

# TRAINING RESTRICTED BOLTZMANN MACHINES WITH AUXILIARY FUNCTION APPROACH

*Hirokazu Kameoka<sup>1),2)</sup> and Norihiro Takamune<sup>1)</sup>*

<sup>1)</sup> Graduate School of Information Science and Technology, The University of Tokyo

<sup>2)</sup> NTT Communication Science Laboratories, Nippon Telegraph and Telephone Corporation

## ABSTRACT

Restricted Boltzmann Machines (RBMs) are neural network models for unsupervised learning, but have recently found a wide range of applications as feature extractors for supervised learning algorithms. They have also received a lot of attention recently after being proposed as building blocks for deep belief networks. The success of these models raises the issue of how best to train them. At present, the most popular training algorithm for RBMs is the Contrastive Divergence (CD) learning algorithm. The aim of this paper is to seek for a new optimization algorithm tailored for training RBMs in the hope of obtaining a faster algorithm than the CD algorithm. We propose deriving a new training algorithm for RBMs based on an auxiliary function approach. Through an experiment on parameter training of an RBM, we confirmed that the present algorithm converged faster and to a better solution than the CD algorithm.

**Index Terms**— Restricted Boltzmann machines, deep belief networks, contrastive divergence learning algorithm, auxiliary function, minorization-maximization

## 1. INTRODUCTION

Restricted Boltzmann machines (RBMs) [1] are probabilistic generative models that can be interpreted as stochastic neural networks. They have found applications in dimensionality reduction, classification, feature learning, etc., and have recently attracted considerable attention after being proposed as building blocks for deep belief networks (DBNs)[2]. Over the last few years, DBNs have been applied with notable success to a wide range of applications including image recognition and speech recognition [3].

A DBN can be constructed by stacking multiple RBMs so that the hidden layer of one RBM becomes the visible layer of another RBM. DBNs can be efficiently trained using greedy layerwise training, in which the RBMs are trained one at a time in a bottom-up fashion. Furthermore, a discriminative neural network for classification or regression can be obtained by adding an output layer on top of the trained RBM or DBN where the additional units represent labels (e.g., class labels) corresponding to the observed data. The initialization

obtained with the layer-wise pretraining of RBMs facilitates the optimization problem of training the discriminative neural network, a procedure known as the global fine-tuning.

As their name implies, RBMs are a special case of Boltzmann machines. As with Boltzmann machines, RBMs consist of visible units and hidden units. The visible units constitute the first layer, which correspond to the components of an observation. The hidden units constitute the second layer, which model dependencies between the components of observations and can be interpreted as non-linear feature detectors. In the RBMs network graph, each unit is connected to all the units in the other layer. However, there are no connections between units in the same layer. This restriction allows for more efficient training algorithms than those applicable for the general class of Boltzmann machines, in particular the gradient-based contrastive divergence (CD) algorithm [4]. The aim of this paper is to seek for a new optimization algorithm tailored for training RBMs in the hope of obtaining a faster algorithm than the CD algorithm.

For many nonlinear optimization problems, parameter estimation algorithms constructed using an auxiliary function have proven to be very effective. The general principle for the parameter estimation scheme using an auxiliary function is referred to as the “auxiliary function approach” (or alternatively the “minorization-maximization (MM) approach” [5]). Note that the auxiliary function approach itself is not an algorithm, but a description of how to construct an optimization algorithm. When applying the auxiliary function approach to a certain optimization problem, the first step is to design an auxiliary function that upper-bounds or lower-bounds the objective function. A parameter estimation algorithm can then be derived using the auxiliary function. An algorithm that consists of iteratively minimizing/maximizing the auxiliary function is guaranteed to converge to a stationary point of the objective function. It should be noted that this concept is adopted in many existing algorithms. For example, the expectation-maximization (EM) algorithm [6] builds a surrogate for a likelihood function of latent variable models by using Jensen’s inequality. It is also well known for its use in devising an algorithm for non-negative matrix factorization [7, 8]. In general, if we can build a tight upper/lower bound function for the objective function of a specific optimization

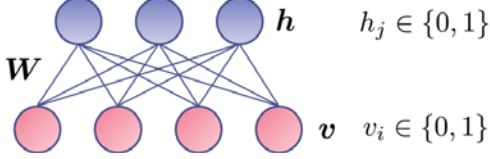


Fig. 1. Graph structure of RBM.

problem, we expect to obtain a fast-converging algorithm. In fact, the authors and colleagues have thus far proposed deriving parameter estimation algorithms based on the auxiliary function approach for various optimization problems, some of which have been proven to be significantly faster than gradient-based methods (e.g., [8–15]). Motivated by these experiences, we propose deriving a new training algorithm for RBMs based on the auxiliary function approach.

