

補助関数法による Gaussian-Bernoulli RBM の 学習アルゴリズムの検討

高宗 典玄^{1,a)} 亀岡 弘和^{1,2,b)}

概要：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おり、その重要な一要素として、制約付きボルツマンマシン (RBM) による pre-training がある。実数の観測データを取り扱うための Gaussian-Bernoulli RBM というモデルがあり、その学習アルゴリズムとして、最急降下法を基とした Contrastive Divergence 法が提案されてきた。そこで、本発表ではその学習問題に対して、経験的に高速で安定に収束する補助関数法による更新アルゴリズムを提案する。小規模な人工データによる実験を行い、その挙動に対して提案法と従来法を比較し議論する。

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] は制約付きボルツマンマシン (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 型制約付きボルツマンマシン

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。このとき、可視層の状態数を I 、可視層の状態を $v = \{v_i\} \in (-\infty, \infty)^I$ 、隠れ層の状態数を J 、隠れ層の状態を $h = \{h_j\} \in \{0, 1\}^J$ 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

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

で定義される。ここで、

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

$$Z(\Theta) = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \sum_h \exp(-E(v, h | \Theta)) dv_1 \cdots dv_I \quad (3)$$

であり、 $\Theta = (b_i^V, b_j^H, W_{ij})$ や (σ_i) は分布パラメータである。ここで、問題の単純化のために (σ_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 $v^{(1)}, \dots, v^{(N)}$ からこのパラメータ Θ を推定することである。

このとき 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。

¹ 東京大学大学院情報理工学系研究科
The University of Tokyo

² 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
NTT Communication Science Laboratories
3-1, Morinosato Wakamiya Atsugi-shi, Kanagawa, 243-0198, Japan

a) takamune@hil.t.u-tokyo.ac.jp

b) kameoka@hil.t.u-tokyo.ac.jp

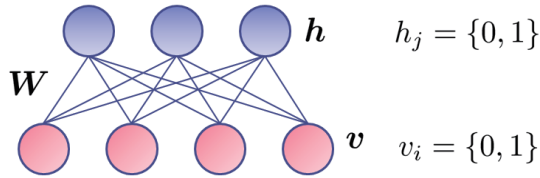


図 1 RBM のグラフ表現

$$J(\Theta) = \frac{1}{N} \sum_n \log p(v^{(n)}|\Theta)$$

$$= \frac{1}{N} \sum_n \log \sum_h p(v^{(n)}, h|\Theta). \quad (4)$$

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 Θ に関して微分すると,

$$\frac{\partial J}{\partial \Theta}(\Theta) = -\frac{1}{N} \sum_n \sum_h p(h|v^{(n)}, \Theta) \frac{\partial E}{\partial \Theta}(v^{(n)}, h|\Theta)$$

$$+ \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \sum_h p(v, h|\Theta) \frac{\partial E}{\partial \Theta}(v, h|\Theta) dv_1 \cdots dv_l$$

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, ϵ は学習率, $\Delta_{\text{cd}} \Theta = \partial J / \partial \Theta$ である.

ここで, h についての和に関しては厳密に計算しようとすると $O(2^J)$ 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,

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

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

という関係がある. ただし,

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

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

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_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 サンプル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.

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

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

ここで, 式 (7), (8) から Gibbs サンプルングは

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

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

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 サンプルン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

3.1 補助関数法 1

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 y を導入して, 任意の x, y で $F(x) \geq F^+(x, y)$ となり, $F(x) = \min_y F^+(x, y)$ となるような下限関数 $F^+(x, y)$ を設計して, $F^+(x, y)$ を x についての最小化と y 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, $F^+(x, y)$ を x についての最小化が容易に行え, かつ $F^+(x, y)$ 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 $F^+(x, y)$ を設計することを目指す.

