
Gaussian-binary Restricted Boltzmann Machines on Modeling Natural Image Statistics

Nan Wang
Institut für Neuroinformatik
Ruhr-Universität Bochum
Bochum, 44780, Germany
nan.wang@ini.rub.de

Jan Melchior
Institut für Neuroinformatik
Ruhr-Universität Bochum
Bochum, 44780, Germany
jan.melchior@ini.rub.de

Laurenz Wiskott
Institut für Neuroinformatik
Ruhr-Universität Bochum
Bochum, 44780, Germany
laurenz.wiskott@ini.rub.de

Abstract

We present a theoretical analysis of Gaussian-binary restricted Boltzmann machines (GRBMs) from the perspective of density models. The key aspect of this analysis is to show that GRBMs can be formulated as a constrained mixture of Gaussians, which gives a much better insight into the model's capabilities and limitations. We show that GRBMs are capable of learning meaningful features both in a two-dimensional blind source separation task and in modeling natural images. Further, we show that reported difficulties in training GRBMs are due to the failure of the training algorithm rather than the model itself. Based on our analysis we are able to propose several training recipes, which allowed successful and fast training in our experiments. Finally, we discuss the relationship of GRBMs to several modifications that have been proposed to improve the model.

1 Introduction

Inspired by the hierarchical structure of the visual cortex, recent studies on probabilistic models used deep hierarchical architectures to learn high order statistics of the data [Karklin and Lewicki(2009), Köster and Hyvärinen(2010)]. One widely used architecture is a deep believe network (DBN), which is usually trained as stacked restricted Boltzmann machines (RBMs) [Hinton and Salakhutdinov(2006), Bengio et al.(2006), Erhan et al.(2010)]. Since the original formulation of RBMs assumes binary input values, the model needs to be modified in order to handle continuous input values. One common way is to replace the binary input units by linear units with independent Gaussian noise, which is known as Gaussian-binary restricted Boltzmann machines (GRBMs) or Gaussian-Bernoulli restricted Boltzmann machines [Krizhevsky(2009), Cho et al.(2011)] first proposed by [Welling et al.(2004)].

The training of GRBMs is known to be difficult, so that several modifications have been proposed to improve the training. [Lee et al.(2007)] used a sparse penalty during training, which allowed them to learn meaningful features from natural image patches. [Krizhevsky(2009)] trained GRBMs on natural images and concluded that the difficulties are mainly due to the existence of high-frequency noise in the images, which further prevents the model from learning the important structures. [Theis et al.(2011)] illustrated that in terms of likelihood estimation GRBMs are already outperformed by simple mixture models. Other researchers focused on improving the model in the

view of generative models [Ranzato et al.(2010), Ranzato and Hinton(2010), Courville et al.(2011), Le Roux et al.(2011)Le Roux, Heess, Shotton, and Winn]. [Cho et al.(2011)] suggested that the failure of GRBMs is due to the training algorithm and proposed some modifications to overcome the difficulties encountered in training GRBMs.

The studies above have shown the failures of GRBMs empirically, but to our knowledge there is no analysis of GRBMs apart from our preliminary work [Wang et al.(2012)], which accounts the reasons behind these failures. In this paper, we extend our work in which we consider GRBMs from the perspective of density models, i.e. how well the model learns the distribution of the data. We show that a GRBM can be regarded as a mixture of Gaussians, which has already been mentioned briefly in previous studies [Bengio(2009), Theis et al.(2011), Courville et al.(2011)] but has gone unheeded. This formulation makes clear that GRBMs are quite limited in the way they can represent data. However we argue that this fact does not necessarily prevent the model from learning the statistical structure in the data. We present successful training of GRBMs both on a two-dimensional blind source separation problem and natural image patches, and that the results are comparable to that of independent component analysis (ICA). Based on our analysis we propose several training recipes, which allowed successful and fast training in our experiments. Finally, we discuss the relationship between GRBMs and above mentioned modifications of the model.

2 Gaussian-binary restricted Boltzmann machines (GRBMs)

2.1 The model

A Boltzmann Machine (BM) is a Markov Random Field with stochastic *visible* and *hidden* units [Smolensky(1986)], which are denoted as $\mathbf{X} := (X_1, \dots, X_M)^T$ and $\mathbf{H} := (H_1, \dots, H_N)^T$, respectively. In general, we use bold letters denote vectors and matrices.

The joint probability distribution is defined as

$$P(\mathbf{X}, \mathbf{H}) := \frac{1}{Z} e^{-\frac{1}{T_0} E(\mathbf{X}, \mathbf{H})}, \quad (1)$$

$$Z := \int \int e^{-\frac{1}{T_0} E(\mathbf{x}, \mathbf{h})} d\mathbf{x} d\mathbf{h} \quad (2)$$

where $E(\mathbf{X}, \mathbf{H})$ denotes an *energy function* as known from statistical physics, which defines the dependence between \mathbf{X} and \mathbf{H} . The temperature parameter T_0 is usually ignored by setting its value to one, but it can play an important role in inference of BMs [Desjardins et al.(2010)]. The *partition function* Z normalizes the probability distribution by integrating over all possible values of \mathbf{X} and \mathbf{H} , which is intractable in most cases. So that in training BMs using gradient descent the partition function is usually estimated using sampling methods. However, even sampling in BMs remains difficult due to the dependencies between all variables.

An RBM is a special case of a BM where the energy function contains no terms combining two different hidden or two different visible units. Viewed as a graphical model, there are no lateral connections within the visible or hidden layer, which results in a bipartite graph. This implies that the hidden units are conditionally independent given the visibles and vice versa, which allows efficient sampling.

The values of the visible and hidden units are usually assumed to be binary, i.e. $X_m, H_n \in \{0, 1\}$. The most common way to extend an RBM to continuous data is a GRBM, which assumes continuous values for the visible units and binary values for the hidden units. Its energy function [Cho et al.(2011), Wang et al.(2012)] is defined as

$$E(\mathbf{X}, \mathbf{H}) := \sum_i^M \frac{(X_i - b_i)^2}{2\sigma^2} - \sum_j^N c_j H_j - \sum_{i,j}^{M,N} \frac{X_i w_{ij} H_j}{\sigma^2} \quad (3)$$

$$= \frac{\|\mathbf{X} - \mathbf{b}\|^2}{2\sigma^2} - \mathbf{c}^T \mathbf{H} - \frac{\mathbf{X}^T \mathbf{W} \mathbf{H}}{\sigma^2}, \quad (4)$$

where $\|\mathbf{u}\|$ denotes the Euclidean norm of \mathbf{u} . In GRBMs the visible units given the hidden values are Gaussian distributed with standard deviation σ . Notice that some authors

[Krizhevsky(2009), Cho et al.(2011), Melchior(2012)] use an independent standard deviation for each visible unit, which comes into account if the data is not whitened [Melchior(2012)].

The conditional probability distribution of the visible given the hidden units is given by

$$P(\mathbf{X}|\mathbf{h}) = \frac{P(\mathbf{X}, \mathbf{h})}{\int P(\mathbf{x}, \mathbf{h}) d\mathbf{x}} \quad (5)$$

$$\stackrel{(1,4)}{=} \frac{e^{\mathbf{c}^T \mathbf{h}} \prod_i^M e^{\frac{x_i \mathbf{w}_{i*}^T \mathbf{h}}{\sigma^2} - \frac{\|x_i - b_i\|^2}{2\sigma^2}}}{\int e^{\mathbf{c}^T \mathbf{h}} \prod_i^M e^{\frac{x_i \mathbf{w}_{i*}^T \mathbf{h}}{\sigma^2} - \frac{\|x_i - b_i\|^2}{2\sigma^2}} d\mathbf{x}} \quad (6)$$

$$= \prod_i^M \frac{e^{\frac{x_i \mathbf{w}_{i*}^T \mathbf{h}}{\sigma^2} - \frac{\|x_i - b_i\|^2}{2\sigma^2}}}{\int e^{\frac{x_i \mathbf{w}_{i*}^T \mathbf{h}}{\sigma^2} - \frac{\|x_i - b_i\|^2}{2\sigma^2}} dx_i} \quad (7)$$

