

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/326463695>

Exploring Helmholtz Machine and Deep Belief Net in the Exponential Family Perspective

Conference Paper · July 2018

CITATION

1

READS

96

2 authors, including:



Yifeng Li

National Research Council Canada

64 PUBLICATIONS 460 CITATIONS

SEE PROFILE

Some of the authors of this publication are also working on these related projects:



Network-Based Prediction [View project](#)

Exploring Helmholtz Machine and Deep Belief Net in the Exponential Family Perspective

Yifeng Li¹ Xiaodan Zhu²

Abstract

Applications of directed deep generative models are still limited by the data types they can model. Here, we formulate the free energy function of exponential family restricted Boltzmann machine, extend the binary Helmholtz machine and deep belief network to the exponential family, and derive the corresponding wake-sleep learning algorithms. We demonstrate that appealing performance can be achieved by the generalized models.

1. Introduction

Generative models with latent variables concern joint distributions of observable measurements and hidden factors. While generative models, such as restricted Boltzmann machine (RBM) (Smolensky, 1986), are often employed to train deep discriminative models (Bengio et al., 2006), they also allow for addressing many more interesting problems. For example, sampling joint distribution $p(\mathbf{x}, \mathbf{h})$ to generate novel data could be useful for a wide variety of problems that may even reach some hard AI tasks, e.g., those involving creativity. A joint distribution also allows us to see association of variables as studied in causality. From the conditional $p(\mathbf{h}|\mathbf{x})$, states of hidden variables can be inferred, which are often used in dimensionality reduction, missing value estimation, prediction, and clustering. From $p(\mathbf{x}|\mathbf{h})$, effects of hidden factors can be evaluated.

Deep generative neural networks inherit this potential versatility and have the distinctive capacity of modeling phenomena and systems of high complexity. For example, undirected generative models, e.g., RBM, have been extended to deep networks, e.g., deep Boltzmann machine (DBM) (Salakhutdinov & Hinton, 2009a) and multi-modal

DBM (MDBM) (Srivastava & Salakhutdinov, 2014), to capture the complex nature of data and association of heterogeneous modalities. Also, directed deep generative models (DGMs), e.g., Helmholtz machine (HM) (Dayan et al., 1995) and deep belief net (DBN) (Hinton et al., 2006; Srivastava & Salakhutdinov, 2012), are used to model complicated (hierarchical or sequential) structures of hidden variables. These DGMs, however, face two major challenges: intractable inference $[p(\mathbf{h}|\mathbf{x})]$ and limited forms of distributions. In undirected DGMs, hidden states could be approximated using variational approximation inference. In directed DGMs, hidden states may be obtained through recognition connections. Nevertheless, applications of existing DGMs are, to a large extent, restricted by the data types they can fit. For example, DBN is often applied to pretrain feedforward networks for classification problems when data are binary or can be scaled to the range of $[0, 1]$. Generalizing DGMs to a larger pool of distributions in a united framework is of interest. Thanks to the exponential family RBM (exp-RBM) (Welling et al., 2005), the extension of binary DBM to the exponential family is straightforward. However, the generalization of binary wake-sleep algorithms (Hinton et al., 1995; 2006) for exponential family HM and DBN is nontrivial and remains a puzzle, thus is revisited here.

In summary, this paper addresses three issues: how to generically define deep generative models? how to derive the stochastic gradients in exponential family wake-sleep algorithms? and how to estimate their (lower bounds of) log-likelihoods? The major contributions of this work are threefold: (1) we discover that the free energy function of exp-RBM and the partition function of its base-rate model can be analytically computed, which enables the estimation of partition functions and (lower bounds of) likelihoods of DGMs; (2) inspired by the concepts of exp-RBMs, we formulate the exponential family HM and DBN (exp-HM and exp-DBN); and (3) we derive the update rules in the wake-sleep algorithms by virtue of an important property of exponential family distributions.

2. Exponential Family RBM

In addition to pretraining DGMs, the conditional distributions of RBM also play an important role in the modelling

¹Digital Technologies Research Centre, National Research Council Canada, Ottawa, Ontario, Canada ²Department of Electrical and Computer Engineering, Queen's University, Kingston, Ontario, Canada. Correspondence to: Yifeng Li <yifeng.li@nrc-nrc.gc.ca>.

of directed deep generative models. Prior to presenting our generalization, we briefly discuss key concepts of exponential family RBM, where we propose that (i) the free energy function can be formulated analytically and (ii) the log-partition function of the base-rate model can be computed analytically, enabling the estimation of partition function in exp-RBM using annealed importance sampling (AIS).

The exponential family includes univariate or multivariate distributions with the following natural parametric form:

$$p(\mathbf{x}) = h(\mathbf{x})e^{\boldsymbol{\theta}^T \mathbf{s}(\mathbf{x}) - A(\boldsymbol{\theta})}, \quad (1)$$

where \mathbf{x} is either a univariate or multivariate random variable, $\boldsymbol{\theta}$ is a vector (or a scalar) of natural parameters, $\mathbf{s}(\mathbf{x})$ is a vector (or a scalar) of sufficient statistics, $A(\boldsymbol{\theta})$ is the log-partition function, and $h(\mathbf{x})$ is the base measure. It has several interesting properties. First, if Eq. (1) is rewritten to $p(\mathbf{x}) = \frac{1}{Z} e^{\boldsymbol{\theta}^T \mathbf{s}(\mathbf{x}) + \log h(\mathbf{x})}$, one will find that $\log Z = A(\boldsymbol{\theta})$, i.e. $A(\boldsymbol{\theta})$ normalizes the distribution:

$$A(\boldsymbol{\theta}) = \log \left[\sum_{\mathbf{x}} h(\mathbf{x}) e^{\boldsymbol{\theta}^T \mathbf{s}(\mathbf{x})} \right]. \quad (2)$$

This property about $A(\boldsymbol{\theta})$ is extremely useful in the estimation of the partition functions of exp-RBMs. Second, the first and second order derivatives of $A(\boldsymbol{\theta})$ equal to the mean and covariance of $\mathbf{s}(\mathbf{x})$, respectively:

$$\frac{\partial A(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \mathbb{E}_{p(\mathbf{x})}[\mathbf{s}(\mathbf{x})]; \quad (3)$$

$$\frac{\partial^2 A(\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} = \text{Cov}_{p(\mathbf{x})}[\mathbf{s}(\mathbf{x})]. \quad (4)$$

$A(\boldsymbol{\theta})$ is convex, as the covariance matrix is positive semi-definite. Eq. (3) is extremely useful in deriving stochastic update rules for exp-HMs (see Section 3.2) and exp-DBNs. For the understanding of this article, we selectively provide the natural and standard forms of a range of univariate member of exponential family in Supplemental Table 1.

Exp-RBM (Welling et al., 2005) is generally formulated as

$$p(\mathbf{x}, \mathbf{h}) = \frac{1}{Z} e^{-E(\mathbf{x}, \mathbf{h})}, \quad (5)$$

where $\mathbf{x} \in \mathbb{R}^M$ is the vector of visible random variables, $\mathbf{h} \in \mathbb{R}^K$ the vector of hidden random variables, $E(\mathbf{x}, \mathbf{h})$ the energy function, and $Z = \sum_{\mathbf{x}} \sum_{\mathbf{h}} e^{-E(\mathbf{x}, \mathbf{h})}$ the partition function. The design of an exp-RBM requires three key components: (1) base distributions of \mathbf{x} and \mathbf{h} , (2) energy function $E(\mathbf{x}, \mathbf{h})$, and (3) conditional distributions $p(\mathbf{x}|\mathbf{h})$ and $p(\mathbf{h}|\mathbf{x})$. They are elaborated below.

First, the base distributions for both \mathbf{x} and \mathbf{h} in the natural forms of exponential class can be defined as:

$$\begin{cases} p(\mathbf{x}) = \prod_{m=1}^M e^{\mathbf{a}_m^T \mathbf{s}_m + \log f_m(x_m) - A_m(\mathbf{a}_m)} \\ p(\mathbf{h}) = \prod_{k=1}^K e^{\mathbf{b}_k^T \mathbf{t}_k + \log g_k(h_k) - B_k(\mathbf{b}_k)}. \end{cases} \quad (6) \quad (7)$$

Here, \mathbf{a}_m and \mathbf{b}_k are vectors/scalars of natural parameters for x_m and h_k , respectively; \mathbf{s}_m and \mathbf{t}_k are the corresponding vectors/scalars of sufficient statistics; $A_m(\mathbf{a}_m)$ and $B_k(\mathbf{b}_k)$ are the corresponding log-partition functions; $f_m(x_m)$ and $g_k(h_k)$ are the corresponding base measures. For example, if $p(\mathbf{x})$ is Gaussian with unknown mean and precision, then $\mathbf{a}_m = [a_m^{(1)}, a_m^{(2)}]^T = [\mu_m \beta_m, -\frac{\beta_m}{2}]^T$ and $\mathbf{s}_m = [s_m^{(1)}, s_m^{(2)}]^T = [x_m, x_m^2]^T$. If $p(\mathbf{h})$ is Bernoulli, then $b_k = \log \frac{p_k}{1-p_k}$ and $t_k = h_k$ are scalars.

Second, the energy function is defined via combining x_m - and h_k -related terms in Eq. (6) and (7) (including base measures), and enabling interactions between \mathbf{s}_m and \mathbf{t}_k :

$$E(\mathbf{x}, \mathbf{h}) = - \sum_{m=1}^M (\mathbf{a}_m^T \mathbf{s}_m + \log f_m(x_m)) - \sum_{k=1}^K (\mathbf{b}_k^T \mathbf{t}_k + \log g_k(h_k)) - \sum_{m=1}^M \sum_{k=1}^K \sum_{r=1}^R \sum_{u=1}^U s_{m,r} w_{m,k,r,u} t_{k,u}, \quad (8)$$

where R and U are the numbers of natural parameters for $p(x_m)$ and $p(h_k)$. The interaction strengths are represented by a tensor $\mathbf{W} \in \mathbb{R}^{M \times K \times R \times U}$. For Gaussian-Bernoulli RBMs (exp-RBMs with Gaussian visible and Bernoulli hidden variables), the weight tensor is of size $M \times K \times 2$, including $\mathbf{W}^{(1)}$ for interactions between \mathbf{x} and \mathbf{h} , and $\mathbf{W}^{(2)}$ for \mathbf{x}^2 and \mathbf{h} . To reduce model complexity, we may consider dropping off some interactions between the sufficient statistics. For instance, in Gaussian-Bernoulli RBMs, we may decide to disregard $\mathbf{W}^{(2)}$, reducing the weights in a tensor to a matrix of size $M \times K$. In fact, we only considered interactions between \mathbf{x} and \mathbf{h} in practice, reducing the interaction term to $\mathbf{x}^T \mathbf{W} \mathbf{h}$.

Third, from the energy function in Eq. (8) and the definition of exponential family in Eq. (1), $p(\mathbf{x}|\mathbf{h})$ can be derived as

$$\begin{aligned} p(\mathbf{x}|\mathbf{h}) &= \frac{p(\mathbf{x}, \mathbf{h})}{p(\mathbf{h})} = e^{\sum_{m=1}^M (\hat{\mathbf{a}}_m^T \mathbf{s}_m + \log f_m(x_m) - A_m(\hat{\mathbf{a}}_m))} \\ &= \prod_{m=1}^M p(x_m|\mathbf{h}, \boldsymbol{\eta}(\hat{\mathbf{a}}_m)), \end{aligned} \quad (9)$$

where $\hat{\mathbf{a}}_m = \mathbf{a}_m + \sum_{k=1}^K \mathbf{W}_{m,k,:} \mathbf{t}_k$, and function $\boldsymbol{\eta}(\hat{\mathbf{a}}_m)$ maps the natural parameters in $\hat{\mathbf{a}}_m$ to the standard forms. As shown in Supplemental Table 1, if $p(x_m|\mathbf{h})$ is Gaussian with unknown mean and precision, we then have $\boldsymbol{\eta}(\hat{\mathbf{a}}_m) = [\hat{\mu}_m, \hat{\beta}_m]^T = [-\frac{\hat{a}_m^{(1)}}{2\hat{a}_m^{(2)}}, -2\hat{a}_m^{(2)}]^T$. Similarly, we can obtain

$$p(\mathbf{h}|\mathbf{x}) = \prod_{k=1}^K p(h_k|\mathbf{x}, \boldsymbol{\eta}(\hat{\mathbf{b}}_k)), \quad (10)$$

where $\hat{\mathbf{b}}_k = \mathbf{b}_k + \sum_{m=1}^M (\mathbf{W}_{m,k,:})^T \mathbf{s}_m$ and $\boldsymbol{\eta}(\hat{\mathbf{b}}_k)$ maps the natural parameters in $\hat{\mathbf{b}}_k$ to the standard ones. For example, if $p(h_k|\mathbf{x})$ is Bernoulli, we have $\boldsymbol{\eta}(\hat{\mathbf{b}}_k) = \hat{p}_k = \sigma(\hat{b}_k)$. Importantly, the conditionals are decomposable and follow the same distributions as the bases, but use posterior parameters to reflect influence of their heterogeneous counterparts.

Stochastic gradient descent can be used to estimate exp-RBMs' parameters which can be represented by $\theta = \{\mathbf{a}, \mathbf{b}, \mathbf{W}\}$ where $\mathbf{a} = \{\mathbf{a}^{(1)}, \dots, \mathbf{a}^{(R)}\}$ and $\mathbf{b} = \{\mathbf{b}^{(1)}, \dots, \mathbf{b}^{(U)}\}$ may respectively include several bias vectors, and $\mathbf{W} = \{\mathbf{W}^{(1,1)}, \dots, \mathbf{W}^{(r,u)}, \dots, \mathbf{W}^{(R,U)}\}$ (but only interactions between \mathbf{x} and \mathbf{h} are considered in practice). The likelihood $p(\mathbf{x})$ can be obtained by marginalizing out \mathbf{h} :

$$p(\mathbf{x}) = \int_{\mathbf{h}'} p(\mathbf{x}, \mathbf{h}') = \frac{1}{Z} e^{-F(\mathbf{x})}, \quad (11)$$

where $F(\mathbf{x}) = -\log \int_{\mathbf{h}'} e^{-E(\mathbf{x}, \mathbf{h}')}$ is the free energy function. The model parameters are estimated by maximizing

$$\log p(\mathbf{x}) = -F(\mathbf{x}) - \log Z. \quad (12)$$

Given N samples, the gradient is computed as

$$\begin{aligned} \Delta_{\theta} &= \frac{1}{N} \sum_{n=1}^N \left(\frac{\partial F(\mathbf{x}_n)}{\partial \theta} - E_{p(\mathbf{x})} \left[\frac{\partial F(\mathbf{x})}{\partial \theta} \right] \right) \\ &= \frac{1}{N} \sum_{n=1}^N \left(E_{p(\mathbf{h}|\mathbf{x}_n)} \left[\frac{\partial E(\mathbf{x}_n, \mathbf{h})}{\partial \theta} \right] - E_{p(\mathbf{x}, \mathbf{h})} \left[\frac{\partial E(\mathbf{x}, \mathbf{h})}{\partial \theta} \right] \right), \end{aligned} \quad (13)$$

where $\frac{\partial E(\mathbf{x}, \mathbf{h})}{\partial \theta}$ w.r.t. \mathbf{a} , \mathbf{b} , and \mathbf{W} , respectively, are:

$$\frac{\partial E(\mathbf{x}, \mathbf{h})}{\partial \mathbf{a}_m^{(r)}} = -s_m^{(r)}; \frac{\partial E(\mathbf{x}, \mathbf{h})}{\partial \mathbf{b}_k^{(u)}} = -t_k^{(u)}; \frac{\partial E(\mathbf{x}, \mathbf{h})}{\partial w_{m,k}^{(r,u)}} = -s_m^{(r)} t_k^{(u)}. \quad (14)$$

The log-likelihood in Eq. (12) may serve as a measure of learning quality, where $\log Z$ must be estimated and we propose that $F(\mathbf{x})$ can be computed analytically:

$$F(\mathbf{x}) = -\zeta(\mathbf{x}) - \sum_{k=1}^K \log \int_{\mathbf{h}_k} e^{-\gamma_k(\mathbf{x}, \mathbf{h}_k)} = -\zeta(\mathbf{x}) - \sum_{k=1}^K B_k(\hat{\mathbf{b}}_k), \quad (15)$$

where $\hat{\mathbf{b}}_k$ and $B_k(\cdot)$ are respectively the posterior parameter and log-partition function of a hidden variable distribution.

