

補助関数法に基づく制約付きボルツマンマシンの学習アルゴリズム

高宗 典玄[†] 亀岡 弘和^{†,††}

[†] 東京大学大学院情報理工学系研究科
〒113-8656 東京都文京区本郷 7-3-1

^{††} 日本電信電話株式会社 NTT コミュニケーション科学基礎研究所
〒243-0198 神奈川県厚木市森の里若宮 3-1
E-mail: †{takamune,kameoka}@hil.t.u-tokyo.ac.jp

あらまし 深層学習の重要な一要素として、レイヤーワイズの pre-training がある。レイヤーワイズの pre-training の一つとして制約付き Boltzmann マシン (RBM) が有名である。RBM には Bernoulli-Bernoulli 型と Gaussian-Bernoulli 型があり、従来法の学習アルゴリズムとして Contrastive Divergence 法が有名である。本発表では Bernoulli-Bernoulli 型、Gaussian-Bernoulli 型の両方の最尤学習アルゴリズムと、最適化規準として新たに最大再構築確率を導入し、その学習アルゴリズムに焦点を当て、経験的に高速で安定に収束する補助関数法による新たな更新アルゴリズムの導出を行う。そして、人工データによる収束性能の比較実験を行い、その挙動に対して議論する。
キーワード 深層学習、制約付き Boltzmann マシン、補助関数法、最大再構築確率学習

Training Algorithm for Restricted Boltzmann Machines Using Auxiliary Function Approach

Norihiro TAKAMUNE[†] and Hirokazu KAMEOKA^{†,††}

[†] Graduate School of Information Science and Technology, The University of Tokyo
Hongo 7-3-1, Bunkyo-ku, Tokyo, 113-8656 Japan

^{††} NTT Communication Science Laboratories, Nippon Telegraph and Telephone Corporation
Morinosato Wakamiya 3-1, Atsugi, Kanagawa, 243-0198, Japan
E-mail: †{takamune,kameoka}@hil.t.u-tokyo.ac.jp

Abstract Layerwise pre-training is one of important elements for deep learning, and Restricted Boltzmann Machines (RBMs) is popular layerwise pre-training method. At present, the most popular training algorithm for RBMs is the Contrastive Divergence (CD) learning algorithm. We propose deriving a new training algorithm based on an auxiliary function approach for RBMs using the likelihood and the reconstruction probability of observations as the optimization criterion. Through an experiment on parameter training of few RBMs, we confirmed that the present algorithm outperformed the CD algorithm in terms of the convergence speed and the reconstruction error when used as an autoencoder.

Key words Deep Learning, Restricted Boltzmann Machine, Auxiliary Function Approach, Maximum Reconstruction Probability Training

1. はじめに

近年、深層学習 (Deep learning) の有効性は画像認識や音声認識をはじめ、様々な分野で示されている。深層学習において、大量のデータの下で学習をいかに効率的に行えるかは重要な課題の一つである。その中で、Hinton らによってレイヤーワイズの pre-training が最終的な Deep Neural Network(DNN) を

学習するときのパラメータの初期値として高い効果が得られることが分かってきた [1], [2] .

レイヤーワイズの pre-training として制約付き Boltzmann マシン (Restricted Boltzmann Machine; RBM) [3] が有名であり、DNN の最下層である実数の観測データを 0 または 1 のバイナリの特徴量に変換する層の学習モデルとして Gaussian-Bernoulli 型 RBM, DNN の中間層としてバイナリからバイナリに変換

する層の学習モデルとして Bernoulli-Bernoulli 型 RBM が存在する。これらを学習するアルゴリズムは Contrastive Divergence(CD) 法 [1], [4] が非常に有名であるが, pre-training にかかる時間と最終的に DNN を学習する時間は pre-training のほうが大きくなる場合があり, pre-training をいかに効率的に行えるかが DNN の全体の学習にかかる計算時間に直結する。

我々の研究室では, これまで様々な音響信号処理問題における最適化問題に対し, 補助関数法と呼ぶ原理に基づく最適化アルゴリズムを導出し, その効果を示してきた (例えば [5])。RBM の学習においても補助関数法に基づく学習則を導出することができれば, DBN の初期学習方式として高い効果を発揮する可能性がある。以上の動機のもと, 本発表では Bernoulli-Bernoulli 型 RBM と Gaussian-Bernoulli 型 RBM の学習問題に焦点を当て, 補助関数法に基づく新しい学習則を提案する。

2. 制約付き Boltzmann マシン [3]

2.1 RBM 学習における目的関数

RBM は, Fig. 1 で示されるように完全 2 部グラフの無向グラフの構造を持ち, 観測される状態を可視層, 背後にある状態を隠れ層と呼ぶ。このとき, 可視層同士や隠れ層同士には結合が無いため “制約付き” と呼ばれる。RBM には, 可視層の状態と隠れ層の状態が 0 または 1 のバイナリ値をとる Bernoulli-Bernoulli 型, 可視層の状態が実数値を取り, 隠れ層の状態が 0 または 1 のバイナリ値をとる Gaussian-Bernoulli 型と呼ばれるものがある。

そこで, まず Bernoulli-Bernoulli 型 RBM について説明する。可視層の状態数を I , 可視層の状態を $\mathbf{v} = \{v_i\} \in \{0, 1\}^I$, 隠れ層の状態数を J , 隠れ層の状態を $\mathbf{h} = \{h_j\} \in \{0, 1\}^J$ とすると可視層の状態と隠れ層の状態を確率変数とした Bernoulli-Bernoulli 型 RBM の同時確率は

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

で定義される。ここで,

$$E(\mathbf{v}, \mathbf{h} | \Theta) = -\sum_i b_i^V v_i - \sum_j b_j^H h_j - \sum_{i,j} W_{ij} v_i h_j, \quad (2)$$

$$Z(\Theta) = \sum_{\mathbf{v}} \sum_{\mathbf{h}} \exp(-E(\mathbf{v}, \mathbf{h} | \Theta)) \quad (3)$$

であり, $\Theta = \{b_i^V, b_j^H, W_{ij}\}$ は分布パラメータである。また, $\sum_{\mathbf{v}}$ や $\sum_{\mathbf{h}}$ は \mathbf{v} や \mathbf{h} のすべてのバイナリパターンに関する和を表す。

ここで, 可視層同士, 隠れ層同士の依存関係が無いことから,

$$p(\mathbf{v} | \mathbf{h}, \Theta) = \prod_i p(v_i | \mathbf{h}, \Theta), \quad (4)$$

$$p(\mathbf{h} | \mathbf{v}, \Theta) = \prod_j p(h_j | \mathbf{v}, \Theta) \quad (5)$$

という関係があり, 式 (1) から

$$p(v_i = 1 | \mathbf{h}, \Theta) = \frac{1}{1 + \exp(-b_i^V - \sum_j W_{ij} h_j)}, \quad (6)$$

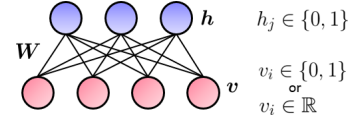


図 1 RBM のグラフ表現

$$p(h_j = 1 | \mathbf{v}, \Theta) = \frac{1}{1 + \exp(-b_j^H - \sum_i W_{ij} v_i)} \quad (7)$$

となる。このように可視層の条件付分布が Bernoulli 分布, 隠れ層の条件付分布が Bernoulli 分布となるため, Bernoulli-Bernoulli 型と呼ばれる。