そこで, 式 (4) による目的関数を考えたとき, Jensen の不等式から

$$J(\Theta) = \frac{1}{N} \sum_n \log \sum_h p(v^{(n)}, h|\Theta)$$

$$\geq \frac{1}{N} \sum_n \sum_h \lambda_{n,h} \log p(v^{(n)}, h|\Theta)$$

$$- \sum_h \lambda_{n,h} \log \lambda_{n,h}. \quad (14)$$

となる. ここで, 等号の成立は

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

である. 式 (14) を整理すると,

$$J(\Theta) \geq -\frac{1}{N} \sum_n \sum_h \lambda_{n,h} E(v^{(n)}, h|\Theta)$$

$$- \log Z(\Theta) - \sum_h \lambda_{n,h} \log \lambda_{n,h} \quad (16)$$

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.

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

ここで, 等号の成立は接点となるので,

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

となる．ここで，式 (2) から $E(\mathbf{v}, \mathbf{h}|\Theta)$ はパラメータ Θ 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．

ここで， $-\exp(\sum_k a_k(\mathbf{v}, \mathbf{h}) \theta_k)$ 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，

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

となり，等号の成立は

$$\begin{aligned} & \forall \beta_k(\mathbf{v}, \mathbf{h}) \in [0, 1], \\ & \sum_k \beta_k(\mathbf{v}, \mathbf{h}) = 1, \end{aligned} \quad (23)$$

$$\begin{aligned} & \alpha_k(\mathbf{v}, \mathbf{h}) = a_k(\mathbf{v}, \mathbf{h}) \theta_k \\ & - \beta_k(\mathbf{v}, \mathbf{h}) \sum_l a_l(\mathbf{v}, \mathbf{h}) \theta_l \end{aligned} \quad (24)$$

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 とできる．更に，新たに補助変数として一反復前の値を表す θ_k^{old} を導入し，補助変数の更新式 (24) を式 (22) に代入し，式 (17) を式 (3), (21) 用いて整理すると，

$$\begin{aligned} & -\log Z(\Theta) \\ & \geq -\int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \sum_{\mathbf{h}} \frac{\exp\left(-\sum_i \frac{v_i^2}{2\sigma_i^2} + \sum_l a_l(\mathbf{v}, \mathbf{h}) \theta_l^{\text{old}}\right)}{\zeta} \\ & \times \sum_k \beta_k \exp\left(\frac{a_k(\mathbf{v}, \mathbf{h}) (\theta_k - \theta_k^{\text{old}})}{\beta_k}\right) dv_1 \cdots dv_I \\ & -\log \zeta + 1 \end{aligned} \quad (25)$$

となる．ここで，式 (1), (18) から式 (25) は

$$\begin{aligned} & -\log Z(\Theta) \\ & \geq -\sum_k \beta_k \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h}|\Theta^{\text{old}}) \\ & \times \exp\left(\frac{a_k(\mathbf{v}, \mathbf{h}) (\theta_k - \theta_k^{\text{old}})}{\beta_k}\right) dv_1 \cdots dv_I \\ & -\log \zeta + 1 \end{aligned} \quad (26)$$

となる．式 (15), (16), (18), (26) 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

$$\begin{aligned} & J^+(\Theta, \bar{\Theta}) \\ & = \frac{1}{N} \sum_n \sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{v}^{(n)}, \Theta^{\text{old}}) \\ & \quad \times \sum_k a_k(\mathbf{v}^{(n)}, \mathbf{h}) \theta_k \\ & - \sum_k \beta_k \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h}|\Theta^{\text{old}}) \\ & \quad \times \exp\left(\frac{a_k(\mathbf{v}, \mathbf{h}) (\theta_k - \theta_k^{\text{old}})}{\beta_k}\right) dv_1 \cdots dv_I \\ & + C(\bar{\Theta}) \end{aligned} \quad (27)$$