$$\stackrel{(11)}{=} \prod_i^M \frac{e^{-\frac{\|x_i - b_i - \mathbf{w}_{i*}^T \mathbf{h}\|^2}{2\sigma^2}}}{\int e^{-\frac{\|x_i - b_i - \mathbf{w}_{i*}^T \mathbf{h}\|^2}{2\sigma^2}} dx_i} \quad (8)$$

$$= \prod_i^M \underbrace{\mathcal{N}(x_i; b_i + \mathbf{w}_{i*}^T \mathbf{h}, \sigma^2)}_{= P(x_i|\mathbf{h})} \quad (9)$$

$$= \mathcal{N}(\mathbf{X}; \mathbf{b} + \mathbf{W}\mathbf{h}, \sigma^2), \quad (10)$$

where \mathbf{w}_{i*} and \mathbf{w}_{*j} denote the i th row and the j th column of the weight matrix, respectively. $\mathcal{N}(x; \mu, \sigma^2)$ denotes a Gaussian distribution with mean μ and variance σ^2 . And $\mathcal{N}(\mathbf{X}; \boldsymbol{\mu}, \sigma^2)$ denotes an isotropic multivariate Gaussian distribution centered at vector $\boldsymbol{\mu}$ with variance σ^2 in all directions. From (7) to (8) we used the relation

$$\begin{aligned} \frac{ax}{\sigma^2} - \frac{(x-b)^2}{2\sigma^2} &= \frac{-x^2 + 2bx + 2ax - b^2}{2\sigma^2} \\ &= \frac{-x^2 + 2bx + 2ax - b^2 + a^2 - a^2 + 2ab - 2ab}{2\sigma^2} \\ &= \frac{-(x-a-b)^2 + a^2 + 2ab}{2\sigma^2}. \end{aligned} \quad (11)$$

The conditional probability distribution of the hidden units given the visibles can be derived as follows

$$P(\mathbf{H}|\mathbf{x}) = \frac{P(\mathbf{x}, \mathbf{H})}{\sum_{\mathbf{h}} P(\mathbf{x}, \mathbf{h})} \quad (12)$$

$$\stackrel{(1,4)}{=} \frac{e^{-\frac{\|\mathbf{x}-\mathbf{b}\|^2}{2\sigma^2}} \prod_j^N e^{\left(c_j + \frac{\mathbf{x}^T \mathbf{w}_{*j}}{\sigma^2}\right) H_j}}{\sum_{\mathbf{h}} e^{-\frac{\|\mathbf{x}-\mathbf{b}\|^2}{2\sigma^2}} \prod_j^N e^{\left(c_j + \frac{\mathbf{x}^T \mathbf{w}_{*j}}{\sigma^2}\right) h_j}} \quad (13)$$

$$= \prod_j^N \frac{e^{\left(c_j + \frac{\mathbf{x}^T \mathbf{w}_{*j}}{\sigma^2}\right) H_j}}{\underbrace{\sum_{h_j} e^{\left(c_j + \frac{\mathbf{x}^T \mathbf{w}_{*j}}{\sigma^2}\right) h_j}}_{= P(H_j|\mathbf{x})}}. \quad (14)$$

$$\Rightarrow P(H_j = 1|\mathbf{x}) = \frac{1}{1 + e^{-\left(c_j + \frac{\mathbf{x}^T \mathbf{w}_{*j}}{\sigma^2}\right)}} \quad (15)$$

$P(\mathbf{H}|\mathbf{x})$ turns out to be a product of independent sigmoid functions, which is a frequently used non-linear activation function in artificial neural networks.

2.2 Maximum likelihood estimation

Maximum likelihood estimation (MLE) is a frequently used technique for training probabilistic models like BMs. In MLE we have a data set $\tilde{\mathcal{X}} = \{\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_L\}$ where the observations $\tilde{\mathbf{x}}_l$ are assumed to be independent and identically distributed (i.i.d.). The goal is to find the optimal parameters $\tilde{\Theta}$ that maximize the likelihood of the data, i.e. maximize the probability that the data is generated by the model [Bishop(2006)]. For practical reasons one often considers the logarithm of the likelihood, which has the same maximum as the likelihood since it is a monotonic function. The log-likelihood is defined as

$$\ln P(\tilde{\mathcal{X}}; \Theta) = \ln \prod_{l=1}^L P(\tilde{\mathbf{x}}_l; \Theta) = \sum_{l=1}^L \ln P(\tilde{\mathbf{x}}_l; \Theta). \quad (16)$$

We use the average log-likelihood per training case denoted by $\hat{\ell}$. For RBMs it is defined as

$$\hat{\ell} := \left\langle \ln P(\tilde{\mathcal{X}}; \Theta) \right\rangle_{\tilde{\mathbf{x}}} = \left\langle \ln \left(\sum_{\mathbf{h}} e^{-E(\tilde{\mathbf{x}}, \mathbf{h})} \right) \right\rangle_{\tilde{\mathbf{x}}} - \ln Z, \quad (17)$$

where $\tilde{\mathbf{x}} \in \tilde{\mathcal{X}}$. And $\langle f(u) \rangle_u$ denotes the expectation of the function $f(u)$ with respect to variable u .

The gradient of the $\hat{\ell}$ turns out to be the difference between the expectations of the energies gradient under the data and model distribution, which is given by

$$\begin{aligned} \frac{\partial \hat{\ell}}{\partial \theta} &\stackrel{(17,2)}{=} \left\langle \sum_{\mathbf{h}} \frac{e^{-E(\tilde{\mathbf{x}}, \mathbf{h})}}{Z} \left(-\frac{\partial E(\tilde{\mathbf{x}}, \mathbf{h})}{\partial \theta} \right) \right\rangle_{\tilde{\mathbf{x}}} - \frac{1}{Z} \sum_{\mathbf{h}} \sum_{\mathbf{x}} e^{-E(\mathbf{x}, \mathbf{h})} \left(-\frac{\partial E(\mathbf{x}, \mathbf{h})}{\partial \theta} \right) \\ &\stackrel{(1)}{=} - \left\langle \sum_{\mathbf{h}} P(\mathbf{h}|\tilde{\mathbf{x}}) \frac{\partial E(\tilde{\mathbf{x}}, \mathbf{h})}{\partial \theta} \right\rangle_{\tilde{\mathbf{x}}} + \left\langle \sum_{\mathbf{h}} P(\mathbf{h}|\mathbf{x}) \frac{\partial E(\mathbf{x}, \mathbf{h})}{\partial \theta} \right\rangle_{\mathbf{x}}. \end{aligned} \quad (18)$$

In practice, a finite set of i.i.d. samples can be used to approximate the expectations in (19). While we can use the training data to estimate the first term, we do not have any i.i.d. samples from the unknown model distribution to estimate the second term. Since we are able to compute the conditional probabilities in RBMs efficiently, Gibbs sampling can be used to generate those samples. But Gibbs-sampling only guarantees to generate samples from the model distribution if we run it infinite long. As this is impossible, a finite number of k sampling steps are used instead. This procedure is known as Contrastive Divergence - k (CD- k) algorithm, in which even $k = 1$ shows good results [Hinton(2002)]. The CD-gradient approximation is given by

$$\frac{\partial \hat{\ell}}{\partial \theta} \approx - \left\langle \sum_{\mathbf{h}} P(\mathbf{h}|\tilde{\mathbf{x}}) \frac{\partial E(\tilde{\mathbf{x}}, \mathbf{h})}{\partial \theta} \right\rangle_{\tilde{\mathbf{x}}} + \left\langle \sum_{\mathbf{h}} P(\mathbf{h}|\mathbf{x}^{(k)}) \frac{\partial E(\mathbf{x}^{(k)}, \mathbf{h})}{\partial \theta} \right\rangle_{\mathbf{x}^{(k)}}, \quad (20)$$

where $\mathbf{x}^{(k)}$ denotes the samples after k steps of Gibbs sampling. The derivatives of the GRBM's energy function with respect to the parameters are given by

$$\frac{\partial E(\mathbf{X}, \mathbf{H})}{\partial \mathbf{b}} = -\frac{\mathbf{X} - \mathbf{b}}{\sigma^2}, \quad (21)$$

$$\frac{\partial E(\mathbf{X}, \mathbf{H})}{\partial \mathbf{c}} = -\mathbf{H}, \quad (22)$$