Using the free energy function, the AIS procedure for binary RBM (Salakhutdinov & Murray, 2008) can be generalized to estimate the log-partition function of exp-RBM, where the intermediate distribution for annealing can be defined as

$$p_t(\mathbf{x}, \mathbf{h}) = \frac{1}{Z_t} e^{(1-\beta_t)E_A(\mathbf{x}, \mathbf{h}^{(A)}) - \beta_t E_B(\mathbf{x}, \mathbf{h}^{(B)})}, \quad (16)$$

where β_t gradually changes from 0 to 1: $0 = \beta_0 < \dots < \beta_t < \dots < \beta_T = 1$, so that $p_0(\mathbf{x}, \mathbf{h})$ corresponds to the base-rate model A and $p_T(\mathbf{x}, \mathbf{h})$ corresponds to the focused model B. The intermediate marginals can be defined as:

$$p_t(\mathbf{x}) = \frac{1}{Z_t} e^{-F_t(\mathbf{x})} = \frac{1}{Z_t} p_t^*(\mathbf{x}), \quad (17)$$

where $p_t^*(\mathbf{x}) = e^{-F_t(\mathbf{x})}$ and the free energy function is

$$\begin{aligned} F_t(\mathbf{x}) &= -(1-\beta_t)\zeta^{(A)}(\mathbf{x}) - \sum_{k=1}^{K_A} B_k((1-\beta_t)\hat{\mathbf{b}}_k^{(A)}) \\ &\quad - \beta_t\zeta^{(B)}(\mathbf{x}) - \sum_{k=1}^{K_B} B_k(\beta_t\hat{\mathbf{b}}_k^{(B)}). \end{aligned} \quad (18)$$

If model B is $\text{RBM}_B(\mathbf{a}, \mathbf{b}, \mathbf{W})$, model A is $\text{RBM}_A(\mathbf{a}, \mathbf{b}, \mathbf{0})$, we unveil that $\log Z_A$ can be analytically computed as

$$\log Z_A = \sum_{m=1}^M A_m(\mathbf{a}_m) + \sum_{k=1}^K B_k(\mathbf{b}_k). \quad (19)$$

where $A_m(\mathbf{a}_m)$ and $B_k(\mathbf{b}_k)$ are the log-partition functions of x_m and h_k , respectively.

The Markov transition $T(\mathbf{x}_t|\mathbf{x}_{t-1})$ can be defined as

$$\begin{cases} p_t((x_t)_m | \mathbf{h}_{t-1}^{(A)}, \mathbf{h}_{t-1}^{(B)}) = p_{\text{vis}}(x | \boldsymbol{\eta}((1-\beta_t)\mathbf{a}_m + \beta_t\hat{\mathbf{a}}_m)) \\ p_t((h_t^{(A)})_k | \mathbf{x}_t) = p_{\text{hid}}(h | \boldsymbol{\eta}((1-\beta_t)\mathbf{b}_k)) \\ p_t((h_t^{(B)})_k | \mathbf{x}_t) = p_{\text{hid}}(h | \boldsymbol{\eta}(\beta_t\hat{\mathbf{b}}_k)). \end{cases}$$

Then, we can sample a sequence: $\{\mathbf{x}_t, \mathbf{h}_t^{(A)}, \mathbf{h}_t^{(B)}\}$ ($0 \leq t \leq T-1$). Hence, the importance weight is computed by

$$w_s = \frac{p_1^*(\mathbf{x}_0)}{p_0^*(\mathbf{x}_0)} \dots \frac{p_T^*(\mathbf{x}_{T-1})}{p_{T-1}^*(\mathbf{x}_{T-1})} = \prod_{t=0}^{T-1} \frac{p_{t+1}^*(\mathbf{x}_t)}{p_t^*(\mathbf{x}_t)}. \quad (20)$$

Thus, the ratio of partition function can be estimated by $\frac{Z_B}{Z_A} \approx w_s$. Using a logarithm version to increase stability (Neal, 2001), we estimate $\log Z_B$ by

$$\log \hat{Z}_B = \log w_s + \log Z_A, \quad (21)$$

where $\log w_s = \sum_{t=0}^{T-1} (F_t(\mathbf{x}_t) - F_{t+1}(\mathbf{x}_t))$ and $\log Z_A$ is computed using Eq. (19). These AIS components for exp-RBMs are summarized in Supplemental Table 3.

3. Exponential Family Helmholtz Machine

3.1. Generic Formulation of Exp-HM

A Helmholtz machine (HM) (Dayan et al., 1995) tackles the intractable inference problem $p(\mathbf{h}|\mathbf{x})$ raised in logistic belief nets (Neal, 1992) by introducing corresponding recognition connections to form an efficient approximate posterior distribution $q(\mathbf{h}|\mathbf{x})$. Thus, HM is a multilayer stochastic neural network with both top-down generative connections and bottom-up recognition connections. However, the original HM model only considers binary visible and hidden variables. Here, we extend it to the exponential family HM (exp-HM), so that a range of discrete and continuous data can be properly modelled.

An exp-HM model with L hidden layers can be generally formulated to a product of conditional distributions:

$$p(\mathbf{x}, \mathbf{h}^{(1)}, \dots, \mathbf{h}^{(L)}) = p(\mathbf{x}|\mathbf{h}^{(1)}) \prod_{l=1}^{L-1} p(\mathbf{h}^{(l)}|\mathbf{h}^{(l+1)})p(\mathbf{h}^{(L)}), \quad (22)$$

where the conditionals and $p(\mathbf{h}^{(L)})$ are driven by the generative connections and belong to the exponential family. The model parameters include generative parameters $\boldsymbol{\theta}^{(G)} = \{\mathbf{a}^{(G)}, \mathbf{b}^{(G,1)}, \dots, \mathbf{b}^{(G,L)}, \mathbf{W}^{(G,1)}, \dots, \mathbf{W}^{(G,L)}\}$

and recognition parameters $\theta^{(R)} = \{\mathbf{b}^{(R,1)}, \dots, \mathbf{b}^{(R,L)}, \mathbf{W}^{(R,1)}, \dots, \mathbf{W}^{(R,L)}\}$, where $\mathbf{a}^{(G)}$ contains the generative bias parameters over the visible variables, $\mathbf{b}^{(G,l)}$ ($\mathbf{b}^{(R,l)}$) contains generative (recognition) bias parameters over the l -th hidden layer, and $\mathbf{W}^{(G,l)}$ ($\mathbf{W}^{(R,l)}$) is the generative (recognition) interaction weight matrix connecting the $(l-1)$ - and l -th hidden layers. The recognition component does not need bias parameters on \mathbf{x} , as \mathbf{x} is observed. By convention, we let $\text{size}(\mathbf{W}^{(R,l)}) = \text{size}(\mathbf{W}^{(G,l)})^T$. Merely for the convenience of discussion, we allow $p(\mathbf{x}|\mathbf{h}^{(1)})$ to follow any suitable exponential family distribution, and restrict $p(\mathbf{h}^{(l)}|\mathbf{h}^{(l+1)})$ and $p(\mathbf{h}^{(L)})$ to be Bernoulli. Writing these distributions in natural forms, we have

$$p(x_m|\mathbf{h}^{(1)}) = e^{\hat{\mathbf{a}}_m^{(G)T} \mathbf{s}_m(x_m) + \log f_m(x_m) - A_m(\hat{\mathbf{a}}_m^{(G)})} \quad (23)$$

$$p(h_k^{(l)}|\mathbf{h}^{(l+1)}) = e^{\hat{b}_k^{(G,l)} h_k^{(l)} + \log g_k^{(l)}(h_k^{(l)}) - B_k^{(l)}(\hat{b}_k^{(G,l)})} \quad (24)$$

$$p(h_k^{(L)}) = e^{\hat{b}_k^{(G,L)} h_k^{(L)} + \log g_k^{(L)}(h_k^{(L)}) - B_k^{(L)}(\hat{b}_k^{(G,L)})}, \quad (25)$$

where $0 \leq l \leq L-1$, $\hat{\mathbf{a}}_m^{(G)}$ hosts all the natural parameters of $p(x_m|\mathbf{h}^{(1)})$, and $\mathbf{s}_m(x_m)$ includes all corresponding natural forms of sufficient statistics for x_m . The posterior bias for the r -th sufficient statistic of x_m is computed as $\hat{a}_m^{(G,r)} = a_m^{(G,r)} + \mathbf{W}_{m,:}^{(G,1)} \mathbf{h}^{(1)}$, if its corresponding sufficient statistic $s_m^{(r)}$ is allowed to interact with $\mathbf{h}^{(1)}$; otherwise $\hat{a}_m^{(G,r)} = a_m^{(G,r)}$. Similarly, the posterior bias in $p(h_k^{(l)}|\mathbf{h}^{(l+1)})$ is computed as $\hat{b}_k^{(G,l)} = b_k^{(G,l)} + \mathbf{W}_{k,:}^{(G,l+1)} \mathbf{h}^{(l+1)}$.

The instances of conditional distribution $p(\mathbf{x}|\mathbf{h}^{(1)})$ for selected distributions in the exponential family are given in Supplemental Table 4. For example, if it follows a Poisson distribution, it is defined as

$$p(\mathbf{x}|\mathbf{h}^{(1)}) = \prod_{m=1}^M \mathcal{PO}(x_m | e^{\hat{a}_m^{(G)}}), \quad (26)$$

where $e^{\hat{a}_m^{(G)}}$ is the mean, and $\hat{a}_m^{(G)} = a_m^{(G)} + \mathbf{W}_{m,:}^{(G,1)} \mathbf{h}^{(1)}$. Conditional $p(\mathbf{h}^{(l)}|\mathbf{h}^{(l+1)})$ follow Bernoulli distributions:

$$p(\mathbf{h}^{(l)}|\mathbf{h}^{(l+1)}) = \prod_{k=1}^{K_l} \mathcal{BE}(h_k | \sigma(\hat{b}_k^{(G,l)})), 1 \leq l \leq L-1 \quad (27)$$

$$p(\mathbf{h}^{(L)}) = \prod_{k=1}^{K_L} \mathcal{BE}(h_k | \sigma(\hat{b}_k^{(G,L)})), \quad (28)$$

where $\sigma(\hat{b}_k^{(G,l)})$ is the success rate.

In HM, the approximation $q(\mathbf{h}|\mathbf{x})$ is defined as

$$q(\mathbf{h}^{(1)}|\mathbf{x}) = \prod_{k=1}^{K_1} \mathcal{BE}(h_k | \sigma(\hat{b}_k^{(R,1)})) \quad (29)$$

$$q(\mathbf{h}^{(l)}|\mathbf{h}^{(l-1)}) = \prod_{k=1}^{K_l} \mathcal{BE}(h_k | \sigma(\hat{b}_k^{(R,l)})), 2 \leq l \leq L. \quad (30)$$

where $\hat{b}_k^{(R,1)} = b_k^{(R,1)} + \mathbf{W}_{k,:}^{(R,1)} \mathbf{x}$ and $\hat{b}_k^{(R,l)} = b_k^{(R,l)} + \mathbf{W}_{k,:}^{(R,l)} \mathbf{h}^{(l-1)}$ for $2 \leq l \leq L$.

3.2. Wake-Sleep Algorithm for Exp-HM

The wake-sleep algorithm (Hinton et al., 1995), originally proposed to learn binary HM, is generalized for exp-HM in this work. Generally speaking, the wake-sleep algorithm is in fact an expectation-maximization (EM) algorithm (Neal & Hinton, 1998) involving a wake phase and a sleep phase, and considers both generative and recognition capabilities. The wake phase is to improve the model to generate a sample close to the training distribution, and the sleep phase is to improve the capability of this model to recognize a fantasy. In the wake phase, values of hidden units are driven by the recognition parameters, and only the generative parameters are updated. In the sleep phase, values of hidden units and visible units are sampled using the generative (and maybe recognition) parameters, and only the recognition parameters are updated. Our wake-sleep algorithm for exp-HM is derived as follows. Note that $q(\mathbf{h}|\mathbf{x})$ [Eq. (29) and (30)] is factorizable and driven by the recognition connections, while $p(\mathbf{h}|\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{h})}{p(\mathbf{x})}$ is nonfactorizable, intractable and driven by the generative model. As an objective, an upper bound of the negated log-likelihood is minimized in HM and exp-HM. According to the variational approximation theory, we know that

$$\log p(\mathbf{x}) = -j(\mathbf{x}) + \text{KL}(q(\mathbf{h}|\mathbf{x})||p(\mathbf{h}|\mathbf{x})), \quad (31)$$

where $j(\mathbf{x}) = -l(\mathbf{x}) = -\text{KL}(p(\mathbf{x}, \mathbf{h})||q(\mathbf{h}|\mathbf{x})) = \int_{\mathbf{h}} q(\mathbf{h}|\mathbf{x}) \log \frac{p(\mathbf{x}, \mathbf{h})}{q(\mathbf{h}|\mathbf{x})}$ is the upper bound of the negated log-likelihood (equivalently $l(\mathbf{x})$ is the variational lower bound of the log-likelihood). The name comes from the fact that $\log p(\mathbf{x}) \geq l(\mathbf{x})$, due to $\text{KL}(q(\mathbf{h}|\mathbf{x})||p(\mathbf{h}|\mathbf{x})) \geq 0$. From this definition of $j(\mathbf{x})$, we have

$$\begin{aligned} j(\mathbf{x}) &= - \int_{\mathbf{h}} q(\mathbf{h}|\mathbf{x}) \log \frac{p(\mathbf{x}, \mathbf{h})}{q(\mathbf{h}|\mathbf{x})} \\ &= - \int_{\mathbf{h}} q(\mathbf{h}|\mathbf{x}) \log p(\mathbf{x}, \mathbf{h}) + \int_{\mathbf{h}} q(\mathbf{h}|\mathbf{x}) \log q(\mathbf{h}|\mathbf{x}) \\ &= -\mathbb{E}_{q(\mathbf{h}|\mathbf{x})}[\log p(\mathbf{x}, \mathbf{h})] - H(q(\mathbf{h}|\mathbf{x})). \end{aligned} \quad (32)$$

where $-\log p(\mathbf{x}, \mathbf{h})$ is the description length of $p(\mathbf{x}, \mathbf{h})$ and $H(q(\mathbf{h}|\mathbf{x}))$ is the functional entropy of the approximate distribution. When $q(\mathbf{h}|\mathbf{x})$ is decomposable, so is $H(q(\mathbf{h}|\mathbf{x}))$. Hence, computing $H(q(\mathbf{h}|\mathbf{x}))$ is technically easy. We can see that minimizing the upper bound $j(\mathbf{x})$ is equivalent to minimizing the expected description length $-\log p(\mathbf{x}, \mathbf{h})$ (also called expected complete data negated log-likelihood) with respect to the approximate posterior distribution and maximizing the functional entropy of the approximate posterior distribution. In general, if $p(\mathbf{x}, \mathbf{h}) = \frac{1}{Z'} e^{-E(\mathbf{x}, \mathbf{h})}$ where $E(\mathbf{x}, \mathbf{h})$ is the internal energy of the system, from a physical perspective, the cost $j(\mathbf{x})$ can also be written as

$$j(\mathbf{x}) = F(\mathbf{x}) = \mathbb{E}_{q(\mathbf{h}|\mathbf{x})}[E(\mathbf{x}, \mathbf{h})] + \log Z' - H(q). \quad (33)$$

Thus, $j(\mathbf{x})$ is called Helmholtz free energy, denoted by $F(\mathbf{x})$ (hence the name of HM). For exp-HM, in the wake phase, $j(\mathbf{x})$ in Eq. (32) is the objective to be minimized.

