

# Restricted Deep Belief Networks for Multi-view Learning

Yoonseop Kang<sup>1</sup> and Seungjin Choi<sup>1,2</sup>

<sup>1</sup> Department of Computer Science

<sup>2</sup> Division of IT Convergence Engineering

Pohang University of Science and Technology

San 31 Hyoja-dong, Nam-gu, Pohang 790-784, Korea

{e0en, seungjin}@postech.ac.kr

**Abstract.** Deep belief network (DBN) is a probabilistic generative model with multiple layers of hidden nodes and a layer of visible nodes, where parameterizations between layers obey harmonium or restricted Boltzmann machines (RBMs). In this paper we present *restricted deep belief network* (RDBN) for multi-view learning, where each layer of hidden nodes is composed of *view-specific* and *shared* hidden nodes, in order to learn individual and shared hidden spaces from multiple views of data. View-specific hidden nodes are connected to corresponding view-specific hidden nodes in the lower-layer or visible nodes involving a specific view, whereas shared hidden nodes follow inter-layer connections without restrictions as in standard DBNs. RDBN is trained using layer-wise contrastive divergence learning. Numerical experiments on synthetic and real-world datasets demonstrate the useful behavior of the RDBN, compared to the multi-wing harmonium (MWH) which is a two-layer undirected model.

## 1 Introduction

Multi-view learning refers to methods for learning from examples that have multiple (independent or dependent) representations, each of which may arise from different views such as modalities or sensors. For instance, in web page classification, one view describes a web page using the text appearing on the document itself, while the other view leads to the anchor text to hyperlinks pointing to this page from other pages [2]. In multimedia mining, images are described by color histograms (one view) and annotated text (the other view), so it is desirable to exploit these two information sources together to boost the performance of image classification [18]. In brain wave classification where EEG data are measured from multiple subjects who undergo the same mental task, each view corresponds to each subject, then learning shared latent spaces provides useful semantic features for EEG classification [11].

Learning latent spaces that capture the relevant information shared by multiple views of data lies at the heart of multi-view learning, especially when views are dependent. One of the oldest but the most popular method is canonical correlation analysis (CCA) [8] which identifies linear relationships between two sets

of observations. The shared Gaussian process latent variable model (sGPLVM) [14] is a nonlinear extension of CCA, in which Gaussian process regression is used to learn common hidden structure shared between corresponding sets of heterogeneous observations. Manifold integration [4] combines similarity matrices (each of which is determined by a specific view) into a compromise matrix that faithfully reflects multiple sensory information.

Most of these methods assume that views are fully independent conditioned on common latent variables, i.e., only shared latent variables are considered to explain the dependency between views. However, in real-world problems, this assumption is not satisfied. Thus it was suggested that separate latent spaces are learned to model the shared characteristics across views and individual components of observations [11,13]. Group nonnegative matrix factorization (GNMF) learn jointly individual latent spaces and shared latent spaces, decomposing non-negative data matrix into a product of two nonnegative factor matrices where the basis matrix is composed of common and individual basis vectors [11]. Factorized orthogonal latent space (FOLS) method [13] factorizes the latent space into shared and private latent spaces, enforcing orthogonality between them. These aforementioned methods have limitations since CCA and GNMF are linear techniques and FOLS requires large memory storage and expensive time complexity to handle operations involving a big Gram matrix, which is not scalable as the number of samples increases.

Multi-wing harmonium (MWH) is a two-layer undirected graphical model designed to handle multiple view data. MWH consists of two or more wings where each wing takes the input associated with a single view, extending harmoniums [16] and restricted Boltzmann machines (RBMs) [5]. Harmoniums were also extended to the exponential family [17]. MWH inherits the advantages of exponential family harmoniums (EFHs) such as easy inference and distributed representations over latent variables (see [17] for other advantages). Deep belief network (DBN) is composed of RBMs with multiple layers of hidden nodes, which is trained efficiently in layer-wise manner based on contrastive divergence [5,6]. See [1] for excellent tutorial on DBNs. Multilayer structure of DBNs allow more complex representation of data than RBMs. However, unlike other Boltzmann machine-based models, MWHs cannot be naturally extended to form a deep network.

In this paper we present *restricted deep belief network* (RDBN) for partially correlated multiple view data. We first present a modification of EFHs where view-specific hidden nodes are restricted to have undirected connections to only visible nodes involving corresponding views whereas shared hidden nodes are connected to all the visible nodes. This model inherits advantages of FOLS and MWH simultaneously. The model can efficiently evaluate latent variables as MWH, while still being capable of modelling shared and private information separately.

Then we stack these modified harmoniums to construct RDBN in such a way that view-specific hidden nodes are connected to corresponding view-specific hidden nodes in the lower-layer and shared hidden nodes are connected to all the

hidden nodes in the lower-layer. We train RDBNs using layer-wise contrastive divergence learning, to learn view-specific and shared hidden spaces from multiple views. Numerical experiments on a synthetic dataset, NORB-small dataset, and ESL photo dataset demonstrated the useful behavior of RDBN, compared to MWHs and its direct multilayer version.

## 2 Related Work

### 2.1 Exponential Family Harmonium

Most of graphical models including with a layer of hidden nodes and other layer of observed nodes are based on a directed graph. These directed two-layer models allows easy sampling of visible layer, and easy handling of latent variables. However, it is often very difficult and time consuming to calculate posterior distribution of hidden nodes given visible nodes. For the tasks that requires fast evaluation of latent variable given a test sample, directed models would not be very effective.

EFH is a two-layer probabilistic graphical model whose probability distributions are constrained to be exponential family distribution. EFH consists of a set of hidden nodes  $\mathbf{h}$ , and a set of visible nodes  $\mathbf{v}$  connected by undirected edges. By incorporating undirected connections, the model calculates posterior distribution of hidden nodes  $p(\mathbf{h}|\mathbf{v})$  much faster than directed models, making it more suitable for the tasks including document searching and automatic image annotation.