## 2. RESTRICTED BOLTZMANN MACHINES

### 2.1. Likelihood function of RBMs

RBMs are undirected generative models that use a layer of hidden variables to model a distribution over visible variables. As depicted in Fig. 1, RBMs can be viewed as a variant of Boltzmann Machines, with the restriction that their neurons must form a bipartite graph structure. This restriction gives the RBM its name.

The standard type of RBM has binary-valued hidden and visible units. A binary RBM with  $J$  hidden units is a parametric model of the joint distribution between a layer of hidden variables  $\mathbf{h} = (h_1, \dots, h_J)^\top \in \{0, 1\}^J$  and the visible variables  $\mathbf{v} = (v_1, \dots, v_I)^\top \in \{0, 1\}^I$ , that is given by

$$p(\mathbf{v}, \mathbf{h} | \Theta) = \frac{\exp(-E(\mathbf{v}, \mathbf{h}; \Theta))}{Z(\Theta)}, \quad (1)$$

where  $E(\mathbf{v}, \mathbf{h}; \Theta)$  denotes a bilinear energy function given by

$$\begin{aligned} E(\mathbf{v}, \mathbf{h}; \Theta) &= -\mathbf{a}^\top \mathbf{v} - \mathbf{b}^\top \mathbf{h} - \mathbf{v}^\top \mathbf{W} \mathbf{h} \\ &= -\sum_i a_i v_i - \sum_j b_j h_j - \sum_{i,j} v_i w_{i,j} h_j, \end{aligned} \quad (2)$$

with parameters  $\Theta = \{\mathbf{W}, \mathbf{a}, \mathbf{b}\}$ .  $\mathbf{W} = (w_{i,j})_{I \times J}$  is a weight matrix whose element  $w_{i,j}$  is associated with the connection between hidden unit  $h_j$  and visible unit  $v_i$ , and  $\mathbf{a} = (a_1, \dots, a_I)^\top$  and  $\mathbf{b} = (b_1, \dots, b_J)^\top$  are bias weight vectors (offsets) associated with the visible units and the hidden units, respectively.  $Z$  is called the “partition function” and is given by the sum of  $\exp(-E(\mathbf{v}, \mathbf{h}; \Theta))$  over all possible pairs of visible and hidden vectors:

$$Z(\Theta) = \sum_{\mathbf{v}} \sum_{\mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h}; \Theta)}. \quad (4)$$

The marginal distribution of  $\mathbf{v}$  is given by summing over all possible hidden vectors:

$$p(\mathbf{v} | \Theta) = \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h} | \Theta) = \frac{1}{Z(\Theta)} \sum_{\mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h}; \Theta)}. \quad (5)$$

For a model of the form (5) with parameters  $\Theta$ , the mean log-likelihood given  $N$  training examples  $\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(N)}$  is thus given by

$$\mathcal{L}(\Theta) = \frac{1}{N} \sum_n \log p(\mathbf{v}^{(n)} | \Theta). \quad (6)$$

Maximum likelihood training of an RBM can be formulated as the problem of maximizing  $\mathcal{L}(\Theta)$  with respect to  $\Theta$ .

### 2.2. Contrastive divergence learning algorithm [2, 4]

Since it is not possible to find the parameters maximizing the mean log-likelihood analytically, one convenient way would be to apply a gradient ascent approach. Although here we consider only the update rule for  $w_{i,j}$ , the update rule for the biases  $a_i, b_j$  can be defined analogously.

The partial derivative of the mean log-likelihood with respect to  $w_{i,j}$  is given by

$$\begin{aligned} \frac{\partial \mathcal{L}(\Theta)}{\partial w_{i,j}} &= -\frac{1}{N} \sum_n \sum_{\mathbf{h}} p(\mathbf{h} | \mathbf{v}^{(n)}, \Theta) v_i^{(n)} h_j \\ &\quad + \sum_{\mathbf{v}} \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h} | \Theta) v_i h_j. \end{aligned} \quad (7)$$

This leads to a simple update rule based on steepest ascent in the log-likelihood:

$$w_{i,j} \leftarrow w_{i,j} + \epsilon \frac{\partial \mathcal{L}(\Theta)}{\partial w_{i,j}}, \quad (8)$$

where  $\epsilon$  is a learning rate. As can be seen from (7), the gradient  $\partial \mathcal{L}(\Theta) / \partial w_{i,j}$  is given by the difference between two expectations: the expectation of the energy gradient under the RBM distribution  $p(\mathbf{v}, \mathbf{h} | \Theta)$  and under the conditional distribution of the hidden variables  $p(\mathbf{h} | \mathbf{v}, \Theta)$  given a set of the training samples. To perform gradient ascent, we need to compute these two terms.

