

# 補助関数法による Gaussian-Bernoulli RBM の学習アルゴリズム\*

高宗 典玄 (東大院・情報理工) 亀岡 弘和 (東大院・情報理工, NTT CS 研)

## 1 はじめに

近年, 深層学習 (Deep learning) の有効性は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, 大量データの下で学習をいかに効率的に行えるかは重要課題の一つである. Deep Neural Network(DNN) の一種である Deep Belief Network(DBN)[1, 2] は制約付き Boltzmann マシン (Restricted Boltzmann Machine; RBM)[3] を多層に積み上げたものと見なせ, 各層の RBM の教師なし学習を順次行っていくことにより初期学習を行う方式が DBN の学習において効果的であることが知られている. RBM の学習アルゴリズムとして, Contrastive Divergence(CD) 法 [1, 4] が非常に有名であるが, RBM の学習をいかに効率的に行えるかが DBN 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]). そこで, 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と考え, Bernoulli-Bernoulli 型の RBM の補助関数法に基づく学習則を導出してきた [6].

Bernoulli-Bernoulli 型の RBM はバイナリデータしか取り扱えないが, 実際の音声や画像といったデータに用いるためには実数を取り扱う必要がある. そこで, 本発表では実数のデータに適用するために Gaussian-Bernoulli 型の 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 Gaussian-Bernoulli 型制約付き Boltzmann マシン

### 2.1 Gaussian-Bernoulli 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, 観測データに該当する状態を可視層, 背後にある特徴量に該当する状態を隠れ層と呼ぶ. 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. このとき, 可視層の状態数を  $I$ , 可視層の状態を  $\mathbf{v} = (v_i) \in (-\infty, \infty)^I$ , 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同時確率分布は

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

で定義される. ここで,

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

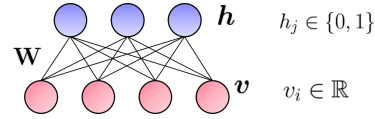


Fig. 1 Gaussian-Bernoulli RBM のグラフ表現

$$Z(\Theta) = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \sum_{\mathbf{h}} \exp(-E(\mathbf{v}, \mathbf{h} | \Theta)) d v_1 \cdots d v_I \quad (3)$$

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  $(\sigma_i)$  は定数として取り扱う.

RBM の学習問題とは, 観測される  $N$  個の可視層のデータ  $\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(N)}$  からこの分布パラメータ  $\Theta$  を推定することである.

このとき RBM の学習問題に対して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.

$$\begin{aligned} J(\Theta) &= \frac{1}{N} \sum_n \log p(\mathbf{v}^{(n)} | \Theta) \\ &= \frac{1}{N} \sum_n \log \sum_{\mathbf{h}} p(\mathbf{v}^{(n)}, \mathbf{h} | \Theta). \end{aligned} \quad (4)$$

ここで,  $\sum_{\mathbf{h}}$  は  $\mathbf{h}$  のとりうるすべての状態に関する和である.

### 2.2 Contrastive Divergence 法 [1, 4]

最急降下法により, 式 (4) 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.  $J(\Theta)$  を  $\Theta$  に関して微分すると,

$$\begin{aligned} \frac{\partial J}{\partial \Theta}(\Theta) &= \\ &= -\frac{1}{N} \sum_n \sum_{\mathbf{h}} p(\mathbf{h} | \mathbf{v}^{(n)}, \Theta) \frac{\partial E}{\partial \Theta}(\mathbf{v}^{(n)}, \mathbf{h} | \Theta) \\ &+ \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h} | \Theta) \\ &\quad \times \frac{\partial E}{\partial \Theta}(\mathbf{v}, \mathbf{h} | \Theta) d v_1 \cdots d v_I \end{aligned} \quad (5)$$

となる. このため, 最急降下法によるパラメータの更新は

$$\Theta \leftarrow \Theta^{\text{old}} + \epsilon \Delta_{\text{cd}} \Theta \quad (6)$$

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

ここで,  $\mathbf{h}$  についての和に関しては厳密に計算しようとする  $O(2^J)$  の計算量となるため, 現実的ではない. しかし, RBM の特徴として可視層同士, 隠れ層同士の依存関係が無いため,

$$p(\mathbf{v} | \mathbf{h}, \Theta) = \prod_i p(v_i | \mathbf{h}, \Theta), \quad (7)$$

$$p(\mathbf{h} | \mathbf{v}, \Theta) = \prod_j p(h_j | \mathbf{v}, \Theta) \quad (8)$$

\* “Training Gaussian-Bernoulli Restricted Boltzmann Machine Using Auxiliary Function Approach” by Takamune Norihiro (Univ. of Tokyo), Kameoka Hirokazu (Univ. of Tokyo, NTT CS Lab.).