To define an EFH, we start from choosing marginal distributions for visible nodes  $\mathbf{v}$  and hidden nodes  $\mathbf{h}$ . As the distributions are constrained to be exponential family distributions, the marginal distributions are defined as below:

$$p(\mathbf{v}) = \prod_i \exp\left\{\sum_a \lambda_{ia} f_{ia}(x_i) - A_i(\{\lambda_{ia}\})\right\}, \quad (1)$$

$$p(\mathbf{h}) = \prod_j \exp\left\{\sum_c \eta_{jc} g_{jc}(h_j) - B_j(\{\eta_{jc}\})\right\}, \quad (2)$$

where  $f_{ia}$  and  $g_{jc}$  are  $a$ , and  $c$ th sufficient statistics of  $v_i$  and  $h_j$ .  $\lambda^x$  and  $\eta$  are parameters, and  $A$  and  $B$  are log partition functions. Note that each  $v_i$ s and  $h_j$  are assumed to be independent to each other.

With this definition, we define joint distribution of  $\mathbf{v}$  and  $\mathbf{h}$  by combining their distributions multiplicatively. The joint distribution is defined by introducing quadratic terms to inter-layer relationship:

$$p(\mathbf{v}, \mathbf{h}|\theta) \propto \exp\left\{\sum_{i,a} \lambda_{ia} f_{ia}(v_i) + \sum_{j,c} \eta_{jc} g_{jc}(h_j) + \sum_{i,j,a,c} W_{ijac} f_{ia}(v_i) g_{jc}(h_j)\right\}. \quad (3)$$

As the model is a bipartite graph, between-layer conditional distributions can be represented as products of distributions of individual nodes. The conditional

distributions are derived as below:

$$p(\mathbf{v}|\mathbf{h}, \theta) = \prod_i \exp\left\{\sum_a \hat{\lambda}_{ia} f_{ia}(v_i) - A_i(\{\hat{\lambda}_{ia}\})\right\}, \quad (4)$$

$$p(\mathbf{h}|\mathbf{v}, \theta) = \prod_j \exp\left\{\sum_c \hat{\eta}_{jc} g_{jc}(h_j) - B_j(\{\hat{\eta}_{jc}\})\right\}, \quad (5)$$

where shifted parameters are

$$\hat{\lambda}_{ia} = \lambda_{ia} + \sum_{j,c} W_{ijac} g_{jc}(h_j), \quad (6)$$

$$\hat{\eta}_{jc} = \eta_{jc} + \sum_{i,a} W_{ijac} f_{ia}(x_i). \quad (7)$$

As the conditional distribution of nodes are independent to each other, sampling posterior distribution of hidden nodes is done by just simply evaluating conditional distribution  $p(\mathbf{h}|\mathbf{v}, \theta)$ . An EFH with binary distributions for  $\mathbf{v}$  and  $\mathbf{h}$  becomes equivalent to a RBM.

## 2.2 Multi-Wing Harmonium

Multi-wing harmonium (MWH) [18] models joint distribution of multi-view data using two-layer graphical model. Given two-view data, MWH uses two sets of visible nodes  $\mathbf{x}$  and  $\mathbf{y}$  connected to hidden nodes  $\mathbf{h}$ . By choosing different distribution for different type of information, the model achieves better representation of data.

Construction of an MWH is similar to the one of EFH. The only difference is that we split visible nodes  $\mathbf{v}$  to two sets  $\mathbf{x}$  and  $\mathbf{y}$ . First we choose marginal distributions for  $\mathbf{x}$ ,  $\mathbf{y}$  and  $\mathbf{h}$  as below:

$$p(\mathbf{x}) = \prod_i \exp\left\{\sum_a \lambda_{ia}^x f_{ia}^x(x_i) - A_i^x(\{\lambda_{ia}^x\})\right\}, \quad (8)$$

$$p(\mathbf{y}) = \prod_k \exp\left\{\sum_b \lambda_{kb}^y f_{kb}^y(y_k) - A_k^y(\{\lambda_{kb}^y\})\right\}, \quad (9)$$

$$p(\mathbf{h}) = \prod_j \exp\left\{\sum_c \eta_{jc} g_{jc}(h_j) - B_j(\{\eta_{jc}\})\right\}, \quad (10)$$

where  $f_{ia}^x$  and  $f_{kb}^y$  are  $a$  and  $b$ th sufficient statistics of  $x_i$  and  $y_k$ .  $\lambda^x$ ,  $\lambda^y$  and  $\eta$  are parameters, and  $A^x$ ,  $A^y$ , and  $B$  are log partition functions. Given marginal distribution, joint distribution of nodes is defined straightforwardly:

$$p(\mathbf{x}, \mathbf{y}, \mathbf{h}|\theta) \propto \exp\left\{\sum_{i,a} \lambda_{ia}^x f_{ia}^x(x_i) + \sum_{k,b} \lambda_{kb}^y f_{kb}^y(y_k) + \sum_{j,c} \eta_{jc} g_{jc}(h_j) + \sum_{i,j,a,c} W_{ijac}^x f_{ia}^x(x_i) g_{jb}(h_j) + \sum_{k,j,b,c} W_{kjbc}^y f_{kb}^y(y_k) g_{jc}(h_j)\right\}. \quad (11)$$