次に Gaussian-Bernoulli 型 RBM について説明する。可視層の状態数, 隠れ層の状態数, 隠れ層の状態は Bernoulli-Bernoulli 型と同様に $I, J, \mathbf{h} = \{h_j\} \in \{0, 1\}^J$ とし, 可視層の状態は実数値をとるので, $\mathbf{v} = \{v_i\} \in (-\infty, \infty)^I$ とおくと可視層の状態と隠れ層の状態を確率変数とした Gaussian-Bernoulli 型 RBM の同時確率分布は

$$p(\mathbf{v}, \mathbf{h} | \Theta) = \frac{\exp(-E(\mathbf{v}, \mathbf{h} | \Theta))}{Z(\Theta)} \quad (8)$$

で定義される。ここで,

$$E(\mathbf{v}, \mathbf{h} | \Theta) = \sum_i \frac{v_i^2}{2\sigma_i^2} - \sum_i \frac{b_i^V v_i}{\sigma_i^2} - \sum_j b_j^H h_j - \sum_{i,j} \frac{W_{ij} v_i h_j}{\sigma_i}, \quad (9)$$

$$Z(\Theta) = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \sum_{\mathbf{h}} \exp(-E(\mathbf{v}, \mathbf{h} | \Theta)) dv_1 \cdots dv_I \quad (10)$$

であり, $\Theta = \{b_i^V, b_j^H, W_{ij}\}$ や $\{\sigma_i\}$ は分布パラメータである。ここで, 問題の簡単化のために本論文では $\{\sigma_i\}$ は定数として取り扱う。

Bernoulli-Bernoulli 型と同様に式 (8) から条件付分布を求めると,

$$p(v_i | \mathbf{h}, \Theta) = \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left(-\frac{(v_i - b_i^V - \sigma_i \sum_j W_{ij} h_j)^2}{2\sigma_i^2}\right), \quad (11)$$

$$p(h_j = 1 | \mathbf{v}, \Theta) = \frac{1}{1 + \exp(-b_j^H - \sum_i W_{ij} v_i)} \quad (12)$$

となり, 可視層の条件付分布が Gauss 分布, 隠れ層の条件付分布が Bernoulli 分布となるため, Gaussian-Bernoulli 型と呼ばれる。

RBM の学習問題とは観測される N 個の可視層のデータ $\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(N)}$ からこれらのパラメータ Θ を推定することである。

このとき, Bernoulli-Bernoulli 型, Gaussian-Bernoulli 型の RBM の学習問題に対するよく用いられる目的関数として, 次の周辺分布の対数尤度関数が挙げられる。

$$J(\Theta) = \frac{1}{N} \sum_n \log p(\mathbf{v}^{(n)} | \Theta) = \frac{1}{N} \sum_n \log \sum_{\mathbf{h}} p(\mathbf{v}^{(n)}, \mathbf{h} | \Theta). \quad (13)$$

2.2 Contrastive Divergence 法 [1], [4]

最急降下法により, 式 (13) で表される周辺分布の対数尤度関

数 $J(\Theta)$ の最大化を考える。

まず, Bernoulli-Bernoulli 型 RBM について, $J(\Theta)$ を Θ に関して微分すると,

$$\begin{aligned} \frac{\partial J}{\partial \Theta}(\Theta) = & -\frac{1}{N} \sum_n \sum_h p(\mathbf{h}|\mathbf{v}^{(n)}, \Theta) \frac{\partial E}{\partial \Theta}(\mathbf{v}^{(n)}, \mathbf{h}|\Theta) \\ & + \sum_v \sum_h p(\mathbf{v}, \mathbf{h}|\Theta) \frac{\partial E}{\partial \Theta}(\mathbf{v}, \mathbf{h}|\Theta) \end{aligned} \quad (14)$$

となる。このため, 最急降下法によるパラメータの更新は

$$\Theta \leftarrow \Theta^{\text{old}} + \epsilon \Delta_{\text{cd}} \Theta \quad (15)$$

となる。ただし, ϵ は学習率, $\Delta_{\text{cd}} \Theta = \partial J / \partial \Theta$ である。

ここで, \mathbf{v} や \mathbf{h} についての和に関しては厳密に計算しようとすると $O(2^I)$ や $O(2^J)$ の計算量となるため, 現実的ではない。しかし, 式 (14) の第 1 項に関しては式 (5) を用いて周辺化を行うことにより, \sum_h を \sum_j にすることが出来る。

しかし, 式 (14) の第 2 項に関してはどうしても計算が困難である。そこで, 式 (14) の第 2 項は同時確率による期待値計算であることに注目すると, 次式のようなギブスサンプリングによる同時確率の近似を用いることで計算量を削減することが考えられる。

$$p(\mathbf{v}, \mathbf{h}|\Theta) \approx \frac{1}{M} \sum_m \delta(\mathbf{v} - \mathbf{v}^{(m)}) p(\mathbf{h}|\mathbf{v}^{(m)}, \Theta) \quad (16)$$

ここで, 式 (4), (5) からギブスサンプリングは

$$h_j^{d-1} \sim p(h_j|\mathbf{v}^{d-1}, \Theta), \quad (17)$$

$$v_i^d \sim p(v_i|\mathbf{h}^{d-1}, \Theta) \quad (18)$$

と容易に行うことが可能である。

式 (15) による更新を式 (16) のギブスサンプリングによる近似で行う手法を CD 法という。

次に Gaussian-Bernoulli 型 RBM について, $J(\Theta)$ を Θ に関して微分すると,

$$\begin{aligned} \frac{\partial J}{\partial \Theta}(\Theta) = & -\frac{1}{N} \sum_n \sum_h p(\mathbf{h}|\mathbf{v}^{(n)}, \Theta) \frac{\partial E}{\partial \Theta}(\mathbf{v}^{(n)}, \mathbf{h}|\Theta) \\ & + \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \sum_h p(\mathbf{v}, \mathbf{h}|\Theta) \frac{\partial E}{\partial \Theta}(\mathbf{v}, \mathbf{h}|\Theta) dv_1 \cdots dv_I \end{aligned} \quad (19)$$

となる。これを計算するには, 第 1 項に関しては Bernoulli-Bernoulli 型と同様に式 (5) を用いて周辺化を行えばよい。また, 第 2 項に関しても Bernoulli-Bernoulli 型と同様に同時確率による期待値計算であるので, 同様にギブスサンプリングによる近似を行うことで近似的に計算が可能となる。よって, Gaussian-Bernoulli 型でも Bernoulli-Bernoulli 型と同様に CD 法によりパラメータ学習を行うことができる。

3. 補助関数法による RBM の学習アルゴリズム

補助関数法による目的関数 $F(x)$ の最大化問題の最適化アルゴリズムは, 補助変数 \bar{x} を導入して, 任意の x, \bar{x} で $F(x) \geq F^+(x, \bar{x})$ となり, $F(x) = \max_{\bar{x}} F^+(x, \bar{x})$ となるような下限関数 $F^+(x, \bar{x})$ を設計して, $F^+(x, \bar{x})$ を x について

の最小化と \bar{x} についての最小化を交互に行うことである。ここで重要なのは, $F^+(x, \bar{x})$ を x についての最小化が容易に行え, かつ $F^+(x, \bar{x})$ が $F(x)$ によくフィットするように設計できるからであるので, 本研究では x の各変数の相互依存をなくすような $F^+(x, \bar{x})$ を設計することを目指す。