Since there are no direct connections between hidden units as well as between visible units in an RBM, hidden units are conditionally independent given the state of the visible variables and vice versa:

$$p(\mathbf{h} | \mathbf{v}, \Theta) = \prod_j p(h_j | \mathbf{v}, \Theta), \quad (9)$$

$$p(\mathbf{v} | \mathbf{h}, \Theta) = \prod_i p(v_i | \mathbf{h}, \Theta). \quad (10)$$

The conditional distributions  $p(h_j | \mathbf{v}, \Theta)$  and  $p(v_i | \mathbf{h}, \Theta)$  can be interpreted as the firing rate of each neuron, given by

$$p(h_j = 1 | \mathbf{v}, \Theta) = \sigma\left(b_j + \sum_i v_i w_{i,j}\right), \quad (11)$$

$$p(v_i = 1 | \mathbf{h}, \Theta) = \sigma\left(a_i + \sum_j w_{i,j} h_j\right), \quad (12)$$

where  $\sigma$  denotes the logistic sigmoid:  $\sigma(x) = 1/(1 + e^{-x})$ . By using these facts, the first term of (7) can be rewritten as

$$\frac{1}{N} \sum_n p(h_i = 1 | \mathbf{v}^{(n)}, \Theta) v_j^{(n)}, \quad (13)$$

which is easy to compute. As for the second term, however, directly calculating the sums that run over all values of  $\mathbf{v}$  and  $\mathbf{h}$  leads to a computational complexity which is in general exponential in the number of variables. To avoid this exponential complexity, one convenient way is to approximate the expectation by samples from the RBM distribution. These samples can be obtained with Gibbs sampling using the conditional distributions (11) and (12). It has been shown that estimates obtained after running the Gibbs chain for just a few steps can be sufficient in practice.

This way of updating the parameters belongs to the family of the contrastive divergence (CD) learning algorithm [4], which has become a standard way of training RBMs.

### 3. AUXILIARY FUNCTION APPROACH FOR TRAINING RBM

#### 3.1. Auxiliary function approach

Here we introduce the general principle of the auxiliary function approach (a.k.a the MM approach).

Let us use  $G(\theta)$  to denote an objective function that we want to maximize with respect to  $\theta$ .  $G^+(\theta, \tilde{\theta})$  is defined as an auxiliary function for  $G(\theta)$  if it satisfies

$$G(\theta) = \max_{\tilde{\theta}} G^+(\theta, \tilde{\theta}). \quad (14)$$

We call  $\tilde{\theta}$  an auxiliary variable. By using  $G^+(\theta, \tilde{\theta})$ ,  $G(\theta)$  can be iteratively increased according to the following theorem:

**Theorem 1.**  $G(\theta)$  is non-decreasing under the updates,  $\theta \leftarrow \operatorname{argmax}_{\theta} G^+(\theta, \tilde{\theta})$  and  $\tilde{\theta} \leftarrow \operatorname{argmax}_{\tilde{\theta}} G^+(\theta, \tilde{\theta})$ .

**Proof:** Let us set  $\theta$  at an arbitrary value  $\theta_\ell$  and define  $\tilde{\theta}_{\ell+1} = \operatorname{argmax}_{\tilde{\theta}} G^+(\theta_\ell, \tilde{\theta})$  and  $\theta_{\ell+1} = \operatorname{argmax}_{\theta} G^+(\theta, \tilde{\theta}_{\ell+1})$ . First, it is obvious that  $G(\theta_\ell) = G^+(\theta_\ell, \tilde{\theta}_{\ell+1})$ . Next, we can confirm that  $G^+(\theta_\ell, \tilde{\theta}_{\ell+1}) \leq G^+(\theta_{\ell+1}, \tilde{\theta}_{\ell+1})$  since  $\theta_{\ell+1}$  is the maximizer of  $G^+(\theta, \tilde{\theta}_{\ell+1})$  with respect to  $\theta$ . By definition, it is obvious that  $G^+(\theta_{\ell+1}, \tilde{\theta}_{\ell+1}) \leq G(\theta_{\ell+1})$  and so we can finally show that  $G(\theta_\ell) \leq G(\theta_{\ell+1})$ .  $\square$

We propose applying this principle for the problem of maximizing the mean log-likelihood  $\mathcal{L}(\Theta)$  of RBM.

#### 3.2. Designing auxiliary function

When applying the auxiliary function approach to a certain maximization problem, the first step is to design an auxiliary function that lower-bounds the objective function. The difficulty in solving the optimization problem of maximizing  $\mathcal{L}(\Theta)$  lies in the nonlinear interdependence of the parameters

due to the sums inside the logarithm function that run over  $\mathbf{v}$  and  $\mathbf{h}$ , and the sums inside the exponential function that run over  $i$  and  $j$ . Since it is preferable to use a function that can be maximized analytically as an auxiliary function, we would like to design an auxiliary function such that the parameters are separated into individual terms.