In the wake phase, the first-order derivative of the objective $j(\mathbf{x})$ [Eq. (32)] w.r.t. the generative parameters is

$$\frac{\partial j(\mathbf{x})}{\partial \theta^{(G)}} = -\mathbb{E}_{q(\mathbf{h}|\mathbf{x})} \left[\frac{\partial \log p(\mathbf{x}, \mathbf{h})}{\partial \theta^{(G)}} \right], \quad (34)$$

where \mathbf{h} is approximated using $q(\mathbf{h}|\mathbf{x})$, and $\frac{\partial \log p(\mathbf{x}, \mathbf{h})}{\partial \theta^{(G)}}$ is crucial. We find that the first-order derivative of $\log p(\mathbf{x}, \mathbf{h})$ w.r.t. $W_{m,k}^{(G,1)}$ can be computed as

$$\begin{aligned} \frac{\partial \log p(\mathbf{x}, \mathbf{h})}{\partial W_{m,k}^{(G,1)}} &= \frac{\partial \log p(x_m | \mathbf{h}^{(1)})}{\partial W_{m,k}^{(G,1)}} \\ &= \frac{\partial \log p(x_m | \mathbf{h}^{(1)})}{\partial \hat{a}_m^{(G,r)}} \frac{\partial \hat{a}_m^{(G,r)}}{\partial W_{m,k}^{(G,1)}} \\ &= (s_m^{(r)} - \mathbb{E}_{p(x_m | \mathbf{h}^{(1)})} [s_m^{(r)}]) h_k^{(1)} \\ &= (x_m - \langle x_m \rangle) h_k^{(1)}, \end{aligned} \quad (35)$$

where in the second line we assume the r -th natural variable interacts with $\mathbf{h}^{(1)}$; in the third line we take advantage of the important property of exponential family from Eq. (3); and finally we only practically allow \mathbf{x} to interact with $\mathbf{h}^{(1)}$. Eq. (35) tells us that the gradient to update $W_{m,k}^{(G,1)}$ equals to the correlation between the observed value of x_m and hidden state $h_k^{(1)}$ minus the correlation between the expected x_m and $h_k^{(1)}$. If $p(x_m | \mathbf{h}^{(1)})$ follows a Bernoulli distribution, then $\langle x_m \rangle = \sigma(\hat{a}_m^{(G)})$, which is consistent with the logistic belief net (Neal, 1992) and the original HM (Hinton et al., 1995). If $p(x_m | \mathbf{h}^{(1)})$ is Gaussian, then $\langle x_m \rangle = -\frac{\hat{a}_m^{(G,1)}}{2a_m^{(G,1)}}$. Likewise, the first order derivative of $\log p(\mathbf{x}, \mathbf{h})$ w.r.t. $a_m^{(G,r)}$ can be computed as

$$\begin{aligned} \frac{\partial \log p(\mathbf{x}, \mathbf{h})}{\partial a_m^{(G,r)}} &= \frac{\partial \log p(x_m | \mathbf{h}^{(1)})}{\partial a_m^{(G,r)}} \\ &= \frac{\partial \log p(x_m | \mathbf{h}^{(1)})}{\partial \hat{a}_m^{(G,r)}} \frac{\partial \hat{a}_m^{(G,r)}}{\partial a_m^{(G,r)}} = s_m^{(r)} - \langle s_m^{(r)} \rangle. \end{aligned} \quad (36)$$

Similarly, the derivative of $\log p(\mathbf{x}, \mathbf{h})$ w.r.t. $W_{k_1,k_2}^{(G,l)}$ ($2 \leq l \leq L$) and $b_k^{(G,l)}$ ($1 \leq l \leq L$) are respectively computed as

$$\frac{\partial \log p(\mathbf{x}, \mathbf{h})}{\partial W_{k_1,k_2}^{(G,l)}} = (h_{k_1}^{(l-1)} - \langle h_{k_1}^{(l-1)} \rangle) h_{k_2}^{(l)}, \quad 2 \leq l \leq L \quad (37)$$

$$\frac{\partial \log p(\mathbf{x}, \mathbf{h})}{\partial b_k^{(G,l)}} = h_k^{(l)} - \langle h_k^{(l)} \rangle, \quad 1 \leq l \leq L, \quad (38)$$

where $\langle h_k^{(l)} \rangle = \sigma(\hat{b}_k^{(G,l)})$ for Bernoulli distributions. In short, we write the gradients in vector and matrix forms:

$$\Delta_{\mathbf{a}^{(G,r)}} = -(\mathbf{s}^{(r)} - \langle \mathbf{s}^{(r)} \rangle) \quad (39)$$

$$\Delta_{\mathbf{W}^{(G,1)}} = -(\mathbf{x} - \langle \mathbf{x} \rangle) \mathbf{h}^{(1)\top} \quad (40)$$

$$\Delta_{\mathbf{W}^{(G,l)}} = -(\mathbf{h}^{(l-1)} - \langle \mathbf{h}^{(l-1)} \rangle) \mathbf{h}^{(l)\top}, \quad 2 \leq l \leq L \quad (41)$$

$$\Delta_{\mathbf{b}^{(G,l)}} = -(\mathbf{h}^{(l)} - \langle \mathbf{h}^{(l)} \rangle), \quad 1 \leq l \leq L, \quad (42)$$

where \mathbf{x} is an actual training sample, and $\mathbf{h} = \{\mathbf{h}^{(1)}, \dots, \mathbf{h}^{(L)}\}$ are sampled from $q(\mathbf{h}|\mathbf{x})$ using the recognition parameters. With these gradients, the generative parameters are updated using rule: $\theta^{(G)} = \theta^{(G)} - \epsilon \Delta_{\theta^{(G)}}$.

In the sleep phase, starting from the top layer, a fantasy can be either (i) unbiasedly sampled from the generative distribution $p(\mathbf{h}^{(L)})$ or (ii) initialized using the states of $\mathbf{h}^{(L)}$ obtained in the wake phase. Values of downstream units are sampled by the generative component. We thus generate a total fantasy, using which the recognition parameters are updated accordingly. The second method (used in our implementation) connects both objectives and is akin to the contrastive divergence (CD) algorithm (Hinton, 2002; Hinton et al., 2006). To remove sampling bias, Gibbs sampling with multiple alternating top-down and bottom-up passes could be adopted, similar to (persistent) CD- k algorithms (Tieleman, 2008). Since the recognition and generative components are structurally opposite, to derive update rules for the recognition parameters, we just need to swap their roles in the objective function [Eq. (32)], resulting in the corresponding gradients:

$$\Delta_{\mathbf{b}^{(R,l)}} = -(\mathbf{h}^{(l)} - \langle \mathbf{h}^{(l)} \rangle), \quad 1 \leq l \leq L \quad (43)$$

$$\Delta_{\mathbf{W}^{(R,1)}} = -(\mathbf{h}^{(1)} - \langle \mathbf{h}^{(1)} \rangle) \mathbf{x}^\top \quad (44)$$

$$\Delta_{\mathbf{W}^{(R,l)}} = -(\mathbf{h}^{(l)} - \langle \mathbf{h}^{(l)} \rangle) \mathbf{h}^{(l-1)\top}, \quad 2 \leq l \leq L, \quad (45)$$

where $\{\mathbf{h}^{(L)}, \dots, \mathbf{h}^{(1)}, \mathbf{x}\} \sim p(\mathbf{x}, \mathbf{h})$; $\langle h_k^{(l)} \rangle = \sigma(\hat{b}_k^{(R,l)})$.

In the algorithm, the wake phase and sleep phase alternately iterate for a pre-specified steps or until the objective $j(\mathbf{x})$ [Eq. (32)] does not reduce dramatically. The gradients for selected instances of exp-HMs are given in Supplemental Table 4. The wake-sleep algorithm with a random initialization may not perform well for deep exp-HMs. Stacked exp-RBMs can be first applied to provide a warm start.

4. Exponential Family DBN

When generating fantasies, one critical issue with HM models is that a point sampled from $p(\mathbf{h}^{(L)})$ may not represent a meaningful object, hence may eventually generate an abnormal \mathbf{x} . To address this issue, DBN adds an RBM on top of an HM, so that a point generated using Gibbs sampling from the RBM is from the domain of interest (Hinton et al., 2006). However, the original DBN was designed for binary variables, limiting its applications in many fields. To overcome it, we generalize it to have visible variables follow any distributions from the exponential family. Although hidden variables can follow any exponential family distributions too, for convenience of derivation, we only discuss Bernoulli hidden variables below. Such an exponential family DBN (exp-DBN) with L hidden layers is defined as

$$\begin{aligned} p(\mathbf{x}, \mathbf{h}) &= p(\mathbf{x} | \mathbf{h}^{(1)}) p(\mathbf{h}^{(1)} | \mathbf{h}^{(2)}) \dots \\ &\quad p(\mathbf{h}^{(L-2)} | \mathbf{h}^{(L-1)}) p(\mathbf{h}^{(L-1)}, \mathbf{h}^{(L)}), \end{aligned} \quad (46)$$

where we denote $\mathbf{h} = \{\mathbf{h}^{(1)}, \dots, \mathbf{h}^{(L)}\}$ and define

$$p(\mathbf{x}|\mathbf{h}^{(1)}) = \prod_{m=1}^M p(x_m|\eta(\hat{\mathbf{a}}_m^{(G)})) \quad (47)$$

$$p(\mathbf{h}^{(l)}|\mathbf{h}^{(l+1)}) = \prod_{k=1}^{K_l} \mathcal{BE}(h_k|\sigma(\hat{b}_k^{(G,l)})), \quad 1 \leq l \leq L-2 \quad (48)$$

$$p(\mathbf{h}^{(L-1)}, \mathbf{h}^{(L)}) = \frac{1}{Z} e^{-E(\mathbf{h}^{(L-1)}, \mathbf{h}^{(L)})}, \quad (49)$$

where $\hat{\mathbf{a}}_m^{(G)}$ and $\hat{b}_k^{(G,l)}$ are defined as in exp-HM and $E(\mathbf{h}^{(L-1)}, \mathbf{h}^{(L)}) = -\mathbf{b}^{(G,L-1)\top} \mathbf{h}^{(L-1)} - \mathbf{b}^{(G,L)\top} \mathbf{h}^{(L)} - \mathbf{h}^{(G,L-1)\top} \mathbf{W}^{(G,L)} \mathbf{h}^{(G,L)}$. The undirected component is formulated by a joint distribution, while the directed component are formulated by conditional distributions. Eq. (47) formulates the exponential family conditional distributions of visible variables. For example, a Bernoulli visible unit has $p(x_m|\mathbf{h}^{(1)}) = \mathcal{BE}(x_m|\sigma(\hat{a}_m^{(G)}))$, where $\hat{a}_m^{(G)} = a_m^{(G)} + \mathbf{W}_{m,:}^{(G,1)} \mathbf{h}^{(1)}$. Other distributions are given in Supplemental Table 4. Same as in exp-HMs, the conditional over a hidden layer follows a Bernoulli distribution which is parameterized by its prior bias and interactions with its upper layer. The top two layers form a Bernoulli-Bernoulli RBM as associative memory. The inference in DBN is also intractable. Same as in exp-HMs, recognition connections define a factorizable approximate distribution:

$$q(\mathbf{h}^{(1 \dots L-1)}|\mathbf{x}) = q(\mathbf{h}^{(1)}|\mathbf{x}) \dots q(\mathbf{h}^{(L-1)}|\mathbf{h}^{(L-2)}), \quad (50)$$

where $\mathbf{h}^{(1 \dots L-1)} = \{\mathbf{h}^{(1)}, \dots, \mathbf{h}^{(L-1)}\}$. In summary, the model parameters in exp-DBN include generative parameters $\theta^{(G)} = \{\mathbf{a}^{(G)}, \mathbf{b}^{(G,1)}, \dots, \mathbf{b}^{(G,L-1)}, \mathbf{b}^{(G,L)}, \mathbf{W}^{(G,1)}, \dots, \mathbf{W}^{(G,L-1)}, \mathbf{W}^{(G,L)}\}$ and recognition parameters $\theta^{(R)} = \{\mathbf{b}^{(R,1)}, \dots, \mathbf{b}^{(R,L-1)}, \mathbf{W}^{(R,1)}, \dots, \mathbf{W}^{(R,L-1)}\}$.

4.1. Model Learning for Exp-DBN

In the original DBN (Hinton et al., 2006), the model parameters are pretrained using binary RBMs (with $\theta^{(R)}$ tied with $\theta^{(G)}$) and then fine-tuned using a wake-sleep (up-down) algorithm (with $\theta^{(R)}$ and $\theta^{(G)}$ untied). We generalize this procedure to exp-DBN, such that its model parameters can be pretrained using exp-RBMs and then fine-tuned using an exponential family wake-sleep algorithm. In the wake phase of this extended algorithm, $\mathbf{h}^{(1 \dots L-1)}$ are first sampled using the corresponding recognition parameters in the exp-HM component. Then, using $\mathbf{h}^{(L-1)}$ as input, (persistent) CD- k sampling (Tieleman, 2008) can be employed in the top exp-RBM and its parameters are updated using the rules for exp-RBM. After that, generative parameters below the top exp-RBM are updated using the same rules as in exp-HMs. In the sleep phase, $\mathbf{h}^{(L-1)}$ sampled during the wake phase is used to generate states of the subsequent layers driven by generative parameters, and the recognition

parameters are updated as in exp-HMs. The gradients to update the model parameters for various distributions are given in Supplemental Table 4.

4.2. Variational Lower Bound of Exp-DBN

As indicated in Eq. (49), the partition function only appears at the top exp-RBM. Inspired by (Salakhutdinov & Murray, 2008), to obtain a variational lower bound of the log-likelihood, $\mathbf{h}^{(L)}$ thus is summed out, leading to

$$p(\mathbf{x}, \mathbf{h}^{(1 \dots L-1)}) = p(\mathbf{x}|\mathbf{h}^{(1)}) \dots p(\mathbf{h}^{(L-2)}|\mathbf{h}^{(L-1)}) \times \frac{1}{Z} e^{-F(\mathbf{h}^{(L-1)})} = \frac{1}{Z} p^*(\mathbf{x}, \mathbf{h}^{(1 \dots L-1)}). \quad (51)$$

The variational lower bound is hence computed as

$$l(\mathbf{x}) = \mathbb{E}_{q(\mathbf{h}^{(1 \dots L-1)}|\mathbf{x})} [\log p^*(\mathbf{x}, \mathbf{h}^{(1 \dots L-1)})] - \log Z + H(q(\mathbf{h}^{(1 \dots L-1)}|\mathbf{x})). \quad (52)$$

The first term can be estimated by Monte Carlo approximation by sampling from $q(\mathbf{h}^{(1 \dots L-1)}|\mathbf{x})$ with \mathbf{x} clamped. The second term $\log Z$ can be estimated using the generalized AIS procedure (Section 2) to the top exp-RBM.

4.3. Applications of Exp-DBN

Certainly, exp-DBN is applicable to pretrain deep discriminative models. More importantly, exp-DBN can be further extended to multi-modal exp-DBNs (see Supplemental Figure 1 for an example) which can be potentially applied in multi-modal learning (Ramachandram & Taylor, 2017; Li et al., 2018) tasks, such as data fusion, integrative classification and clustering, multi-label learning, transfer learning, missing value estimation, information retrieval, machine translation, autonomous navigation, and Internet of things. Our future work will focus on these applications. In our experiments, we designed a two-modal exp-DBN to showcase some potentials of our exponential family generalization.