3.1 最尤学習

ここでは最尤学習, つまり, 目的関数として式 (13) で表される対数尤度を用い, これを最大化することを考える。

まず, Bernoulli-Bernoulli 型 RBM について考えると, Jensen の不等式から

$$\begin{aligned} J(\Theta) &= \frac{1}{N} \sum_n \log \sum_h p(\mathbf{v}^{(n)}, \mathbf{h}|\Theta) \\ &\geq \frac{1}{N} \sum_n \sum_h \lambda_{n,h} \log p(\mathbf{v}^{(n)}, \mathbf{h}|\Theta) \\ &\quad - \sum_h \lambda_{n,h} \log \lambda_{n,h}. \end{aligned} \quad (20)$$

となる。ここで, 等号の成立は

$$\lambda_{n,h} = p(\mathbf{h}|\mathbf{v}^{(n)}, \Theta) \quad (21)$$

のときである。式 (20) を整理すると,

$$\begin{aligned} J(\Theta) &\geq -\frac{1}{N} \sum_n \sum_h \lambda_{n,h} E(\mathbf{v}^{(n)}, \mathbf{h}|\Theta) \\ &\quad - \log Z(\Theta) - \sum_h \lambda_{n,h} \log \lambda_{n,h} \end{aligned} \quad (22)$$

となる。ここで, この式の第 2 項について考えると, 負の対数関数は凸関数であるので, 接線の方程式を用いて下から抑えることが出来る。

$$-\log Z(\Theta) \geq -\frac{Z(\Theta)}{\zeta} - \log \zeta + 1. \quad (23)$$

ここで, 等号の成立は接点となるので,

$$\zeta = Z(\Theta) \quad (24)$$

となる。ここで, 式 (2) から $E(\mathbf{v}, \mathbf{h}|\Theta)$ はパラメータ Θ に対し線形であるので,

$$a_k(\mathbf{v}, \mathbf{h}) = \begin{cases} v_i & (k = i) \\ h_j & (k = I + j) \\ v_i h_j & (k = I + J + J \times (i - 1) + j) \end{cases}, \quad (25)$$

$$\theta_k = \begin{cases} b_i^V & (k = i) \\ b_j^H & (k = I + j) \\ W_{ij} & (k = I + J + J \times (i - 1) + j) \end{cases} \quad (26)$$

とおくと, 式 (2) は

$$E(\mathbf{v}, \mathbf{h}|\Theta) = -\sum_k a_k(\mathbf{v}, \mathbf{h}) \theta_k \quad (27)$$

と表すことができる。このとき, $-\exp(\sum_k a_k(\mathbf{v}, \mathbf{h}) \theta_k)$ に対して, 複素 NMF [5] で用いられている Jensen の不等式を用いた補助関数の設計法を用いると,

$$\begin{aligned}
& - \exp\left(\sum_k a_k(\mathbf{v}, \mathbf{h})\theta_k\right) \\
& \geq - \sum_k \beta_k(\mathbf{v}, \mathbf{h}) \exp\left(\frac{a_k(\mathbf{v}, \mathbf{h})\theta_k - \alpha_k(\mathbf{v}, \mathbf{h})}{\beta_k(\mathbf{v}, \mathbf{h})}\right)
\end{aligned} \quad (28)$$

となり，等号の成立は

$$\forall \beta_k(\mathbf{v}, \mathbf{h}) \in [0, 1], \sum_k \beta_k(\mathbf{v}, \mathbf{h}) = 1, \quad (29)$$

$$\alpha_k(\mathbf{v}, \mathbf{h}) = a_k(\mathbf{v}, \mathbf{h})\theta_k - \beta_k(\mathbf{v}, \mathbf{h}) \sum_l a_l(\mathbf{v}, \mathbf{h})\theta_l \quad (30)$$

となる．ここで， $\beta_k(\mathbf{v}, \mathbf{h})$ は任意に設計できるので， \mathbf{v}, \mathbf{h} に依存しない定数 β_k とする．また，補助関数法は補助変数と従来変数を交互に更新する手法であるので，一反復前の従来変数 θ^{old} を新たな補助変数として導入すると，補助変数の更新式の従来変数 θ を θ^{old} に置き換えて代入することができるので，式 (28) に補助変数の更新式 (30) を代入し，式 (23) を式 (3), (27) 用いて整理すると，

$$\begin{aligned}
- \log Z(\Theta) & \geq - \sum_{\mathbf{v}} \sum_{\mathbf{h}} \frac{\exp(\sum_l a_l(\mathbf{v}, \mathbf{h})\theta_l^{\text{old}})}{\zeta} \\
& \times \sum_k \beta_k \exp\left(\frac{a_k(\mathbf{v}, \mathbf{h})(\theta_k - \theta_k^{\text{old}})}{\beta_k}\right) \\
& - \log \zeta + 1
\end{aligned} \quad (31)$$

となる．ここで，式 (1), (24) から式 (31) は

$$\begin{aligned}
- \log Z(\Theta) & \geq - \sum_k \beta_k \sum_{\mathbf{v}} \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h} | \Theta^{\text{old}}) \\
& \times \exp\left(\frac{a_k(\mathbf{v}, \mathbf{h})(\theta_k - \theta_k^{\text{old}})}{\beta_k}\right) \\
& - \log \zeta + 1
\end{aligned} \quad (32)$$

となる．式 (21), (22), (24), (32) より， $\bar{\Theta} = \{\Theta^{\text{old}}, \beta_k\}$ とおいたとき， $J(\Theta)$ の下限関数 $J^+(\Theta, \bar{\Theta})$ は

$$\begin{aligned}
J^+(\Theta, \bar{\Theta}) & = \frac{1}{N} \sum_n \sum_{\mathbf{h}} p(\mathbf{h} | \mathbf{v}^{(n)}, \Theta^{\text{old}}) \sum_k a_k(\mathbf{v}^{(n)}, \mathbf{h})\theta_k \\
& - \sum_k \beta_k \sum_{\mathbf{v}} \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h} | \Theta^{\text{old}}) \\
& \times \exp\left(\frac{a_k(\mathbf{v}, \mathbf{h})(\theta_k - \theta_k^{\text{old}})}{\beta_k}\right) \\
& + C(\bar{\Theta})
\end{aligned} \quad (33)$$

と定義できる．ただし $C(\bar{\Theta})$ は Θ に対して定数の項である．

ここで，式 (33) を θ_k について微分すると，

$$\begin{aligned}
\frac{\partial J^+}{\partial \theta_k}(\Theta, \bar{\Theta}) & = \frac{1}{N} \sum_n \sum_{\mathbf{h}} p(\mathbf{h} | \mathbf{v}^{(n)}, \Theta^{\text{old}}) a_k(\mathbf{v}^{(n)}, \mathbf{h}) \\
& - \sum_{\mathbf{v}} \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h} | \Theta^{\text{old}}) a_k(\mathbf{v}, \mathbf{h}) \\
& \times \exp\left(\frac{a_k(\mathbf{v}, \mathbf{h})(\theta_k - \theta_k^{\text{old}})}{\beta_k}\right)
\end{aligned} \quad (34)$$

となる．ここで，式 (25) から Bernoulli-Bernoulli 型 RBM においては $a_k(\mathbf{v}, \mathbf{h}) \in \{0, 1\}$ であるため，式 (34) の第二項は

$$\begin{aligned}
& \sum_{\mathbf{v}} \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h} | \Theta^{\text{old}}) a_k(\mathbf{v}, \mathbf{h}) \exp\left(\frac{a_k(\mathbf{v}, \mathbf{h})(\theta_k - \theta_k^{\text{old}})}{\beta_k}\right) \\
& = \exp\left(\frac{\theta_k - \theta_k^{\text{old}}}{\beta_k}\right) \sum_{\mathbf{v}} \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h} | \Theta^{\text{old}}) a_k(\mathbf{v}, \mathbf{h})
\end{aligned} \quad (35)$$