First, let us focus on the sum over  $\mathbf{h}$  inside the logarithm function in  $\mathcal{L}(\Theta)$ . We can invoke the following inequality to construct a lower bound function of  $\mathcal{L}(\Theta)$ .

**Lemma 1** (Jensen's inequality for concave functions with non-negative arguments). *For any concave function  $f$ , any non-negative values  $x_1, \dots, x_K \in [0, \infty)$  and any positive weights  $\lambda_1, \dots, \lambda_K \in (0, 1)$  that sums to unity (i.e.,  $\sum_k \lambda_k = 1$ ), we have*

$$f\left(\sum_k x_k\right) \geq \sum_k \lambda_k f\left(\frac{x_k}{\lambda_k}\right), \quad (15)$$

where equality holds if and only if

$$\lambda_k = \frac{x_k}{\sum_m x_m}. \quad (16)$$

Since the logarithm function is a concave function, and of course the value of  $p(\mathbf{v}^{(n)}, \mathbf{h} | \Theta)$  is always non-negative, we can use the above inequality:

$$\begin{aligned} \mathcal{L}(\Theta) &= \frac{1}{N} \sum_n \log \sum_{\mathbf{h}} p(\mathbf{v}^{(n)}, \mathbf{h} | \Theta) \\ &\geq \frac{1}{N} \sum_n \sum_{\mathbf{h}} \lambda_n(\mathbf{h}) \log \frac{p(\mathbf{v}^{(n)}, \mathbf{h} | \Theta)}{\lambda_n(\mathbf{h})}, \end{aligned} \quad (17)$$

where  $\lambda_n(\mathbf{h})$  is a positive weight that sums to unity, such that  $\sum_{\mathbf{h}} \lambda_n(\mathbf{h}) = 1$ . Here, equality holds when

$$\lambda_n(\mathbf{h}) = \frac{p(\mathbf{v}^{(n)}, \mathbf{h} | \Theta)}{\sum_{\mathbf{h}'} p(\mathbf{v}^{(n)}, \mathbf{h}' | \Theta)} = p(\mathbf{h} | \mathbf{v}^{(n)}, \Theta). \quad (18)$$

Thus, (17) can be rewritten as

$$\begin{aligned} \mathcal{L}(\Theta) &\geq -\frac{1}{N} \sum_n \sum_{\mathbf{h}} \lambda_n(\mathbf{h}) E(\mathbf{v}^{(n)}, \mathbf{h}; \Theta) \\ &\quad - \log Z(\Theta) - \sum_{\mathbf{h}} \lambda_n(\mathbf{h}) \log \lambda_n(\mathbf{h}). \end{aligned} \quad (19)$$

It is still impossible to obtain a closed-form update rule for  $\Theta$  from this lower bound function because of the second term. In the second term, we notice that the sum is taken over  $\mathbf{v}$  and  $\mathbf{h}$  inside the negative logarithm function. We can use the following inequality to construct a looser lower bound.

**Lemma 2.** *For any differentiable convex function  $f$  and any real numbers  $x$  and  $\zeta$ , we have*

$$f(x) \geq f(\zeta) + f'(\zeta)(x - \zeta), \quad (20)$$

where equality holds if and only if  $x = \zeta$ .

Since the negative logarithm function is a differentiable convex function, we can apply this inequality to  $-\log Z(\Theta)$ :

$$-\log Z(\Theta) \geq -\log \zeta - \frac{1}{\zeta}(Z(\Theta) - \zeta), \quad (21)$$

where  $\zeta$  is a positive number. Here, equality holds when

$$\zeta = Z(\Theta). \quad (22)$$

A lower bound function can thus be constructed by replacing  $-\log Z(\Theta)$  with the right-hand side of (21). However, it is still impossible to obtain a closed-form update rule for  $\Theta$  from this lower bound function because of the sums that run over  $i$  and  $j$  inside the negative exponential function in  $Z(\Theta)$ . Since the negative exponential function is a concave function, it may appear that Lemma 1 can be applied in a similar way that we obtain the first inequality. However, since the argument of the negative exponential function, i.e.,  $-E(\mathbf{v}, \mathbf{h}; \Theta)$ , can possibly take on negative values, the inequality given in Lemma 1 cannot be applied in this case. Instead, here we introduce the following inequality (which we have previously employed for a different optimization problem [11]).

