathbf{v}) h_i = \mathbb{E}[h_i | \mathbf{v}]$$

## Partial derivatives related to the oRBM and the iRBM

Recall equation (19) representing the free energy of the oRBM:

$$F(\mathbf{v}) = \ln \left( \sum_{z=1}^K e^{-F(\mathbf{v}, z)} \right)$$

where

$$F(\mathbf{v}, z) = -\mathbf{v}^T \mathbf{b}^v - \sum_{i=1}^z (\text{soft}_+(\mathbf{W}_i \cdot \mathbf{v} + b_i^h) - \beta \text{soft}_+(b_i^h))$$

The partial derivatives of  $F(\mathbf{v}, z)$  w.r.t.  $W_{ij}$ ,  $b_i^h$  and  $b_j^v$  are similar to equations (30), (31) and (32) from the RBM and are respectively given by

$$\begin{aligned}\frac{\partial F(\mathbf{v}, z)}{\partial W_{ij}} &= -\sum_{k=1}^z \frac{\partial \text{soft}_+(\mathbf{W}_k \cdot \mathbf{v} + b_k^h)}{\partial W_{ij}} \\ &= -\sum_{k=1}^z \sigma(\mathbf{W}_k \cdot \mathbf{v} + b_k^h) \frac{\partial \mathbf{W}_k \cdot \mathbf{v}}{\partial W_{ij}} \\ &= -H(z - i) \sigma(\mathbf{W}_i \cdot \mathbf{v} + b_i^h) v_j\end{aligned}\tag{33}$$

$$\begin{aligned}\frac{\partial F(\mathbf{v}, z)}{\partial b_i^h} &= -\sum_{k=1}^z \frac{\partial \text{soft}_+(\mathbf{W}_k \cdot \mathbf{v} + b_k^h) - \beta \text{soft}_+(b_k^h)}{\partial b_i^h} \\ &= -\sum_{k=1}^z (\sigma(\mathbf{W}_k \cdot \mathbf{v} + b_k^h) - \beta \sigma(b_k^h)) \frac{\partial b_k^h}{\partial b_i^h} \\ &= -H(z - i) (\sigma(\mathbf{W}_i \cdot \mathbf{v} + b_i^h) - \beta \sigma(b_i^h))\end{aligned}\tag{34}$$

$$\frac{\partial F(\mathbf{v}, z)}{\partial b_j^v} = -\frac{\mathbf{v}^T \mathbf{b}^v}{\partial b_j^v} = -v_j\tag{35}$$

with the Heaviside step function denoted as

$$H(n) = \begin{cases} 0, & n < 0 \\ 1, & n \geq 0 \end{cases}.$$

Then, the partial derivatives of  $F(\mathbf{v})$  w.r.t.  $W_{ij}$ ,  $b_i^h$  and  $b_j^v$  are obtained respec-

tively as follows:

$$\begin{aligned}
\frac{\partial F(\mathbf{v})}{\partial W_{ij}} &= \sum_{z=1}^K \frac{e^{-F(\mathbf{v},z)}}{\sum_{z'=1}^K e^{-F(\mathbf{v},z')}} \frac{\partial F(\mathbf{v},z)}{\partial W_{ij}} \\
&= - \sum_{z=1}^K P(z|\mathbf{v}) H(z-i) \sigma(\mathbf{W}_i \cdot \mathbf{v} + b_i^h) v_j \\
&= - \sum_{z=1}^K H(z-i) P(z|\mathbf{v}) P(h_i = 1|\mathbf{v}, z) v_j \\
&= - \sum_{z=1}^K \sum_{h_i \in \{0,1\}} H(z-i) h_i P(z|\mathbf{v}) P(h_i|\mathbf{v}, z) v_j \\
&= - \mathbb{E}_{h_i, z} [H(z-i) h_i | \mathbf{v}] v_j \tag{36}
\end{aligned}$$