という関係がある．ただし，

$$p(v_i | \mathbf{h}, \Theta) = \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left(-\frac{(v_i - b_i^V - \sigma_i \sum_j W_{ij} h_j)^2}{2\sigma_i^2}\right), \quad (9)$$

$$p(h_j = 1 | \mathbf{v}, \Theta) = \frac{1}{1 + \exp(-b_j^H - \sum_i W_{ij} v_i)} \quad (10)$$

である．そこで，式 (5) の第 1 項に関しては周辺化を行うことにより， $\sum_{\mathbf{h}}$  を  $\sum_j$  にすることが出来る．

しかし，式 (5) 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 (5) 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 Gibbs サンプリ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．

$$p(\mathbf{v}, \mathbf{h} | \Theta) \approx \frac{1}{M} \sum_m \delta(\mathbf{v} - \mathbf{v}^{(m)}) p(\mathbf{h} | \mathbf{v}^{(m)}, \Theta) \quad (11)$$

ここで， $\mathbf{v}^{(m)}$  は Gibbs サンプリグによりサンプリグされた値であり， $M$  はその個数である．本稿では， $M = N$  とし，各  $\mathbf{v}^{(m)}$  は対応する  $\mathbf{v}^{(n)}$  を初期値としてサンプリグされたものを用いる．ここで，式 (7)，(8) から Gibbs サンプリグは

$$h_j^{d-1} \sim p(h_j | \mathbf{v}^{d-1}, \Theta), \quad (12)$$

$$v_i^d \sim p(v_i | \mathbf{h}^{d-1}, \Theta) \quad (13)$$

と容易に行うことが可能である．

式 (6) による更新を式 (11) の Gibbs サンプリグによる近似で行う手法を CD 法という．

### 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とは，補助変数  $\mathbf{y}$  を導入して，任意の  $x$  について  $F(x) = \max_{\mathbf{y}} F^+(x, \mathbf{y})$  を満たすような下限関数  $F^+(x, \mathbf{y})$  を設計して， $F^+(x, \mathbf{y})$  を  $x$  についての最大化と  $\mathbf{y}$  についての最大化を交互に行うことである．ここで重要なのは， $F^+(x, \mathbf{y})$  を  $x$  についての最小化が容易に行え，かつ  $F^+(x, \mathbf{y})$  が  $F(x)$  によくフィットするように設計できるかであるので，本研究では  $x$  の各変数の相互依存をなくすような  $F^+(x, \mathbf{y})$  を設計することを目指す．

#### 3.1 最尤学習

そこでまず，式 (4) を目的関数として考えたとき，Jensen の不等式から

$$\begin{aligned} J(\Theta) &= \frac{1}{N} \sum_n \log \sum_{\mathbf{h}} p(\mathbf{v}^{(n)}, \mathbf{h} | \Theta) \\ &\geq \frac{1}{N} \sum_n \sum_{\mathbf{h}} \lambda_{n, \mathbf{h}} \log p(\mathbf{v}^{(n)}, \mathbf{h} | \Theta) \\ &\quad - \sum_{\mathbf{h}} \lambda_{n, \mathbf{h}} \log \lambda_{n, \mathbf{h}} \end{aligned} \quad (14)$$

となる．ここで，等号の成立は

$$\lambda_{n, \mathbf{h}} = p(\mathbf{h} | \mathbf{v}^{(n)}, \Theta) \quad (15)$$

である．式 (14) を整理すると，

$$\begin{aligned} J(\Theta) &\geq -\frac{1}{N} \sum_n \sum_{\mathbf{h}} \lambda_{n, \mathbf{h}} E(\mathbf{v}^{(n)}, \mathbf{h} | \Theta) \\ &\quad - \log Z(\Theta) - \sum_{\mathbf{h}} \lambda_{n, \mathbf{h}} \log \lambda_{n, \mathbf{h}} \end{aligned} \quad (16)$$