となる．よって， $\partial J^+ / \partial \theta_k = 0$ を解析的に求めることができ， $\Delta_{\text{af}} \theta_k$ を

$$\begin{aligned}
\Delta_{\text{af}} \theta_k & = \log \frac{1}{N} \sum_n \sum_{\mathbf{h}} p(\mathbf{h} | \mathbf{v}^{(n)}, \Theta^{\text{old}}) a_k(\mathbf{v}^{(n)}, \mathbf{h}) \\
& - \log \sum_{\mathbf{v}} \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h} | \Theta^{\text{old}}) a_k(\mathbf{v}, \mathbf{h})
\end{aligned} \quad (36)$$

と定義すると， θ_k の更新式は，

$$\theta_k \leftarrow \theta_k^{\text{old}} + \beta_k \Delta_{\text{af}} \theta_k \quad (37)$$

となる．ここで，式 (36) の右辺の第 2 項は 2.2 節で言及したように厳密に計算することは困難である．そこで，CD 法と同様にギブスサンプリングによる \mathbf{v} の周辺確率の近似を行うことで，式 (37) の計算が可能となる．

さて，式 (15), (37) を比較すると，式 (37) の更新式はあたかも学習率が β_k であるかのような形となっている．ところが，式 (29) よりパラメータ数が増えると β_k の平均的な値は小さくなるため，収束が遅くなることが予想される．そこで，収束を速くするために以下の近似を考える．

$$\beta'_k \leftarrow \beta_k. \quad (38)$$

このとき， $\gamma \in [0, 1]$ ならば， $\beta_k \in [0, 1]$ を満たすので β'_k は β_k よりも大きくなる．そのため，式 (37) による更新が速くなることが期待できる．ただし， β'_k は式 (29) を満たさないため，補助関数法における収束性は保証されないことには注意されたい．

以上より Bernoulli-Bernoulli 型 RBM の補助関数法による最尤学習アルゴリズムは，次の 1)~3) を反復することである．1) 補助変数 $\bar{\Theta}$ を求める．2) ギブスサンプリングで $p(\mathbf{v}, \mathbf{h} | \Theta^{\text{old}})$ の近似値を求める．3) 式 (37) の更新式でパラメータを更新．

次に，Gaussian-Bernoulli 型 RBM について考える．式 (1) と式 (8) は同じ構造をしており，式 (2) と式 (9) を比較すると， $\{\sigma_i\}$ を定数だとしているので，学習すべきパラメータは Bernoulli-Bernoulli 型と同じ $\Theta = \{b_i^V, b_j^H, W_{ij}\}$ であり，Bernoulli-Bernoulli 型と同じく $E(\mathbf{v}, \mathbf{h} | \Theta)$ において線形和で表現されているため，Bernoulli-Bernoulli 型と同じ手順で補助関数が設計可能である．よって， $J(\Theta)$ の下限関数 $J^+(\Theta, \bar{\Theta})$ は

$$\begin{aligned}
J^+(\Theta, \bar{\Theta}) & = \frac{1}{N} \sum_n \sum_{\mathbf{h}} p(\mathbf{h} | \mathbf{v}^{(n)}, \Theta^{\text{old}}) \sum_k a_k(\mathbf{v}^{(n)}, \mathbf{h})\theta_k \\
& - \sum_k \beta_k \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h} | \Theta^{\text{old}}) \\
& \times \exp\left(\frac{a_k(\mathbf{v}, \mathbf{h})(\theta_k - \theta_k^{\text{old}})}{\beta_k}\right) d\mathbf{v}_1 \cdots d\mathbf{v}_I \\
& + C(\bar{\Theta})
\end{aligned} \quad (39)$$

と設計できる．ただし $C(\bar{\Theta})$ は Θ に対して定数の項であり，

$$a_k(\mathbf{v}, \mathbf{h}) = \begin{cases} \frac{v_i}{\sigma_i} & (k = i) \\ h_j & (k = I + j) \\ \frac{v_i h_j}{\sigma_i} & (k = I + J + J \times (i - 1) + j) \end{cases}, \quad (40)$$

$$\theta_k = \begin{cases} b_i^V & (k = i) \\ b_j^H & (k = I + j) \\ W_{ij} & (k = I + J + J \times (i - 1) + j) \end{cases}, \quad (41)$$

また、 β_k は $\beta_k \in [0, 1]$ 、 $\sum_k \beta_k = 1$ を満たす任意定数である。ここで、式 (39) を θ_k について微分すると、

$$\begin{aligned} \frac{\partial J^+}{\partial \theta_k}(\Theta, \bar{\Theta}) &= \frac{1}{N} \sum_n \sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{v}^{(n)}, \Theta^{\text{old}}) a_k(\mathbf{v}^{(n)}, \mathbf{h}) \\ &\quad - \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h}|\Theta^{\text{old}}) a_k(\mathbf{v}, \mathbf{h}) \\ &\quad \times \exp\left(\frac{a_k(\mathbf{v}, \mathbf{h})(\theta_k - \theta_k^{\text{old}})}{\beta_k}\right) dv_1 \cdots dv_I \end{aligned} \quad (42)$$

となる。ここで、式 (42) の右辺の第 2 項は 2.2 節で言及し、また Bernoulli-Bernoulli 型に対する補助関数法によるアルゴリズムと同様に厳密に計算することは困難である。そこで、CD 法などと同様に Gibbs サンプルングによる \mathbf{v} の周辺確率の近似を行い、同時確率を式 (16) で近似すると式 (42) は

$$\begin{aligned} \frac{\partial J^+}{\partial \theta_k}(\Theta, \bar{\Theta}) &= \frac{1}{N} \sum_n \sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{v}^{(n)}, \Theta^{\text{old}}) a_k(\mathbf{v}^{(n)}, \mathbf{h}) \\ &\quad - \frac{1}{M} \sum_m \sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{v}^{(m)}, \Theta^{\text{old}}) a_k(\mathbf{v}^{(m)}, \mathbf{h}) \\ &\quad \times \exp\left(\frac{a_k(\mathbf{v}^{(m)}, \mathbf{h})(\theta_k - \theta_k^{\text{old}})}{\beta_k}\right) \end{aligned} \quad (43)$$

となる。ここで、 $a_k(\mathbf{v}, \mathbf{h}) \in \{0, 1\}$ になる場合を除いて $\partial J^+ / \partial \theta_k = 0$ は解析的に解くことは困難である。Bernoulli-Bernoulli 型の RBM では常に $a_k(\mathbf{v}, \mathbf{h}) \in \{0, 1\}$ となっているが、Gaussian-Bernoulli 型の RBM では θ_k が b_j^H のみ $a_k(\mathbf{v}, \mathbf{h}) \in \{0, 1\}$ となる。しかし、 $\partial^2 J^+ / \partial \theta_k^2$ は常に負となるため、Newton 法により安定して最大点を求めることが期待できる。また、Bernoulli-Bernoulli 型に対する補助関数と同型の補助関数を用いているので、 β_k は同様の役割を担っていると考えられる。そのため、Bernoulli-Bernoulli 型で議論したように、収束を速くするために式 (38) による近似を考える。

以上より Gaussian-Bernoulli 型 RBM の補助関数法による最尤学習アルゴリズムは、次の 1)~3) を反復することである。1) 補助変数 $\bar{\Theta}$ を求める。2) Gibbs サンプルングで $p(\mathbf{v}, \mathbf{h}|\Theta^{\text{old}})$ の近似値を求める。3) $\partial J^+ / \partial \theta_k = 0$ となるようにパラメータを更新 (一部 Newton 法を用いる)。