$$\frac{\partial E(\mathbf{X}, \mathbf{H})}{\partial \mathbf{W}} = -\frac{\mathbf{X}\mathbf{H}^T}{\sigma^2}, \quad (23)$$

$$\frac{\partial E(\mathbf{X}, \mathbf{H})}{\partial \sigma} = -\frac{\|\mathbf{X} - \mathbf{b}\|^2}{\sigma^3} + \frac{2\mathbf{X}^T\mathbf{W}\mathbf{H}}{\sigma^3}, \quad (24)$$

and the corresponding gradient approximations (20) become

$$\frac{\partial \hat{\ell}}{\partial \mathbf{b}} \approx \left\langle \frac{\tilde{\mathbf{x}} - \mathbf{b}}{\sigma^2} \right\rangle_{\tilde{\mathbf{x}}} - \left\langle \frac{\mathbf{x}^{(k)} - \mathbf{b}}{\sigma^2} \right\rangle_{\mathbf{x}^{(k)}}, \quad (25)$$

$$\frac{\partial \hat{\ell}}{\partial \mathbf{c}} \approx \langle P(\mathbf{h} = \mathbf{1} | \tilde{\mathbf{x}}) \rangle_{\tilde{\mathbf{x}}} - \langle P(\mathbf{h} = \mathbf{1} | \mathbf{x}^{(k)}) \rangle_{\mathbf{x}^{(k)}}, \quad (26)$$

$$\frac{\partial \hat{\ell}}{\partial \mathbf{w}} \approx \left\langle \frac{\tilde{\mathbf{x}} P(\mathbf{h} = \mathbf{1} | \tilde{\mathbf{x}})^T}{\sigma^2} \right\rangle_{\tilde{\mathbf{x}}} - \left\langle \frac{\mathbf{x}^{(k)} P(\mathbf{h} = \mathbf{1} | \mathbf{x}^{(k)})^T}{\sigma^2} \right\rangle_{\mathbf{x}^{(k)}}, \quad (27)$$

$$\begin{aligned} \frac{\partial \hat{\ell}}{\partial \sigma} \approx & \left\langle \frac{\|\tilde{\mathbf{x}} - \mathbf{b}\|^2 - 2\tilde{\mathbf{x}}^T \mathbf{W} P(\mathbf{h} = \mathbf{1} | \tilde{\mathbf{x}})}{\sigma^3} \right\rangle_{\tilde{\mathbf{x}}} \\ & - \left\langle \frac{\|\mathbf{x}^{(k)} - \mathbf{b}\|^2 - 2\mathbf{x}^{(k)T} \mathbf{W} P(\mathbf{h} = \mathbf{1} | \mathbf{x}^{(k)})}{\sigma^3} \right\rangle_{\mathbf{x}^{(k)}}, \end{aligned} \quad (28)$$

where $P(\mathbf{h} = \mathbf{1} | \mathbf{x}) := (P(h_1 = 1 | \mathbf{x}), \dots, P(h_N = 1 | \mathbf{x}))^T$, i.e. $P(\mathbf{h} = \mathbf{1} | \mathbf{x})$ denotes a vector of probabilities.

2.3 The marginal probability distribution of the visible units

From the perspective of density estimation, the performance of the model can be assessed by examining how well the model estimates the data distribution. We therefore take a look at the model's marginal probability distribution of the visible units, which can be formalized as a product of experts (PoE) or as a mixture of Gaussians (MoG)¹.

2.3.1 In the Form of Product of Experts

We derive the marginal probability distribution of the visible units $P(\mathbf{X})$ by factorizing the joint probability distribution over the hidden units.

$$P(\mathbf{X}) = \sum_{\mathbf{h}} P(\mathbf{X}, \mathbf{h}) \quad (29)$$

$$\stackrel{(1,4)}{=} \frac{1}{Z} e^{-\frac{\|\mathbf{x} - \mathbf{b}\|^2}{2\sigma^2}} \prod_j \sum_{h_j} e^{c_j + \frac{\mathbf{x}^T \mathbf{w}_{*j}}{\sigma^2} h_j} \quad (30)$$

$$h_j \in \{0,1\} \quad \frac{1}{Z} \prod_j \left(e^{-\frac{\|\mathbf{x} - \mathbf{b}\|^2}{2N\sigma^2}} + e^{c_j + \frac{\mathbf{x}^T \mathbf{w}_{*j}}{\sigma^2} - \frac{\|\mathbf{x} - \mathbf{b}\|^2}{2N\sigma^2}} \right) \quad (31)$$

$$\stackrel{(11)}{=} \frac{1}{Z} \prod_j \left(e^{-\frac{\|\mathbf{x} - \mathbf{b}\|^2}{2N\sigma^2}} + e^{\frac{\|\mathbf{b} + N\mathbf{w}_{*j}\|^2 - \|\mathbf{b}\|^2}{2N\sigma^2} + c_j - \frac{\|\mathbf{x} - \mathbf{b} - N\mathbf{w}_{*j}\|^2}{2N\sigma^2}} \right) \quad (32)$$

$$\begin{aligned} &= \frac{1}{Z} \prod_j \left(\sqrt{2\pi N\sigma^2} \right)^M \left[\mathcal{N}(\mathbf{X}; \mathbf{b}, N\sigma^2) \right. \\ & \quad \left. + e^{\frac{\|\mathbf{b} + N\mathbf{w}_{*j}\|^2 - \|\mathbf{b}\|^2}{2N\sigma^2} + c_j} \mathcal{N}(\mathbf{X}; \mathbf{b} + N\mathbf{w}_{*j}, N\sigma^2) \right] \end{aligned} \quad (33)$$

$$=: \frac{1}{Z} \prod_j p_j(\mathbf{X}). \quad (34)$$

Equation (34) illustrates that $P(\mathbf{X})$ can be written as a product of N factors, referred to as a product of experts [Hinton(2002)]. Each expert $p_j(\mathbf{X})$ consists of two isotropic Gaussians with the same variance $N\sigma^2$. The first Gaussian is placed at the visible bias \mathbf{b} . The second Gaussian is shifted

¹Some part of this analysis has been previously reported by [Freund & Haussler(1992)]. Thanks to the anonymous reviewer for pointing out this coincidence.

relative to the first one by N times the weight vector \mathbf{w}_{*j} and scaled by a factor that depends on \mathbf{w}_{*j} and \mathbf{b} . Every hidden unit leads to one expert, each mode of which corresponds to one state of the corresponding hidden unit. Figure 1 (a) and (b) illustrate $P(\mathbf{X})$ of a GRBM-2-2 viewed as a PoE, where GRBM- M - N denotes a GRBM with M visible and N hidden units.

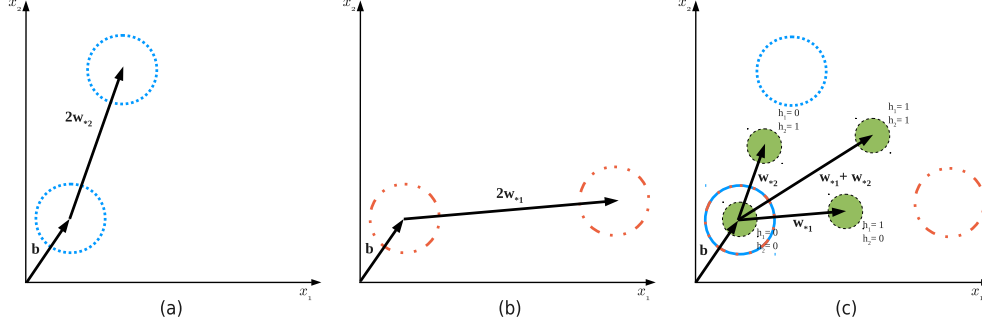


Figure 1: Illustration of a GRBM-2-2 as a PoE and MoG, in which arrows indicate the roles of the visible bias vector and the weight vectors. (a) and (b) visualize the two experts of the GRBM. The red (dotted) and blue (dashed) circles indicate the center of two Gaussians in each expert. (c) visualizes the components in the GRBM. Denoted by the green (filled) circles, the four components are the results of the product of the two experts. Notice how each component sits right between a red (dotted) and a blue (dashed) circle.