The conditional distributions are derived as below:

$$p(\mathbf{x}|\mathbf{h}, \theta) = \prod_i \exp\left\{\sum_a \hat{\lambda}_{ia}^x f_{ia}^x(x_i) - A_i^x(\{\hat{\lambda}_{ia}^x\})\right\}, \quad (12)$$

$$p(\mathbf{y}|\mathbf{h}, \theta) = \prod_k \exp\left\{\sum_b \hat{\lambda}_{kb}^y f_{kb}^y(y_k) - A_k^y(\{\hat{\lambda}_{kb}^y\})\right\}, \quad (13)$$

$$p(\mathbf{h}|\mathbf{x}, \mathbf{y}, \theta) = \prod_j \exp\left\{\sum_c \hat{\eta}_{jc} g_{jc}(h_j) - B_j(\{\hat{\eta}_{jc}\})\right\}, \quad (14)$$

with shifted parameters

$$\hat{\lambda}_{ia}^x = \lambda_{ia}^x + \sum_{j,c} W_{ijac}^x g_{jc}(h_j), \quad (15)$$

$$\hat{\lambda}_{kb}^y = \lambda_{kb}^y + \sum_{j,c} W_{kjbc}^y g_{jc}(h_j), \quad (16)$$

$$\hat{\eta}_{jc} = \eta_{jc} + \sum_{i,a} W_{ijac}^x f_{ia}^x(x_i) + \sum_{k,b} W_{kjbc}^y f_{jb}^y(y_j). \quad (17)$$

The model can be extended to the case of more than two sets of visible nodes. With a single set of visible nodes, the model shrinks down to a EFH or RBM. To enhance discriminative performance of MWH, labeled version [19] and large-margin approach [3] were also proposed.

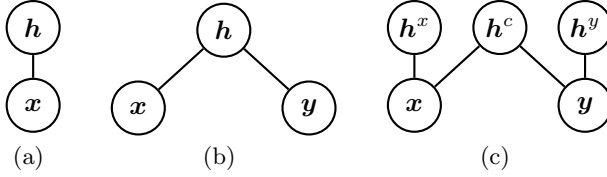
### 3 Restricted DBNs

#### 3.1 Multi-view Harmonium

Many multi-view algorithms are based on a weak assumption that all views are completely correlated. MWH also assumes that views are completely independent when the values of hidden nodes are given. However, the views are often incompletely correlated on real-world datasets. On these datasets, MWH will mix view-specific information with shared information and fail to obtain optimal representation of data. To overcome this problem, we need to model the view-specific information and shared information separately.

Recent multi-view learning algorithms including work of Salzman et al.[13] learns from partially independent multi-view data sets by separating view-specific latent variables from latent variables shared among views.

However, the limitation of existing models is evident. For example, as the FOLS framework requires a Gram matrix, so the model consumes memory storage proportional to  $N^2$ , where  $N$  is the number of training samples. The evaluation also requires computing time of complexity  $O(N)$ , as we need to calculate kernel function for every training sample and a test sample. Moreover, existing models including FOLS and group NMF requires additional parameters to control orthogonality between view-specific and shared latent spaces, and these parameters should be selected by going through computationally expensive procedures including cross-validation or by hand.



**Fig. 1.** Graphical models of (a) EFH, (b) MWH, and (c) multi-view harmonium. Repetitions of variables are not denoted for simplicity.

Taking advantage of fast learning and inference of harmonium models and the idea of separating shared and common latent variables, we extend MWH to devise a new model named *multi-view harmonium*.

This model incorporates view-specific hidden nodes  $\mathbf{h}^x$  and  $\mathbf{h}^y$  in addition to common hidden nodes  $\mathbf{h}^c$ , to model view-specific, uncorrelated information.  $\mathbf{h}^x$  and  $\mathbf{h}^y$  are only connected to their corresponding set of visible nodes  $\mathbf{x}$  and  $\mathbf{y}$  with connection weights  $\mathbf{U}^x$  and  $\mathbf{U}^y$ . Graphical representations show the difference between other models and our model (Fig. 1). Joint probability of visible and hidden nodes is as below:

$$\begin{aligned}
 p(\mathbf{x}, \mathbf{y}, \mathbf{h}^x, \mathbf{h}^y, \mathbf{h}^c | \theta) \propto & \\
 \exp \left\{ \sum_{i,a} \lambda_{ia}^x f_{ia}^x(x_i) + \sum_{k,b} \lambda_{kb}^y f_{kb}^y(y_k) + \sum_{j,c} \eta_{jc}^c g^c(h_j^c) + \sum_{m,d} \eta_{md}^x g^x(h_m^x) \right. & \\
 + \sum_{n,e} \eta_{ne}^y g^y(h_n^y) + \sum_{i,j,a,c} W_{ijac}^x f_{ia}^x(x_i) g_{jb}(h_j) + \sum_{k,j,b,c} W_{kjbc}^y f_{kb}^y(y_k) g_{jc}(h_j) & \\
 \left. + \sum_{i,m,a,d} U_{imad}^x f_{ia}^x(x_i) g_{md}^x(h_m^x) + \sum_{k,n,b,e} U_{knbe}^y f_{kb}^y(y_k) g_{ne}^y(h_n^y) \right\}. & \quad (18)
 \end{aligned}$$

Conditional distributions can easily be derived from this joint distribution.

Depending on the type of data, we can choose appropriate distributions from exponential-family. To handle continuous-valued inputs, one can use Gaussian distribution for visible nodes  $\mathbf{x}$  [7], where the conditional probability is defined as:

$$p(x_i | \mathbf{h}^x, \mathbf{h}^c, \theta) = \mathcal{N}(x_i | \mathbf{W}_i^x \mathbf{h}^c + \mathbf{U}_i^x \mathbf{h}^x + b_i^x, 1), \quad (19)$$

assuming  $x_i$  has zero mean and unit variance. Rectified linear units(ReLU) for hidden nodes is also helpful in handling continuous values [12], where value of hidden node  $h_j$  given visible node  $\mathbf{x}$  is defined as:

$$h_j = \max(0, \sum_i W_{ij} x_i + h_j + \mathcal{N}(0, 1)). \quad (20)$$

When modeling term occurrence on bag-of-words representation, we can use Poisson distribution:

$$p(x_i | \mathbf{h}^x, \mathbf{h}^c, \theta) = \text{Poisson}(x_i | \exp(\alpha_i + \mathbf{W}_i^x \mathbf{h}^x + \mathbf{U}_i^x \mathbf{h}^c)). \quad (21)$$