3.2 最大再構築確率学習

次に、RBM の学習アルゴリズムとして、今度は目的関数が他と異なるものを導出する。RBM を学習する際は、可視層に入力されたデータが隠れ層で変換された特徴量を更に可視層に変換するときに元の入力データに近くなるように学習が進むという特徴がある。これは、変換された特徴量が入力データの情報をなるべく保持しようと学習されていると考えられる。そこで、本研究ではこの特徴を陽に記述した目的関数を導入するこ

とを考え、次のような可視層に観測データが来たときに、ギブス サンプルングを 1 回行ったときに元の観測データが再現される確率の対数を考える。

$$J_r(\Theta) = \frac{1}{N} \sum_n \log \sum_{\mathbf{h}} p(\mathbf{v}^{(n)}|\mathbf{h}, \Theta) p(\mathbf{h}|\mathbf{v}^{(n)}, \Theta). \quad (44)$$

この元の観測データが再現される確率を再構築確率と呼び、これを最大化する最大再構築確率学習について補助関数法を用いて学習アルゴリズムを設計することを目指す。

この式を直接最大化するのは困難であるので、この式に対して、

$$\mathbf{h}^{(n)} \sim p(\mathbf{h}|\mathbf{v}^{(n)}, \Theta) \quad (45)$$

でサンプルングした値を元に

$$J_r(\Theta) \approx \frac{1}{N} \sum_n \log p(\mathbf{v}^{(n)}|\mathbf{h}^{(n)}, \Theta) p(\mathbf{h}^{(n)}|\mathbf{v}^{(n)}, \Theta) \quad (46)$$

と近似することを考える。この右辺を \tilde{J}_r とおく。

まず、Bernoulli-Bernoulli 型 RBM について \tilde{J}_r の補助関数の設計を考える。式 (4), (5), Bernoulli-Bernoulli 型 RBM の条件付分布が Bernoulli 分布となることより、

$$\begin{aligned} \tilde{J}_r &= \frac{1}{N} \sum_n \left\{ \sum_i v_i^{(n)} \log q_{ni}^V + \sum_i (1 - v_i^{(n)}) \log(1 - q_{ni}^V) \right. \\ &\quad \left. + \sum_j h_j^{(n)} \log q_{nj}^H + \sum_j (1 - h_j^{(n)}) \log(1 - q_{nj}^H) \right\} \end{aligned} \quad (47)$$

となる。ただし、 $q_{ni}^V = p(v_i^{(n)} = 1|\mathbf{h}^{(n)}, \Theta)$ 、 $q_{nj}^H = p(h_j^{(n)} = 1|\mathbf{v}^{(n)}, \Theta)$ である。よって、式 (6), (7) より式 (47) を整理すると、

$$\begin{aligned} \tilde{J}_r &= -\frac{1}{N} \sum_n \left\{ \sum_i v_i^{(n)} \log(1 + \exp(-f_{ni}^V)) \right. \\ &\quad \left. + \sum_i (1 - v_i^{(n)}) \log(1 + \exp f_{ni}^V) \right. \\ &\quad \left. + \sum_j h_j^{(n)} \log(1 + \exp(-f_{nj}^H)) \right. \\ &\quad \left. + \sum_j (1 - h_j^{(n)}) \log(1 + \exp f_{nj}^H) \right\} \end{aligned} \quad (48)$$

となる。ただし、 $f_{ni}^V = b_i^V + \sum_j W_{ij} h_j^{(n)}$ 、 $f_{nj}^H = b_j^H + \sum_i W_{ij} v_i^{(n)}$ である。ここで、式 (48) のそれぞれの項は負の対数関数であり、その引数の項をみると、式 (23) ~ (32) と同様の論理で下限関数が設計できることが分かる。

よって、下限関数 $\tilde{J}_r^+(\Theta, \bar{\Theta})$ は

$$\begin{aligned} \tilde{J}_r^+(\Theta, \bar{\Theta}) &= -\frac{1}{N} \left\{ \sum_i v_i^{(n)} (1 - \hat{q}_{ni}^V) \xi_{ni}^V + \sum_i (1 - v_i^{(n)}) \hat{q}_{ni}^V \eta_{ni}^V \right. \\ &\quad \left. + \sum_j h_j^{(n)} (1 - \hat{q}_{nj}^H) \xi_{nj}^H + \sum_j (1 - h_j^{(n)}) \hat{q}_{nj}^H \eta_{nj}^H \right\} \\ &\quad + C(\bar{\Theta}) \end{aligned} \quad (49)$$

となる。ただし、 $\hat{q}_{ni}^V = p(v_i^{(n)} = 1|\mathbf{h}^{(n)}, \Theta^{\text{old}})$ 、 $\hat{q}_{nj}^H =$

$$p(h_j^{(n)} = 1 | \mathbf{v}^{(n)}, \Theta^{\text{old}}), \xi_{ni}^V = \beta_{i0}^V e^{-\hat{b}_i^V} + \sum_j \beta_{ij}^V e^{-\hat{W}_{nij}^V}, \eta_{ni}^V = \beta_{i0}^V e^{\hat{b}_i^V} + \sum_j \beta_{ij}^V e^{\hat{W}_{nij}^V}, \xi_{nj}^H = \beta_{0j}^H e^{-\hat{b}_j^H} + \sum_i \beta_{ij}^H e^{-\hat{W}_{nij}^H}, \eta_{nj}^H = \beta_{0j}^H e^{\hat{b}_j^H} + \sum_i \beta_{ij}^H e^{\hat{W}_{nij}^H}, \hat{b}_i^V = \frac{b_i^V - b_i^{V, \text{old}}}{\beta_{i0}^V}, \hat{b}_j^H = \frac{b_j^H - b_j^{H, \text{old}}}{\beta_{0j}^H}, \hat{W}_{nij}^V = \frac{W_{ij} - W_{ij}^{\text{old}}}{\beta_{ij}^V} h_j^{(n)}, \hat{W}_{nij}^H = \frac{W_{ij} - W_{ij}^{\text{old}}}{\beta_{ij}^H} v_i^{(n)}, \text{であり, } \beta_{i0}^V, \beta_{0j}^H, \beta_{ij}^H \text{ は}$$

$$\begin{aligned} \beta_{i0}^V, \beta_{ij}^V, \beta_{0j}^H, \beta_{ij}^H &\in [0, 1], \\ \beta_{i0}^V + \sum_j \beta_{ij}^V &= 1, \quad \forall i, \quad \beta_{0j}^H + \sum_i \beta_{ij}^H = 1, \quad \forall j, \end{aligned} \quad (50)$$

を満たす任意定数である。また、 $\bar{\Theta} = \{\Theta^{\text{old}}, \beta_{i0}^V, \beta_{ij}^V, \beta_{0j}^H, \beta_{ij}^H\}$ であり、 $C(\bar{\Theta})$ は Θ に対して定数の項である。

ここで、 $\partial \tilde{J}_r^+ / \partial \Theta = 0$ を解くことについて考えると、 b_i^V, b_j^H については解析的に解くことが出来るが、 W_{ij} については解析的に解くことが出来ない。しかし、これは、 $\partial^2 \tilde{J}_r^+ / \partial W_{ij}^2 < 0$ であるので、Newton 法を用いて効率的に求めることが出来る。