$$\begin{aligned}
\frac{\partial F(\mathbf{v})}{\partial b_i^h} &= \sum_{z=1}^K P(z|\mathbf{v}) \frac{\partial F(\mathbf{v},z)}{\partial b_i^h} \\
&= - \sum_{z=1}^K P(z|\mathbf{v}) H(z-i) (\sigma(\mathbf{W}_i \cdot \mathbf{v} + b_i^h) - \beta \sigma(b_i^h)) \\
&= - \sum_{z=1}^K H(z-i) P(z|\mathbf{v}) (P(h_i = 1|\mathbf{v}, z) - \beta \sigma(b_i^h)) \\
&= - \sum_{z=1}^K H(z-i) P(z|\mathbf{v}) (-\beta \sigma(b_i^h)(1 - P(h_i = 1|\mathbf{v}, z)) + (1 - \beta \sigma(b_i^h))P(h_i = 1|\mathbf{v}, z)) \\
&= - \sum_{z=1}^K H(z-i) P(z|\mathbf{v}) ((0 - \beta \sigma(b_i^h))P(h_i = 0|\mathbf{v}, z) + (1 - \beta \sigma(b_i^h))P(h_i = 1|\mathbf{v}, z)) \\
&= - \sum_{z=1}^K \sum_{h_i \in \{0,1\}} H(z-i) (h_i - \beta \sigma(b_i^h)) P(z|\mathbf{v}) P(h_i|\mathbf{v}, z) \\
&= - \mathbb{E}_{h_i, z} [H(z-i) (h_i - \beta \sigma(b_i^h)) | \mathbf{v}] \tag{37}
\end{aligned}$$

$$\frac{\partial F(\mathbf{v})}{\partial b_j^y} = - \frac{\mathbf{v}^T \mathbf{b}^y}{\partial b_j^y} = -v_j$$

Observe that in equations (36) and (37),  $\sum_{z=1}^K P(z|\mathbf{v}) H(z-i)$  corresponds to  $P(z \geq i|\mathbf{v})$ . This then translates to  $(1 - \mathbf{cdf}(z|\mathbf{v}))$  when deriving the gradients as shown in equations (23) and (24).

## B Convergence of the partition function for the iRBM

We will show that the partition function  $Z$  of the iRBM is finite. To do so, we take the limit of  $K \rightarrow \infty$  of equation (16):

$$\begin{aligned} Z &= \sum_{\mathbf{v} \in \mathcal{V}} \prod_{z=1}^{\infty} \sum_{\mathbf{h} \in \mathcal{H}_z} e^{-E(\mathbf{v}, \mathbf{h}, z)} \\ &= \sum_{\mathbf{v} \in \mathcal{V}} \prod_{z=1}^{\infty} e^{-F(\mathbf{v}, z)} \\ &= \sum_{\mathbf{v} \in \mathcal{V}} Z(\mathbf{v}) \end{aligned} \tag{38}$$

Since the sum over all  $\mathbf{v} \in \mathcal{V}$  is finite and we know from equation (27) that  $Z(\mathbf{v})$  is finite, then  $Z$  is also finite.

## References

- Bastien, F., Lamblin, P., Pascanu, R., Bergstra, J., Goodfellow, I. J., Bergeron, A., Bouchard, N., and Bengio, Y. (2012). Theano: new features and speed improvements. Deep Learning and Unsupervised Feature Learning NIPS 2012 Workshop.
- Bergstra, J., Breuleux, O., Bastien, F., Lamblin, P., Pascanu, R., Desjardins, G., Turian, J., Warde-Farley, D., and Bengio, Y. (2010). Theano: a CPU and GPU math expression compiler. In *Proceedings of the Python for Scientific Computing Conference (SciPy)*.
- Cho, K., Raiko, T., and Ilin, A. (2013). Enhanced gradient for training restricted Boltzmann machines. *Neural computation*, 25:805–31.
- Cueto, M. A., Morton, J., and Sturmfels, B. (2010). Geometry of the restricted Boltzmann machine. In Viana, M. A. G. and Wynn, H. P., editors, *Algebraic methods in statistics and probability II, AMS Special Session*, volume 516 of *Contemporary Mathematics*. American Mathematical Society.
- Dahl, G. E., Adams, R. P., and Larochelle, H. (2012). Training Restricted Boltzmann Machines on Word Observations. In *Proceedings of the 29th International Conference on Machine Learning (ICML 2012)*, pages 679–686.
- Desjardins, G., Courville, A., Bengio, Y., Vincent, P., and Delalleau, O. (2010). Parallel Tempering for Training of Restricted Boltzmann Machines. *AISTATS*, 9:145–152.
- Duchi, J., Hazan, E., and Singer, Y. (2011). Adaptive subgradient methods for online learning and stochastic optimization. *J. Mach. Learn. Res.*, 12:2121–2159.
- Fahlman, S. E. and Lebiere, C. (1990). The cascade-correlation learning architecture. In Touretzky, D., editor, *Advances in Neural Information Processing Systems 2 (NIPS 1989)*, pages 524–532. Morgan-Kaufmann.
- Gregor, K., Mnih, A., and Wierstra, D. (2014). Deep AutoRegressive Networks. In *International Conference on Machine Learning*, pages 1242–1250.
- Hinton, G. (2002). Training products of experts by minimizing contrastive divergence. *Neural computation*.
- Larochelle, H. and Murray, I. (2011). The Neural Autoregressive Distribution Estimator. volume 15, pages 29–37. AISTATS.
- Marlin, B. M., Swersky, K., Chen, B., and de Freitas, N. Inductive Principles for Restricted Boltzmann Machine Learning. *Proc. Intl. Conference on Artificial Intelligence and Statistics*, pages 305–306.
- Montufar, G. and Ay, N. (2010). Refinements of Universal Approximation Results for Deep Belief Networks and Restricted Boltzmann Machines. *Neural Computation*, (2010):1–12.