2.3.2 In the Form of Mixture of Gaussians

Using Bayes' theorem, the marginal probability of \mathbf{X} can also be formalized as:

$$P(\mathbf{X}) = \sum_{\mathbf{h}} P(\mathbf{X}|\mathbf{h}) P(\mathbf{h}) \quad (35)$$

$$= \sum_{\mathbf{h}} \mathcal{N}(\mathbf{X}; \mathbf{b} + \mathbf{W}\mathbf{h}, \sigma^2) \frac{(\sqrt{2\pi\sigma^2})^M}{Z} e^{\mathbf{c}^T \mathbf{h} + \frac{\|\mathbf{b} + \mathbf{W}\mathbf{h}\|^2 - \|\mathbf{b}\|^2}{2\sigma^2}} \quad (36)$$

$$= \underbrace{\frac{(\sqrt{2\pi\sigma^2})^M}{Z}}_{P(\mathbf{h}; \mathbf{h} \in \mathcal{H}_0)} \mathcal{N}(\mathbf{X}; \mathbf{b}, \sigma^2) + \sum_{j=1}^N \underbrace{\frac{(\sqrt{2\pi\sigma^2})^M}{Z} e^{\frac{\|\mathbf{b} + \mathbf{w}_{*j}\|^2 - \|\mathbf{b}\|^2}{2\sigma^2} + c_j}}_{P(\mathbf{h}_j; \mathbf{h}_j \in \mathcal{H}_1)} \mathcal{N}(\mathbf{X}; \mathbf{b} + \mathbf{w}_{*j}, \sigma^2) + \sum_{j=1}^{N-1} \sum_{k>j}^N \underbrace{\frac{(\sqrt{2\pi\sigma^2})^M}{Z} e^{\frac{\|\mathbf{b} + \mathbf{w}_{*j} + \mathbf{w}_{*k}\|^2 - \|\mathbf{b}\|^2}{2\sigma^2} + c_j + c_k}}_{P(\mathbf{h}_{jk}; \mathbf{h}_{jk} \in \mathcal{H}_2)} \mathcal{N}(\mathbf{X}; \mathbf{b} + \mathbf{w}_{*j} + \mathbf{w}_{*k}, \sigma^2) + \dots, \quad (37)$$

where \mathcal{H}_k denotes the set of all possible binary vectors with exactly k ones and $M - k$ zeros respectively. As an example, $\sum_{j=1}^{N-1} \sum_{k>j}^N P(\mathbf{h}_{jk} : \mathbf{h}_{jk} \in \mathcal{H}_2) = \sum_{\mathbf{h} \in \mathcal{H}_2} P(\mathbf{h})$ sums over the probabilities of all binary vectors having exactly two entries set to one. $P(\mathbf{H})$ in (36) is derived as

follows

$$P(\mathbf{H}) = \int P(\mathbf{x}, \mathbf{H}) d\mathbf{x} \quad (38)$$

$$\stackrel{(1,4)}{=} \frac{1}{Z} \int e^{\mathbf{c}^T \mathbf{H}} \prod_i^M e^{\frac{x_i \mathbf{w}_{i*}^T \mathbf{H}}{\sigma^2} - \frac{\|x_i - b_i\|^2}{2\sigma^2}} d\mathbf{x} \quad (39)$$

$$= \frac{e^{\mathbf{c}^T \mathbf{H}}}{Z} \prod_i^M \int e^{\frac{x_i \mathbf{w}_{i*}^T \mathbf{H}}{\sigma^2} - \frac{\|x_i - b_i\|^2}{2\sigma^2}} dx_i \quad (40)$$

$$\stackrel{(11)}{=} \frac{e^{\mathbf{c}^T \mathbf{H}}}{Z} \prod_i^M \left(e^{\frac{(b_i + \mathbf{w}_{i*}^T \mathbf{H})^2 - b_i^2}{2\sigma^2}} \int e^{\frac{\|x_i - b_i - \mathbf{w}_{i*}^T \mathbf{H}\|^2}{2\sigma^2}} dx_i \right) \quad (41)$$

$$= \frac{e^{\mathbf{c}^T \mathbf{H}}}{Z} \left(\sqrt{2\pi\sigma^2} \right)^M e^{\sum_i^M \frac{(b_i + \mathbf{w}_{i*}^T \mathbf{H})^2 - b_i^2}{2\sigma^2}} \quad (42)$$

$$= \frac{\left(\sqrt{2\pi\sigma^2} \right)^M}{Z} e^{\mathbf{c}^T \mathbf{H} + \frac{\|\mathbf{b} + \mathbf{W}\mathbf{H}\|^2 - \|\mathbf{b}\|^2}{2\sigma^2}} \quad (43)$$

Since the form in (37) is similar to a mixture of isotropic Gaussians, we follow its naming convention. Each Gaussian distribution is called a *component* of the model distribution, which is exactly the conditional probability of the visible units given a particular state of the hidden units. As well as in MoGs, each component has a *mixing coefficient*, which is the marginal probability of the corresponding state and can also be viewed as the prior probability of picking the corresponding component. The total number of components in a GRBM is 2^N , which is exponential in the number of hidden units, see Figure 1 (c) for an example.

The locations of the components in a GRBM are not independent of each other as it is the case in MoGs. They are centered at $\mathbf{b} + \mathbf{W}\mathbf{h}$, which is the vector sum of the visible bias and selected weight vectors. The selection is done by the corresponding entries in \mathbf{h} taking the value one. This implies that only the $M + 1$ components that sum over exactly one or zero weights can be placed and scaled independently. We name them first order components and the anchor component respectively. All $2^N - M - 1$ higher order components are then determined by the choice of the anchor and first order components. This indicates that GRBMs are constrained MoGs with isotropic components.

3 Experiments

3.1 Two-dimensional blind source separation

The general presumption in the analysis of natural images is that they can be considered as a mixture of independent super-Gaussian sources [Bell and Sejnowski(1997)], but see [Zetsche and Rohrbein(2001)] for an analysis of remaining dependencies. In order to be able to visualize how GRBMs model natural image statistics, we use a mixture of two independent Laplacian distributions as a toy example.

The independent sources $\mathbf{s} = (s_1, s_2)^T$ are mixed by a random mixing matrix \mathbf{A} yielding

$$\tilde{\mathbf{x}}' = \mathbf{A}\mathbf{s}, \quad (44)$$

where $p(s_i) = \frac{e^{-\sqrt{2}|s_i|}}{\sqrt{2}}$. It is common to whiten the data (see Section 4.1), resulting in

$$\tilde{\mathbf{x}} = \mathbf{V}\tilde{\mathbf{x}}' = \mathbf{V}\mathbf{A}\mathbf{s}, \quad (45)$$

where $\mathbf{V} = \left\langle \tilde{\mathbf{x}}' \tilde{\mathbf{x}}'^T \right\rangle^{-\frac{1}{2}}$ is the whitening matrix calculated with principle component analysis (PCA). Through all this paper, we used the whitened data.

In order to assess the performance of GRBMs in modeling the data distribution, we ran the experiments for 200 times and calculated the $\hat{\ell}$ for test data analytically. For comparison, we also

calculated the $\hat{\ell}$ over the test data for ICA², an isotropic two-dimensional Gaussian distribution and the true data distribution³. The results are presented in Table 1, which confirm the conclusion of [Theis et al.(2011)] that GRBMs are not as good as ICA in terms of $\hat{\ell}$.

Table 1: Comparison of $\hat{\ell}$ between different models

	$\hat{\ell} \pm \text{std}$
Gaussian	-2.8367 ± 0.0086
GRBM	-2.8072 ± 0.0088
ICA	-2.7382 ± 0.0091
data distribution	-2.6923 ± 0.0092

To illustrate how GRBMs model the statistical structure of the data, we looked at the probability distributions of the 200 trained GRBMs. About half of them (110 out of 200) recovered the independent components, see Figure 2 (a) as an example. This can be further illustrated by plotting the Amari errors⁴ between the true unmixing matrix \mathbf{A}^{-1} and estimated model matrices, i.e. the unmixing matrix of ICA and the weight matrix of the GRBM, as shown in Figure 3. One can see that these 110 GRBMs estimated the unmixing matrix quite well, although GRBMs are not as good as ICA. This is due to the fact that the weight vectors in GRBMs are not restricted to be orthogonal as in ICA.