**Lemma 3** (Jensen's inequality for concave functions with real arguments). *For any concave function  $f$ , any real numbers  $x_1, \dots, x_K \in \mathbb{R}$ , any positive weights  $\beta_1, \dots, \beta_K \in (0, 1)$  that sums to unity (i.e.,  $\sum_k \beta_k = 1$ ), and any real numbers  $\alpha_1, \dots, \alpha_K \in \mathbb{R}$  that sums to zero (i.e.,  $\sum_k \alpha_k = 0$ ), we have*

$$f\left(\sum_k x_k\right) \geq \sum_k \beta_k f\left(\frac{x_k - \alpha_k}{\beta_k}\right), \quad (23)$$

where equality holds if and only if

$$\alpha_k = x_k - \beta_k \sum_m x_m. \quad (24)$$

By using this inequality, we obtain

$$\begin{aligned} & -\exp\left(\sum_i a_i v_i + \sum_j b_j h_j + \sum_{i,j} v_i w_{i,j} h_j\right) \\ & \geq -\sum_i \beta_i^a \exp\left(\frac{a_i v_i - \alpha_i^a(\mathbf{v}, \mathbf{h})}{\beta_i^a}\right) \\ & \quad -\sum_j \beta_j^b \exp\left(\frac{b_j h_j - \alpha_j^b(\mathbf{v}, \mathbf{h})}{\beta_j^b}\right) \\ & \quad -\sum_{i,j} \beta_{i,j}^w \exp\left(\frac{v_i w_{i,j} h_j - \alpha_{i,j}^w(\mathbf{v}, \mathbf{h})}{\beta_{i,j}^w}\right), \end{aligned} \quad (26)$$

where  $\alpha_i^a(\mathbf{v}, \mathbf{h})$ ,  $\alpha_j^b(\mathbf{v}, \mathbf{h})$  and  $\alpha_{i,j}^w(\mathbf{v}, \mathbf{h})$  are arbitrary real numbers that sum to zero, i.e.,  $\sum_i \alpha_i^a(\mathbf{v}, \mathbf{h}) + \sum_j \alpha_j^b(\mathbf{v}, \mathbf{h}) + \sum_{i,j} \alpha_{i,j}^w(\mathbf{v}, \mathbf{h}) = 0$ , and  $\beta_{i,j}^a$ ,  $\beta_{i,j}^b$  and  $\beta_{i,j}^w$  are arbitrary positive numbers that sum to one, i.e.,  $\sum_i \beta_i^a + \sum_j \beta_j^b + \sum_{i,j} \beta_{i,j}^w = 1$ . Here, equality holds when

$$\alpha_i^a(\mathbf{v}, \mathbf{h}) = a_i v_i + \beta_i^a E(\mathbf{v}, \mathbf{h}; \Theta),$$

$$\begin{aligned} \alpha_j^b(\mathbf{v}, \mathbf{h}) &= b_j h_j + \beta_j^b E(\mathbf{v}, \mathbf{h}; \Theta), \\ \alpha_{i,j}^w(\mathbf{v}, \mathbf{h}) &= v_i w_{i,j} h_j + \beta_{i,j}^w E(\mathbf{v}, \mathbf{h}; \Theta). \end{aligned} \quad (27)$$

We can obtain a slightly looser lower bound function for  $\mathcal{L}(\Theta)$  by further replacing  $-e^{-E(\mathbf{v}, \mathbf{h}; \Theta)}$  with the right-hand side of (26). It is important to note that this lower bound function is given in such a way that the parameters are separated in individual terms, thus allowing us to derive closed-form update equations for the parameters.

To sum up, we obtain the following theorem:

**Theorem 2** (Auxiliary function for RBMs). *Define  $\mathcal{L}^+$  as*

$$\begin{aligned} \mathcal{L}^+(\Theta, \tilde{\Theta}) &= -\frac{1}{N} \sum_n \sum_{\mathbf{h}} \lambda_n(\mathbf{h}) E(\mathbf{v}^{(n)}, \mathbf{h}; \Theta) \\ & \quad -\frac{1}{\zeta} \sum_{\mathbf{v}} \sum_{\mathbf{h}} G(\mathbf{v}, \mathbf{h}; \Theta) \\ & \quad -\log \zeta + 1 - \sum_{\mathbf{h}} \lambda_n(\mathbf{h}) \log \lambda_n(\mathbf{h}), \end{aligned} \quad (28)$$

where

$$\begin{aligned} G(\mathbf{v}, \mathbf{h}; \Theta) &= \sum_i \beta_i^a \exp\left(\frac{a_i v_i - \alpha_i^a(\mathbf{v}, \mathbf{h})}{\beta_i^a}\right) \\ & \quad + \sum_j \beta_j^b \exp\left(\frac{b_j h_j - \alpha_j^b(\mathbf{v}, \mathbf{h})}{\beta_j^b}\right) \\ & \quad + \sum_{i,j} \beta_{i,j}^w \exp\left(\frac{v_i w_{i,j} h_j - \alpha_{i,j}^w(\mathbf{v}, \mathbf{h})}{\beta_{i,j}^w}\right), \end{aligned} \quad (29)$$

with  $\tilde{\Theta} = \{\lambda, \zeta, \alpha\}$ ,  $\lambda = \{\lambda_n(\mathbf{h}) | n=1, \dots, N, \mathbf{h} \in \{0, 1\}^J\}$  and  $\alpha = \{\alpha_i^a(\mathbf{v}, \mathbf{h}), \alpha_j^b(\mathbf{v}, \mathbf{h}), \alpha_{i,j}^w(\mathbf{v}, \mathbf{h}) | i=1, \dots, I, j=1, \dots, J, \mathbf{v} \in \{0, 1\}^I, \mathbf{h} \in \{0, 1\}^J\}$ . Then,  $\mathcal{L}^+(\Theta, \tilde{\Theta})$  is an auxiliary function for  $\mathcal{L}(\Theta)$ .