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 (27) を θ_k について微分すると，

$$\begin{aligned} & \frac{\partial J^+}{\partial \theta_k}(\Theta, \bar{\Theta}) \\ & = \frac{1}{N} \sum_n \sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{v}^{(n)}, \Theta^{\text{old}}) a_k(\mathbf{v}^{(n)}, \mathbf{h}) \\ & - \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h}|\Theta^{\text{old}}) a_k(\mathbf{v}, \mathbf{h}) \\ & \quad \times \exp\left(\frac{a_k(\mathbf{v}, \mathbf{h}) (\theta_k - \theta_k^{\text{old}})}{\beta_k}\right) dv_1 \cdots dv_I \end{aligned} \quad (28)$$

となる．ここで，式 (28) 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 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) で近似すると式 (28) は

$$\begin{aligned} & \frac{\partial J^+}{\partial \theta_k}(\Theta, \bar{\Theta}) \\ & = \frac{1}{N} \sum_n \sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{v}^{(n)}, \Theta^{\text{old}}) a_k(\mathbf{v}^{(n)}, \mathbf{h}) \\ & - \frac{1}{M} \sum_m \sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{v}^{(m)}, \Theta^{\text{old}}) a_k(\mathbf{v}^{(m)}, \mathbf{h}) \\ & \quad \times \exp\left(\frac{a_k(\mathbf{v}^{(m)}, \mathbf{h}) (\theta_k - \theta_k^{\text{old}})}{\beta_k}\right) \end{aligned} \quad (29)$$

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 [6] 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，[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以下の近似を考える．

$$\beta_k' \leftarrow \beta_k^\gamma. \quad (30)$$

ただし， β_k' は式 (23) 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1つ目の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 補助関数法 2

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尤度としたもので、

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

と表される。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、

$$\mathbf{h}^{(n)} \sim p(\mathbf{h}|\mathbf{v}^{(n)}, \Theta) \quad (32)$$

でサンプリグした値を元に

$$J_r(\Theta) \approx \frac{1}{N} \sum_n \log p(\mathbf{v}^{(n)}|\mathbf{h}^{(n)}, \Theta) p(\mathbf{h}^{(n)}|\mathbf{v}^{(n)}, \Theta) \quad (33)$$

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、

$$\begin{aligned} \tilde{J}_r(\Theta) = & -\frac{1}{N} \sum_n \left\{ \sum_i \frac{(v_i^{(n)} - f_{ni}^V)^2}{2\sigma_i^2} \right. \\ & + \sum_i \log \sigma_i + \frac{I}{2} \log(2\pi) \\ & - \sum_j \left(h_j^{(n)} - \frac{1}{2} \right) f_{nj}^H \\ & \left. + \sum_j \log(\exp(f_{nj}^H/2) + \exp(-f_{nj}^H/2)) \right\} \end{aligned} \quad (34)$$

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、

$$\begin{aligned} & \log(\exp(x) + \exp(-x)) \\ & \leq \frac{\tanh(\kappa)}{2\kappa} x^2 - \frac{\kappa \tanh(\kappa)}{2} + \log(2 \cosh(\kappa)) \end{aligned} \quad (35)$$

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

$$\kappa = \pm x \quad (36)$$

である。この不等式を用いると式(34)は

$$\begin{aligned} \tilde{J}_r(\Theta) \geq & -\frac{1}{N} \sum_n \left\{ \sum_i \frac{(f_{ni}^V)^2}{2\sigma_i^2} \right. \\ & + \sum_j \frac{\tanh(\kappa_{nj})}{8\kappa_{nj}} (f_{nj}^H)^2 \\ & - \sum_i \frac{v_i^{(n)}}{\sigma_i^2} f_{ni}^V \\ & - \sum_j \left(h_j^{(n)} - \frac{1}{2} \right) f_{nj}^H \\ & + \sum_i \left(\frac{(v_i^{(n)})^2}{2\sigma_i^2} + \log \sigma_i \right) + \frac{I}{2} \log(2\pi) \\ & \left. + \sum_j \left(\log(\cosh(\kappa_{nj})) - \frac{\kappa_{nj} \tanh(\kappa_{nj})}{8} \right) \right\} \end{aligned} \quad (37)$$