We train multi-view harmonium by maximizing log-likelihood. Given the joint distribution  $p(\mathbf{x}, \mathbf{y}, \mathbf{h}^x, \mathbf{h}^y, \mathbf{h}^c | \theta)$  we can integrate out the hidden units to obtain likelihood function:

$$p(\mathbf{x}, \mathbf{y} | \theta) \propto \exp \left\{ \sum_{i,a} \lambda_{ia}^x f_{ia}^x(x_i) + \sum_{k,b} \lambda_{kb}^y f_{kb}^y(y_k) - \sum_j B^c(\{\hat{\eta}_{jc}^c\}) - \sum_n B^x(\{\hat{\eta}_{md}^x\}) - \sum_n B^y(\{\hat{\eta}_{ne}^y\}) \right\}, \quad (22)$$

where  $B^x$ ,  $B^y$  and  $B^c$  are log-partition functions of marginal distribution of hidden variables  $\mathbf{h}^c$ ,  $\mathbf{h}^x$  and  $\mathbf{h}^y$  and the shifted parameters are

$$\hat{\eta}_{jc}^c = \eta_{jc}^c + \sum_{i,a} W_{ijac}^x f_{ia}^x(x_i) + \sum_{k,b} W_{kjbc}^y f_{kb}^y(y_k), \quad (23)$$

$$\hat{\eta}_{md}^x = \eta_{md}^x + \sum_{i,a} U_{imad}^x f_{ia}^x(x_i), \quad (24)$$

$$\hat{\eta}_{ne}^y = \eta_{ne}^y + \sum_{k,n,b,e} U_{knbe}^y f_{kb}^y(y_k). \quad (25)$$

To maximize expected log-likelihood over data distribution  $\mathcal{L} = \langle \log p(\mathbf{x}, \mathbf{y} | \theta) \rangle_{\bar{p}}$ , we can use gradient ascent to update parameters  $\theta$ . The derivatives on parameters related to  $\mathbf{x}$  are given as follows:

$$\frac{\partial \mathcal{L}}{\partial W_{ijac}^x} = \langle f_{ia}(x_i) B'_{jc}(\hat{\eta}_{jc}^c) \rangle_{\bar{p}} - \langle f_{ia}(x_i) B'_{jc}(\hat{\eta}_{jc}^c) \rangle_p, \quad (26)$$

$$\frac{\partial \mathcal{L}}{\partial U_{imad}^x} = \langle f_{ia}(x_i) B'_{md}(\hat{\eta}_{md}^x) \rangle_{\bar{p}} - \langle f_{ia}(x_i) B'_{md}(\hat{\eta}_{md}^x) \rangle_p, \quad (27)$$

$$\frac{\partial \mathcal{L}}{\partial \lambda_{ia}^x} = \langle f_{ia}^x(x_i) \rangle_{\bar{p}} - \langle f_{ia}^x(x_i) \rangle_p, \quad (28)$$

$$\frac{\partial \mathcal{L}}{\partial \eta_{jc}^c} = \langle B'_{jc}(\hat{\eta}_{jc}^c) \rangle_{\bar{p}} - \langle B'_{jc}(\hat{\eta}_{jc}^c) \rangle_p, \quad (29)$$

$$\frac{\partial \mathcal{L}}{\partial \eta_{md}^x} = \langle B'_{md}(\hat{\eta}_{md}^x) \rangle_{\bar{p}} - \langle B'_{md}(\hat{\eta}_{md}^x) \rangle_p. \quad (30)$$

where  $\langle \cdot \rangle_{\bar{p}}$  and  $\langle \cdot \rangle_p$  are expectation over data distribution and model distribution, and  $B'_{jc} = \partial B_{jc} / \partial \eta_{jc}^c$  and  $B'_{md} = \partial B_{md}^x / \partial \eta_{md}^x$  are partial derivatives of log-partition functions. Derivatives related to  $\mathbf{y}$  can be derived similarly. To calculate the expectations from data and model distribution, we need samples from the distributions. However, as we cannot directly sample from the joint distribution of visible nodes and hidden nodes, we use Gibbs sampling and sample from conditional distributions in each Gibbs steps. Values of each node can be sampled from conditional distribution in a single step when other layer's nodes are given.

However, we need infinite steps of Gibbs sampling to calculate exact model distribution. To avoid the problem, we may approximate log-likelihood by using contrastive divergence (CD) learning instead [5]. The key idea of CD learning is

---

**Algorithm 1.** Training multi-view harmonium using CD learning with mini-batch size  $K$  and  $M$  Gibbs steps

---

```

1: Input:  $\mathcal{D} = \{(\mathbf{x}_t, \mathbf{y}_t)\}$ 
2: Output:  $\theta = \{W^x, W^y, U^x, U^y, \eta^x, \eta^y, \eta^c, \lambda^x, \lambda^y\}$ 
3: procedure CDLEARNING( $\mathcal{D}$ )
4:   Randomly initialize parameters  $\theta$ .
5:   repeat
6:     Pick  $K$  samples as a mini-batch
7:     repeat
8:       Sample  $\tilde{\mathbf{h}}^x \sim p(\mathbf{h}^x | \mathbf{x}_t, \theta)$ ,  $\tilde{\mathbf{h}}^y \sim p(\mathbf{h}^y | \mathbf{y}_t, \theta)$ .
9:       Sample  $\tilde{\mathbf{h}}^c \sim p(\mathbf{h}^c | \mathbf{x}_t, \mathbf{y}_t, \theta)$ .
10:      repeat
11:        Sample  $\mathbf{x} \sim p(\mathbf{x} | \mathbf{h}^x, \mathbf{h}^c, \theta)$ .
12:        Sample  $\mathbf{y} \sim p(\mathbf{y} | \mathbf{h}^y, \mathbf{h}^c, \theta)$ .
13:        Sample  $\mathbf{h}^x \sim p(\mathbf{h}^x | \mathbf{x}, \theta)$ .
14:        Sample  $\mathbf{h}^y \sim p(\mathbf{h}^y | \mathbf{y}, \theta)$ .
15:        Sample  $\mathbf{h}^c \sim p(\mathbf{h}^c | \mathbf{x}, \mathbf{y}, \theta)$ .
16:      until  $M$  steps
17:    until for all  $(\mathbf{x}_t, \mathbf{y}_t)$  in mini-batch.
18:    update  $\theta$  using gradient ascent on  $\mathcal{L}$ .
19:  until  $\theta$  is converged
20: end procedure

```

---

to initialize Gibbs chain to the training data, and to run just a few Gibbs steps instead of infinite steps to approximate model distribution. Even running just a single Gibbs step is known to work well in practice.