5. Experiments

5.1. Modelling Image Data

To demonstrate the potential of these exponential family models for image generation, we first designed and learned exp-HM, exp-DBN, and two-modal exp-DBN models on the MNIST (LeCun et al., 1998) and Fashion-MNIST (Xiao et al., 2017) data sets respectively. Each data set contains 60,000 training and 10,000 test images of size 28×28 . Bernoulli and Multinoulli distributions were assumed respectively for the image (pixel intensities in range $[0,1]$) and class (categories from 0 to 9) modalities in the two-modal exp-DBN. All models were pretrained using exp-RBMs and fine-tuned using our generalized wake-sleep algorithms. The structures of these models and the generated images are

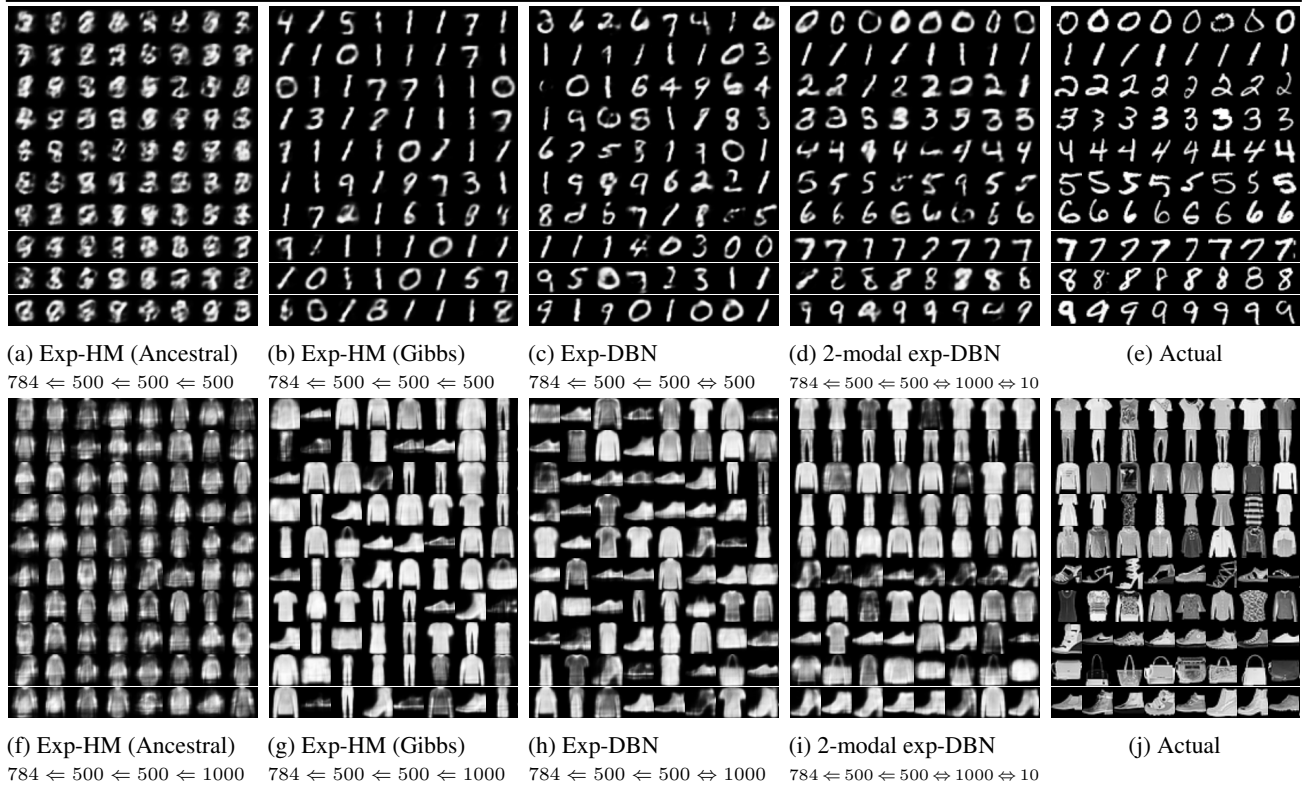


Figure 1. Actual and generated images on MNIST (top) and Fashion-MNIST (bottom) data sets.

given in Figure 1. First, from Figures 1a and 1f, we see that exp-HM did not generate good images using a traditional ancestral sampling, indicating a hidden state sampled from the generative distribution of the top hidden layer is unlikely from the domain of interest when its dimensionality is large. Second, we uncover that exp-HM’s capacity of generating fantasies can be improved by a Gibbs sampling using multiple alternating top-down and bottom-up passes (see Figures 1b and 1g). This discovery was not reported before. Third, Figures 1c and 1h corroborate the early finding in (Hinton et al., 2006) that, by adding an undirected network on top of an exp-HM, exp-DBN could generate images of good quality, the majority of which look similar to the actual examples in Figures 1e and 1j. Fourth, as show in Figures 1d and 1i, by fixing the class modality, the two-modal exp-DBN could take attention to generating images of specified classes.

5.2. Modelling Text Data

To further show the importance of proper distribution in modelling, we applied a two-modal exp-DBN on the 20 Newsgroups data (Lang, 1995). The 11,269 training and 7,505 test documents from 20 topics were transformed to word-count profiles. Stop-words and words appear in less than 60 training samples were discarded, leaving 4,031 words for modelling. The network has one modality for word counts with a hidden layer of 500 Bernoulli units, one modality for class labels without any hidden layer, and a joint hidden layer with 520 Bernoulli units. We used

Bernoulli, Poisson, and Multinomial distributions, separately, for word counts, and Multinoulli distribution for class labels. Each profile was normalized to have 1,000 word counts for Poisson and Multinomial models, but binarized for Bernoulli models. For comparison, we extended exp-RBM to a two-modal exp-DBM with same network architecture as the two-modal exp-DBN. Both types of generative models were pretrained using exp-RBMs. The two-modal exp-DBNs were fine-tuned by our generalized wake-sleep algorithms, while the two-modal exp-DBMs were fine-tuned by a modified algorithm originally proposed in (Salakhutdinov & Hinton, 2009a). A multilayer perceptron with same network skeleton but RELU units was also applied as a discriminative model on the tf-idf normalized data.

The log-partition functions and variational lower bounds of these generative models were estimated and shown in Table 1. We see that log-partition functions can be estimated stably with tiny variances in the exponential family framework, and the log-likelihoods of Poisson and Multinomial models are roughly similar (note: we cannot compare their log-likelihoods with these of the Bernoulli models, because they were experimented on different formats of the data). After training, fantasies were drawn using Gibbs sampling from the two-modal exp-DBN with Poisson distribution. The top 15 words and class label of some cherry-picked generations are given in Table 2, which implies that the model learned useful information and generated meaningful samples. In these generative models, class labels of test samples can

be treated as missing data and inferred using a mean-field method. Table 1 also shows the test accuracies after pre-training and after fine-tuning, respectively. First, we find that the Poisson models outperform Multinoulli models which are named replicated softmax models in (Salakhutdinov & Hinton, 2009b). The models based on conventionally used Bernoulli distributions obtained non-competitive results. Second, the two-modal exp-DBNs achieved better accuracies than their two-modal exp-DBM rivals. Moreover, although discriminative models often claim victories over generative models in classification tasks, we find that with appropriate distributions, the directed and undirected generative models under investigation may perform better than them which are less advantageous in handling data types. Of course, classification performance can be further improved if sequential dependency is properly modelled.

Table 1. Log-partition functions, variational lower bounds, and prediction accuracies on 20 Newsgroups dataset. BE: Bernoulli, MU: Multinoulli, PO: Poisson, MN: Multinomial.

	Vis. Type	Data	log Z (STD)	$l(\mathbf{x})$	Acc. pretrain	Acc. fine-tune
MDBN	BE+MU	Train	3105.74(± 0.35)	-2691.43	35.03%	38.81%
		Test	3106.05(± 0.46)	-2694.58		
	PO+MU	Train	1469.18(± 0.44)	-3976.12	68.46%	69.30%
		Test	1468.72(± 0.39)	-4273.88		
	MN+MU	Train	1617.29(± 0.64)	-3897.86	66.80 %	68.27%
		Test	1616.88(± 0.63)	-4216.70		
MDBM	BE+MU	Train	3028.27(± 2.61)	-1754.63	27.24%	40.58%
		Test	3028.28(± 1.43)	-1757.03		
	PO+MU	Train	2261.53(± 1.76)	-3772.34	67.31%	66.84%
		Test	2260.72(± 1.76)	-4043.79		
	MN+MU	Train	9761.72(± 0.96)	-3873.40	63.94 %	63.96%
		Test	9762.22(± 0.76)	-4108.18		
Multilayer perceptron with RELU units					58.53%	

Table 2. Top 15 words and classes of samples generated by the Poisson-Multinoulli two-modal exp-DBN.

\mathbf{x}	\mathbf{y}
government system writes key clipper encryption keys don	sci.crypt
nsa article escrow secure chip people time	
article writes people killed gun fire apr control country children weapons government time guns police	talk.politics.guns
sale offer university email includes time condition windows call contact shipping interested sell computer original	misc.forsale
division games toronto boston hockey st team cup pittsburgh series detroit play chicago game montreal	rec.sport.hockey
card bus motherboard ram controller board isa friend machine pc slot scsi work dx cards	comp.sys.mac.hardware

6. Conclusion and Discussion

We discuss RBM’s free energy function, HM and DBN from the exponential family perspective. We derive the corresponding wake-sleep algorithms in a general framework. Our experiments show that exp-DBN and exp-HM with Gibbs sampling are superior to exp-HM with ancestral sampling in generating samples. With appropriate distributions, the two-modal exp-DBN outperformed the two-modal exp-DBM and a discriminative model on a text data. We clarify that, after finishing this work, we identified that a model, named deep exponential family (DEF), is independently explored in (Ranganath et al., 2015). DEF and our work share similar motivation, but complement each other. The differences are explained below. (1) DEF was a result

of hierarchical chaining of exponential family distributions, while our work was inspired by existing and our new findings of exp-RBMs. In the DEF framework, only means of nodes within current layer are affected by their parents. But in our framework, any parameters of nodes in current layer can be modified by their parents in theory. (2) We generalize the wake-sleep algorithm for directed DGM, which is not trivial because strict derivation of generic update rules requires wise application of several properties of exp-RBM and exponential family distributions. DEF resorts to an algorithm similar to neural variational inference and learning (NVIL) (Mnih & Gregor, 2014), a REINFORCE algorithm. (3) We find that the free energy function of exp-RBM is a linear combination of log-partition functions of hidden variable, and the log-partition function of a base-rate exp-RBM can be beautifully computed as a summation of log-partition functions of individual variables. This discovery enables us to easily estimate log-partition functions of exp-RBM, exp-DBN and exp-DBM using AIS, and thus their (lower bounds of) likelihoods. This is not addressed in DEF.

Current DGMs with inference networks (including ours) are partially united under the umbrella of variational inference, using either separate or joint objectives for generative and recognition components. Specific learning and inference methods, such as wake-sleep algorithms (Bornschein & Bengio, 2015), NVIL based algorithms (Mnih & Gregor, 2014; Ranganath et al., 2015) and variational autoencoder (VAE) based algorithms (Rezende et al., 2014; Kingma & Welling, 2014), mainly differ in how the inference component is trained. However, these methods could be transformable under mild conditions. For example, if we apply mean-field approximation to alternating hidden layers in exp-HM, we would obtain a VAE-like model which uses deterministic networks to regress parameters of stochastic layers. Moreover, level of novelty and quality within generated samples vary among models, depending on how the generative and inference networks coordinate (Burda et al., 2016).

Our investigation indicates that ancestral sampling is not always effective. Association and (or) attention may be needed in both directed and undirected generative networks for novelty generation. While celebrating limited successes in current DGMs [including generative adversarial network (GAN) (Goodfellow et al., 2014)], methodological engineering may only partially solve the puzzle, we thus need to think beyond this niche. From the cognitive science perspective, innovation comes from interaction and combination of cognits, units of any knowledge or concept representation in the cerebral cortex (Fuster, 2003). Many existing DGMs require justification in cognitive neuroscience. Certainly, inspiration from hierarchies of perceptual, associative and executive networks would create a new generation of DGMs which might enhance the statistical foundation and drive us closer towards the long-dreamed strong AI.

References

- Bengio, Y., Lamblin, P., Popovici, D., and Larochelle, H. Greedy layer-wise training of deep networks. In *Advances in Neural Information Processing Systems*, pp. 153–160, 2006.
- Bornschein, J. and Bengio, Y. Reweighted wake-sleep. In *International Conference on Learning Representations*, 2015.
- Burda, Y., Grosse, R., and Salakhutdinov, R. Importance weighted autoencoders. In *International Conference on Learning Representations*, 2016.
- Dayan, P., Hinton, G.E., Neal, R., and Zemel, R.S. The Helmholtz machine. *Neural Computation*, 7:1022–1037, 1995.
- Fuster, J.M. *Cortex and Mind*. Oxford University Press, 2003.
- Goodfellow, I.J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. Generative adversarial networks. In *Advances in Neural Information Processing Systems*, pp. 2672–2680, 2014.
- Hinton, G. Training products of experts by minimizing contrastive divergence. *Neural Computation*, 14(8):1771–1800, 2002.
- Hinton, G., Dayan, P., Frey, B., and Neal, R. The wake-sleep algorithm for unsupervised neural networks. *Science*, 268: 1558–1161, 1995.
- Hinton, G.E., Osindero, S., and Teh, Y. A fast learning algorithm for deep belief nets. *Neural Computation*, 18: 1527–1554, 2006.
- Kingma, D.P. and Welling, M. Auto-encoding variational Bayes. In *International Conference on Learning Representations*, 2014.
- Lang, Ken. Newsweeder: Learning to filter netnews. In *International Conference on Machine Learning*, pp. 331–339, 1995. <http://qwone.com/~jason/20NewsGroups>.
- LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Li, Y., Wu, F.X., and Ngom, A. A review on machine learning principles for multi-view biological data integration. *Briefings in Bioinformatics*, 29(2):325–340, 2018.
- Mnih, A. and Gregor, K. Neural variational inference and learning in belief networks. In *International Conference on Machine Learning*, pp. II–1791–II–1799, 2014.
- Neal, R. M. and Hinton, G.E. A view of the EM algorithm that justifies incremental, sparse, and other variants. In Jordan, M.I. (ed.), *Learning in Graphical Models*, Adaptive Computation and Machine Learning series, chapter 11, pp. 355–368. MIT, Cambridge, MA, 1998.
- Neal, R.M. Connectionist learning of belief networks. *Artificial Intelligence*, 56:71–113, 1992.
- Neal, R.M. Annealed importance sampling. *Statistics and Computing*, 11:125–139, 2001.
- Ramachandram, D. and Taylor, G.W. Deep multimodal learning: A survey on recent advances and trends. *IEEE Signal Processing Magazine*, 34(6):96–108, 2017.
- Ranganath, R., Tang, L., Charlin, L., and Blei, D.M. Deep exponential families. In *International Conference on Artificial Intelligence and Statistics*, pp. 762–771, 2015.
- Rezende, D. J., Mohamed, S., and Wierstra, D. Stochastic backpropagation and approximate inference in deep generative models. In *International Conference on Machine Learning*, pp. II–1278–II–1286, 2014.
- Salakhutdinov, R. and Hinton, G. Deep Boltzmann machine. In *International Conference on Artificial Intelligence and Statistics*, pp. 448–455, 2009a.
- Salakhutdinov, R. and Hinton, G. Replicated softmax: An undirected topic model. In *Advances in Neural Information Processing Systems*, pp. 1607–1614, 2009b.
- Salakhutdinov, R. and Murray, I. On the quantitative analysis of deep belief networks. In *International Conference on Machine Learning*, pp. 872–879, 2008.
- Smolensky, P. Information processing in dynamical systems: Foundations of harmony theory. In Rumelhart, D.E. and McClelland, J.L. (eds.), *Parallel Distributed Processing: Explorations in the Microstructure of Cognition, Vol. 1: Foundation*, chapter 6, pp. 194–281. MIT, Cambridge, MA, 1986.
- Srivastava, N. and Salakhutdinov, R. Learning representations for multimodal data with deep belief nets. In *International Conference on Machine Learning Workshop on Representation Learning*, 2012.
- Srivastava, N. and Salakhutdinov, R. Multimodal learning with deep Boltzmann machines. *Journal of Machine Learning Research*, 15:2949–2980, 2014.
- Tieleman, T. Training restricted Boltzmann machines using approximations to the likelihood gradient. In *International Conference on Machine Learning*, pp. 1064–1071, 2008.