For the remaining 90 GRBMs, the two weight vectors pointed to the opposite direction as shown in Figure 2 (b). Accordingly, these GRBMs failed to estimate the unmixing matrix, but in terms of density estimation these solutions have the same quality as the orthogonal ones. Thus all the 200 GRBMs were able to learn the statistical structures in the data and model the data distribution pretty well.

For comparison, we plotted the probability distribution of a learned GRBM with four hidden units, see Figure 2 (c), in which GRBMs can always find the two independent components correctly.

To further show how the components contribute to the model distribution, we randomly chose one of the 110 GRBMs and calculated the mixing coefficients of the anchor and the first order components, as shown in Table 2. The large mixing coefficient for the anchor component indicates that the model will most likely reach hidden states in which none of the hidden units are activated. In general, the more activated hidden units a state has, the less likely it will be reached, which leads naturally to a sparse representation of the data.

Table 2: The mixing coefficients of a successfully-trained GRBM-2-2, GRBM-2-4 and an MoG-3.

	$\sum_{\mathbf{h} \in \mathcal{H}_0} P(\mathbf{h})$	$\sum_{\mathbf{h} \in \mathcal{H}_1} P(\mathbf{h})$	$\sum_{\mathbf{h} \in \mathcal{H}_2} P(\mathbf{h})$	$\sum_{\mathbf{h} \in \mathcal{H}_3} P(\mathbf{h})$	$\sum_{\mathbf{h} \in \mathcal{H}_4} P(\mathbf{h})$
GRBM-2-2	0.9811	0.0188	$7.8856e-05$	-	-
GRBM-2-4	0.9645	0.0352	$3.4366e-04$	$1.2403e-10$	$6.9977e-18$
MoG-3	0.9785	0.0215	-	-	-

²For the fast ICA algorithm [Hyvärinen(1999)] we used for training, the $\hat{\ell}$ for super Gaussian sources can also be assessed analytically by $\hat{\ell} = -\left\langle \sum_{j=1}^N \ln 2 \cosh^2 \mathbf{w}_{*j}^T \tilde{\mathbf{x}}_l \right\rangle_{\tilde{\mathbf{x}}_l} + \ln |\det \mathbf{W}|$.

³As we know the true data distribution, the exact $\hat{\ell}$ can be calculated by $\hat{\ell} = -\sqrt{2} \left\langle |\mathbf{u}_{1*} \tilde{\mathbf{x}}_l| + |\mathbf{u}_{2*} \tilde{\mathbf{x}}_l| \right\rangle_{\tilde{\mathbf{x}}_l} - \ln 2 + \ln |\det \mathbf{U}|$, where $\mathbf{U} = (\mathbf{V}\mathbf{A})^{-1}$.

⁴The Amari error [Bach and Jordan(2002)] between two matrices \mathbf{A} and \mathbf{B} is defined as: $\frac{1}{2N} \left(\sum_{i=1}^N \sum_{j=1}^N \frac{|(\mathbf{A}\mathbf{B}^{-1})_{ij}|}{\max_k |(\mathbf{A}\mathbf{B}^{-1})_{ik}|} + \frac{|(\mathbf{A}\mathbf{B}^{-1})_{ij}|}{\max_k |(\mathbf{A}\mathbf{B}^{-1})_{kj}|} \right) - 1$.

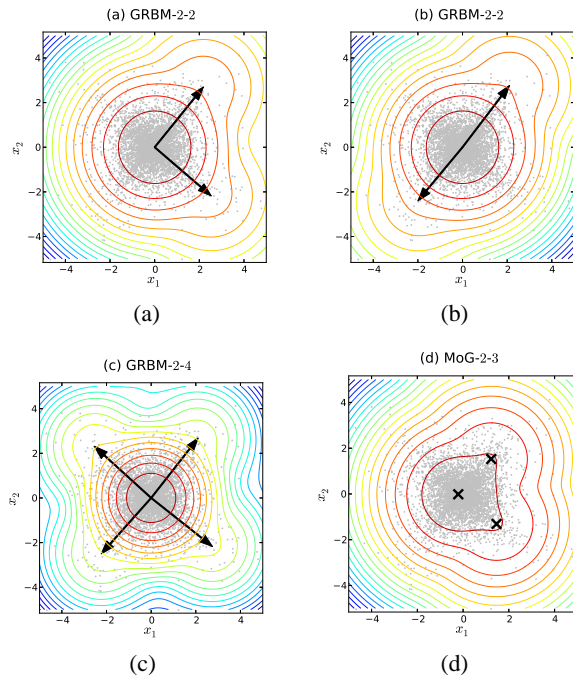


Figure 2: Illustration of the log-probability densities. The data is plotted in blue dots. (a)GRBM-2-2 which learned two independent components. (b)GRBM-2-2 which learned one independent component with opposite directions. (c)GRBM-2-4. (d)An isotropic MoG with three components. The arrows indicate the weight vectors of GRBM, while the crosses denote the means of the MoG components. Comparing (a) and (d), the contribution of the second order component is so insignificant that the probability distribution of the GRBM with four components is almost the same as the MoG with only three components.

The dominance of $\sum_{\mathbf{h} \in \mathcal{H}_0} P(\mathbf{h})$ and all $\sum_{\mathbf{h} \in \mathcal{H}_1} P(\mathbf{h})$ can also be seen in Figure 2 by comparing a GRBM-2-2 (a) with a two dimensional MoG having three isotropic components denoted by MoG-2-3 (d). Although the MoG-2-3 has one component fewer than the GRBM-2-2, their probability distributions are almost the same.

3.2 Natural image patches

In contrast to random images, natural images have a common underlying structure which could be used to code them more efficiently than with a pixel-wise representation. [Olshausen and Field(1996)] showed that sparse coding is such an efficient coding scheme and that it is in addition a biological plausible model for the simple cells in the primary visual cortex. [Bell and Sejnowski(1997)] showed that the independent components provide a comparable representation for natural images. We now want to test empirically whether GRBMs generate such biological plausible results like sparse coding and ICA.

We used the `imlog` natural image Database of [van Hateren and van der Schaaf(1998)] and randomly sampled 70000 grey scale image patches with a size of 14×14 pixels. The data was whitened using Zero-phase Component Analysis (ZCA), afterwards it was divided into 40000 training and 30000 testing image patches. We followed the training recipes mentioned in Section 4, since training a GRBM on natural image patches is not a trivial task.

In Figure 4, we show the learned weight vectors namely features or filters, which can be regarded as receptive fields of the hidden units. They are fairly similar to the filters learned by ICA [Bell and Sejnowski(1997)]. Similar to the 2D experiment, we calculated the anchor and first order mixing coefficients, as shown in Table 3. The coefficients are much smaller compared to the anchor and first order coefficients of the GRBMs in the two dimensional case. However, they are

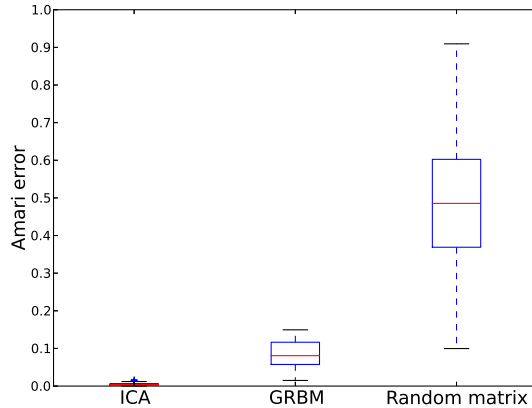


Figure 3: The Amari errors between the real unmixing matrix and the estimations from ICA and the 110 GRBMs. The box extends from the lower to the upper quantile values of the data, with a line at the median. The whiskers extend from the box to show the range of the reliable data points. The outlier points are marked by “+”. As a base line, the amari errors between the real unmixing matrices and random matrices are provided.

