

## 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\*

高宗典玄, 石原達馬 (東大院・情報理工研), 亀岡弘和 (東大院・情報理工, NTT CS 研)

## 1 はじめに

近年, 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. Deep learning 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]). 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 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 制約付きボルツマンマシン

## 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. このとき, 可視層の状態数を  $I$ , 可視層の状態を  $\mathbf{v} = \{v_i\} \in \{0, 1\}^I$ , 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.

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.

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

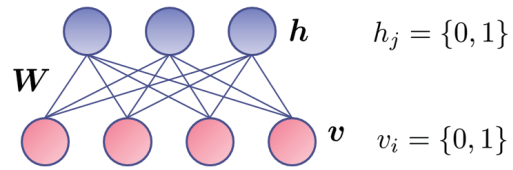


Fig. 1 RBM のグラフ表現

## 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,

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

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,  $\mathbf{v}$  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, 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,

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

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

という関係がある. ただし,

$$p(v_i = 1 | \mathbf{h}, \Theta) = \frac{1}{1 + \exp(-b_i^V - \sum_j W_{ij} h_j)}, \quad (9)$$

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

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.

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

ここで, 式 (7), (8) からギブスサンプリングは

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

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

\* “Training Restricted Boltzmann Machine with Auxiliary Function Approach” by Takamune Norihiro, Ishihara Tatsuma (Univ. of Tokyo), Kameoka Hirokazu (Univ. of Tokyo, NTT CS Lab.).

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。

### 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るように  $F^+(x, y)$  を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

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

となる。ここで、等号の成立は

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

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、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、

$$\zeta \leftarrow Z(\Theta^{\text{old}}) \quad (18)$$

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

$$E(\mathbf{v}, \mathbf{h} | \Theta) = -\sum_k a_k(\mathbf{v}, \mathbf{h}) \theta_k \quad (19)$$

とおく。ここで注意すべきことは  $a_k \in \{0, 1\}$  である。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、

$$\begin{aligned} &-\exp\left(\sum_k a_k(\mathbf{v}, \mathbf{h}) \theta_k\right) \\ &\geq -\sum_k \beta_k(\mathbf{v}, \mathbf{h}) \\ &\quad \times \exp\left(\frac{a_k(\mathbf{v}, \mathbf{h}) \theta_k - \alpha_k(\mathbf{v}, \mathbf{h})}{\beta_k(\mathbf{v}, \mathbf{h})}\right) \end{aligned} \quad (20)$$

となり、等号の成立は

$$\begin{aligned} &\forall \beta_k(\mathbf{v}, \mathbf{h}) \in [0, 1], \\ &\sum_k \beta_k(\mathbf{v}, \mathbf{h}) = 1, \end{aligned} \quad (21)$$

$$\begin{aligned} \alpha_k(\mathbf{v}, \mathbf{h}) &\leftarrow a_k(\mathbf{v}, \mathbf{h}) \theta_k^{\text{old}} \\ &\quad - \beta_k(\mathbf{v}, \mathbf{h}) \sum_l a_l(\mathbf{v}, \mathbf{h}) \theta_l^{\text{old}} \end{aligned} \quad (22)$$

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  $\beta_k$  とし、式 (20) に補助変数の更新式 (22) を代入し、式 (17) を式 (3), (19) 用いて整理すると、

$$\begin{aligned} &-\log Z(\Theta) \\ &\geq -\sum_{\mathbf{v}} \sum_{\mathbf{h}} \frac{\exp(\sum_l a_l(\mathbf{v}, \mathbf{h}) \theta_l^{\text{old}})}{\zeta} \\ &\quad \times \sum_k \beta_k \exp\left(\frac{a_k(\mathbf{v}, \mathbf{h}) (\theta_k - \theta_k^{\text{old}})}{\beta_k}\right) \\ &\quad - \log \zeta + 1 \end{aligned} \quad (23)$$

となる。ここで、式 (1), (18) から式 (23) は

$$\begin{aligned} &-\log Z(\Theta) \\ &\geq -\sum_k \beta_k \sum_{\mathbf{v}} \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h} | \Theta^{\text{old}}) \\ &\quad \times \exp\left(\frac{a_k(\mathbf{v}, \mathbf{h}) (\theta_k - \theta_k^{\text{old}})}{\beta_k}\right) \\ &\quad - \log \zeta + 1 \end{aligned} \quad (24)$$

となる。式 (15), (16), (18), (24) 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

$$\begin{aligned} J^+(\Theta, \bar{\Theta}) &= \frac{1}{N} \sum_n \sum_{\mathbf{h}} p(\mathbf{h} | \mathbf{v}^{(n)}, \Theta^{\text{old}}) \\ &\quad \times \sum_k a_k(\mathbf{v}^{(n)}, \mathbf{h}) \theta_k \\ &\quad - \sum_k \beta_k \sum_{\mathbf{v}} \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h} | \Theta^{\text{old}}) \\ &\quad \times \exp\left(\frac{a_k(\mathbf{v}, \mathbf{h}) (\theta_k - \theta_k^{\text{old}})}{\beta_k}\right) \\ &\quad + C(\bar{\Theta}) \end{aligned} \quad (25)$$

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 (25) を  $\theta_k$  について微分すると、