### 3.3. Update equations

Now we can derive the update equations for the model parameters by using the auxiliary function defined above. As stated in Theorem 1, we can iteratively increase the value of  $\mathcal{L}(\Theta)$  by performing the updates

$$\tilde{\Theta} \leftarrow \underset{\tilde{\Theta}}{\operatorname{argmax}} \mathcal{L}^+(\Theta, \tilde{\Theta}), \quad (30)$$

$$\Theta \leftarrow \underset{\Theta}{\operatorname{argmax}} \mathcal{L}^+(\Theta, \tilde{\Theta}). \quad (31)$$

(30) is given explicitly as (18), (22) and (27), respectively. (31) can be obtained by solving  $\partial \mathcal{L}^+ / \partial a_i = 0$ ,  $\partial \mathcal{L}^+ / \partial b_j = 0$  and  $\partial \mathcal{L}^+ / \partial w_{i,j} = 0$ . For example,  $\partial \mathcal{L}^+ / \partial w_{i,j}$  is given by

$$\begin{aligned} \frac{\partial \mathcal{L}^+(\Theta)}{\partial w_{i,j}} &= \frac{1}{N} \sum_n \sum_{\mathbf{h}} \lambda_n(\mathbf{h}) v_i^{(n)} h_j \\ & \quad - \frac{1}{\zeta} \sum_{\mathbf{v}} \sum_{\mathbf{h}} v_i h_j \exp\left(\frac{v_i w_{i,j} h_j - \alpha_{i,j}^w(\mathbf{v}, \mathbf{h})}{\beta_{i,j}^w}\right). \end{aligned} \quad (32)$$

Now, by substituting (30) into (32),  $\partial\mathcal{L}^+/\partial w_{i,j}$  becomes

$$\frac{\partial\mathcal{L}^+(\Theta)}{\partial w_{i,j}} = \frac{1}{N} \sum_n \sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{v}^{(n)}, \Theta') v_i^{(n)} h_j \quad (33)$$

$$- \sum_{\mathbf{v}} \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h}|\Theta') v_i h_j \exp\left(\frac{v_i h_j (w_{i,j} - w'_{i,j})}{\beta_{i,j}^w}\right),$$

where  $w'_{i,j}$  and  $\Theta'$  denote the estimates of  $w_{i,j}$  and  $\Theta$  updated at the previous iteration. Since  $v_i h_j \in \{0, 1\}$ , we can write the second term of (33) as

$$\sum_{\mathbf{v}} \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h}|\Theta') v_i h_j \exp\left(\frac{v_i h_j (w_{i,j} - w'_{i,j})}{\beta_{i,j}^w}\right)$$

$$= \exp\left(\frac{w_{i,j} - w'_{i,j}}{\beta_{i,j}^w}\right) \sum_{\mathbf{v}} \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h}|\Theta') v_i h_j. \quad (34)$$

This eventually allows us to solve  $\partial\mathcal{L}^+/\partial w_{i,j} = 0$ :

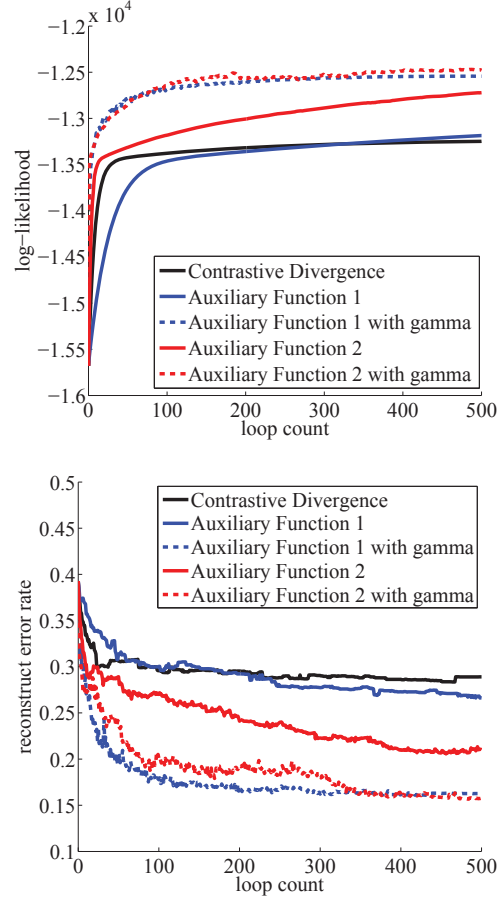
$$w_{i,j} = w'_{i,j} + \beta_{i,j}^w \log \frac{\frac{1}{N} \sum_n \sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{v}^{(n)}, \Theta') v_i^{(n)} h_j}{\sum_{\mathbf{v}} \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h}|\Theta') v_i h_j}. \quad (35)$$