となる．ここで，式 (16) 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．

$$-\log Z(\Theta) \geq -\frac{Z(\Theta)}{\zeta} - \log \zeta + 1. \quad (17)$$

ここで，等号の成立は接点となるので，

$$\zeta = Z(\Theta) \quad (18)$$

となる．

さて，式 (2) から  $E(\mathbf{v}, \mathbf{h} | \Theta)$  はパラメータ  $\Theta$  の各要素に対し線形であるので，

$$a_k(\mathbf{v}, \mathbf{h}) = \begin{cases} \frac{v_i}{\sigma_i^2} & (k = i) \\ h_j & (k = I + j) \\ \frac{v_i h_j}{\sigma_i} & (k = I + J + J \times (i - 1) + j) \end{cases}, \quad (19)$$

$$\theta_k = \begin{cases} b_i^V & (k = i) \\ b_j^H & (k = I + j) \\ W_{ij} & (k = I + J + J \times (i - 1) + j) \end{cases} \quad (20)$$

とおくと，式 (2) は

$$E(\mathbf{v}, \mathbf{h} | \Theta) = \sum_i \frac{v_i^2}{2\sigma_i^2} - \sum_k a_k(\mathbf{v}, \mathbf{h}) \theta_k \quad (21)$$

と表すことができる．そこで， $-\exp(\cdot)$  は凹関数であるので，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，

$$\begin{aligned} &-\exp\left(\sum_k a_k(\mathbf{v}, \mathbf{h}) \theta_k\right) \\ &\geq -\sum_k \beta_k(\mathbf{v}, \mathbf{h}) \exp\left(\frac{a_k(\mathbf{v}, \mathbf{h}) \theta_k - \alpha_k(\mathbf{v}, \mathbf{h})}{\beta_k(\mathbf{v}, \mathbf{h})}\right) \end{aligned} \quad (22)$$

となり，等号の成立は

$$\alpha_k(\mathbf{v}, \mathbf{h}) = a_k(\mathbf{v}, \mathbf{h}) \theta_k - \beta_k(\mathbf{v}, \mathbf{h}) \sum_l a_l(\mathbf{v}, \mathbf{h}) \theta_l, \quad (23)$$

ただし， $\beta_k(\mathbf{v}, \mathbf{h})$  は

$$\beta_k(\mathbf{v}, \mathbf{h}) \in [0, 1], \quad \sum_k \beta_k(\mathbf{v}, \mathbf{h}) = 1 \quad (24)$$

を満たす任意定数である．そこで，計算の簡単のため  $\beta_k(\mathbf{v}, \mathbf{h})$  を  $\mathbf{v}, \mathbf{h}$  に依存しない定数  $\beta_k$  とする．更に，新たに補助変数として一反復前の値を表す  $\theta_k^{\text{old}}$  を導入し，補助変数の更新式 (23) を式 (22) に代入し，式 (17) を式 (3)，(21) 用いて整理し，式 (1)，(18) を用いて書き換えると，

$$\begin{aligned} &-\log Z(\Theta) \\ &\geq -\sum_k \beta_k \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h} | \Theta^{\text{old}}) \\ &\quad \times \exp\left(\frac{a_k(\mathbf{v}, \mathbf{h})(\theta_k - \theta_k^{\text{old}})}{\beta_k}\right) dv_1 \cdots dv_I \\ &\quad - \log \zeta + 1 \end{aligned} \quad (25)$$

となる．式 (15)，(16)，(18)，(25) 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

$$\begin{aligned} &J^+(\Theta, \bar{\Theta}) \\ &= \frac{1}{N} \sum_n \sum_{\mathbf{h}} p(\mathbf{h} | \mathbf{v}^{(n)}, \Theta^{\text{old}}) \sum_k a_k(\mathbf{v}^{(n)}, \mathbf{h}) \theta_k \\ &\quad - \sum_k \beta_k \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h} | \Theta^{\text{old}}) \end{aligned} \quad (26)$$