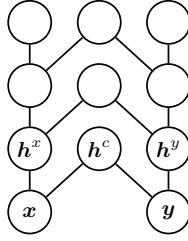
Training harmoniums using CD learning requires a gradient calculated over whole training set, and this is often a time consuming task. Instead, we can pick mini-batches from training set, and calculate gradient over mini-batches to save time and space required for training. Summary of procedure for training multi-view harmonium using CD learning is shown in Algorithm 1.

### 3.2 Restricted DBN

Introducing view-specific hidden nodes brings our model capability of modeling partial correlation, but view-specific hidden nodes also enables us to easily extend our model to form a deep network. Values of view specific hidden layer nodes will be fed as data for upper layer of deep network (Fig. 2). By forming deep network in this manner, the information considered to be "uncorrelated" by limited representation power of current layer will be sent to upper layer and view-to-view correlation will be further analyzed.

Training a deep network using back-propagation is known to be often inefficient. Therefore, training an RDBN is done in greedy, layer-wise manner [6]. We train each layer of multi-view harmoniums using algorithm 1, starting from the bottom layer. After training each layer is finished, a harmonium on the next layer uses samples of view-specific hidden nodes from the trained harmonium





**Fig. 2.** The graphical model of RDBN. View-specific hidden nodes act as new visible nodes for upper layer of the deep network model.

as its training set. Training of RDBN ends when the procedure reaches the top layer.

### 3.3 Inferring One View from the Other

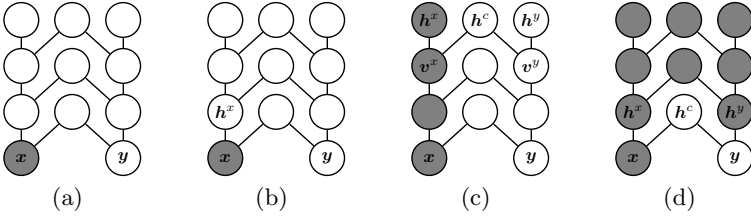
Inferring values of unobserved views from other observed views on RDBN is not as easy as inferring the values from two-layer MWH. Instead of using Variational approximation, we choose to perform Gibbs sampling in a systematic manner. The procedure can be divided into three steps, and each step is executed layer by layer (Fig. 3).

1. In the first step, we start from the bottom layer. Given observed values  $\mathbf{x}$ , the hidden nodes of observed view  $\mathbf{h}^x$  are only nodes that whose probability can be exactly calculated. Therefore we sample the hidden nodes  $\mathbf{h}^x$  and proceed to the upper layer to repeat this procedure (Fig. 3-b).
2. When the first step reaches the top layer, we run Gibbs sampling for top two layers, while nodes of observed views  $\mathbf{v}^x, \mathbf{h}^x$  are fixed to values determined on the first step. In other words, we sample from  $p(\mathbf{v}^y, \mathbf{h}^c, \mathbf{h}^y | \mathbf{v}^x, \mathbf{h}^x, \theta)$ . Then we get the values for top two layers of unobserved views  $\mathbf{h}^c, \mathbf{v}^y$ . We use average of repeatedly sampled values for the unobserved nodes (Fig. 3-c).
3. Finally, we move to lower layers. Given observed nodes  $\mathbf{x}$  and  $\mathbf{h}^x$ , and hidden units of unobserved views sampled in previous steps  $\mathbf{h}^y$ , we sample from  $p(\mathbf{y}, \mathbf{h}^c | \mathbf{x}, \mathbf{h}^x, \mathbf{h}^y, \theta)$ . Repeating the process until we reach bottom layer gives us the values of unobserved views (Fig. 3-d).

When common hidden nodes  $\mathbf{h}^c$  of bottom layer of deep network are binary logistic nodes, we can skip Gibbs sampling and directly evaluate  $p(\mathbf{y} | \mathbf{x}, \theta)$  up to a normalization constant by rearranging terms:

$$\begin{aligned}
 p(\mathbf{y} | \mathbf{x}, \theta) \propto \prod_j \left[ 1 + \exp \left\{ \sum_{i,a} (W_{ija}^x f_a^x(x_i) + \lambda_{ia}^x f_a^x(x_i)) \right. \right. \\
 \left. \left. + \sum_{k,b} ((W_{kjb}^y + U_{knb}^y) f_b^y(y_k) + \lambda_{kb}^y f_b^y(y_k)) + c_j^c \right\} \right], \quad (31)
 \end{aligned}$$

For other layers of RDBN, we use Gibbs sampling instead. Empirically, about 10 times of Gibbs sampling for each layer gave a fair result.



**Fig. 3.** Inferring procedure of RDBN. Variables with known values at current step is marked as a shaded circle.

## 4 Numerical Experiments

In this section, we compare our model to MWH. We also compare our method with FOLS-GPLVM [13], and shared GPLVM(SGPLVM) [14], which are directed, non-parametric multi-view latent variable models based on Gaussian process latent variable model [9]. We considered CCA as a baseline too. On a synthetic dataset, we show the effect of modelling view-specific information. Then we run experiments on widely-used NORB-small dataset, and ESL photo dataset in a view-to-view conversion task to numerically compare performance of RDBN and MWH and GPLVM-based models. We used Gibbs sampling for inference on the graphical models.