Welling, M., Rosen-zvi, M., and Hinton, G. Exponential family harmoniums with an application to information retrieval. In *Advances in Neural Information Processing Systems*, pp. 1481–1488, 2005.

Xiao, Han, Rasul, Kashif, and Vollgraf, Roland. Fashion-MNIST: A novel image dataset for benchmarking machine learning algorithms, 2017.

Supplemental Figures and Tables for “Exploring Helmholtz Machine and Deep Belief Net in the Exponential Family Perspective”

Yifeng Li^{*1} and Xiaodan Zhu^{†2}

¹Digital Technologies Research Centre, National Research Council Canada, Ottawa, Ontario, K1A 0R6, Canada

²Department of Electrical and Computer Engineering, Queen’s University, Kingston, Ontario, K7L 3N6 Canada

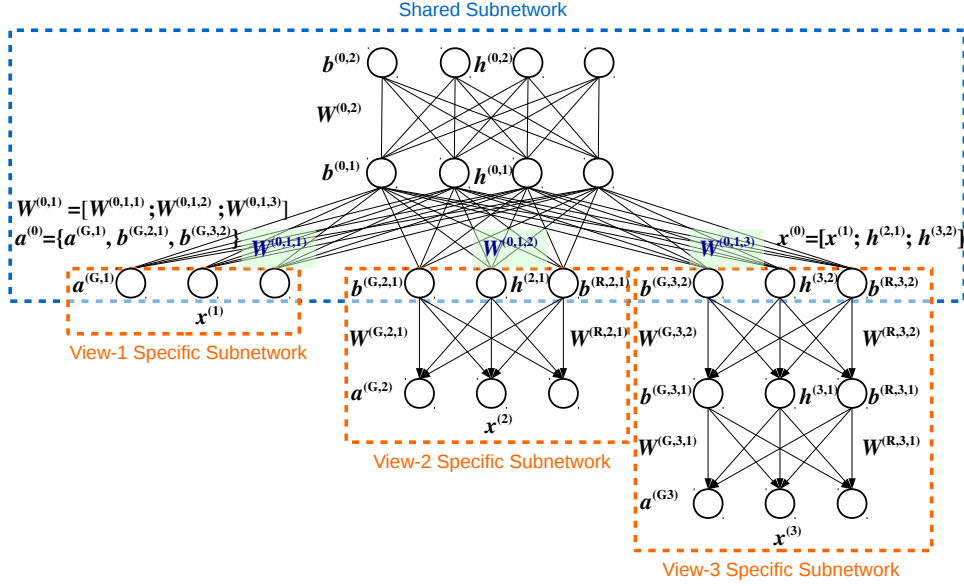


Figure 1: A schematic example of exponential family multi-modal deep belief net (exp-MDBN). It integrates three modalities/views whose visible layers can respectively follow any types of distributions from the exponential family. Modality 1 has a trivial structure without any hidden layers, and may be used for class labels. The subnetworks of modalities 2 and 3 are exp-HMs. The shared/joint subnetwork is a DBM, functioning as a deep associative memory.

^{*}Corresponding author. E-mail address: yifeng.li@nrc-cnrc.gc.ca yifeng.li.cn@gmail.com

[†]Email address: xiaodan.zhu@queensu.ca

Table 1: Selected members of the exponential family.

Distribution	Standard Form	Natural Form	$\theta(\eta)$	$\eta(\theta)$	$s(x)$	$h(x)$	$A(\eta)$	$A(\theta)$	Sufficient Statistics
Gaussian	$p(x \mu, \lambda^{-1}) = \frac{\sqrt{\lambda}}{\sqrt{2\pi}} e^{-\frac{\lambda}{2}(x-\mu)^2}$, $\lambda > 0$	$p(x \theta) = e^{\theta_1 s_1 + \theta_2 s_2 - A(\theta)}$, $\theta_2 < 0$	$\begin{bmatrix} \mu\lambda \\ -\frac{\lambda}{2} \end{bmatrix}$	$\begin{bmatrix} -\frac{\theta_1}{2\theta_2} \\ -\frac{\theta_1}{2\theta_2} \end{bmatrix}$	$\begin{bmatrix} x \\ x^2 \end{bmatrix}$	1	$\frac{1}{2} \log \frac{2\pi}{\lambda} + \frac{\mu^2 \lambda}{2}$	$\frac{1}{2} \log \frac{\pi}{-\theta_2} - \frac{\theta_1^2}{4\theta_2^2}$	$E[x] = \mu = -\frac{\theta_1}{2\theta_2}$ $\text{Var}[x] = \frac{1}{\lambda} = -\frac{1}{2\theta_2}$ $E[x^2] = \mu^2 + \frac{1}{\lambda} = \frac{\theta_1^2}{4\theta_2^2} + \frac{1}{-2\theta_2}$
Gaussian (fix precision λ)	$p(x \mu, \lambda^{-1}) = \frac{\sqrt{\lambda}}{\sqrt{2\pi}} e^{-\frac{\lambda}{2}(x-\mu)^2}$, $\lambda > 0$	$p(x \theta) = (\frac{\sqrt{\lambda}}{\sqrt{2\pi}} e^{-\frac{\lambda}{2}x^2}) e^{\theta s - A(\theta)}$	μ	θ	λx	$\frac{\sqrt{\lambda}}{\sqrt{2\pi}} e^{-\frac{\lambda}{2}x^2}$	$\frac{\mu^2 \lambda}{2}$	$\frac{\theta^2 \lambda}{2}$	$E[x] = \mu = \frac{\theta}{\lambda}$ $\text{Var}[x] = \frac{1}{\lambda}$ $E[x^2] = \mu^2 + \frac{1}{\lambda} = \theta^2 + \frac{1}{\lambda}$
Gaussian (fix precision λ)	$p(x \mu, \lambda^{-1}) = \frac{\sqrt{\lambda}}{\sqrt{2\pi}} e^{-\frac{\lambda}{2}(x-\mu)^2}$, $\lambda > 0$	$p(x \theta) = (\frac{\sqrt{\lambda}}{\sqrt{2\pi}} e^{-\frac{\lambda}{2}x^2}) e^{\theta s - A(\theta)}$	$\mu\lambda$	$\frac{\theta}{\lambda}$	x	$\frac{\sqrt{\lambda}}{\sqrt{2\pi}} e^{-\frac{\lambda}{2}x^2}$	$\frac{\mu^2 \lambda}{2}$	$\frac{\theta^2}{2\lambda}$	$E[x] = \mu = \frac{\theta}{\lambda}$ $\text{Var}[x] = \frac{1}{\lambda}$ $E[x^2] = \mu^2 + \frac{1}{\lambda} = \frac{\theta^2}{\lambda} + \frac{1}{\lambda}$
Poisson	$p(x \lambda) = \frac{e^{-\lambda}}{x!} \lambda^x$, $\lambda > 0$, $x \geq 0$	$p(x \theta) = \frac{1}{x!} e^{\theta s - A(\theta)}$	$\log \lambda$	e^θ	x	$\frac{1}{x!}$	λ	e^θ	$E[x] = \lambda = e^\theta$ $\text{Var}[x] = \lambda = e^\theta$
Bernoulli	$p(x p) = p^x (1-p)^{1-x}$, $p \in (0, 1)$, $x \in \{0, 1\}$	$p(x \theta) = e^{\theta s - A(\theta)}$	$\log \frac{p}{1-p}$	$\frac{e^\theta}{1+e^\theta} = \sigma(\theta)$	x	1	$-\log(1-p)$	$-\log(1 - \sigma(\theta)) = \log(1 + e^\theta)$	$E[x] = p = \sigma(\theta)$ $\text{Var}[x] = p(1-p) = \sigma(\theta)(1 - \sigma(\theta))$
Binomial (fix number of trials n)	$p(x n, p) = \binom{n}{x} p^x (1-p)^{n-x}$, $p \in (0, 1)$, $x \geq 0$	$p(x \theta) = \frac{n!}{x!(n-x)!} e^{\theta s - A(\theta)}$	$\log \frac{p}{1-p}$	$\frac{e^\theta}{1+e^\theta} = \sigma(\theta)$	x	$\frac{n!}{x!(n-x)!}$	$-n \log(1-p)$	$-n \log(1 - \sigma(\theta))$	$E[x] = np = n\sigma(\theta)$ $\text{Var}[x] = np(1-p) = n\sigma(\theta)(1 - \sigma(\theta))$
Negative binomial (fix number of successes k , success rate p)	$p(x k, p) = \binom{x+k-1}{k-1} p^k (1-p)^x$, $p \in (0, 1)$, $x \geq 0$	$p(x \theta) = \frac{(x+k-1)!}{(k-1)!x!} e^{\theta s - A(\theta)}$, $\theta < 0$	$\log(1-p)$	$1 - e^\theta$	x	$\frac{(x+k-1)!}{(k-1)!x!}$	$-k \log p$	$-k \log(1 - e^\theta)$	$E[x] = k \frac{1-p}{p} = k \frac{e^\theta}{1-e^\theta}$ $\text{Var}[x] = k \frac{1-p}{p^2} = k \frac{e^\theta}{(1-e^\theta)^2}$
Multinoulli	$p(x_1, \dots, x_M p_1, \dots, p_M) = p_1^{x_1} \dots p_M^{x_M}$, $p_m \geq 0$, $\sum_{m=1}^M p_m = 1$, $x_m \in \{0, 1\}$, $\sum_{m=1}^M x_m = 1$	$p(x \theta) = \sum_{m=1}^M e^{\theta_1 s_m + \dots + \theta_M s_M - \log C}$, where $C = \sum_{m=1}^M e^{\theta_m}$	$\begin{bmatrix} \log p_1 + \log C \\ \vdots \\ \log p_M + \log C \end{bmatrix}$	$\begin{bmatrix} \frac{\theta_1}{C} \\ \vdots \\ \frac{\theta_M}{C} \end{bmatrix}$	$\begin{bmatrix} x_1 \\ \vdots \\ x_M \end{bmatrix}$	1	$\log C$	$\log C$	$E[x_m = 1] = p_m = \frac{e^{\theta_m}}{C}$ $\text{Var}[x_m = 1] = p_m(1 - p_m) = \frac{e^{\theta_m}}{C} - \frac{\theta_m}{C}$
Multinomial (fix number of trials n)	$p(x_1, \dots, x_M n, p_1, \dots, p_M) = \frac{n!}{\prod_{m=1}^M x_m!} p_1^{x_1} \dots p_M^{x_M}$, $p_m \geq 0$, $\sum_{m=1}^M p_m = 1$, $x_m \in \{0, 1, \dots, n\}$, $\sum_{m=1}^M x_m = n$	$p(x \theta) = \frac{n!}{\prod_{m=1}^M x_m!} e^{\theta_1 s_m + \dots + \theta_M s_M - n \log C}$, where $C = \sum_{m=1}^M e^{\theta_m}$	$\begin{bmatrix} \log p_1 + \log C \\ \vdots \\ \log p_M + \log C \end{bmatrix}$	$\begin{bmatrix} \frac{\theta_1}{C} \\ \vdots \\ \frac{\theta_M}{C} \end{bmatrix}$	$\begin{bmatrix} x_1 \\ \vdots \\ x_M \end{bmatrix}$	$\frac{n!}{\prod_{m=1}^M x_m!}$	$n \log C$	$n \log C$	$E[x_m] = np_m = n \frac{e^{\theta_m}}{C}$ $\text{Var}[x] = np_m(1 - p_m) = n \frac{e^{\theta}}{C} - e^{\theta}$

Table 2: Selected instances of exp-RBM.