$$\times \exp\left(\frac{a_k(\mathbf{v}, \mathbf{h})(\theta_k - \theta_k^{\text{old}})}{\beta_k}\right) dv_1 \cdots dv_I$$

$$+ C(\bar{\Theta})$$

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

ここで，式 (26) の右辺の第 2 項は  $p(\mathbf{v}, \mathbf{h}|\Theta^{\text{old}})$  による期待値の計算になっており，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 Gibbs サンプリングによる  $\mathbf{v}$  の周辺確率の近似を行い，同時確率を式 (11) で近似をした後， $\theta_k$  について微分すると，

$$\frac{\partial J^+}{\partial \theta_k}(\Theta, \bar{\Theta})$$

$$= \frac{1}{N} \sum_n \sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{v}^{(n)}, \Theta^{\text{old}}) a_k(\mathbf{v}^{(n)}, \mathbf{h})$$

$$- \frac{1}{M} \sum_m \sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{v}^{(m)}, \Theta^{\text{old}}) a_k(\mathbf{v}^{(m)}, \mathbf{h})$$

$$\times \exp\left(\frac{a_k(\mathbf{v}^{(m)}, \mathbf{h})(\theta_k - \theta_k^{\text{old}})}{\beta_k}\right)$$
(27)

となる．しかし， $a_k(\mathbf{v}, \mathbf{h}) \in \{0, 1\}$  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 [6] が，Gaussian-Bernoulli 型の RBM では  $\theta_k$  が  $b_j^{\text{H}}$  のときのみ  $a_k(\mathbf{v}, \mathbf{h}) \in \{0, 1\}$  となる．ところが， $\partial^2 J^+/\partial \theta_k^2$  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，[6] と同型の補助関数を用いているので，収束を速くするために同様に以下の近似を考える．

$$\beta'_k \leftarrow \beta_k^\gamma. \quad (28)$$

ただし， $\beta'_k$  は式 (24) 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 サンプリングで  $p(\mathbf{v}, \mathbf{h}|\Theta^{\text{old}})$  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，補助関数法を用いた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．ここで用いる目的関数は可視層に観測データが来たときに，Gibbs サンプリングを 1 回行ったときに元の観測データが再構築される確率の対数で，

$$J_r(\Theta) = \frac{1}{N} \sum_n \log \sum_{\mathbf{h}} p(\mathbf{v}^{(n)}|\mathbf{h}, \Theta) p(\mathbf{h}|\mathbf{v}^{(n)}, \Theta)$$
(29)

と表される．式 (29) を直接最大化するのは困難であるので，この式に対して，

$$\mathbf{h}^{(n)} \sim p(\mathbf{h}|\mathbf{v}^{(n)}, \Theta) \quad (30)$$

でサンプリングした値を元に

$$J_r(\Theta)$$

$$\approx \frac{1}{N} \sum_n \log p(\mathbf{v}^{(n)}|\mathbf{h}^{(n)}, \Theta) p(\mathbf{h}^{(n)}|\mathbf{v}^{(n)}, \Theta) \quad (31)$$

と近似することを考える．この右辺を  $\tilde{J}_r(\Theta)$  とおくと，式 (7)~(10) より，

$$\tilde{J}_r(\Theta) = -\frac{1}{N} \sum_n \left\{ \sum_i \frac{(v_i^{(n)} - f_{ni}^{\text{V}})^2}{2\sigma_i^2} + \sum_i \log \sigma_i \right.$$

$$+ \frac{I}{2} \log(2\pi) - \sum_j \left( h_j^{(n)} - \frac{1}{2} \right) f_{nj}^{\text{H}} \quad (32)$$

$$\left. + \sum_j \log(\cosh(f_{nj}^{\text{H}}/2)) + J \log 2 \right\}$$

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

$$\text{ここで，}$$

$$\log(\cosh(x))$$

$$\leq \frac{\tanh(\kappa)}{2\kappa} x^2 - \frac{\kappa \tanh(\kappa)}{2} + \log(\cosh(\kappa)) \quad (33)$$