### 4.1 Synthetic Example

To show the effectiveness of additional hidden nodes, we performed an experiment taken from recent work of Salzman et al. [13]. We constructed a synthetic partially correlated multi-view dataset. To illustrate common and view-specific latent variables, we used sinusoidal functions of  $\mathbf{t}$  with different phases and frequencies:

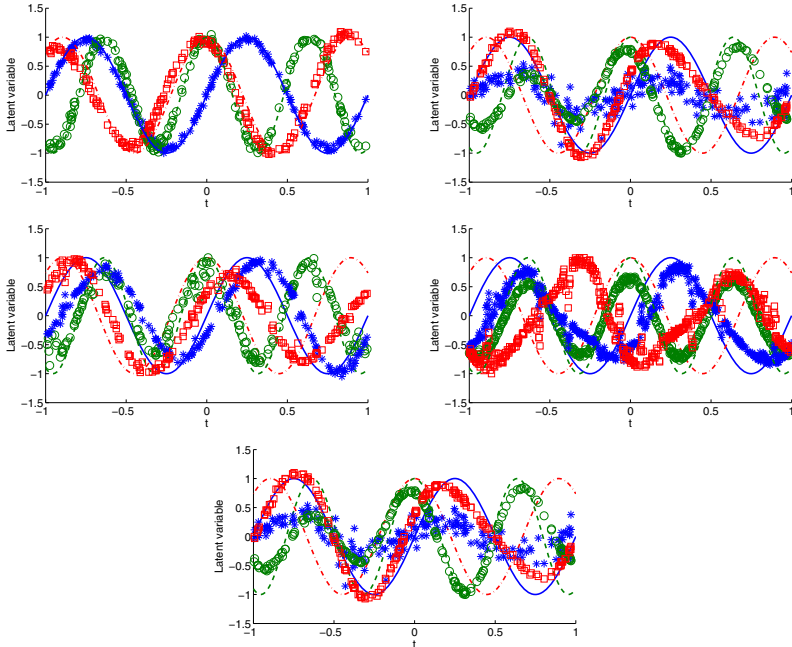
$$\begin{aligned} \mathbf{x} &= \sin(2\pi\mathbf{t}), & \mathbf{z}^1 &= \cos(\pi\pi\mathbf{t}), & \mathbf{z}^2 &= \sin(\sqrt{5}\pi\mathbf{t}), \\ \mathbf{m}^1 &= [\mathbf{x}, \mathbf{z}^1], & \mathbf{m}^2 &= [\mathbf{x}, \mathbf{z}^2]. \end{aligned}$$

We randomly projected  $\mathbf{m}^1$  and  $\mathbf{m}^2$  to 20 dimensional space and added independent Gaussian noise of variance 0.01 and correlated noise  $0.02 \sin(3.6\pi\mathbf{t})$  to obtain our final multi-view synthetic dataset.

A multi-view harmonium, which is a two-layer RDBN, with 1 common hidden node and 1 view-specific hidden node is used, and we used a MWH and CCA for comparison. 3 hidden nodes for MWH and projection to 3-dimensional space for CCA was chosen for fair comparison. We used Gaussian nodes and ReLU for visible and hidden nodes.

Training was done using 4000 training samples and we calculated hidden node activations from 500 test samples. We checked the correspondence between ground-truth latent variables, and obtained hidden node activations.

MWH, SGPLVM and CCA failed to infer latent variables correctly. No hidden node activations corresponded to latent variables used to generate data. In the



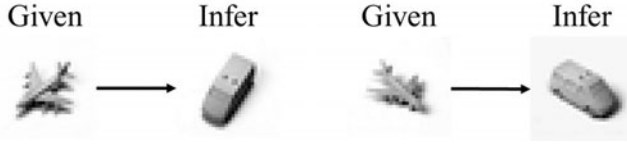
**Fig. 4.** Latent variables generated by multi-view harmonium, MWH, FOLS, Shared GPLVM and CCA (starting from top-left to bottom-right). Outputs of methods are normalized to have zero mean and maximum absolute value to be 1. Latent variables used to generate data are depicted as lines: Blue solid line:  $\mathbf{x}$ , green dotted line:  $\mathbf{z}^1$ , and red dash-dotted line:  $\mathbf{z}^2$ . Each dimension of outputs are depicted as markers with different colors and shapes.

other hand, our model, multi-view harmonium, found common and view-specific latent variables by its hidden node activations 4. Common hidden nodes and view-specific hidden node activations exactly corresponded to the common and view-specific latent variables. Although FOLS-GPLVM also discovered latent variables correctly, the result was not as accurate as the result of our model. Moreover, our model was able to separate common and view-specific information without any help of additional scheme, while FOLS had to minimize mutual information explicitly.

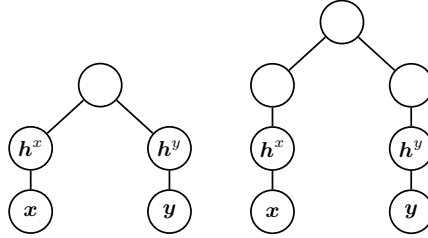
## 4.2 Object Conversion on NORB-Small

Training set of NORB-small dataset [10] contains 24,300 stereo images of 5 generic categories of objects on a white background, imaged under various lighting conditions and camera angles. Each category contains 5 instances of objects. For example, the category 'car' contains pictures of 5 different cars.