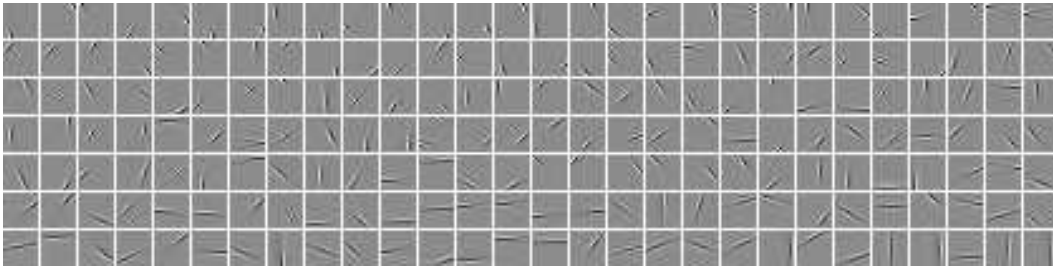


Figure 4: Illustration of 196 learned filters of a GRBM-196-196. The plot has been ordered from left to right and from top to bottom by the increasing average activation level of the corresponding hidden units.

still significantly large, considering that the total number of components in this case is 2^{196} . Similar to the two-dimensional experiments, the more activated hidden units a state has, the less likely it will be reached, which leads naturally to a sparse representation. To support this statement, we plotted the histogram of the number of activated hidden units per training sample, as shown in Figure 5.

Table 3: The mixing coefficients of GRBMs-196-196 per component (the Partition function was estimated using AIS).

	$\sum_{\mathbf{h} \in \mathcal{H}_0} P(\mathbf{h})$	$\sum_{\mathbf{h} \in \mathcal{H}_1} P(\mathbf{h})$	$\sum_{\mathbf{h} \in \mathcal{H} \setminus \{\mathcal{H}_0 \cup \mathcal{H}_1\}} P(\mathbf{h})$
GRBM-196-196	0.04565	0.00070	0.95365

We also examined the results of GRBMs in the over-complete case, i.e. GRBM-196-588. There is no prominent difference of the filters compared to the complete case shown in Figure 4. To further compare the filters in the complete and over-complete case, we estimated the spatial frequency, location and orientation for all filters in the spatial and frequency domains, see Figure 6 and Figure 7 respectively. This is achieved by fitting a Gabor function of the form used

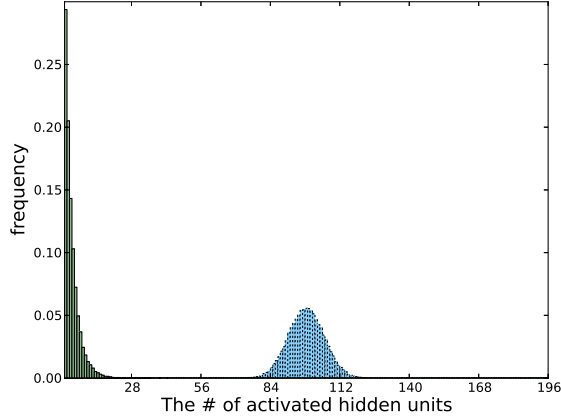


Figure 5: The histogram of the number of activated hidden units per training sample. The histograms before and after training are plotted in blue (dotted) and in green (solid), respectively.

by [Lewicki and Olshausen(1999)]. Note that the additional filters in the over-complete case increase the variety of spatial frequency, location and orientation.

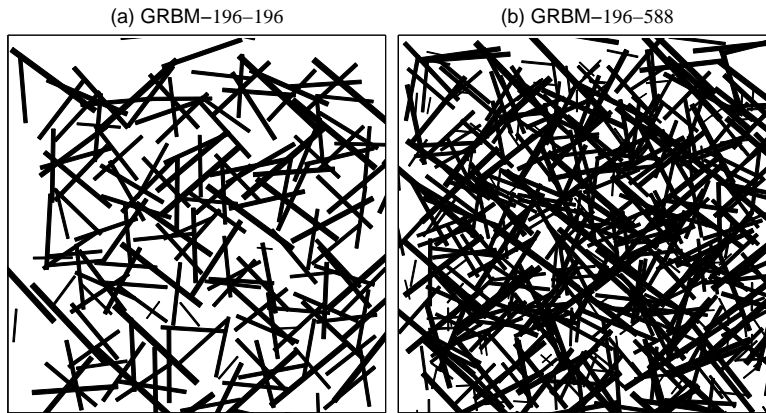


Figure 6: The spatial layout and size of the filters, which are described by the position and size of the bars. Each bar denotes the center position and the orientation of a fitted Gabor function within 14×14 grid. The thickness and length of each bar are proportional to its spatial-frequency bandwidth.

4 Successful Training of GRBMs on Natural Images

The training of GRBMs has been reported to be difficult [Krizhevsky(2009), Cho et al.(2011)]. Based on our analysis we are able to propose some recipes which should improve the success and speed of training GRBMs on natural image patches. Some of them do not depend on the data distribution and should therefore improve the training in general.

4.1 Preprocessing of the Data

The preprocessing of the data is important especially if the model is highly restricted like GRBMs. Whitening is a common preprocessing step for natural images. It removes the first and second order statistics from the data, so that it has zero mean and unit variance in all directions. This allows train-

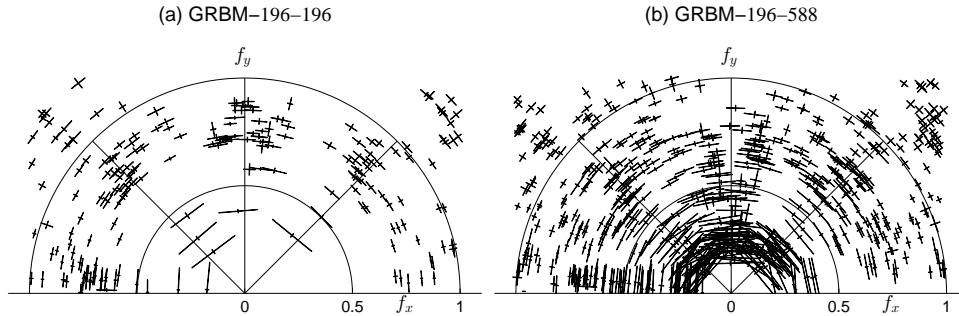


Figure 7: A polar plot of frequency tuning and orientation of the learned filters. The crosshairs describe the selectivity of the filters, which is given by the 1/16-bandwidth in spatial-frequency and orientation, [Lewicki and Olshausen(1999)].

ing algorithms to focus on higher order statistics like kurtosis, which is assumed to play an important role in natural image representations [Olshausen and Field(1996), Hyvärinen et al.(2001)].

The components of GRBMs are isotropic Gaussians, so that the model would use several components for modeling covariances. But the whitened data has a spherical covariance matrix so that the distribution can be modelled already fairly well by a single component. The other components can then be used to model higher order statistics, so that we claim that whitening is also an important preprocessing step for GRBMs.

4.2 Parameter Initialization

The initial choice of model parameters is important for optimization process. Using prior knowledge about the optimization problem can help to derive an initialization, which can improve the speed and success of the training.

For GRBMs we know from the analysis above that the anchor component, which is placed at the visible bias, represents most of the whitened data. Therefore it is reasonable in practice to set the visible bias to the data's mean.

Learning the right scaling is usually very slow since the weights and biases determine both the position and scaling of the components. In the final stage of training GRBMs on whitened natural images, the first components are scaled down extremely compared to the anchor component. Therefore, it will usually speed up the training process if we initialize the parameters so that the first order scaling factors are already very small. Considering equation (37), we are able to set a specific first order scaling factor by initializing the hidden bias to

$$c_j = -\frac{\|\mathbf{b} + \mathbf{w}_{*j}\|^2 - \|\mathbf{b}\|^2}{2\sigma^2} + \ln \tau_j, \quad (46)$$

so that the scaling is determined by τ_j , which should ideally be chosen close to the unknown final scaling factors. In practice, the choice of 0.01 showed good performance in most cases. The learning rate for the hidden bias can then be set much smaller than the learning rate for the weights.

According to [Bengio(2010)], the weights should be initialized to $w_{ij} \sim U\left(-\frac{\sqrt{6}}{\sqrt{N+M}}, \frac{\sqrt{6}}{\sqrt{N+M}}\right)$, where $U(a, b)$ is the uniform distribution in the interval $[a, b]$. In our experience, this works better than the commonly used initialization to small Gaussian-distributed random values.