また、 $\beta_{i0}^V, \beta_{ij}^V, \beta_{0j}^H, \beta_{ij}^H$ が最尤学習における補助関数法を用いた学習アルゴリズムと同様に学習率を担っていると考えられるため、これらについても γ 乗を行う近似を考える。

以上より Bernoulli-Bernoulli 型 RBM の補助関数法による最大再構築確率学習アルゴリズムは、次の 1)~3) を反復することである。1) サンプリングにより $h_j^{(n)}$ を求める。2) 補助変数 $\bar{\Theta}$ を求める。3) $\partial \tilde{J}_r^+ / \partial \Theta = 0$ から求まる更新式でパラメータを更新。

次に、Gaussian-Bernoulli 型 RBM について補助関数の設計を考える。式 (46) に式 (4), (5), (11), (12) を代入すると、

$$\begin{aligned} \tilde{J}_r(\Theta) = & -\frac{1}{N} \sum_n \left\{ \sum_i \frac{(v_i^{(n)} - f_{ni}^V)^2}{2\sigma_i^2} + \sum_i \log \sigma_i \right. \\ & + \frac{I}{2} \log(2\pi) - \sum_j \left(h_j^{(n)} - \frac{1}{2} \right) f_{nj}^H \\ & \left. + \sum_j \log(\exp(f_{nj}^H/2) + \exp(-f_{nj}^H/2)) \right\} \end{aligned} \quad (51)$$

となる。ただし、 $f_{ni}^V = b_i^V + \sigma_i \sum_j W_{ij} h_j^{(n)}$, $f_{nj}^H = b_j^H + \sum_i W_{ij} v_i^{(n)} / \sigma_i$ である。

ここで、

$$\begin{aligned} \log(\exp(x) + \exp(-x)) \\ \leq \frac{\tanh(\kappa)}{2\kappa} x^2 - \frac{\kappa \tanh(\kappa)}{2} + \log(2 \cosh(\kappa)) \end{aligned} \quad (52)$$

という不等式を考える。(右辺) - (左辺) が連続関数で、その極小値が $x = \pm \kappa$ で 0 になり、かつ $x \rightarrow \pm \infty$ で正に発散することからこの不等式が成立することが分かる。また、等号の成立は

$$\kappa = \pm x \quad (53)$$

のときである。この不等式を用いると式 (51) は

$$\begin{aligned} \tilde{J}_r(\Theta) \geq & -\frac{1}{N} \sum_n \left\{ \sum_i \frac{(f_{ni}^V)^2}{2\sigma_i^2} + \sum_j \frac{\tanh(\kappa_{nj})}{8\kappa_{nj}} (f_{nj}^H)^2 \right. \\ & - \sum_i \frac{v_i^{(n)}}{\sigma_i^2} f_{ni}^V - \sum_j \left(h_j^{(n)} - \frac{1}{2} \right) f_{nj}^H \\ & + \sum_i \left(\frac{(v_i^{(n)})^2}{2\sigma_i^2} + \log \sigma_i \right) + \frac{I}{2} \log(2\pi) \\ & \left. + \sum_j \left(\log(\cosh(\kappa_{nj})) - \frac{\kappa_{nj} \tanh(\kappa_{nj})}{8} \right) \right\} \end{aligned} \quad (54)$$

と下から押さえられる。

ここで、 $f_{ni}^V = b_i^V + \sigma_i \sum_j W_{ij} h_j^{(n)}$, $f_{nj}^H = b_j^H + \sum_i W_{ij} v_i^{(n)} / \sigma_i$ であるので、 $-(f_{ni}^V)^2$ と $-(f_{nj}^H)^2$ に対して複素 NMF [5] で用いられている Jensen の不等式を用いた補助関数の設計法を用いると

$$-(f_{ni}^V)^2 \geq -\frac{(b_i^V - \alpha_{ni0}^V)^2}{\beta_{ni0}^V} - \sum_j \frac{(\sigma_i W_{ij} h_j^{(n)} - \alpha_{nij}^V)^2}{\beta_{nij}^V}, \quad (55)$$

$$-(f_{nj}^H)^2 \geq -\frac{(b_j^H - \alpha_{n0j}^H)^2}{\beta_{n0j}^H} - \sum_i \frac{(W_{ij} v_i^{(n)} / \sigma_i - \alpha_{nij}^H)^2}{\beta_{nij}^H} \quad (56)$$

となり、等号の成立は

$$\forall \beta_{ni0}^V, \beta_{nij}^V \in [0, 1], \beta_{ni0}^V + \sum_j \beta_{nij}^V = 1, \quad (57)$$

$$\forall \beta_{n0j}^H, \beta_{nij}^H \in [0, 1], \beta_{n0j}^H + \sum_i \beta_{nij}^H = 1, \quad (58)$$

$$\alpha_{ni0}^V = b_i^V - \beta_{ni0}^V f_{ni}^V, \alpha_{nij}^V = \sigma_i W_{ij} h_j^{(n)} - \beta_{nij}^V f_{ni}^V, \quad (59)$$

$$\alpha_{n0i}^H = b_j^H - \beta_{n0j}^H f_{nj}^H, \alpha_{nij}^H = \frac{W_{ij} v_i^{(n)}}{\sigma_i} - \beta_{nij}^H f_{nj}^H \quad (60)$$

のときである。

よって、式 (54)、式 (55)、式 (56) から下限関数 $\tilde{J}_r^+(\Theta, \bar{\Theta})$ は

$$\begin{aligned} \tilde{J}_r^+(\Theta, \bar{\Theta}) = & -\frac{1}{N} \sum_n \left\{ \sum_i \frac{(b_i^V - \alpha_{ni0}^V)^2}{2\beta_{ni0}^V \sigma_i^2} \right. \\ & + \sum_i \sum_j \frac{(\sigma_i W_{ij} h_j^{(n)} - \alpha_{nij}^V)^2}{2\beta_{nij}^V \sigma_i^2} \\ & + \sum_j \frac{\tanh(\kappa_{nj})}{8\kappa_{nj} \beta_{n0j}^H} (b_j^H - \alpha_{n0j}^H)^2 \\ & + \sum_j \sum_i \frac{\tanh(\kappa_{nj})}{8\kappa_{nj} \beta_{nij}^H} (W_{ij} v_i^{(n)} / \sigma_i - \alpha_{nij}^H)^2 \\ & - \sum_i \frac{v_i^{(n)}}{\sigma_i^2} \left(b_i^V + \sigma_i \sum_j W_{ij} h_j^{(n)} \right) \\ & \left. - \sum_j \left(h_j^{(n)} - \frac{1}{2} \right) \left(b_j^H + \sum_i W_{ij} v_i^{(n)} / \sigma_i \right) \right\} + C(\bar{\Theta}) \end{aligned} \quad (61)$$

となる。ただし、

$$\bar{\Theta} = \{\kappa_{nj}, \alpha_{ni0}^V, \alpha_{nij}^V, \alpha_{n0j}^H, \alpha_{nij}^H, \beta_{ni0}^V, \beta_{nij}^V, \beta_{n0j}^H, \beta_{nij}^H\} \quad (62)$$