Model	Energy Function	Conditional	Gradient
Bernoulli-Bernoulli	$E(\mathbf{a}, \mathbf{h}) = -\mathbf{a}^T \mathbf{x} - \mathbf{b}^T \mathbf{h} - \mathbf{x}^T \mathbf{W} \mathbf{h}$, where $\mathbf{a} = [\log \frac{p_1}{1-p_1}, \dots, \log \frac{p_M}{1-p_M}]^T$, $\mathbf{b} = [\log \frac{p_1}{1-p_1}, \dots, \log \frac{p_K}{1-p_K}]$	$p(\mathbf{a} \mathbf{h}) = \prod_{m=1}^M \mathcal{B}\mathcal{E}(a_m \sigma(\hat{a}_m)), \hat{a}_m = a_m + \mathbf{W}_{m,:} \mathbf{h}$ $p(\mathbf{h} \mathbf{a}) = \prod_{k=1}^K \mathcal{B}\mathcal{E}(h_k \sigma(\hat{b}_k)), \hat{b}_k = b_k + (\mathbf{W}_{:,k})^T \mathbf{a}$	$\Delta \mathbf{a} \approx \frac{1}{N} \sum_{n=1}^N (-\mathbf{a}_n) - \frac{1}{S} \sum_{s=1}^S (-\mathbf{a}_s)$ $\Delta \mathbf{b} \approx \frac{1}{N} \sum_{n=1}^N (-\mathbf{h}_n) - \frac{1}{S} \sum_{s=1}^S (-\mathbf{h}_s)$ $\Delta \mathbf{W} \approx \frac{1}{N} \sum_{n=1}^N (-\mathbf{a}_n \mathbf{h}_n^T) - \frac{1}{S} \sum_{s=1}^S (-\mathbf{a}_s \mathbf{h}_s^T)$
Gaussian-Bernoulli	$E(\mathbf{a}, \mathbf{h}) = -\mathbf{a}^{(1)T} \mathbf{x} - \mathbf{a}^{(2)T} \mathbf{x}^{*2} - \mathbf{b}^T \mathbf{h} - \mathbf{x}^T \mathbf{W} \mathbf{h}$, where $\mathbf{a}^{(1)} = [\mu_1 \lambda_1, \dots, \mu_M \lambda_M]^T$, $\mathbf{a}^{(2)} = [-\frac{\lambda_1}{2}, \dots, -\frac{\lambda_M}{2}]^T$, $\mathbf{a}^{(2)} < \mathbf{0}$, $\mathbf{b} = [\log \frac{p_1}{1-p_1}, \dots, \log \frac{p_K}{1-p_K}]$	$p(\mathbf{a} \mathbf{h}) = \prod_{m=1}^M \mathcal{G}\mathcal{S}(a_m -\frac{\hat{a}_m}{2\sigma_m^2}, (-2\sigma_m^2)^{-1})$, $\hat{a}_m = a_m^{(1)} + \mathbf{W}_{m,:} \mathbf{h}$ $p(\mathbf{h} \mathbf{a}) = \prod_{k=1}^K \mathcal{B}\mathcal{E}(h_k \sigma(\hat{b}_k))$, $\hat{b}_k = b_k + (\mathbf{W}_{:,k})^T \mathbf{a}$	$\Delta \mathbf{a}^{(1)} \approx \frac{1}{N} \sum_{n=1}^N (-\mathbf{a}_n) - \frac{1}{S} \sum_{s=1}^S (-\mathbf{a}_s)$ $\Delta \mathbf{a}^{(2)} \approx \frac{1}{N} \sum_{n=1}^N (-\mathbf{a}_n^{*2}) - \frac{1}{S} \sum_{s=1}^S (-\mathbf{a}_s^{*2})$ $\Delta \mathbf{b} \approx \frac{1}{N} \sum_{n=1}^N (-\mathbf{h}_n) - \frac{1}{S} \sum_{s=1}^S (-\mathbf{h}_s)$ $\Delta \mathbf{W} \approx \frac{1}{N} \sum_{n=1}^N (-\mathbf{a}_n \mathbf{h}_n^T) - \frac{1}{S} \sum_{s=1}^S (-\mathbf{a}_s \mathbf{h}_s^T)$
Gaussian-Bernoulli fix precision λ , model 1	$E(\mathbf{a}, \mathbf{h}) = -\mathbf{a}^T (\lambda * \mathbf{x}) - \sum_{m=1}^M (\log \frac{\sqrt{\lambda_m}}{\sqrt{2\pi}} - \frac{\lambda_m x_m^2}{2}) - \mathbf{b}^T \mathbf{h} - (\lambda * \mathbf{x})^T \mathbf{W} \mathbf{h}$, where $\mathbf{a} = [\mu_1, \dots, \mu_M]^T$, $\mathbf{b} = [\log \frac{p_1}{1-p_1}, \dots, \log \frac{p_K}{1-p_K}]$	$p(\mathbf{a} \mathbf{h}) = \prod_{m=1}^M \mathcal{G}\mathcal{S}(a_m \hat{a}_m, \lambda_m)$, $\hat{a}_m = a_m + \mathbf{W}_{m,:} \mathbf{h}$ $p(\mathbf{h} \mathbf{a}) = \prod_{k=1}^K \mathcal{B}\mathcal{E}(h_k \sigma(\hat{b}_k))$, $\hat{b}_k = b_k + (\mathbf{W}_{:,k})^T (\lambda * \mathbf{a})$	$\Delta \mathbf{a} \approx \frac{1}{N} \sum_{n=1}^N (-\lambda * \mathbf{a}_n) - \frac{1}{S} \sum_{s=1}^S (-\lambda * \mathbf{a}_s)$ $\Delta \mathbf{b} \approx \frac{1}{N} \sum_{n=1}^N (-\mathbf{h}_n) - \frac{1}{S} \sum_{s=1}^S (-\mathbf{h}_s)$ $\Delta \mathbf{W} \approx \frac{1}{N} \sum_{n=1}^N (-\lambda * \mathbf{a}_n) \mathbf{h}_n^T - \frac{1}{S} \sum_{s=1}^S (-\lambda * \mathbf{a}_s) \mathbf{h}_s^T$
Gaussian-Bernoulli fix precision λ , model 2	$E(\mathbf{a}, \mathbf{h}) = -\mathbf{a}^T \mathbf{x} - \sum_{m=1}^M (\log \frac{\sqrt{\lambda_m}}{\sqrt{2\pi}} - \frac{\lambda_m x_m^2}{2}) - \mathbf{b}^T \mathbf{h} - \mathbf{x}^T \mathbf{W} \mathbf{h}$, where $\mathbf{a} = [\mu_1 \lambda_1, \dots, \mu_M \lambda_M]^T$, $\mathbf{b} = [\log \frac{p_1}{1-p_1}, \dots, \log \frac{p_K}{1-p_K}]$	$p(\mathbf{a} \mathbf{h}) = \prod_{m=1}^M \mathcal{G}\mathcal{S}(a_m \frac{\hat{a}_m}{\lambda_m}, \lambda_m)$, $\hat{a}_m = a_m + \mathbf{W}_{m,:} \mathbf{h}$ $p(\mathbf{h} \mathbf{a}) = \prod_{k=1}^K \mathcal{B}\mathcal{E}(h_k \sigma(\hat{b}_k))$, $\hat{b}_k = b_k + (\mathbf{W}_{:,k})^T \mathbf{a}$	$\Delta \mathbf{a} \approx \frac{1}{N} \sum_{n=1}^N (-\mathbf{a}_n) - \frac{1}{S} \sum_{s=1}^S (-\mathbf{a}_s)$ $\Delta \mathbf{b} \approx \frac{1}{N} \sum_{n=1}^N (-\mathbf{h}_n) - \frac{1}{S} \sum_{s=1}^S (-\mathbf{h}_s)$ $\Delta \mathbf{W} \approx \frac{1}{N} \sum_{n=1}^N (-\mathbf{a}_n \mathbf{h}_n^T) - \frac{1}{S} \sum_{s=1}^S (-\mathbf{a}_s \mathbf{h}_s^T)$
Poisson-Bernoulli	$E(\mathbf{a}, \mathbf{h}) = -\mathbf{a}^T \mathbf{x} + \sum_{m=1}^M \log(x_m!) - \mathbf{b}^T \mathbf{h} - \mathbf{x}^T \mathbf{W} \mathbf{h}$, where $\mathbf{a} = [\log \lambda_1, \dots, \log \lambda_M]^T$, $\mathbf{b}_k = [\log \frac{p_1}{1-p_1}, \dots, \log \frac{p_K}{1-p_K}]^T$	$p(\mathbf{a} \mathbf{h}) = \prod_{m=1}^M \mathcal{P}\mathcal{O}(a_m e^{\hat{a}_m})$, $\hat{a}_m = a_m + \mathbf{W}_{m,:} \mathbf{h}$ $p(\mathbf{h} \mathbf{a}) = \prod_{k=1}^K \mathcal{B}\mathcal{E}(h_k \sigma(\hat{b}_k))$, $\hat{b}_k = b_k + (\mathbf{W}_{:,k})^T \mathbf{a}$	$\Delta \mathbf{a} \approx \frac{1}{N} \sum_{n=1}^N (-\mathbf{a}_n) - \frac{1}{S} \sum_{s=1}^S (-\mathbf{a}_s)$ $\Delta \mathbf{b} \approx \frac{1}{N} \sum_{n=1}^N (-\mathbf{h}_n) - \frac{1}{S} \sum_{s=1}^S (-\mathbf{h}_s)$ $\Delta \mathbf{W} \approx \frac{1}{N} \sum_{n=1}^N (-\mathbf{a}_n \mathbf{h}_n^T) - \frac{1}{S} \sum_{s=1}^S (-\mathbf{a}_s \mathbf{h}_s^T)$
Poisson-Binomial	$E(\mathbf{a}, \mathbf{h}) = -\mathbf{a}^T \mathbf{x} + \sum_{m=1}^M \log(x_m!) - \mathbf{b}^T \mathbf{h} - \sum_{k=1}^K \log \frac{n!}{h_k!(n-h_k)!} - \mathbf{x}^T \mathbf{W} \mathbf{h}$, where $\mathbf{a} = [\log \lambda_1, \dots, \log \lambda_M]^T$, $\mathbf{b} = [\log \frac{p_1}{1-p_1}, \dots, \log \frac{p_K}{1-p_K}]^T$	$p(\mathbf{a} \mathbf{h}) = \prod_{m=1}^M \mathcal{P}\mathcal{O}(a_m e^{\hat{a}_m})$, $\hat{a}_m = a_m + \mathbf{W}_{m,:} \mathbf{h}$ $p(\mathbf{h} \mathbf{a}) = \prod_{k=1}^K \mathcal{B}\mathcal{N}(n_k, \sigma(\hat{b}_k))$, $\hat{b}_k = b_k + (\mathbf{W}_{:,k})^T \mathbf{a}$	$\Delta \mathbf{a} \approx \frac{1}{N} \sum_{n=1}^N (-\mathbf{a}_n) - \frac{1}{S} \sum_{s=1}^S (-\mathbf{a}_s)$ $\Delta \mathbf{b} \approx \frac{1}{N} \sum_{n=1}^N (-\mathbf{h}_n) - \frac{1}{S} \sum_{s=1}^S (-\mathbf{h}_s)$ $\Delta \mathbf{W} \approx \frac{1}{N} \sum_{n=1}^N (-\mathbf{a}_n \mathbf{h}_n^T) - \frac{1}{S} \sum_{s=1}^S (-\mathbf{a}_s \mathbf{h}_s^T)$
Negative-Binomial-Bernoulli	$E(\mathbf{a}, \mathbf{h}) = -\mathbf{a}^T \mathbf{x} - \sum_{m=1}^M \log \frac{(sm-1)!}{(sm-1)!x_m!} - \mathbf{b}^T \mathbf{h} - \mathbf{x}^T \mathbf{W} \mathbf{h}$, where $\mathbf{a} = [\log p_1, \dots, \log p_M]^T$, $\mathbf{a} < \mathbf{0}$ $\mathbf{b} = [\log \frac{p_1}{1-p_1}, \dots, \log \frac{p_K}{1-p_K}]^T$	$p(\mathbf{a} \mathbf{h}) = \prod_{m=1}^M \mathcal{N}\mathcal{B}(a_m b_m, 1 - e^{\hat{a}_m})$, $\hat{a}_m = a_m + \mathbf{W}_{m,:} \mathbf{h}$, $\hat{a}_m < 0$ $p(\mathbf{h} \mathbf{a}) = \prod_{k=1}^K \mathcal{B}\mathcal{E}(\sigma(\hat{b}_k))$, $\hat{b}_k = b_k + (\mathbf{W}_{:,k})^T \mathbf{a}$	$\Delta \mathbf{a} \approx \frac{1}{N} \sum_{n=1}^N (-\mathbf{a}_n) - \frac{1}{S} \sum_{s=1}^S (-\mathbf{a}_s)$ $\Delta \mathbf{b} \approx \frac{1}{N} \sum_{n=1}^N (-\mathbf{h}_n) - \frac{1}{S} \sum_{s=1}^S (-\mathbf{h}_s)$ $\Delta \mathbf{W} \approx \frac{1}{N} \sum_{n=1}^N (-\mathbf{a}_n \mathbf{h}_n^T) - \frac{1}{S} \sum_{s=1}^S (-\mathbf{a}_s \mathbf{h}_s^T)$
Multinoulli-Bernoulli	$E(\mathbf{a}, \mathbf{h}) = -\sum_{m=1}^M \mathbf{a}^{(m)T} \mathbf{x}^{(m)} - \mathbf{b}^T \mathbf{h} - \sum_{m=1}^M (\mathbf{a}^{(m)})^T \mathbf{W}^{(m)} \mathbf{h}$, where $\mathbf{a}^{(m)} = [\log p_1^{(m)} + \log C^{(m)}, \dots, \log p_{C^{(m)}}^{(m)} + \log C^{(m)}]^T$, $C^{(m)} = \sum_{c'=1}^{C_m} \exp(a_{c'}^{(m)})$, $\mathbf{b} = [\log \frac{p_1}{1-p_1}, \dots, \log \frac{p_K}{1-p_K}]^T$, $\mathbf{W}^{(m)} \in \mathbb{R}^{C_m \times K}$	$p(\mathbf{a}^{(m)} \mathbf{h}) = \mathcal{M}\mathcal{U}(\mathbf{a}^{(m)} \mathbf{p}^{(m)})$, $\hat{a}_c^{(m)} = a_c^{(m)} + \mathbf{W}_{c,:}^{(m)} \mathbf{h}$, $\mathbf{p}^{(m)} = [\frac{\exp(\hat{a}_1^{(m)})}{\sum_{c'=1}^{C_m} \exp(\hat{a}_{c'}^{(m)})}, \dots, \frac{\exp(\hat{a}_{C_m}^{(m)})}{\sum_{c'=1}^{C_m} \exp(\hat{a}_{c'}^{(m)})}]^T$ $p(\mathbf{h} \mathbf{a}) = \prod_{k=1}^K \mathcal{B}\mathcal{E}(\sigma(\hat{b}_k))$, $\hat{b}_k = b_k + \sum_{m=1}^M (\mathbf{W}_{:,k}^{(m)})^T \mathbf{a}^{(m)}$	$\Delta \mathbf{a}^{(m)} \approx \frac{1}{N} \sum_{n=1}^N (-\mathbf{a}_n^{(m)}) - \frac{1}{S} \sum_{s=1}^S (-\mathbf{a}_s^{(m)})$ $\Delta \mathbf{b} \approx \frac{1}{N} \sum_{n=1}^N (-\mathbf{h}_n) - \frac{1}{S} \sum_{s=1}^S (-\mathbf{h}_s)$ $\Delta \mathbf{W} \approx \frac{1}{N} \sum_{n=1}^N (-\mathbf{a}_n^{(m)} \mathbf{h}_n^T) - \frac{1}{S} \sum_{s=1}^S (-\mathbf{a}_s^{(m)} \mathbf{h}_s^T)$
Multinomial-Bernoulli	$E(\mathbf{a}, \mathbf{h}) = -\mathbf{a}^T \mathbf{x} - \log \frac{n!}{\prod_{m=1}^M x_m!} - \mathbf{b}^T \mathbf{h} - \mathbf{x}^T \mathbf{W} \mathbf{h}$, where $\mathbf{a} = [\log p_1 + \log C, \dots, \log p_M + \log C]^T$, $\mathbf{b} = [\log \frac{p_1}{1-p_1}, \dots, \log \frac{p_K}{1-p_K}]^T$	$p(\mathbf{a} \mathbf{h}) = \mathcal{M}\mathcal{N}(\mathbf{a} \mathbf{n}, \mathbf{p})$, $\hat{a}_m = a_m + \mathbf{W}_{m,:} \mathbf{h}$, $\hat{p}_m = \frac{e^{\hat{a}_m}}{\sum_{m'=1}^M e^{\hat{a}_{m'}}}$ $p(\mathbf{h} \mathbf{a}) = \prod_{k=1}^K \mathcal{B}\mathcal{E}(\sigma(\hat{b}_k))$, $\hat{b}_k = b_k + (\mathbf{W}_{:,k})^T \mathbf{a}$	$\Delta \mathbf{a} \approx \frac{1}{N} \sum_{n=1}^N (-\mathbf{a}_n) - \frac{1}{S} \sum_{s=1}^S (-\mathbf{a}_s)$ $\Delta \mathbf{b} \approx \frac{1}{N} \sum_{n=1}^N (-\mathbf{h}_n) - \frac{1}{S} \sum_{s=1}^S (-\mathbf{h}_s)$ $\Delta \mathbf{W} \approx \frac{1}{N} \sum_{n=1}^N (-\mathbf{a}_n \mathbf{h}_n^T) - \frac{1}{S} \sum_{s=1}^S (-\mathbf{a}_s \mathbf{h}_s^T)$

Table 3: Estimation of log-partition functions of exp-RBMs.