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

$$\begin{aligned} -(f_{ni}^V)^2 \geq & -\frac{(b_i^V - \alpha_{ni0}^V)^2}{\beta_{ni0}^V} \\ & - \sum_j \frac{(\sigma_i W_{ij} h_j^{(n)} - \alpha_{nij}^V)^2}{\beta_{nij}^V}, \end{aligned} \quad (38)$$

$$\begin{aligned} -(f_{nj}^H)^2 \geq & -\frac{(b_j^H - \alpha_{n0j}^H)^2}{\beta_{n0j}^H} \\ & - \sum_i \frac{(W_{ij} v_i^{(n)}/\sigma_i - \alpha_{nij}^H)^2}{\beta_{nij}^H} \end{aligned} \quad (39)$$

となり、等号の成立は

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

$$\begin{aligned} & \forall \beta_{n0j}^H, \beta_{nij}^H \in [0, 1], \\ & \beta_{n0j}^H + \sum_i \beta_{nij}^H = 1, \end{aligned} \quad (41)$$

$$\begin{aligned} \alpha_{ni0}^V & = b_i^V - \beta_{ni0}^V f_{ni}^V, \\ \alpha_{nij}^V & = \sigma_i W_{ij} h_j^{(n)} - \beta_{nij}^V f_{ni}^V, \end{aligned} \quad (42)$$

$$\begin{aligned} \alpha_{n0i}^H & = b_j^H - \beta_{n0i}^H f_{nj}^H, \\ \alpha_{nij}^H & = \frac{W_{ij} v_i^{(n)}}{\sigma_i} - \beta_{nij}^H f_{nj}^H \end{aligned} \quad (43)$$

である。

よって、式(37)、式(38)、式(39)から下限関数 $\tilde{J}_r^+(\Theta, \bar{\Theta})$ は

$$\begin{aligned}
\tilde{J}_r^+(\Theta, \bar{\Theta}) = & -\frac{1}{N} \sum_n \left\{ \sum_i \frac{(b_i^V - \alpha_{ni0}^V)^2}{2\beta_{ni0}^V \sigma_i^2} \right. \\
& + \sum_i \sum_j \frac{(\sigma_i W_{ij} h_j^{(n)} - \alpha_{nij}^V)^2}{2\beta_{nij}^V \sigma_i^2} \\
& + \sum_j \frac{\tanh(\kappa_{nj})}{8\kappa_{nj} \beta_{n0j}^H} (b_j^H - \alpha_{n0j}^H)^2 \\
& + \sum_j \sum_i \frac{\tanh(\kappa_{nj})}{8\kappa_{nj} \beta_{nij}^H} (W_{ij} v_i^{(n)} / \sigma_i - \alpha_{nij}^H)^2 \quad (44) \\
& - \sum_i \frac{v_i^{(n)}}{\sigma_i^2} \left(b_i^V + \sigma_i \sum_j W_{ij} h_j^{(n)} \right) \\
& - \sum_j \left(h_j^{(n)} - \frac{1}{2} \right) \left(b_j^H + \sum_i W_{ij} v_i^{(n)} / \sigma_i \right) \left. \right\} \\
& + C(\bar{\Theta})
\end{aligned}$$

となる。ただし、

$$\bar{\Theta} = (\kappa_{nj}, \alpha_{ni0}^V, \alpha_{nij}^V, \alpha_{n0j}^H, \alpha_{nij}^H, \beta_{ni0}^V, \beta_{nij}^V, \beta_{n0j}^H, \beta_{nij}^H) \quad (45)$$