4.3 Gradient Restriction and Choices of the Hyperparameters

The choice of the hyper-parameters has an significant impact on the speed and success of training GRBMs. For successful training in an acceptable number of updates, the learning rate needs to be sufficiently big. Otherwise the learning process becomes too slow or the algorithm converges to a local optimum where all components are placed in the data's mean. But if the learning rate is chosen too big, the gradient can easily diverge resulting in a number overflow of the weights. This effect

Method	Time per epoch in s
CD-1	2.1190
PCD-1	2.1348
CD-10	10.8052
PCD-10	10.8303
PT-10	21.4855

Table 4: Comparison of the CPU time for training a GRBM with different methods.

becomes even more crucial as the model dimensionality increases, so that a GRBM with 196 visible and 1000 hidden units diverges already for a learning rate of 0.001.

We therefore propose restricting the weight gradient column norms $\nabla w_{:j}$ to a meaningful size to prevent divergence. Since we know that the components are placed in the region of data, there is no need for a weight norm to be bigger than twice the maximal data norm. Consequently, this natural bound also holds for the gradient and can in practice be chosen even smaller. It allows to choose big learning rates even for very large models and therefore enables fast and stable training. In practice, one should restrict the norm of the update matrix rather than the gradient matrix to also restrict the effects of the momentum term and etc.

Since the components are placed on the data they are naturally restricted, which makes the use of a weight decay useless or even counter productive since we want the weights to grow up to a certain norm. Thus we do recommend not to use a weight decay regularization.

A momentum term adds a percentage of the old gradient to the current gradient which leads to a more robust behavior especially for small batch-sizes. In the early stage of training the gradient usually varies a lot, a large momentum can therefore be used to prevent the weights from converging to zero. In the late stage however, it can also prevent convergence so that in practice a momentum of 0.9 that will be reduced to zero in the final stage of training is recommended.

4.4 Training Method

Using the gradient approximation, RBMs are usually trained as described in Section 2.2. The quality of the approximation highly depends on the set of samples used for estimating the model expectation, which should ideally be i.i.d. But Gibbs sampling usually has a low mixing rate, which means that the samples tend to stay close to the previously presented samples. Therefore, a few steps of Gibbs sampling commonly leads to a biased approximation of the gradient. In order to increase the mixing rate [Tieleman(2008)] suggested to use a persistent Markov chain for drawing samples from the model distribution, which is referred as persistent Contrastive Divergence (PCD). [Desjardins et al.(2010)] proposed to use parallel tempering (PT), which selects samples from a persistent Markov chain with a different scaling of the energy function. In particular, [Cho et al.(2011)] analyzed PT algorithm for training GRBMs and proposed a modified version of PT.

In our experiments all methods above lead to meaningful features and comparable $\hat{\ell}$, but differ in convergence speed as shown in Figure 8. As for PT, we used original algorithm [Desjardins et al.(2010)] together with weight restrictions and temperatures from 0.1 to 1 with step-size 0.1. Although, PT has a better performance than CD, it has also a much higher computational cost as shown in Table 4.

5 Discussion

The difficulties of using GRBMs for modeling natural images have been reported by several authors [Krizhevsky(2009), Bengio et al.(2006)] and various modifications have been proposed to address this problem.

[Ranzato and Hinton(2010)] analyzed the problem from the view of generative models and argued that the failure of GRBMs is due to the model’s focus on predicting the mean intensity of each pixel rather than the dependence between pixels. To model the covariance matrices at the same time, they proposed the mean-covariance RBM (mCRBM). In addition to the conventional hidden units \mathbf{h}^m ,

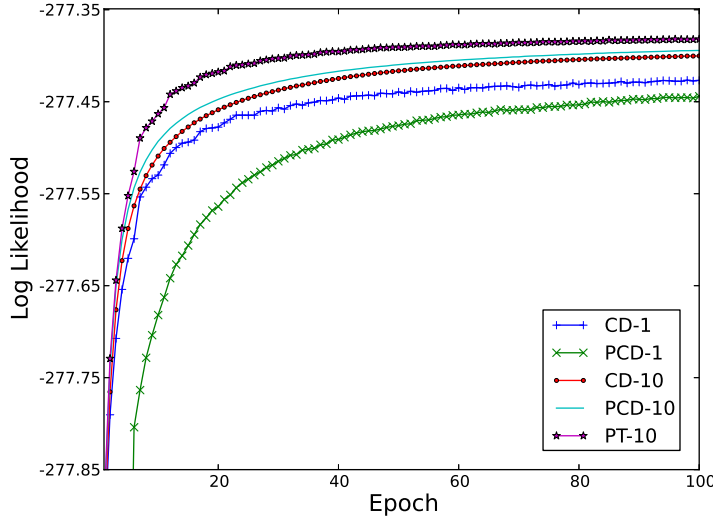


Figure 8: Evolution of the $\hat{\ell}$ of a GRBM 196-16 on the whitened natural image dataset for CD, PCD using a k of 1, 10 each and PT with 10 temperatures. The learning curves are the average over 40 trials. The learning rate was 0.1, an initial momentum term of 0.9 was multiplied with 0.9 after each fifth epoch, the gradient was restricted to one hundredth of the maximal data norm (0.48), no weight decay was used.

there is a group of hidden units \mathbf{h}^c dedicated to model the covariance between the visible units. From the view of density models, mcRBMs can be regarded as improved GRBMs such that the additional hidden units are used to depict the covariances. The conditional probabilities of mcRBM are given by

$$P(\mathbf{X}|\mathbf{h}^m, \mathbf{h}^c) = \mathcal{N}(\mathbf{X}; \mathbf{\Sigma} \mathbf{W} \mathbf{h}^m, \mathbf{\Sigma}), \quad (47)$$

where $\mathbf{\Sigma} = (\mathbf{C} \text{diag}(\mathbf{P}\mathbf{h}^c) \mathbf{C}^T)^{-1}$ [Ranzato and Hinton(2010)]. By comparing (47) and (9), it can be seen that the components of mcRBM can have a covariance matrix that is not restricted to be diagonal as it is the case for GRBMs.

From the view of generative models another explanation for the failure of GRBMs is provided by [Courville et al.(2011)]. Although they agree with the poor ability of GRBMs in modeling covariances, [Courville et al.(2011)] argue that the deficiency is due to the binary nature of the hidden units. In order to overcome this limitation, they developed the spike-and-slab RBM (ssRBM), which splits each binary hidden unit into a binary spike variable h_j and a real valued slab variable s_j . The conditional probability of visible units is given by

$$P(\mathbf{X}|\mathbf{s}, \mathbf{h}, \|\mathbf{X}\|^2 < R) = \frac{1}{B} \mathcal{N}\left(\mathbf{X}; \mathbf{\Lambda}^{-1} \sum_{j=1}^N \mathbf{w}_{*j} s_j h_j, \mathbf{\Lambda}^{-1}\right), \quad (48)$$

where $\mathbf{\Lambda}$ is a diagonal matrix and B is determined by integrating the Gaussian $\mathcal{N}(\mathbf{X}; \mathbf{\Lambda}^{-1} \sum_{j=1}^N \mathbf{w}_{*j} s_j h_j, \mathbf{\Lambda}^{-1})$ over the ball $\|\mathbf{X}\|^2 < R$ [Courville et al.(2011)]. In contrast to the conditional probability of GRBMs (9), \mathbf{w}_{*j} in (48) is scaled by the continuous variable s_j , which implies that the components can be shifted along their weight vectors.

We have shown that GRBMs are capable of modeling natural image patches and that the reported failures are due to the training procedure. [Lee et al.(2007)] showed also that GRBMs could learn meaningful filters by using a sparse penalty. But this penalty changes the objective function and introduced a new hyper-parameter.

[Cho et al.(2011)] addressed these training difficulties, by proposing a modification of PT and an adaptive learning rate. However, we claim that the reported difficulties of training GRBMs with PT

are due to the mentioned gradient divergence problem. With gradient restriction we were able to overcome the problem and train GRBMs with normal PT successfully.

6 Conclusion

In this paper, we provide a theoretical analysis of GRBM and showed that its product of experts formulation can be rewritten as a constrained mixture of Gaussians. This representation gives a much better insight into the capabilities and limitations of the model. We use two-dimensional blind source separation task as a toy problem to demonstrate how GRBMs model the data distribution. In our experiments, GRBMs were capable of learning meaningful features both in the toy problem and in modeling natural images.