To evaluate view-to-view conversion performance, we constructed a multi-view dataset by taking pairs of images with same lighting and camera angles, but different category of objects (airplanes - cars), that means each sample in



**Fig. 5.** Graphical representation of view to view conversion task on NORB-small dataset



**Fig. 6.** Graphical model of 3 and 4-layer MWHs used in experiments on NORB-small dataset and ESL dataset

the dataset contains a image of airplane and a car, taken with same condition. 1200 image pairs were used as a training set and 3600 image pairs were used as a test set. After training, model was asked to generate a car image with same imaging condition from a given airplane image (Fig. 5). Images were resized to the size  $32 \times 32$ , and pixels on each position of an image were normalized to have zero mean and variance 1. The variance of pixel values was calculated across whole training set on the same position on images. The root mean squared error (RMSE) was computed:

$$\text{RMSE} = \sqrt{\frac{1}{N_{\text{samples}}} \sum_t \frac{1}{N_{\text{pixels}}} \|\mathbf{x}_t^{\text{truth}} - \mathbf{x}_t^{\text{rec}}\|^2}, \quad (32)$$

where  $N_{\text{samples}}$  and  $N_{\text{pixels}}$  are the number of samples and of pixels, respectively, and  $\mathbf{x}_t^{\text{truth}}$  and  $\mathbf{x}_t^{\text{rec}}$  are ground-truth and reconstructed images, respectively.

For our experiment, we constructed 2, 3, and 4-layer RDBN. For two layers from the bottom, we used Gaussian and ReLU, to handle continuous valued data. Logistic binary nodes were used for remaining two layers.

We also constructed a deep network model for MWH by attaching harmoniums under each views of MWH for a fair comparison with RDBN. For example, 4-layer MWH was constructed by attaching additional two layers of Harmoniums under each view of a MWH (Fig. 6). Inferring an unobserved view was done by the procedure below.

1. Sample nodes of observed view layer by layer from the bottom.
2. Infer the unobserved view of MWH on the top layer.
3. Sample unobserved nodes, from top to bottom layer.

**Table 1.** Averaged RMS reconstruction errors of RDBN, MWH models and random reconstruction in the view-to-view conversion task on NORB-small dataset

Method	RMSE	Method	RMSE
4-layer RDBN	<b>0.0923</b>	4-layer MWH	0.1379
3-layer RDBN	0.0936	3-layer MWH	0.1377
2-layer RDBN	0.0958	2-layer MWH	0.0973
FOLS-GPLVM	0.2489	SGPLVM	0.2089
Random	0.5205		

We also added additional hidden nodes to MWH to ensure same number of parameters are used. For example, if a RDBN had 200 view-specific hidden nodes and 800 common hidden nodes, each view of corresponding layer of a MWH had  $200 + 800 = 1000$  hidden nodes. We also tried 2-layer and 3-layer RDBN to see the effect of number of layers. Each network was trained 200 iterations with learning rate of 0.001. Instead of full batch learning, we used mini-batches with 100 samples. SGPLVM and FOLS-GPLVM were trained with 200 iterations, and their latent dimension was automatically decided by a heuristic method of Neil Lawrence’s software. As a baseline, we used images with random pixel values as reconstruction images and compared the reconstruction errors with the results of RDBNs and MWH models.

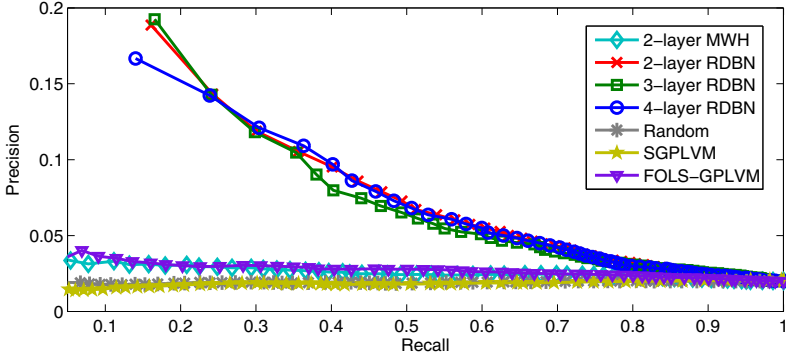
The experimental result showed remarkable difference between RDBN and MWH algorithm. 4-layer RDBN showed the smallest error than any configuration of MWH, FOLS-GPLVM and SGPLVM. Moreover, adding each additional layer to RDBN reduced reconstruction error, while adding layer to MWH damaged the performance of the model (Table 1).

### 4.3 Image Annotation on ESL Photo Dataset

To examine the capability of our deep model to find nontrivial relation between views, we applied our model to the task of image annotation. The experiment was done using ESL photo dataset [15], which contains 3464 photos collected from Flickr, annotated with 59 different tags. We calculated autocorrelogram of images with the setting of radius 6 and 32 colors on HSV color space. By this process we obtained a 192-dimensional feature vector for each image. We also calculated TF-IDF from our tag data. We used 2000 images for training, and 1044 images testing.

We compared 2, 3, 4-layer RDBNs, 2-layer MWH, sGPLVM, and FOLS-GPLVM. 30, 40, 60 hidden units were used for each hidden layer of RDBNs, and GPLVM models were trained with 200 iterations. Image autocorrelogram and TF-IDF of tags were assigned to each view of the models. On visible layer of the models, we used Gaussian nodes for images and Poisson nodes for tags. Each model was trained 500 iterations with learning rate 0.01. Again, we used mini-batches with 100 samples in training. We also tried annotating images with random tags as a baseline.

Precision-recall curve and mean average precision (MAP) were used as evaluation measure. MAP is defined as the mean of average precision (AP) over test



**Fig. 7.** Precision-recall curves for RDBNs, 2-layer MWH, SGPLVM, FOLS-GPLVM and random annotation. Vertical axis of the plot is trimmed for better comparison between results of RDBNs with different number of layers.

**Table 2.** The mean average precision of RDBNs, a MWH, SGPLVM, FOLS-GPLVM and random annotation in image annotation task on ESL dataset.