Model	$F(\mathbf{w})$	$F_t(\mathbf{w}) \ln p_t(\mathbf{w})$	$T(\mathbf{w}_t, \mathbf{w}_{t-1})$	$\log Z_A$
Bernoulli-Bernoulli	$-\mathbf{a}^T \mathbf{w}$ $-\sum_{k=1}^K \log(1 + e^{b_k + \mathbf{W}^T_{k,:} \mathbf{w}})$	$-\mathbf{a}^T \mathbf{w}$ $-\sum_{k=1}^K \left(\log(1 + e^{(1-\beta_t)b_k} + \log(1 + e^{\beta_t(b_k + \mathbf{W}^T_{k,:} \mathbf{w})}) \right)$	$p_t(\mathbf{w}_t) p_t(\mathbf{w}_{t-1}) = \mathcal{B}\mathcal{E}((h_t^{(A)})_k \sigma((1-\beta_t)a_m + \beta_t(am + \mathbf{W}^T_{m,:} h_{t-1}^{(B)})))$ $p_t((h_t^{(A)})_k \mathbf{w}_t) = \mathcal{B}\mathcal{E}((h_t^{(A)})_k \sigma((1-\beta_t)b_k))$ $p_t((h_t^{(B)})_k \mathbf{w}_t) = \mathcal{B}\mathcal{E}((h_t^{(B)})_k \sigma(\beta_t(b_k + \mathbf{W}^T_{k,:} \mathbf{w}_t)))$	$\sum_{m=1}^M \log(1 + e^{am}) + \sum_{k=1}^K \log(1 + e^{b_k})$
Gaussian-Bernoulli	$-\mathbf{a}^{(1)T} \mathbf{w} - \mathbf{a}^{(2)T} \mathbf{w}^{\ast 2}$ $-\sum_{k=1}^K \log(1 + e^{b_k + \mathbf{W}^T_{k,:} \mathbf{w}})$	$-\mathbf{a}^{(1)T} \mathbf{w} - \mathbf{a}^{(2)T} \mathbf{w}^{\ast 2}$ $-\sum_{k=1}^K \left(\log(1 + e^{(1-\beta_t)b_k} + \log(1 + e^{\beta_t(b_k + \mathbf{W}^T_{k,:} \mathbf{w})}) \right)$	$p_t(\mathbf{w}_t) p_t(\mathbf{w}_{t-1}) = \mathcal{G}\mathcal{S}((x_t)_m -\frac{(1-\beta_t)a_m + \beta_t(am + \mathbf{W}^T_{m,:} h_{t-1}^{(B)})}{a_m}, (-2a_m)^{-1})$ $p_t((h_t^{(A)})_k \mathbf{w}_t) = \mathcal{B}\mathcal{E}((h_t^{(A)})_k \sigma((1-\beta_t)b_k))$ $p_t((h_t^{(B)})_k \mathbf{w}_t) = \mathcal{B}\mathcal{E}((h_t^{(B)})_k \sigma(\beta_t(b_k + \mathbf{W}^T_{k,:} \mathbf{w}_t)))$	$\sum_{m=1}^M \left(\frac{1}{2} \log \frac{\pi}{-a_m} - \frac{a_m^2}{4a_m} \right) + \sum_{k=1}^K \log(1 + e^{b_k})$
Gaussian-Bernoulli fix precision, model 1	$-\mathbf{a}^T (\boldsymbol{\lambda} \ast \mathbf{w}) - \sum_{m=1}^M \left(\log \frac{\sqrt{\lambda_m}}{\sqrt{2\pi}} - \frac{\lambda_m}{2} x_m^2 \right)$ $-\sum_{k=1}^K \log(1 + e^{b_k + \mathbf{W}^T_{k,:} (\boldsymbol{\lambda} \ast \mathbf{w})})$	$-\mathbf{a}^T (\boldsymbol{\lambda} \ast \mathbf{w}) - \sum_{m=1}^M \left(\log \frac{\sqrt{\lambda_m}}{\sqrt{2\pi}} - \frac{\lambda_m}{2} x_m^2 \right)$ $-\sum_{k=1}^K \left(\log(1 + e^{(1-\beta_t)b_k} + \log(1 + e^{\beta_t(b_k + \mathbf{W}^T_{k,:} (\boldsymbol{\lambda} \ast \mathbf{w}))}) \right)$	$p_t(\mathbf{w}_t) p_t(\mathbf{w}_{t-1}) = \mathcal{G}\mathcal{S}((x_t)_m (1-\beta_t)a_m + \beta_t(am + \mathbf{W}^T_{m,:} h_{t-1}^{(B)}), \lambda_m)$ $p_t((h_t^{(A)})_k \mathbf{w}_t) = \mathcal{B}\mathcal{E}((h_t^{(A)})_k \sigma((1-\beta_t)b_k))$ $p_t((h_t^{(B)})_k \mathbf{w}_t) = \mathcal{B}\mathcal{E}((h_t^{(B)})_k \sigma(\beta_t(b_k + \mathbf{W}^T_{k,:} (\boldsymbol{\lambda} \ast \mathbf{w}_t))))$	$\sum_{m=1}^M \left(\frac{\lambda_m}{2} a_m^2 \right) + \sum_{k=1}^K \log(1 + e^{b_k})$
Gaussian-Bernoulli fix precision, model 2	$-\mathbf{a}^T \mathbf{w} - \sum_{m=1}^M \left(\log \frac{\sqrt{\lambda_m}}{\sqrt{2\pi}} - \frac{\lambda_m}{2} x_m^2 \right)$ $-\sum_{k=1}^K \log(1 + e^{b_k + \mathbf{W}^T_{k,:} \mathbf{w}})$	$-\mathbf{a}^T \mathbf{w} - \sum_{m=1}^M \left(\log \frac{\sqrt{\lambda_m}}{\sqrt{2\pi}} - \frac{\lambda_m}{2} x_m^2 \right)$ $-\sum_{k=1}^K \left(\log(1 + e^{(1-\beta_t)b_k} + \log(1 + e^{\beta_t(b_k + \mathbf{W}^T_{k,:} \mathbf{w})}) \right)$	$p_t(\mathbf{w}_t) p_t(\mathbf{w}_{t-1}) = \mathcal{G}\mathcal{S}((x_t)_m \frac{(1-\beta_t)a_m + \beta_t(am + \mathbf{W}^T_{m,:} h_{t-1}^{(B)})}{\lambda}, \lambda_m)$ $p_t((h_t^{(A)})_k \mathbf{w}_t) = \mathcal{B}\mathcal{E}((h_t^{(A)})_k \sigma((1-\beta_t)b_k))$ $p_t((h_t^{(B)})_k \mathbf{w}_t) = \mathcal{B}\mathcal{E}((h_t^{(B)})_k \sigma(\beta_t(b_k + \mathbf{W}^T_{k,:} \mathbf{w}_t)))$	$\sum_{m=1}^M \left(\frac{a_m^2}{2\lambda_m} \right) + \sum_{k=1}^K \log(1 + e^{b_k})$
Poisson-Bernoulli	$-\mathbf{a}^T \mathbf{w} + \sum_{m=1}^M \log(x_m!)$ $-\sum_{k=1}^K \log(1 + e^{b_k + \mathbf{W}^T_{k,:} \mathbf{w}})$	$-\mathbf{a}^T \mathbf{w} + \sum_{m=1}^M \log(x_m!)$ $-\sum_{k=1}^K \left(\log(1 + e^{(1-\beta_t)b_k} + \log(1 + e^{\beta_t(b_k + \mathbf{W}^T_{k,:} \mathbf{w})}) \right)$	$p_t(\mathbf{w}_t) p_t(\mathbf{w}_{t-1}) = \mathcal{B}\mathcal{E}((h_t^{(A)})_k \sigma((1-\beta_t)b_k))$ $p_t((h_t^{(B)})_k \mathbf{w}_t) = \mathcal{B}\mathcal{E}((h_t^{(B)})_k \sigma(\beta_t(b_k + \mathbf{W}^T_{k,:} \mathbf{w}_t)))$	$\sum_{m=1}^M e^{am} + \sum_{k=1}^K \log(1 + e^{b_k})$
Poisson-Binomial	$-\mathbf{a}^T \mathbf{w} + \sum_{m=1}^M \log(x_m!)$ $-\sum_{k=1}^K n_k \log(1 + e^{b_k + \mathbf{W}^T_{k,:} \mathbf{w}})$	$-\mathbf{a}^T \mathbf{w} + \sum_{m=1}^M \log(x_m!)$ $-\sum_{k=1}^K \left(n_k \log(1 + e^{(1-\beta_t)b_k} + \log(1 + e^{\beta_t(b_k + \mathbf{W}^T_{k,:} \mathbf{w})}) \right)$	$p_t(\mathbf{w}_t) p_t(\mathbf{w}_{t-1}) = \mathcal{P}\mathcal{O}((x_t)_m e^{(1-\beta_t)a_m + \beta_t(am + \mathbf{W}^T_{m,:} h_{t-1}^{(B)})})$ $p_t((h_t^{(A)})_k \mathbf{w}_t) = \mathcal{B}\mathcal{E}((h_t^{(A)})_k \sigma((1-\beta_t)b_k))$ $p_t((h_t^{(B)})_k \mathbf{w}_t) = \mathcal{B}\mathcal{E}((h_t^{(B)})_k \sigma(\beta_t(b_k + \mathbf{W}^T_{k,:} \mathbf{w}_t)))$	$\sum_{m=1}^M e^{am} + \sum_{k=1}^K n_k \log(1 + e^{b_k})$
Negative-Binomial- Bernoulli	$-\mathbf{a}^T \mathbf{w} - \sum_{m=1}^M \log \frac{(x_m + sm - 1)!}{(sm-1)! x_m!}$ $-\sum_{k=1}^K \log(1 + e^{b_k + \mathbf{W}^T_{k,:} \mathbf{w}})$	$-\mathbf{a}^T \mathbf{w} - \sum_{m=1}^M \log \frac{(x_m + sm - 1)!}{(sm-1)! x_m!}$ $-\sum_{k=1}^K \left(n_k \log(1 + e^{(1-\beta_t)b_k} + \log(1 + e^{\beta_t(b_k + \mathbf{W}^T_{k,:} \mathbf{w})}) \right)$	$p_t(\mathbf{w}_t) p_t(\mathbf{w}_{t-1}) = \mathcal{NB}((x_t)_m sm, 1 - e^{(1-\beta_t)a_m + \beta_t(am + \mathbf{W}^T_{m,:} h_{t-1}^{(B)})})$ $p_t((h_t^{(A)})_k \mathbf{w}_t) = \mathcal{B}\mathcal{E}((h_t^{(A)})_k \sigma((1-\beta_t)b_k))$ $p_t((h_t^{(B)})_k \mathbf{w}_t) = \mathcal{B}\mathcal{E}((h_t^{(B)})_k \sigma(\beta_t(b_k + \mathbf{W}^T_{k,:} \mathbf{w}_t)))$	$\sum_{m=1}^M -s_m \log(1 - e^{am}) + \sum_{k=1}^K \log(1 + e^{b_k})$
Multinomial- Bernoulli	$-\sum_{m=1}^M \mathbf{a}^{(m)T} \mathbf{w}^{(m)}$ $-\sum_{k=1}^K \log(1 + e^{b_k + \sum_{m=1}^M (\mathbf{W}^{(m)T}_{k,:} \mathbf{w}^{(m)})})$	$-\sum_{m=1}^M \mathbf{a}^{(m)T} \mathbf{w}^{(m)}$ $-\sum_{k=1}^K \left(\log(1 + e^{(1-\beta_t)b_k} + \log(1 + e^{\beta_t(b_k + \mathbf{W}^T_{k,:} \mathbf{w})}) \right)$ $+\log(1 + e^{\beta_t(b_k + \sum_{m=1}^M (\mathbf{W}^{(m)T}_{k,:} \mathbf{w}^{(m)})})$	$p_t(\mathbf{w}_t) p_t(\mathbf{w}_{t-1}) = \mathcal{ML}(\mathbf{w}_t \left[\frac{e^{(1-\beta_t)a_c^{(m)} + \beta_t(c_c^{(m)} h_{t-1}^{(m)})}}{\sum_{c'=1}^C e^{(1-\beta_t)a_{c'}^{(m)} + \beta_t(c_{c'}^{(m)} h_{t-1}^{(m)})}} \right])$ $p_t((h_t^{(A)})_k \mathbf{w}_t) = \mathcal{B}\mathcal{E}((h_t^{(A)})_k \sigma((1-\beta_t)b_k))$ $p_t((h_t^{(B)})_k \mathbf{w}_t) = \mathcal{B}\mathcal{E}((h_t^{(B)})_k \sigma(\beta_t(b_k + \sum_{m=1}^M (\mathbf{W}^{(m)T}_{k,:} \mathbf{w}_t^{(m)})))$	$\sum_{m=1}^M \log \sum_{c'=1}^C e^{c' a_c^{(m)}} + \sum_{k=1}^K \log(1 + e^{b_k})$
Multinomial- Bernoulli	$-\mathbf{a}^T \mathbf{w} - \log \frac{M^{n1}}{\prod_{m=1}^M x_m!}$ $-\sum_{k=1}^K \log(1 + e^{b_k + \mathbf{W}^T_{k,:} \mathbf{w}})$	$-\mathbf{a}^T \mathbf{w} - \log \frac{M^{n1}}{\prod_{m=1}^M x_m!}$ $-\sum_{k=1}^K \left(\log(1 + e^{(1-\beta_t)b_k} + \log(1 + e^{\beta_t(b_k + \mathbf{W}^T_{k,:} \mathbf{w})}) \right)$	$p_t(\mathbf{w}_t) p_t(\mathbf{w}_{t-1}) = \mathcal{MN}(\mathbf{w}_t n, \left[\frac{\sum_{m'=1}^M e^{(1-\beta_t)a_{m'} + \beta_t(a_{m'} + \mathbf{W}^T_{m',:} h_{t-1}^{(B)})}}{\sum_{m'=1}^M e^{(1-\beta_t)a_{m'} + \beta_t(a_{m'} + \mathbf{W}^T_{m',:} h_{t-1}^{(B)})}} \right])$ $p_t((h_t^{(A)})_k \mathbf{w}_t) = \mathcal{B}\mathcal{E}((h_t^{(A)})_k \sigma((1-\beta_t)b_k))$ $p_t((h_t^{(B)})_k \mathbf{w}_t) = \mathcal{B}\mathcal{E}((h_t^{(B)})_k \sigma(\beta_t(b_k + \mathbf{W}^T_{k,:} \mathbf{w}_t)))$	$n \log \sum_{m=1}^M e^{a_m} + \sum_{k=1}^K \log(1 + e^{b_k})$

Table 4: Selected instances of exp-HM and exp-DBN.