Similarly to the CD approach, the numerator inside the logarithm function of (35) can be easily computed since it can be written as (13). However, the denominator involves the sums that run over all values of  $\mathbf{v}$  and  $\mathbf{h}$ . To avoid directly calculating these sums, one convenient way is to approximate the expectation by samples from the RBM distribution, as with the CD approach. As mentioned in 2.2, these samples can be obtained with Gibbs sampling using the conditional distributions (11) and (12). The update rule for the biases  $a_i, b_j$  can be derived analogously.

It is interesting to compare the proposed update rule (35) with the update rule (8) of the CD algorithm. While each step of the CD algorithm moves the parameters in the direction given by the difference of the two expectations, the proposed update rule moves the parameters in the direction given by the difference of the logarithms of the two expectations. From this perspective,  $\beta_{i,j}^w$  can be thought of as the learning rate. Since  $\beta_{i,j}^w$  must satisfy the sum-to-unity constraint, the larger the number of the hidden and visible variables becomes, the smaller  $\beta_{i,j}^w$  becomes in average, thus slowing the convergence speed of the algorithm. We can thus expect to accelerate the algorithm by replacing  $\beta_{i,j}^w$  with  $\beta_{i,j}^w \gamma$ , setting  $\gamma$  at a reasonably small value at the early stage of the algorithm and moving it towards 1 as the iteration proceeds.

#### 4. EXPERIMENTS

We conducted an experiment to compare the convergence speeds of the present algorithm and the CD algorithm. In order to evaluate the values of the likelihoods, We used a



**Fig. 2.** Evolution of the mean log-likelihood (top) and the reconstruction error (bottom) during training of a small scale RBM using the CD algorithm (black solid line) and the proposed algorithms (red and blue) with and without the  $\gamma$  scheduling (solid and dashed line), respectively. “Auxiliary 1” refers to the the present algorithm and “Auxiliary 2” refers to the algorithm proposed in [16].

small scale RBM with  $I = 10$  visible units and  $J = 8$  hidden units. We also compared the present algorithm with another algorithm we proposed in [16] (“auxiliary function 1” and “auxiliary function 2” hereafter).

The experimental conditions were as follows: 2000 randomly-generated binary data were used as the training examples. All the algorithms were run for  $T = 500$  iterations. The Gibbs sampling was run for 1 iteration at each iterative step.  $\beta = \{\beta_i^a, \beta_j^b, \beta_{i,j}^w | i = 1, \dots, I, j = 1, \dots, J\}$  was set uniformly, i.e.,

$$\beta_i^a = \beta_j^b = \beta_{i,j}^w = \frac{1}{I + J + IJ}, \quad (36)$$

for all  $i$  and  $j$ . The learning rate  $\epsilon$  of the CD algorithm and the value of  $\gamma$  were scheduled in the following way: At step  $t$  of the algorithms, we updated  $\epsilon$  and  $\gamma$  at

$$\epsilon_t = \epsilon_{\text{init}} \left( \frac{\epsilon_{\text{end}}}{\epsilon_{\text{init}}} \right)^{\frac{t-1}{T-1}}, \quad (37)$$

$$\gamma_t = \gamma_{\text{init}} \left( \frac{\gamma_{\text{end}}}{\gamma_{\text{init}}} \right)^{\frac{t-1}{T-1}}, \quad (38)$$

respectively, where we set  $\epsilon_{\text{init}} = 1$ ,  $\epsilon_{\text{end}} = 0.1$ ,  $\gamma_{\text{init}} = 0.1$  and  $\gamma_{\text{end}} = 1$ . In addition to the log-likelihood given by (6), we used the reconstruction error rate  $E$  for the measure for comparisons:

$$E = \frac{1}{NI} \sum_n \sum_i (v_i^{(n)} - \hat{v}_i^{(n)})^2, \quad (39)$$

where

$$\hat{\mathbf{v}}^{(n)} = \underset{\mathbf{v}}{\operatorname{argmax}} p(\mathbf{v} | \hat{\mathbf{h}}^{(n)}, \Theta), \quad (40)$$

$$\hat{\mathbf{h}}^{(n)} = \underset{\mathbf{h}}{\operatorname{argmax}} p(\mathbf{h} | \mathbf{v}^{(n)}, \Theta). \quad (41)$$