であり、 $C(\bar{\Theta})$ は $\bar{\Theta}$ 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

$$b_i^V \leftarrow \frac{\sum_n (\alpha_{ni0}^V / \beta_{ni0}^V + v_i^{(n)})}{\sum_n 1 / \beta_{ni0}^V}, \quad (46)$$

$$b_j^H \leftarrow \frac{\sum_n (\tanh(\kappa_{nj}) \alpha_{n0j}^H / 4\kappa_{nj} \beta_{n0j}^H + h_j^{(n)} - 1/2)}{\sum_n \tanh(\kappa_{nj}) / 4\kappa_{nj} \beta_{n0j}^H}, \quad (47)$$

$$W_{ij} \leftarrow \frac{\sum_n \left(\frac{\alpha_{nij}^V}{\beta_{nij}^V \sigma_i} + \frac{\tanh(\kappa_{nj}) \alpha_{nij}^H}{4\kappa_{nj} \beta_{nij}^H} + \frac{v_i^{(n)} (2h_j^{(n)} - 1/2)}{\sigma_i} \right)}{\sum_n \left(\frac{h_j^{(n)}}{\beta_{nij}^V \sigma_i} + \frac{\tanh(\kappa_{nj}) v_i^{(n)}}{4\kappa_{nj} \beta_{nij}^H \sigma_i} \right)} \quad (48)$$

となる。ここで、 $\beta_{ni0}^V, \beta_{nij}^V, \beta_{n0j}^H, \beta_{nij}^H$ は式 (40), (41) を満たす任意定数であるので、 n によらない定数 $\beta_{ni0}^V, \beta_{ij}^V, \beta_{0j}^H, \beta_{ij}^H$ 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 (36), (42), (43) を代入して整理すると

$$b_i^V \leftarrow b_i^{V,old} + \frac{\beta_{i0}^V}{N} \sum_n (v_i^{(n)} - f_{ni}^{V,old}), \quad (49)$$

$$b_j^H \leftarrow b_j^{H,old} + \beta_{0j}^H \frac{\sum_n (h_j^{(n)} - q_{nj})}{\sum_n (q_{nj} - 1/2) / f_{nj}^{H,old}}, \quad (50)$$

$$\begin{aligned}
W_{ij} \leftarrow & W_{ij}^{old} \quad (51) \\
& + \frac{\sum_n \left((v_i^{(n)} - f_{ni}^{V,old}) h_j^{(n)} + v_i^{(n)} (h_j^{(n)} - q_{nj}) \right)}{\sum_n \left(\frac{h_j^{(n)} \sigma_i}{\beta_{nij}^V} + \frac{(v_i^{(n)})^2 (q_{nj} - 1/2)}{\beta_{nij}^H \sigma_i f_{nj}^{H,old}} \right)}
\end{aligned}$$

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

$$\bar{\Theta} = (b_i^{V,old}, b_j^{H,old}, W_{ij}^{old}, \beta_{ni0}^V, \beta_{nij}^V, \beta_{n0j}^H, \beta_{nij}^H) \quad (52)$$

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1つ目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。

以上より補助関数法による3つ目の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 (49)~(51) の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、可視層の状態数 $I = 20$, 隠れ層の状態数 $J = 15$ という非常に小さな系で実験を行った。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 $N = 5000$ となる。また、 σ_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、つまり、

$$\beta_k = \frac{1}{I + J + IJ}, \quad (53)$$

$$\beta_{ni0}^V = \beta_{nij}^V = \frac{1}{1 + J}, \quad (54)$$

$$\beta_{n0j}^H = \beta_{nij}^H = \frac{1}{1 + I} \quad (55)$$

とし、CD法の学習率 ϵ や式 (30) に示す γ のスケジューリングを t 回目の反復のとき

$$\epsilon(t) = \epsilon_{init} \left(\frac{\epsilon_{end}}{\epsilon_{init}} \right)^{\frac{t-1}{T-1}}, \quad (56)$$

$$\gamma(t) = \gamma_{init} \left(\frac{\gamma_{end}}{\gamma_{init}} \right)^{\frac{t-1}{T-1}} \quad (57)$$

とした。このとき、 $\epsilon_{init} = 0.001$, $\epsilon_{end} = 0.0001$ とし、補助関数法1に関しては $\gamma_{init} = 0.03$, $\gamma_{init} = 0.01$ とし、補助関数法2に関しては $\gamma_{init} = 1$, $\gamma_{init} = 1$ とした。また、Gibbs 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、補助関数法1による計算時間も補助関数法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 (4) に示す対数尤度がどのように遷移したかを Fig. 2 に示し、以下のように定義した再構築誤差 $e_{reconsnt}$ 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. 3 に示す。