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

$$\kappa = \pm x \quad (34)$$

である．また，係数が負の 2 次関数は凹関数であるので，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ことができるため，紙面の都合上，詳細は割愛するが，式 (33) と Jensen の不等式から下限関数  $\tilde{J}_r^+(\Theta, \bar{\Theta})$  は

$$\tilde{J}_r^+(\Theta, \bar{\Theta}) = -\frac{1}{N} \sum_n \left\{ \sum_i \frac{(b_i^{\text{V}} - \alpha_{ni0}^{\text{V}})^2}{2\beta_{ni0}^{\text{V}} \sigma_i^2} \right.$$

$$+ \sum_i \sum_j \frac{(\sigma_i W_{ij} h_j^{(n)} - \alpha_{nij}^{\text{V}})^2}{2\beta_{nij}^{\text{V}} \sigma_i^2}$$

$$+ \sum_j \frac{\tanh(\kappa_{nj})}{8\kappa_{nj} \beta_{n0j}^{\text{H}}} (b_j^{\text{H}} - \alpha_{n0j}^{\text{H}})^2$$

$$+ \sum_j \sum_i \frac{\tanh(\kappa_{nj})}{8\kappa_{nj} \beta_{nij}^{\text{H}}} \left( W_{ij} v_i^{(n)}/\sigma_i - \alpha_{nij}^{\text{H}} \right)^2 \quad (35)$$

$$- \sum_i \frac{v_i^{(n)}}{\sigma_i^2} \left( b_i^{\text{V}} + \sigma_i \sum_j W_{ij} h_j^{(n)} \right)$$

$$\left. - \sum_j \left( h_j^{(n)} - \frac{1}{2} \right) \left( b_j^{\text{H}} + \sum_i W_{ij} v_i^{(n)}/\sigma_i \right) \right\}$$

$$+ C(\bar{\Theta})$$

となる．ただし，

$$\bar{\Theta} = (\kappa_{nj}, \alpha_{ni0}^{\text{V}}, \alpha_{nij}^{\text{V}}, \alpha_{n0j}^{\text{H}}, \alpha_{nij}^{\text{H}}, \beta_{ni0}^{\text{V}}, \beta_{nij}^{\text{V}}, \beta_{n0j}^{\text{H}}, \beta_{nij}^{\text{H}}) \quad (36)$$

であり， $C(\bar{\Theta})$  は  $\Theta$  に対して定数の項である．また，等号の成立は，

$$\kappa_{nj} = \frac{f_{nj}^{\text{H}}}{2}, \quad (37)$$

$$\alpha_{ni0}^{\text{V}} = b_i^{\text{V}} - \beta_{ni0}^{\text{V}} f_{ni}^{\text{V}}, \quad (38)$$

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

$$\alpha_{n0i}^{\text{H}} = b_j^{\text{H}} - \beta_{ni0}^{\text{H}} f_{nj}^{\text{H}}, \quad (39)$$

$$\alpha_{nij}^{\text{H}} = \frac{W_{ij} v_i^{(n)}}{\sigma_i} - \beta_{nij}^{\text{H}} f_{nj}^{\text{H}}$$

であり,  $\beta_{ni0}^V, \beta_{nij}^V, \beta_{n0j}^H, \beta_{nij}^H$  は

$$\begin{aligned} \beta_{ni0}^V, \beta_{nij}^V, \beta_{n0j}^H, \beta_{nij}^H &\in [0, 1], \\ \beta_{ni0}^V + \sum_j \beta_{nij}^V &= 1, \quad \beta_{n0j}^H + \sum_i \beta_{nij}^H = 1 \end{aligned} \quad (40)$$

を満たす任意定数である.

よって,  $b_i^V, b_j^H, W_{ij}$  の更新式は  $\partial \tilde{J}_r^+ / \partial \Theta = 0$  により求めることができる. このとき,  $\beta_{i0}^V, \beta_{ij}^V, \beta_{0j}^H, \beta_{ij}^H$  もまた, 補助関数法による最尤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  $\gamma$  乗にする近似を考える.

以上より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) を反復することとなる. 1) サンプルングにより  $h_j^{(n)}$  を求める. 2) 補助変数  $\Theta$  を求める. 3)  $\partial \tilde{J}_r^+ / \partial \Theta = 0$  から得られる更新式でパラメータを更新.