- Murray, I. and Salakhutdinov, R. (2009). Evaluating probabilities under high-dimensional latent variable models. In *Advances in neural information processing systems*, pages 1137–1144.
- Nair, V. and Hinton, G. (2008). Implicit Mixtures of Restricted Boltzmann Machines. *NIPS*, pages 1–8.
- Orbanz, P. and Teh, Y. W. (2010). Bayesian Nonparametric Models. In *Encyclopedia of Machine Learning*. Springer.
- Raiko, T. and Bengio, Y. (2014). Iterative Neural Autoregressive Distribution. In *Advances in neural information processing systems*, pages 1–9.
- Ranzato, M., Krizhevsky, A., and Hinton, G. E. (2010). Factored 3-way restricted boltzmann machines for modeling natural images. *Journal of Machine Learning Research - Proceedings Track*, 9:621–628.
- Rippel, O., Gelbart, M., and Adams, R. (2014). Learning Ordered Representations with Nested Dropout. *arXiv preprint arXiv:1402.0915*.
- Salakhutdinov, R., Mnih, A., and Hinton, G. (2007). Restricted Boltzmann machines for collaborative filtering. In *Proceedings of the 24th International Conference on Machine Learning (ICML 2007)*, pages 791–798, New York, NY, USA. ACM.
- Salakhutdinov, R. and Murray, I. (2008). On the quantitative analysis of Deep Belief Networks. In McCallum, A. and Roweis, S., editors, *Proceedings of the 25th Annual International Conference on Machine Learning (ICML 2008)*, pages 872–879. Omnipress.
- Swersky, K., Sutskever, I., Tarlow, D., Zemel, R. S., Salakhutdinov, R. R., and Adams, R. P. (2012). Cardinality Restricted Boltzmann Machines. In *Advances in Neural Information Processing Systems*, pages 3293–3301.
- Taylor, G. W., Hinton, G. E., and Roweis, S. T. (2011). Two distributed-state models for generating high-dimensional time series. *Journal of Machine Learning Research*, 12:1025–1068.
- Tieleman, T. (2008). Training restricted Boltzmann machines using approximations to the likelihood gradient. *ICML; Vol. 307*, page 7.
- Uria, B., Murray, I., and Larochelle, H. (2014). A Deep and Tractable Density Estimator. In *International Conference on Machine Learning*, page 9.
- Welling, M., Zemel, R. S., and Hinton, G. E. (2003). Self supervised boosting. In Becker, S., Thrun, S., and Obermayer, K., editors, *Advances in Neural Information Processing Systems 15 (NIPS 2002)*, pages 681–688. MIT Press.

Zhou, G., Sohn, K., and Lee, H. (2012). Online incremental feature learning with denoising autoencoders. In Lawrence, N. D. and Girolami, M. A., editors, *Proceedings of the Fifteenth International Conference on Artificial Intelligence and Statistics (AISTATS-12)*, volume 22, pages 1453–1461.