Fig. 2 shows the evolution of the log-likelihoods (6) and the reconstruction error rate (39) obtained with all the algorithms. With the MATLAB implementation, the execution time per iteration for the CD algorithm and the proposed algorithms was almost the same. For the present algorithm, the log-likelihood converged faster when updating  $\gamma$  than when fixing it at 1, as expected. As Fig. 2 shows, the present algorithm increased the log-likelihood and decreased the reconstruction error rate of the RBM faster than the CD algorithm. It is interesting to note that the present algorithm converged to a greater value of the log-likelihood than the CD algorithm. This implies that the present algorithm had a higher ability to avoid getting trapped into local maxima.

## 5. CONCLUSIONS

This paper proposed deriving a new training algorithm for RBMs based on the auxiliary function concept. While each step of the CD algorithm moves the parameters in the direction given by the difference of the two expectations, the proposed update rule moves the parameters in the direction given by the difference of the logarithms of the two expectations. Through an experiment on parameter training of an RBM, we confirmed that the present algorithm converged faster and to a better solution than the CD algorithm.

## 6. REFERENCES

- [1] P. Smolensky, "Information processing in dynamical systems: Foundations of harmony theory," in *Parallel Distributed Processing: Explorations in the Microstructure of Cognition*, D.E. Rumelhart and J.L. McClelland, Eds. MIT Press, 1986.
- [2] G. E. Hinton, S. Osindero, and Y. W. Teh, "A fast learning algorithm for deep belief nets," *Neural Computation*, vol. 18, no. 7, 2006.
- [3] G. E. Hinton, L. Deng, D. Yu, G. E. Dahl, A. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, T. N. Sainath, and B. Kingsbury, "Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups," *IEEE Signal Processing Magazine*, vol. 29, no. 6, pp. 82–97, 2012.
- [4] G. E. Hinton, "Training products of experts by minimizing contrastive divergence," *Neural Computation*, vol. 14, no. 8, 2002.
- [5] J. De Leeuw and W. J. Heiser, "Convergence of correction matrix algorithms for multidimensional scaling," in *Geometric representations of relational data*, J. C. Lingoes, E. E. Roskam, and I. Borg, Eds. Ann Arbor, MI: Mathesis Press, 1977.
- [6] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood from incomplete data via the EM algorithm," *J. of Royal Statistical Society Series B*, vol. 39, 1977.
- [7] D. D. Lee and H. S. Seung, "Algorithms for non-negative matrix factorization," in *Adv. Neural Information and Processing Systems (NIPS)*. 2001, pp. 556–562, MIT Press.
- [8] M. Nakano, H. Kameoka, J. Le Roux, Y. Kitano, N. Ono, and S. Sagayama, "Convergence-guaranteed multiplicative algorithms for non-negative matrix factorization with beta-divergence," in *Proc. 2010 IEEE International Workshop on Machine Learning for Signal Processing (MLSP 2010)*, 2010, pp. 283–288.
- [9] H. Kameoka, *Statistical Approach to Multipitch Analysis*, Ph.D. thesis, The University of Tokyo, 2007.
- [10] H. Kameoka, N. Ono, and S. Sagayama, "Auxiliary function approach to parameter estimation of constrained sinusoidal model for monaural speech separation," in *Proc. 2008 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP2008)*, 2008, pp. 29–32.
- [11] H. Kameoka, N. Ono, K. Kashino, and S. Sagayama, "Complex NMF: A new sparse representation for acoustic signals," in *Proc. 2009 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP2009)*, 2009, pp. 3437–3440.
- [12] H. Kameoka, T. Nakatani, and T. Yoshioka, "Robust speech dereverberation based on non-negativity and sparse nature of speech spectrograms," in *Proc. 2009 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP2009)*, 2009, pp. 45–48.
- [13] N. Ono, "Stable and fast update rules for independent vector analysis based on auxiliary function technique," in *Proc. 2011 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA2011)*, 2011, pp. 261–264.
- [14] N. Yasuraoka, H. Kameoka, T. Yoshioka, and H. G. Okuno, "I-divergence-based dereverberation method with auxiliary function approach," in *Proc. 2011 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP2011)*, 2011, pp. 369–372.
- [15] H. Sawada, H. Kameoka, S. Araki, and N. Ueda, "Efficient algorithms for multichannel extensions of Itakura-Saito non-negative matrix factorization," in *Proc. 2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP2012)*, 2012, pp. 261–264.
- [16] N. Takamune and H. Kameoka, "Maximum reconstruction probability training of restricted Boltzmann machines with auxiliary function approach," in *Proc. 2014 IEEE International Workshop on Machine Learning for Signal Processing (MLSP 2014)*, 2014, submitted.