であり、 $C(\bar{\Theta})$ は Θ に対して定数の項である。

よって、 b_i^V, b_j^H, W_{ij} の更新式は $\partial \tilde{J}_r / \partial \Theta = 0$ より

$$b_i^V \leftarrow \frac{\sum_n (\alpha_{ni0}^V / \beta_{ni0}^V + v_i^{(n)})}{\sum_n 1 / \beta_{ni0}^V}, \quad (63)$$

$$b_j^H \leftarrow \frac{\sum_n (\tanh(\kappa_{nj}) \alpha_{n0j}^H / 4\kappa_{nj} \beta_{n0j}^H + h_j^{(n)} - 1/2)}{\sum_n \tanh(\kappa_{nj}) / 4\kappa_{nj} \beta_{n0j}^H}, \quad (64)$$

$$W_{ij} \leftarrow \frac{\sum_n \left(\frac{\alpha_{nij}^V}{\beta_{nij}^V \sigma_i} + \frac{\tanh(\kappa_{nj}) \alpha_{nij}^H}{4\kappa_{nj} \beta_{nij}^H} + \frac{v_i^{(n)} (2h_j^{(n)} - 1/2)}{\sigma_i} \right)}{\sum_n \left(\frac{h_j^{(n)}}{\beta_{nij}^V \sigma_i} + \frac{\tanh(\kappa_{nj}) v_i^{(n)}}{4\kappa_{nj} \beta_{nij}^H \sigma_i} \right)} \quad (65)$$

となる。ここで、 $\beta_{ni0}^V, \beta_{nij}^V, \beta_{n0j}^H, \beta_{nij}^H$ は式 (57), (58) を満たす任意定数であるので、 n によらない定数 $\beta_{i0}^V, \beta_{ij}^V, \beta_{0j}^H, \beta_{ij}^H$ とする。そして、新たに補助変数として一反復前の値を表す $b_i^{V,old}, b_j^{H,old}, W_{ij}^{old}$ を導入し、補助変数の更新式 (53), (59), (60) を代入して整理すると

$$b_i^V \leftarrow b_i^{V,old} + \frac{\beta_{i0}^V}{N} \sum_n (v_i^{(n)} - f_{ni}^{V,old}), \quad (66)$$

$$b_j^H \leftarrow b_j^{H,old} + \beta_{0j}^H \frac{\sum_n (h_j^{(n)} - q_{nj})}{\sum_n (q_{nj} - 1/2) / f_{nj}^{H,old}}, \quad (67)$$

$$W_{ij} \leftarrow W_{ij}^{old} + \frac{\sum_n ((v_i^{(n)} - f_{ni}^{V,old}) h_j^{(n)} + v_i^{(n)} (h_j^{(n)} - q_{nj}))}{\sum_n \left(\frac{h_j^{(n)} \sigma_i}{\beta_{ij}^V} + \frac{(v_i^{(n)})^2 (q_{nj} - 1/2)}{\beta_{ij}^H \sigma_i f_{nj}^{H,old}} \right)} \quad (68)$$

となる。ただし、 $f_{ni}^{V,old} = b_i^{V,old} + \sigma_i \sum_j W_{ij}^{old} h_j^{(n)}$, $f_{nj}^{H,old} = b_j^{H,old} + \sum_i W_{ij}^{old} v_i^{(n)} / \sigma_i$, $q_{nj} = 1 / (1 + \exp(-f_{nj}^{H,old}))$ である。また、このときの補助変数は

$$\bar{\Theta} = \{b_i^{V,old}, b_j^{H,old}, W_{ij}^{old}, \beta_{i0}^V, \beta_{ij}^V, \beta_{0j}^H, \beta_{ij}^H\} \quad (69)$$

となる。

この更新式を見てみると、 $\beta_{i0}^V, \beta_{ij}^V, \beta_{0j}^H, \beta_{ij}^H$ は他の補助関数法を用いた学習アルゴリズムと同様に学習率を担っていると考えられるため、これらについても γ 乗にする近似を考える。

以上より Gaussian-Bernoulli 型 RBM の補助関数法による最大再構築確率学習アルゴリズムは、次の 1)~3) を反復することである。1) サンプリングにより $h_j^{(n)}$ を求める。2) 補助変数 $\bar{\Theta}$ を求める。3) 式 (66)~(68) の更新式でパラメータを更新。

4. 収束性能の比較実験

2., 3. 章で説明した各学習アルゴリズムがどのような挙動を示すかについて、まず、Bernoulli-Bernoulli 型 RBM について実験を行った。また、これ以降、提案した補助関数法に基づく最尤学習アルゴリズムを提案法 1、補助関数法に基づく最大再構築確率学習アルゴリズムを提案法 2 と呼ぶこととする。最初に可視層の状態数 $I = 15$ 、隠れ層の状態数 $J = 10$ という非常に小さな系で実験を行い、次に可視層の状態数 $I = 1000$ 、隠れ層の状態数 $J = 800$ で実験を行った。小さな状態数で行うのは、尤度を厳密に計算するには状態数に対して指数的な計算量が必

要となるため、尤度による収束の比較を行いたいからである。このとき、可視層に入力するバイナリデータは乱数で 100 個生成し、生成したそれぞれに対し、80% のノイズをかけたものを 100 個づつ用意した。つまり、入力するデータ数は $N = 10000$ となる。また、学習の反復回数 T は 3000 回とし、各パラメータの初期値は $[-1, 1]$ の一様乱数から生成し、すべてのアルゴリズムで共通の初期値とした。

次に、 $\beta_k, \beta_{i0}^V, \beta_{ij}^V, \beta_{0j}^H, \beta_{ij}^H$ は一様、つまり、

$$\beta_k = \frac{1}{I + J + IJ}, \beta_{i0}^V = \beta_{ij}^V = \frac{1}{1 + J}, \beta_{0j}^H = \beta_{ij}^H = \frac{1}{1 + I} \quad (70)$$

とし、CD 法の学習率 ϵ や式 (38) に示す γ のスケジューリングを t 回目の反復のとき

$$\epsilon(t) = \epsilon_{init} \left(\frac{\epsilon_{end}}{\epsilon_{init}} \right)^{\frac{t-1}{T-1}}, \gamma(t) = \gamma_{init} \left(\frac{\gamma_{end}}{\gamma_{init}} \right)^{\frac{t-1}{T-1}} \quad (71)$$

とした。このとき、 $J = 15$ かつ $I = 10$ とした実験のときは $\epsilon_{init} = 0.1, \epsilon_{end} = 0.01$ 、提案法 1 に関して $\gamma_{init} = 0.1, \gamma_{init} = 1$ 、提案法 2 に関して $\gamma_{init} = 0.1, \gamma_{init} = 1$ とし、 $J = 1000$ かつ $I = 800$ とした実験のときは $\epsilon_{init} = 0.1, \epsilon_{end} = 0.01$ 、提案法 1 に関しては $\gamma_{init} = 0.3, \gamma_{init} = 0.5$ 、提案法 2 に関しては $\gamma_{init} = 0.5, \gamma_{init} = 1.0$ とした。また、ギブスサンプリングの回数を 1 回、Newton 法を用いる場合はその反復数を 1 回とした。

MATLAB による実装により計算時間を計ったところ、CD 法と比べ、提案法 1 による計算時間も提案法 2 による計算時間もおおそ同じくらいとなった。また、各学習アルゴリズムにより式 (13) に示す対数尤度がどのように遷移したかを Fig. 2(a) に示し、以下のように定義した再構築誤差 $e_{reconst}$ が各学習アルゴリズムによりどのように遷移したかを Fig. 2(b), 3 に示す。

$$\bar{h}^{(n)} = \arg \max_{\mathbf{h}} p(\mathbf{h} | \mathbf{v}^{(n)}, \Theta), \quad (72)$$