Model	$p(\mathbf{a} \mathbf{h}^{(1)})$	Gradient for Exp-HMs	Gradient for Exp-DBNs
Bernoulli-Bernoulli	$\prod_{m=1}^M \mathcal{BE}(x_m \sigma(\hat{a}_m)), \hat{a}_m^{(G)} = a_m^{(G)} + \mathbf{W}_{m,:}^{(G,1)} \mathbf{h}^{(1)}$	$\Delta_{\mathbf{a}}(G,1) \approx -(\mathbf{a} - \langle \mathbf{a} \rangle), \langle \mathbf{a} \rangle = \sigma(\hat{\mathbf{a}}^{(G)}), \hat{\mathbf{a}}^{(G)} = \mathbf{a}^{(G)} + \mathbf{W}^{(G,1)} \mathbf{h}^{(1)}$ $\Delta_{\mathbf{b}}(G,l) \approx -(\mathbf{h}^{(l)} - \langle \mathbf{h}^{(l)} \rangle), \langle \mathbf{h}^{(l)} \rangle = \sigma(\hat{\mathbf{b}}^{(G,l)}), \hat{\mathbf{b}}^{(G,l)} = \mathbf{b}_k^{(G,l)} + \mathbf{W}^{(G,l+1)} \mathbf{h}^{(l+1)}, 1 \leq l \leq L-2$ $\Delta_{\mathbf{b}}(G,L) \approx -(\mathbf{h}^{(L)} - \langle \mathbf{h}^{(L)} \rangle), \langle \mathbf{h}^{(L)} \rangle = \sigma(\hat{\mathbf{b}}^{(G,L)})$ $\Delta_{\mathbf{W}}(G,1) \approx -(\mathbf{a} - \langle \mathbf{a} \rangle) \mathbf{h}^{(1)\top}$ $\Delta_{\mathbf{W}}(G,l) \approx -(\mathbf{h}^{(l-1)} - \langle \mathbf{h}^{(l-1)} \rangle) \mathbf{h}^{(l)\top}, 2 \leq l \leq L$	$\Delta_{\mathbf{a}}(G) \approx -(\mathbf{a} - \langle \mathbf{a} \rangle), \langle \mathbf{a} \rangle = \sigma(\hat{\mathbf{a}}^{(G)}), \hat{\mathbf{a}}^{(G)} = \mathbf{a}^{(G)} + \mathbf{W}^{(G,1)} \mathbf{h}^{(1)}$ $\Delta_{\mathbf{b}}(G,l) \approx -(\mathbf{h}^{(l)} - \langle \mathbf{h}^{(l)} \rangle), \langle \mathbf{h}^{(l)} \rangle = \sigma(\hat{\mathbf{b}}^{(G,l)}), \hat{\mathbf{b}}^{(G,l)} = \mathbf{b}_k^{(G,l)} + \mathbf{W}^{(G,l+1)} \mathbf{h}^{(l+1)}, 1 \leq l \leq L-2$ $\Delta_{\mathbf{W}}(G,1) \approx -(\mathbf{a} - \langle \mathbf{a} \rangle) \mathbf{h}^{(1)\top}$ $\Delta_{\mathbf{W}}(G,l) \approx -(\mathbf{h}^{(l-1)} - \langle \mathbf{h}^{(l-1)} \rangle) \mathbf{h}^{(l)\top}, 2 \leq l \leq L-1$ $\Delta_{\mathbf{b}}(G,L-1) \approx -(\mathbf{h}^{(G,L-1)} - \hat{\mathbf{h}}_k^{(G,L-1)}, \{\hat{\mathbf{h}}_k^{(G,L-1)}, \hat{\mathbf{h}}_k^{(G,L)}\} \sim p(\mathbf{h}^{(G,L-1)}, \mathbf{h}^{(G,L)})$ $\Delta_{\mathbf{b}}(G,L) \approx -(\mathbf{h}^{(G,L)} - \hat{\mathbf{h}}_k^{(G,L)})$ $\Delta_{\mathbf{W}}(G,L) \approx -(\mathbf{h}^{(L-1)} \mathbf{h}^{(L)\top} - \hat{\mathbf{h}}^{(L-1)} \hat{\mathbf{h}}^{(L)\top})$
Gaussian-Bernoulli	$\prod_{m=1}^M \mathcal{GS}(x_m -\frac{\hat{a}_m^{(G,1)}}{2\sigma_m^{(G,2)}}, (-2\sigma_m^{(2)})^{-1}), \hat{a}_m^{(G,1)} = a_m^{(G,1)} + \mathbf{W}_{m,:}^{(G,1)} \mathbf{h}^{(1)}$	$\Delta_{\mathbf{a}}(G,1) \approx -(\mathbf{a} - \langle \mathbf{a} \rangle), \langle \mathbf{a} \rangle = -\frac{\hat{\mathbf{a}}^{(G,1)}}{2\sigma_m^{(G,2)}}, \hat{\mathbf{a}}^{(G,1)} = \mathbf{a}^{(G,1)} + \mathbf{W}^{(G,1)} \mathbf{h}^{(1)}$ $\Delta_{\mathbf{a}}(G,2) \approx -(\mathbf{a}^{*2} - \langle \mathbf{a}^{*2} \rangle), \langle \mathbf{a}^{*2} \rangle = \frac{(\hat{\mathbf{a}}^{(G,1)})^2}{4(\sigma_m^{(G,2)})^2} - \frac{1}{2\sigma_m^{(G,2)}}$ $\Delta_{\mathbf{W}}(G,1) \approx -(\mathbf{a} - \langle \mathbf{a} \rangle) \mathbf{h}^{(1)\top}$ $\Delta_{\mathbf{b}}(R,1) \approx -(\mathbf{h}^{(1)} - \langle \mathbf{h}^{(1)} \rangle), \langle \mathbf{h}^{(1)} \rangle = \sigma(\hat{\mathbf{b}}^{(R,1)}), \hat{\mathbf{b}}^{(R,1)} = \mathbf{b}^{(R,1)} + \mathbf{W}^{(R,1)} \mathbf{a}$ $\Delta_{\mathbf{W}}(R,1) \approx -(\mathbf{h}^{(1)} - \langle \mathbf{h}^{(1)} \rangle) \mathbf{a}^\top$ the rest is the same as in Bernoulli-Bernoulli HM	$\Delta_{\mathbf{a}}(G,1) \approx -(\mathbf{a} - \langle \mathbf{a} \rangle), \langle \mathbf{a} \rangle = -\frac{\hat{\mathbf{a}}^{(G,1)}}{2\sigma_m^{(G,2)}}, \hat{\mathbf{a}}^{(G,1)} = \mathbf{a}^{(G,1)} + \mathbf{W}^{(G,1)} \mathbf{h}^{(1)}$ $\Delta_{\mathbf{a}}(G,2) \approx -(\mathbf{a}^{*2} - \langle \mathbf{a}^{*2} \rangle), \langle \mathbf{a}^{*2} \rangle = \frac{(\hat{\mathbf{a}}^{(G,1)})^2}{4(\sigma_m^{(G,2)})^2} - \frac{1}{2\sigma_m^{(G,2)}}$ $\Delta_{\mathbf{W}}(G,1) \approx -(\mathbf{a} - \langle \mathbf{a} \rangle) \mathbf{h}^{(1)\top}$ $\Delta_{\mathbf{b}}(R,1) \approx -(\mathbf{h}^{(1)} - \langle \mathbf{h}^{(1)} \rangle), \langle \mathbf{h}^{(1)} \rangle = \sigma(\hat{\mathbf{b}}^{(R,1)}), \hat{\mathbf{b}}^{(R,1)} = \mathbf{b}^{(R,1)} + \mathbf{W}^{(R,1)} \mathbf{a}$ $\Delta_{\mathbf{W}}(R,1) \approx -(\mathbf{h}^{(1)} - \langle \mathbf{h}^{(1)} \rangle) \mathbf{a}^\top$ the rest is the same as in Bernoulli-Bernoulli DBN
Poisson-Bernoulli	$\prod_{m=1}^M \mathcal{PO}(x_m e^{\hat{a}_m^{(G)}}), \hat{a}_m^{(G)} = a_m^{(G)} + \mathbf{W}_{m,:}^{(G,1)} \mathbf{h}^{(1)}$	$\Delta_{\mathbf{a}}(G) \approx -(\mathbf{a} - \langle \mathbf{a} \rangle), \langle \mathbf{a} \rangle = e^{\hat{\mathbf{a}}^{(G)}}, \hat{\mathbf{a}}^{(G)} = \mathbf{a}^{(G)} + \mathbf{W}^{(G,1)} \mathbf{h}^{(1)}$ $\Delta_{\mathbf{W}}(G,1) \approx -(\mathbf{a} - \langle \mathbf{a} \rangle) \mathbf{h}^{(1)\top}$ $\Delta_{\mathbf{b}}(R,1) \approx -(\mathbf{h}^{(1)} - \langle \mathbf{h}^{(1)} \rangle), \langle \mathbf{h}^{(1)} \rangle = \sigma(\hat{\mathbf{b}}^{(R,1)}), \hat{\mathbf{b}}^{(R,1)} = \mathbf{b}^{(R,1)} + \mathbf{W}^{(R,1)} \mathbf{a}$ $\Delta_{\mathbf{W}}(R,1) \approx -(\mathbf{h}^{(1)} - \langle \mathbf{h}^{(1)} \rangle) \mathbf{a}^\top$ the rest is the same as in Bernoulli-Bernoulli HM	$\Delta_{\mathbf{a}}(G) \approx -(\mathbf{a} - \langle \mathbf{a} \rangle), \langle \mathbf{a} \rangle = e^{\hat{\mathbf{a}}^{(G)}}, \hat{\mathbf{a}}^{(G)} = \mathbf{a}^{(G)} + \mathbf{W}^{(G,1)} \mathbf{h}^{(1)}$ $\Delta_{\mathbf{W}}(G,1) \approx -(\mathbf{a} - \langle \mathbf{a} \rangle) \mathbf{h}^{(1)\top}$ $\Delta_{\mathbf{b}}(R,1) \approx -(\mathbf{h}^{(1)} - \langle \mathbf{h}^{(1)} \rangle), \langle \mathbf{h}^{(1)} \rangle = \sigma(\hat{\mathbf{b}}^{(R,1)}), \hat{\mathbf{b}}^{(R,1)} = \mathbf{b}^{(R,1)} + \mathbf{W}^{(R,1)} \mathbf{a}$ $\Delta_{\mathbf{W}}(R,1) \approx -(\mathbf{h}^{(1)} - \langle \mathbf{h}^{(1)} \rangle) \mathbf{a}^\top$ the rest is the same as in Bernoulli-Bernoulli DBN
Negative-Binomial-Bernoulli	$\prod_{m=1}^M \mathcal{NB}(x_m s_m, 1 - e^{\hat{a}_m^{(G)}}), \hat{a}_m^{(G)} = a_m^{(G)} + \mathbf{W}_{m,:}^{(G,1)} \mathbf{h}^{(1)}$	$\Delta_{\mathbf{a}}(G) \approx -(\mathbf{a} - \langle \mathbf{a} \rangle), \langle \mathbf{a} \rangle = sm \frac{e^{\hat{\mathbf{a}}^{(G)}}}{1 - e^{\hat{\mathbf{a}}^{(G)}}}, \hat{\mathbf{a}}^{(G)} = \mathbf{a}^{(G)} + \mathbf{W}^{(G,1)} \mathbf{h}^{(1)}$ $\Delta_{\mathbf{W}}(G,1) \approx -(\mathbf{a} - \langle \mathbf{a} \rangle) \mathbf{h}^{(1)\top}$ $\Delta_{\mathbf{b}}(R,1) \approx -(\mathbf{h}^{(1)} - \langle \mathbf{h}^{(1)} \rangle), \langle \mathbf{h}^{(1)} \rangle = \sigma(\hat{\mathbf{b}}^{(R,1)}), \hat{\mathbf{b}}^{(R,1)} = \mathbf{b}^{(R,1)} + \mathbf{W}^{(R,1)} \mathbf{a}$ $\Delta_{\mathbf{W}}(R,1) \approx -(\mathbf{h}^{(1)} - \langle \mathbf{h}^{(1)} \rangle) \mathbf{a}^\top$ the rest is the same as in Bernoulli-Bernoulli HM	$\Delta_{\mathbf{a}}(G) \approx -(\mathbf{a} - \langle \mathbf{a} \rangle), \langle \mathbf{a} \rangle = sm \frac{e^{\hat{\mathbf{a}}^{(G)}}}{1 - e^{\hat{\mathbf{a}}^{(G)}}}, \hat{\mathbf{a}}^{(G)} = \mathbf{a}^{(G)} + \mathbf{W}^{(G,1)} \mathbf{h}^{(1)}$ $\Delta_{\mathbf{W}}(G,1) \approx -(\mathbf{a} - \langle \mathbf{a} \rangle) \mathbf{h}^{(1)\top}$ $\Delta_{\mathbf{b}}(R,1) \approx -(\mathbf{h}^{(1)} - \langle \mathbf{h}^{(1)} \rangle), \langle \mathbf{h}^{(1)} \rangle = \sigma(\hat{\mathbf{b}}^{(R,1)}), \hat{\mathbf{b}}^{(R,1)} = \mathbf{b}^{(R,1)} + \mathbf{W}^{(R,1)} \mathbf{a}$ $\Delta_{\mathbf{W}}(R,1) \approx -(\mathbf{h}^{(1)} - \langle \mathbf{h}^{(1)} \rangle) \mathbf{a}^\top$ the rest is the same as in Bernoulli-Bernoulli DBN
Multinomial-Bernoulli	$p(\mathbf{a}_c^{(m)} \mathbf{h}^{(1)}) = \frac{\mathcal{ML}(\mathbf{a}^{(m)} \hat{\mathbf{p}}^{(m)})}{\hat{a}_c^{(G,m)} = \sigma_c^{(G,m)} + \mathbf{W}_{c,:}^{(G,m)}}, \hat{a}_m^{(G)} = a_m^{(G)} + \mathbf{W}_{m,:}^{(G,1)} \mathbf{h}^{(1)}, \hat{\mathbf{p}}^{(m)} = [\frac{\exp(\hat{a}_1^{(G,m)})}{\sum_{c'=1}^{C_m} \exp(\hat{a}_{c'}^{(G,m)}), \dots, \frac{\exp(\hat{a}_{C_m}^{(G,m)})}{\sum_{c'=1}^{C_m} \exp(\hat{a}_{c'}^{(G,m)})}]^\top$	$\Delta_{\mathbf{a}}(G,m) \approx -(\mathbf{a}^{(m)} - \langle \mathbf{a}^{(m)} \rangle), \langle \mathbf{a}^{(m)} \rangle = \hat{\mathbf{p}}^{(G,m)}, \hat{\mathbf{a}}^{(G,m)} = \mathbf{a}^{(G,m)} + \mathbf{W}^{(G,1,m)} \mathbf{h}^{(1)},$ $\hat{\mathbf{p}}^{(m)} = [\frac{\exp(\hat{a}_1^{(G,m)})}{\sum_{c'=1}^{C_m} \exp(\hat{a}_{c'}^{(G,m)}), \dots, \frac{\exp(\hat{a}_{C_m}^{(G,m)})}{\sum_{c'=1}^{C_m} \exp(\hat{a}_{c'}^{(G,m)})}]^\top$ $\Delta_{\mathbf{W}}(G,1,m) \approx -(\mathbf{a}^{(m)} - \langle \mathbf{a}^{(m)} \rangle) \mathbf{h}^{(1)\top}$ $\Delta_{\mathbf{b}}(R,1) \approx -(\mathbf{h}^{(1)} - \langle \mathbf{h}^{(1)} \rangle), \langle \mathbf{h}^{(1)} \rangle = \sigma(\hat{\mathbf{b}}^{(R,1)}), \hat{\mathbf{b}}^{(R,1)} = \mathbf{b}^{(R,1)} + \sum_{m'=1}^M \mathbf{W}^{(R,1,m')} \mathbf{a}^{(m')\top}$ $\Delta_{\mathbf{W}}(R,1,m) \approx -(\mathbf{h}^{(1)} - \langle \mathbf{h}^{(1)} \rangle) \mathbf{a}^{(m)\top}$ the rest is the same as in Bernoulli-Bernoulli HM	$\Delta_{\mathbf{a}}(G,m) \approx -(\mathbf{a}^{(m)} - \langle \mathbf{a}^{(m)} \rangle), \langle \mathbf{a}^{(m)} \rangle = \hat{\mathbf{p}}^{(G,m)}, \hat{\mathbf{a}}^{(G,m)} = \mathbf{a}^{(G,m)} + \mathbf{W}^{(G,1,m)} \mathbf{h}^{(1)},$ $\hat{\mathbf{p}}^{(m)} = [\frac{\exp(\hat{a}_1^{(G,m)})}{\sum_{c'=1}^{C_m} \exp(\hat{a}_{c'}^{(G,m)}), \dots, \frac{\exp(\hat{a}_{C_m}^{(G,m)})}{\sum_{c'=1}^{C_m} \exp(\hat{a}_{c'}^{(G,m)})}]^\top$ $\Delta_{\mathbf{W}}(G,1,m) \approx -(\mathbf{a}^{(m)} - \langle \mathbf{a}^{(m)} \rangle) \mathbf{h}^{(1)\top}$ $\Delta_{\mathbf{b}}(R,1) \approx -(\mathbf{h}^{(1)} - \langle \mathbf{h}^{(1)} \rangle), \langle \mathbf{h}^{(1)} \rangle = \sigma(\hat{\mathbf{b}}^{(R,1)}), \hat{\mathbf{b}}^{(R,1)} = \mathbf{b}^{(R,1)} + \sum_{m'=1}^M \mathbf{W}^{(R,1,m')} \mathbf{a}^{(m')\top}$ $\Delta_{\mathbf{W}}(R,1,m) \approx -(\mathbf{h}^{(1)} - \langle \mathbf{h}^{(1)} \rangle) \mathbf{a}^{(m)\top}$ the rest is the same as in Bernoulli-Bernoulli DBN