#### 4 動作確認実験

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, 人工データによる比較実験を行った. 可視層の状態数  $I = 1024$ , 隠れ層の状態数  $J = 512$  であり, 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  $N = 5000$  となる. また,  $\sigma_i$  は一様に 1 とし, 学習の反復回数  $T$  は 20000 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,  $\beta_k, \beta_{i0}^V, \beta_{ij}^V, \beta_{0j}^H, \beta_{ij}^H$  は一様とし, CD 法の学習率  $\epsilon$  や式 (28) に示す  $\gamma$  のスケジューリングを  $t$  回目の反復のとき

$$\epsilon(t) = \epsilon_{\text{init}} \left( \frac{\epsilon_{\text{end}}}{\epsilon_{\text{init}}} \right)^{\frac{t-1}{T-1}}, \quad (41)$$

$$\gamma(t) = \gamma_{\text{init}} \left( \frac{\gamma_{\text{end}}}{\gamma_{\text{init}}} \right)^{\frac{t-1}{T-1}} \quad (42)$$

とした. このとき,  $\epsilon_{\text{init}} = 0.001$ ,  $\epsilon_{\text{end}} = 0.0001$  とし, 補助関数法による最尤学習アルゴリズムに関しては  $\gamma_{\text{init}} = 0.03$ ,  $\gamma_{\text{end}} = 0.05$  とし, 補助関数法の最大再構築確率学習に関しては  $\gamma_{\text{init}} = 0.9$ ,  $\gamma_{\text{end}} = 1$  とした. また, Gibbs サンプルングの回数を 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, 補助関数法による最尤学習アルゴリズムによる計算時間も補助関数法の最大再構築確率学習アルゴリズムによる計算時間も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 (4) に示す対数尤度については厳密に計算することが困難であるので, 代わりに以下のように定義した再構築誤差  $e_{\text{reconst}}$  が各学習アルゴリズムによりどのように遷移したのかを Fig. 2 に示す.

$$\bar{\mathbf{h}}^{(n)} = \arg \max_{\mathbf{h}} p(\mathbf{h} | \mathbf{v}^{(n)}, \Theta), \quad (43)$$

$$\bar{\mathbf{v}}^{(n)} = \arg \max_{\mathbf{v}} p(\mathbf{v} | \bar{\mathbf{h}}^{(n)}, \Theta), \quad (44)$$

$$e_{\text{reconst}} = \frac{1}{NI} \sum_n \sum_i (v_i^{(n)} - \bar{v}_i^{(n)})^2. \quad (45)$$

$\gamma$  による学習率の近似をしなかったアルゴリズムは, 特に最尤学習の場合は収束が遅く, それに対し,  $\gamma$  に

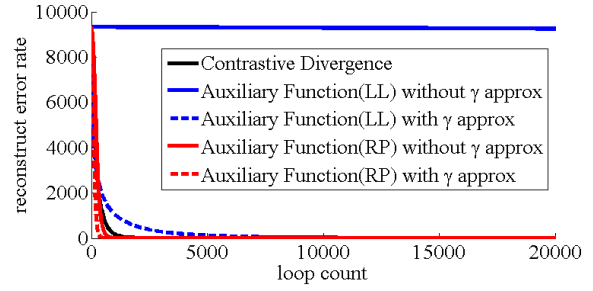


Fig. 2 各学習アルゴリズムにおける反復毎の再構築誤差. 黒い実線は CD 法を表し, 青, 赤の実線と破線はそれぞれ提案手法の  $\gamma$  による学習率の近似をしなかったものとしたものを表す.

よる学習率の近似をしたアルゴリズムは最尤学習, 最大再構築確率学習とも速く収束するという挙動が観測された. 特に最大再構築確率学習は  $\gamma$  の近似を行わなくても CD 法よりも速く誤差は小さくなるという結果が得られた.

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  $\gamma$  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, Gibbs サンプル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.

#### 5 まとめ

本稿では, Gaussian-Bernoulli 型の RBM の学習アルゴリズムとして, 補助関数法を用いて新たに 2 つの学習アルゴリズムを導出した. そして, 人工データによる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, 多層に重ねて深層学習を行ったときにどのような挙動を示すかの観察や実データへの適用が挙げられる.

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

#### 参考文献

- [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.
- [6] 高宗 典玄, 他. "補助関数法による制約付きボルツマンマシンの学習アルゴリズムの検討," 音講論 (春), 2014, No. 1-5-4.