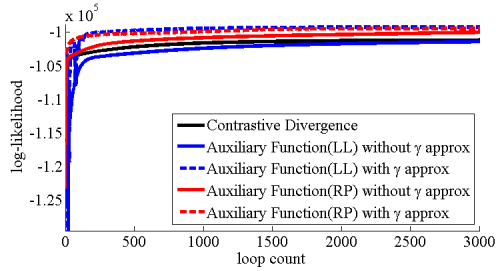
$$\bar{v}^{(n)} = \arg \max_{\mathbf{v}} p(\mathbf{v} | \bar{\mathbf{h}}^{(n)}, \Theta), \quad (73)$$

$$e_{reconst} = \frac{1}{NI} \sum_n \sum_i (v_i^{(n)} - \bar{v}_i^{(n)}). \quad (74)$$

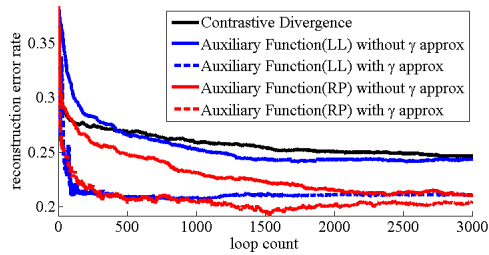
予想したように γ による学習率の近似をしなかったアルゴリズムは収束が遅く、それに対し、 γ による学習率の近似をしたアルゴリズムは非常に速く収束するという挙動が観測された。また、 γ の加速を行うことにより、CD 法よりも速くなる可能性を示す結果が得られた。

本来、補助関数法は設計パラメータが少ないという利点があるが、このときは γ という設計パラメータが生じてしまう。しかし、補助関数法の原理より、ギブスサンプリングの近似が十分ならば $\gamma = 1$ のときは安定して収束するので、 $\gamma = 1$ に近づくようなスケジューリングをすれば良いことから、CD 法の学習率のスケジューリングより設計の指針がはっきりしていると考えられる。

さらに、興味深いことは、本来、目的関数が式 (13) とは異なるものから出発した学習アルゴリズムである最大再構築確率



(a) 対数尤度



(b) 再構築誤差

図2 Bernoulli-Bernoulli 型 RBM で $J = 15$, $I = 10$ における各学習アルゴリズムの反復毎の対数尤度 (a) と再構築誤差 (b) . 黒い実線は CD 法を表し, 青, 赤の実線と破線はそれぞれ提案手法の γ による学習率の近似をしなかったものとしたものを表す .

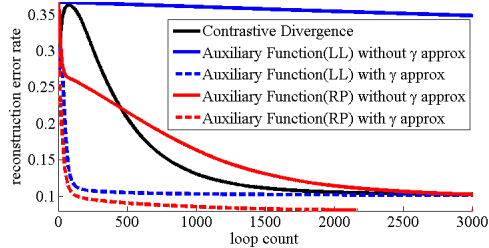


図3 Bernoulli-Bernoulli 型 RBM で $J = 1000$, $I = 800$ における各学習アルゴリズムの反復毎の再構築誤差 . 黒い実線は CD 法を表し, 青, 赤の実線と破線はそれぞれ提案手法の γ による学習率の近似をしなかったものとしたものを表す .

学習アルゴリズムが対数尤度に対して増加傾向にあることである . また, このアルゴリズムは γ による学習の加速を行うと, 始めの数反復において他のアルゴリズムより対数尤度の上昇が見られた .

次に Gaussian-Bernoulli 型 RBM について, 各学習どのような挙動をするのか, 可視層の状態数 $I = 1000$, 隠れ層の状態数 $J = 800$ で実験を行った . このとき, 可視層に入力するデータは平均 0, 標準偏差 100 の正規分布し従う乱数で 50 個生成し, 生成したそれぞれに対し, 平均 0, 標準偏差 0.1 の正規分布に従うノイズを足し合わせたものを 100 個づつ用意した . つまり, 入力するデータ数は $N = 10000$ となる . また, σ_i は一様に 1 とし, 学習の反復回数 T は 3000 回とし, 各パラメータの初期値は $[-1, 1]$ の一様乱数から生成し, すべてのアルゴリズムで共通の初期値とした . また, Bernoulli-Bernoulli 型と同様に任意定数の各 β をそれぞれ一様にし, CD 法の学習率 ϵ や式 (38) に示す γ のスケジューリングを式 (71) に従うよ

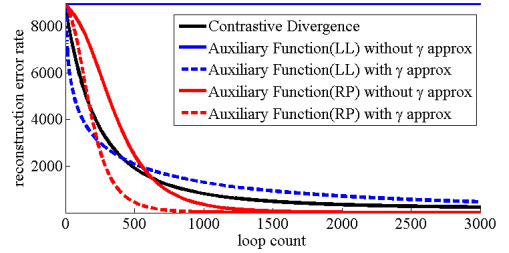


図4 Gaussian-Bernoulli 型 RBM で $J = 1000$, $I = 800$ における各学習アルゴリズムの反復毎の再構築誤差 . 黒い実線は CD 法を表し, 青, 赤の実線と破線はそれぞれ提案手法の γ による学習率の近似をしなかったものとしたものを表す .

うにし, それぞれ, $\epsilon_{\text{init}} = 0.001$, $\epsilon_{\text{end}} = 0.0001$ とし, 提案法 1 に関しては $\gamma_{\text{init}} = 0.03$, $\gamma_{\text{init}} = 0.05$ とし, 提案法 2 に関しては $\gamma_{\text{init}} = 0.9$, $\gamma_{\text{init}} = 1$ とし, Gibbs サンプリングの回数, Newton 法を用いる場合はその反復数もそれぞれ 1 回とした .

MATLAB による実装により計算時間を計ったところ, CD 法と比べ, 提案法 1 による計算時間も提案法 2 による計算時間もおおよそ同じくらいとなった . また, 各学習アルゴリズムにより式 (72) に示す再構築誤差がどのように遷移したのかを Fig. 4 に示す .

Bernoulli-Bernoulli 型と同様に, γ による学習率の近似をしなかったアルゴリズムは収束が遅く, それに対し, γ による学習率の近似をしたアルゴリズムは非常に速く収束するという挙動が観測された . また, γ の加速を行うことにより, CD 法よりも速くなる可能性を示す結果が得られた .

5. まとめ

本稿では, Bernoulli-Bernoulli 型 RBM, Gaussian-Bernoulli 型 RBM の学習アルゴリズムとして, それぞれ補助関数法を用いて新たに最尤学習アルゴリズムを導出した . また, それぞれのモデルに対して, 新たに最大再構築確率学習を提案し, その学習アルゴリズムを導出した . そして, 動作確認実験を通して, 既存手法と同等以上の性能を見込めることが確認できた . 今後の課題として, 実データに用いた場合や多層に重ねて Deep learning を行ったときにどのような挙動を示すかの観察が挙げられる .

謝辞 本研究は JSPS 科研費 26730100 の助成を受けたものです .

文献

- [1] Hinton, G. E., et al. "A fast learning algorithm for deep belief nets," *Neural Computation*, 2006, 18.7, pp. 1527-1554.
- [2] Bengio, Y., et al. "Greedy layer-wise training of deep networks," *NIPS*, 2007, 19: 153.
- [3] Smolensky, P. "Information processing in dynamical systems: Foundations of harmony theory," MIT Press, 1986, pp. 194-281.
- [4] Hinton, G. E. "Training Products of Experts by Minimizing Contrastive Divergence," *Neural Computation*, 2002, 14.8, pp. 1771-1800.
- [5] Kameoka, H., et al. "Complex NMF: A new sparse representation for acoustic signals," In *Proc. of ICASSP*, 2009, p. 3437-3440.