In both cases, the results are comparable to that of ICA. But in contrast to ICA the features are not restricted to be orthogonal and can form an over-complete representation. However, the success of training GRBMs highly depends on the training setup, for which we proposed several recipes based on the theoretical analysis. Some of them can be further generalized to other datasets or directly applied like the gradient restriction.

In our experience, maximizing the $\hat{\ell}$ does not imply good features and vice versa. Prior knowledge about the data distribution will be beneficial in the modeling process. For instance, our recipes are based on the prior knowledge of the natural image statistics, which is center peaked and has heavy tails. It will be an interesting topic to integrate prior knowledge of the data distribution into the model rather than starting modeling from scratch.

Considering the simplicity and easiness of training with our proposed recipe, we believe that GRBMs provide a possible way for modeling natural images. Since GRBMs are usually used as first layer in deep belief networks, the successful training of GRBMs should therefore improve the performance of the whole network.

References

- [Bach and Jordan(2002)] F. R. Bach and M. I. Jordan. Kernel independent component analysis. *Journal of Machine Learning Research*, 3:1–48, 2002.
- [Bell and Sejnowski(1997)] A. Bell and T. Sejnowski. The “independent components” of natural scenes are edge filters. *Vision Research*, 37(23):3327–3338, 12 1997. ISSN 00426989.
- [Bengio(2010)] X. Glorot Bengio. Understanding the difficulty of training deep feedforward neural networks. *AISTATS*, 2010.
- [Bengio(2009)] Y. Bengio. Learning deep architectures for AI. *Foundation and Trends in Machine Learning*, 2(1):1–127, 2009. Also published as a book. Now Publishers, 2009.
- [Bengio et al.(2006)] Y. Bengio, P. Lamblin, D. Popovici, and H. Larochelle. Greedy layer-wise training of deep networks. In *Proceedings of the Conference on Neural Information Processing Systems*, pages 153–160, 2006.
- [Bishop(2006)] C. M. Bishop. *Pattern Recognition and Machine Learning*, chapter 8, pages 359–422. Springer, Secaucus, NJ, USA, 2006. ISBN 0387310738.
- [Cho et al.(2011)] K. Cho, A. Ilin, and T. Raiko. Improved learning of gaussian-bernoulli restricted boltzmann machines. In *Proceedings of the International Conference on Artificial Neural Networks*, 10–17, 2011.
- [Courville et al.(2011)] A. C. Courville, J. Bergstra, and Y. Bengio. A spike and slab restricted boltzmann machine. *Journal of Machine Learning Research*, 15:233–241, 2011.
- [Desjardins et al.(2010)] G. Desjardins, A. Courville, Y. Bengio, P. Vincent, and O. Delalleau. Parallel tempering for training of restricted boltzmann machines. In *Proceedings of the International conference on Artificial Intelligence and Statistics*, 2010.
- [Erhan et al.(2010)] D. Erhan, Y. Bengio, A. C. Courville, P. Manzagol, P. Vincent, and Y. Bengio. Why does unsupervised pre-training help deep learning? *Journal of Machine Learning Research*, 11:625–660, 2010.
- [Freund & Haussler(1992)] Y. Freund, and D. Haussler(1992) Unsupervised learning of distributions of binary vectors using two layer networks. *Proceedings of the Conference on Neural Information Processing Systems*, 912 – 919.
- [Hinton and Salakhutdinov(2006)] G. E. Hinton and R. Salakhutdinov. Reducing the dimensionality of data with neural networks. *Science*, 313(5786):504–507, 7 2006.
- [Hinton(2002)] G. E. Hinton. Training products of experts by minimizing contrastive divergence. *Neural Computation*, 14:1771–1800, 8 2002. ISSN 0899-7667.
- [Hyvärinen(1999)] A. Hyvärinen. Fast and robust fixed-point algorithms for independent component analysis. *IEEE Transactions on Neural Networks*, 10(3):626–634, 1999.
- [Hyvärinen et al.(2001)] A. Hyvärinen, J. Karhunen, and E. Oja. *Independent component analysis*. John Wiley & Sons, USA, 2001. ISBN 0-471-40540-X.
- [Karklin and Lewicki(2009)] Y. Karklin and M.S. Lewicki. Emergence of complex cell properties by learning to generalize in natural scenes. *Nature*, 457:83–86, 1 2009.
- [Köster and Hyvärinen(2010)] U. Köster and A. Hyvärinen. A two-layer model of natural stimuli estimated with score matching. *Neural Computation*, 22(9):2308–2333, 2010.
- [Krizhevsky(2009)] A. Krizhevsky. Learning multiple layers of features from tiny images. Master’s thesis, University of Toronto, Toronto, 4 2009.
- [Le Roux et al.(2011)] Le Roux, Heess, Shotton, and Winn] N. Le Roux, N. Heess, J. Shotton, and J. Winn. Learning a generative model of images by factoring appearance and shape. *Neural Computation*, 23(3):593–650, 12 2011.
- [Lee et al.(2007)] H. Lee, C. Ekanadham, and A. Y. Ng. Sparse deep belief net model for visual area v2. In *Proceedings of the 20th Conference on Neural Information Processing Systems*. MIT Press, 2007.
- [Lewicki and Olshausen(1999)] M. S. Lewicki and B. A. Olshausen. Probabilistic framework for the adaptation and comparison of image codes. *Journal of the Optical Society of America*, 16 (7):1587–1601, 7 1999.

- [Melchior(2012)] J. Melchior. Learning natural image statistics with gaussian-binary restricted boltzmann machines. Master’s thesis, ET-IT Dept., Univ. of Bochum, Germany, 2012.
- [Olshausen and Field(1996)] B. A. Olshausen and D. J. Field. Natural image statistics and efficient coding. *Networks*, 7(2):333–339, 5 1996. ISSN 0954-898X.
- [Ranzato and Hinton(2010)] M. Ranzato and G. E. Hinton. Modeling pixel means and covariances using factorized third-order boltzmann machines. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 2551–2558, 2010.
- [Ranzato et al.(2010)] M. Ranzato, A. Krizhevsky, and G. E. Hinton. Factored 3-way restricted boltzmann machines for modeling natural images. *Journal of Machine Learning Research*, 9: 621–628, 2010.
- [Smolensky(1986)] P. Smolensky. Information processing in dynamical systems: foundations of harmony theory. In *Parallel distributed processing: explorations in the microstructure of cognition*, pages 194–281. MIT Press, Cambridge, MA, USA, 1986. ISBN 0-262-68053-X.
- [Theis et al.(2011)] L. Theis, S. Gerwinn, F. Sinz, and M. Bethge. In all likelihood, deep belief is not enough. *Journal of Machine Learning Research*, 12:3071–3096, 11 2011.
- [Tieleman(2008)] T. Tieleman. Training restricted boltzmann machines using approximations to the likelihood gradient. In *Proceedings of the 25th Annual International Conference on Machine Learning*, pages 1064–1071, New York, NY, USA, 2008. ACM. ISBN 978-1-60558-205-4.
- [van Hateren and van der Schaaf(1998)] J. H. van Hateren and A. van der Schaaf. Independent component filters of natural images compared with simple cells in primary visual cortex. *Proceedings of the Royal Society*, 265(1394):359–366, 1998. PMC1688904.
- [Wang et al.(2012)] N. Wang, J. Melchior, and L. Wiskott. An analysis of gaussian-binary restricted boltzmann machines for natural image. In Michel Verleysen, editor, *Proceedings of the 20th European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning*, pages 287–292, Bruges, Belgium, Arpil 2012.
- [Welling et al.(2004)] M. Welling, M. Rosen-Zvi, and G. E. Hinton. Exponential family harmoniums with an application to information retrieval. In *Proceedings of the 17th Conference on Neural Information Processing Systems*. MIT Press, 12 2004.
- [Zetsche and Rohrbein(2001)] C. Zetsche and F. Rohrbein. Nonlinear and extra-classical receptive field properties and the statistics of natural scenes. *Network Comp Neural Sys*, 12:331–350, Aug 2001.