Method	MAP	Method	MAP
4-layer RDBN	<b>0.0620</b>	2-layer MWH	0.0251
3-layer RDBN	0.0598	SGPLVM	0.0180
2-layer RDBN	0.0565	FOLS-GPLVM	0.0267
Random	0.0190		

samples:

$$\text{MAP} = \frac{\sum_t \text{AP}(\mathbf{x}_t)}{N_{\text{samples}}} = \frac{1}{N_{\text{samples}}} \sum_t \sum_{j=1}^{N_{\text{tags}}} \text{Precision}(\mathbf{x}_t, j), \quad (33)$$

where  $N_{\text{samples}}$  and  $N_{\text{tags}}$  are the number of samples and of tags, respectively, and  $\text{Precision}(\mathbf{x}_t, j)$  denotes the precision when top  $j$  tags are chosen, given the query image  $\mathbf{x}_t$ . MAP is also calculated by area under precision-recall curve. We used this alternative definition of MAP to calculate the results reported here.

Precision-recall curve shows that RDBN algorithms clearly outperformed other methods, showing the advantage of modelling view-specific information separately. Although the difference was not very significant, deep models showed higher average precision than 2-layer RDBN, showing the benefit of using deep models. In contrast, average precision of 2-layer MWH and SGPLVM were not very distinguishable from the result of random annotation (Table 2). Precision-recall curve also shows the difference between RDBNs and other models (Fig. 7).

## 5 Conclusions

In this paper, we have proposed the harmonium-based model that uses view-specific hidden nodes for multi-view data, to capture the view-specific

information of the data separated from shared information among views. Moreover, our model can be naturally extended to deep network to model more complex relationship between views by using view-specific hidden nodes as inputs to next layer. We have demonstrated the effectiveness of our approach by comparing our model to MWH, and their directed, non-parametric counterparts including SGPLVM and FOLS-GPLVM. the existing harmonium-based multi-view model, on various experiments on synthetic and real-world examples. And we showed significant improvement over other methods in the tested examples. In the future, we plan to investigate the modifications of RDBN model such as conditional RDBN to directly model conditional distribution between views. We also plan to extend our model to use label information for supervised and semi-supervised learning. Another future topic is extension of our model to time-series data or other structured data.

**Acknowledgments.** This work was supported by the Converging Research Center Program funded by the Ministry of Education, Science, and Technology (No. 2010K001171), NIPA ITRC Support Program (NIPA-2011-C1090-1131-0009), NIPA-MSRA Creative IT/SW Research Project, and NRF WCU Program (R31-10100).

## References

1. Bengio, Y.: Learning deep architectures for AI. *Foundations and Trends<sup>®</sup> in Machine Learning* 2(1), 1–127 (2009)
2. Blum, A., Mitchell, T.: Combining labeled and unlabeled data with co-training. In: *Proceedings of the Annual Conference on Learning Theory (COLT)*, Madison, Wisconsin (1998)
3. Chen, N., Zhu, J., Xing, E.P.: Predictive subspace learning for multi-view data: A large margin approach. In: *Advances in Neural Information Processing Systems (NIPS)*, vol. 23. MIT Press, Cambridge (2010)
4. Choi, H., Choi, S., Choe, Y.: Manifold integration with Markov random walk. In: *Proceedings of the AAAI National Conference on Artificial Intelligence (AAAI)*, Chicago, IL (2008)
5. Hinton, G.E.: Training products of experts by minimizing contrastive divergence. *Neural Computation* 14(8), 1771–1800 (2002)
6. Hinton, G.E., Osindero, S., Teh, Y.W.: A fast learning algorithm for deep belief nets. *Neural Computation* 18(7), 1527–1554 (2006)
7. Hinton, G.E., Salakhutdinov, R.: Reducing the dimensionality of data with neural networks. *Science* 313(5786), 504–507 (2006)
8. Hotelling, H.: Relations between two sets of variates. *Biometrika* 28, 312–377 (1936)
9. Lawrence, N.: Gaussian process latent variable models for visualisation of high dimensional data. In: *Advances in Neural Information Processing Systems (NIPS)*, vol. 16. MIT Press, Cambridge (2004)
10. LeCun, Y., Huang, F.J., Bottou, L.: Learning methods for generic object recognition with invariance to pose and lighting. In: *Proceedings of the IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, Washington, DC (2004)

11. Lee, H., Choi, S.: Group nonnegative matrix factorization for EEG classification. In: Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS), Clearwater Beach, Florida (2009)
12. Nair, V., Hinton, G.E.: Rectified linear units improve restricted Boltzmann machines. In: Proceedings of the International Conference on Machine Learning (ICML), Haifa, Israel (2010)
13. Salzmann, M., Ek, C.H., Urtasun, R., Darrell, T.: Factorized orthogonal latent spaces. In: Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS), Sardinia, Italy (2010)
14. Shon, A.P., Grochow, K., Hertzmann, A., Rao, R.P.N.: Learning shared latent structure for image synthesis and robotic imitation. In: Advances in Neural Information Processing Systems (NIPS), vol. 17. MIT Press, Cambridge (2006)
15. Sinha, P., Jaini, R.: Classification and annotation of digital photos using optical context data. In: Proceedings of the ACM International Conference on Image and Video Retrieval (CIVR), Niagara Falls, Canada (2008)
16. Smolensky, P.: Information processing in dynamical systems: Foundations of harmony theory. In: Rumelhart, D.E., McClelland, J.L. (eds.) *Parallel Distributed Processing: Explorations in the Microstructures of Cognition: Foundations*, vol. 1, pp. 194–281. MIT Press, Cambridge (1986)
17. Welling, M., Rosen-Zvi, M., Hinton, G.: Exponential family harmoniums with an application to information retrieval. In: Advances in Neural Information Processing Systems (NIPS), vol. 17. MIT Press, Cambridge (2005)
18. Xing, E.P., Yan, R., Hauptmann, A.G.: Mining associated text and images with dual-wing harmonium. In: Proceedings of the Annual Conference on Uncertainty in Artificial Intelligence (UAI), Edinburgh, UK (2005)
19. Yang, J., Liu, Y., Xing, E.P., Hauptmann, A.G.: Harmonium models for semantic video representation and classification. In: Proceedings of the SIAM International Conference on Data Mining (SDM), Minneapolis, MN (2007)