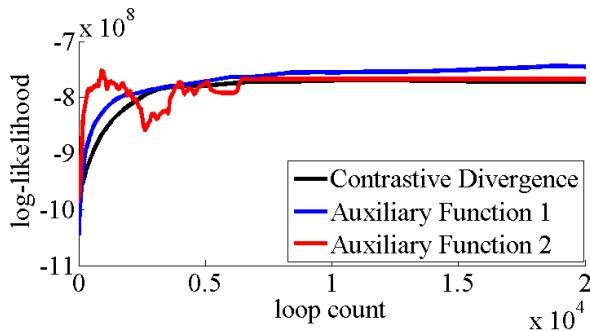


図 2 各学習アルゴリズムにおける反復毎の対数尤度．黒い実線は CD 法を表し，青，赤の実線と破線はそれぞれ提案手法を表す．

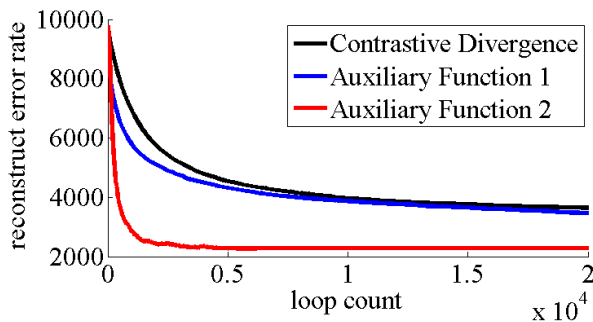


図 3 各学習アルゴリズムにおける反復毎の再構築誤差．黒い実線は CD 法を表し，青，赤の実線と破線はそれぞれ提案手法を表す．

$$\bar{h}^{(n)} = \arg \max_{\mathbf{h}} p(\mathbf{h} | \mathbf{v}^{(n)}, \Theta), \quad (58)$$

$$\bar{v}^{(n)} = \arg \max_{\mathbf{v}} p(\mathbf{v} | \bar{h}^{(n)}, \Theta), \quad (59)$$

$$e_{\text{reconst}} = \frac{1}{NI} \sum_n \sum_i (v_i^{(n)} - \bar{v}_i^{(n)}). \quad (60)$$

対数尤度による比較をみると，補助関数法 1 は CD 法よりも速く尤度が増大するが，補助関数法 2 は目的関数が対数尤度でないことが起因しているのか初期の反復において尤度の上下が激しくなるといった結果になった．また，再構築誤差による比較を見ると，補助関数法 1 は CD 法よりも若干速く誤差は小さくなり，補助関数法 2 にいたっては大幅に速く誤差は小さくなり，収束先も誤差がより小さいところに収束しそうな結果が得られた．

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，Gibbs 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．

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，小規模の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，可視層と隠れ層の状態数が多くなったときや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が挙げられる．

謝辞 議論に参加してくださった東大院・情報理工の石原達馬氏に深く感謝する．

参考文献

- [1] Hinton, G. E., Osindero, S. and Teh, Y. W.: "A fast learning algorithm for deep belief nets," *Neural Computation*, 2006, 18.7, pp. 1527-1554.
- [2] Bengio, Y., Lamblin, P., Popovici, D. and Larochelle, H.: "Greedy layer-wise training of deep networks," *Advances in neural information processing systems*, 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., Ono, N., Kashino, K. and Sagayama, S.: "Complex NMF: A new sparse representation for acoustic signals," *Acoustics, Speech and Signal Processing*, 2009. ICASSP 2009. IEEE International Conference on, 2009, p. 3437-3440.
- [6] 高宗 典玄, 石原 達馬, 亀岡 弘和: "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," *日本音響学会 2014 年春季研究発表会講演論文集*, 2014, No. 1-5-4.