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

となる。ここで、 $a_k(\mathbf{v}, \mathbf{h}) \in \{0, 1\}$  である事実を用

いと、式 (26) の第二項は

$$\begin{aligned} & \sum_{\mathbf{v}} \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h} | \Theta^{\text{old}}) a_k(\mathbf{v}, \mathbf{h}) \\ & \quad \times \exp\left(\frac{a_k(\mathbf{v}, \mathbf{h})(\theta_k - \theta_k^{\text{old}})}{\beta_k}\right) \\ & = \exp\left(\frac{\theta_k - \theta_k^{\text{old}}}{\beta_k}\right) \\ & \quad \times \sum_{\mathbf{v}} \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h} | \Theta^{\text{old}}) a_k(\mathbf{v}, \mathbf{h}) \end{aligned} \quad (27)$$

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

$$\begin{aligned} \Delta_{\text{af}} \theta_k & = \log \frac{1}{N} \sum_n \sum_{\mathbf{h}} p(\mathbf{h} | \mathbf{v}^{(n)}, \Theta^{\text{old}}) a_k(\mathbf{v}^{(n)}, \mathbf{h}) \\ & \quad - \log \sum_{\mathbf{v}} \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h} | \Theta^{\text{old}}) a_k(\mathbf{v}, \mathbf{h}) \end{aligned} \quad (28)$$

と定義すると、 $\theta_k$  の更新式は、

$$\theta_k \leftarrow \theta_k^{\text{old}} + \beta_k \Delta_{\text{af}} \theta_k \quad (29)$$

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 (29) の計算が可能となる．

さて、式 (6)、(29) を比較すると、式 (29) の更新式はあたかも学習率が  $\beta_k$  であるかのような形となっている．ところが、式 (21) よりパラメータ数が増えると  $\beta_k$  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．

$$\beta'_k \leftarrow \beta_k^\gamma \quad (30)$$

このとき、 $\gamma \in [0, 1]$  ならば、 $\beta_k \in [0, 1]$  を満たすので  $\beta'_k$  は  $\beta_k$  よりも大きくなる．そのため、式 (29) による更新が速くなることが期待できる．ただし、 $\beta'_k$  は式 (21) 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  $\Theta$  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 (29) の更新式でパラメータを更新．

### 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、ギブスサンプリングを 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、

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

でサンプリングした値を元に

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

と近似することを考える．この右辺を  $\tilde{J}_r$  とおくと、式 (7)、(8) より、

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

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 (9)、(10) より式 (34) を整理すると、

$$\begin{aligned} \tilde{J}_r & = -\frac{1}{N} \sum_n \left\{ \sum_i v_i^{(n)} \log (1 + \exp(-f_{ni}^V)) \right. \\ & \quad + \sum_i (1 - v_i^{(n)}) \log (1 + \exp(f_{ni}^V)) \\ & \quad + \sum_j h_j^{(n)} \log (1 + \exp(-f_{nj}^H)) \\ & \quad \left. + \sum_j (1 - h_j^{(n)}) \log (1 + \exp(f_{nj}^H)) \right\} \end{aligned} \quad (35)$$

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 (35) 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 (17) ~ (24) 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

$$\begin{aligned} \tilde{J}_r^+(\Theta, \bar{\Theta}) & = -\frac{1}{N} \left\{ \sum_i v_i^{(n)} (1 - \hat{q}_{ni}^V) \xi_{ni}^V \right. \\ & \quad + \sum_i (1 - v_i^{(n)}) \hat{q}_{ni}^V \eta_{ni}^V \\ & \quad + \sum_j h_j^{(n)} (1 - \hat{q}_{nj}^H) \xi_{nj}^H \\ & \quad \left. + \sum_j (1 - h_j^{(n)}) \hat{q}_{nj}^H \eta_{nj}^H \right\} + C(\bar{\Theta}) \end{aligned} \quad (36)$$

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

$$\begin{aligned} & \beta_{i0}^V, \beta_{ij}^V, \beta_{0j}^H, \beta_{ij}^H \in [0, 1], \\ & \beta_{i0}^V + \sum_j \beta_{ij}^V = 1, \quad \forall i, \\ & \beta_{0j}^H + \sum_i \beta_{ij}^H = 1, \quad \forall j, \end{aligned} \quad (37)$$

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  $\Theta$  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。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  $\gamma$  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)  $\partial \tilde{J}_r^+ / \partial \Theta = 0$  から求ま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、可視層の状態数  $I = 10$ 、隠れ層の状態数  $J = 8$  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乱数で20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個ずつ用意した。つまり、入力するデータ数は  $N = 2000$  となる。また、学習の反復回数  $T$  は500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、

$$\beta_k = \frac{1}{I + J + IJ}, \quad (38)$$

$$\beta_{i0}^V = \beta_{ij}^V = \frac{1}{1 + J}, \quad (39)$$

$$\beta_{0j}^H = \beta_{ij}^H = \frac{1}{1 + I} \quad (40)$$

とし、CD法の学習率  $\epsilon$  や式(30)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、 $\epsilon_{\text{init}} = 1, \epsilon_{\text{end}} = 0.1, \gamma_{\text{init}} = 0.1, \gamma_{\text{end}} = 1$  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、補助関数法1による学習時間はおおよそ同じくらいであり、補助関数法2による学習時間はおおよそ3/4倍となった。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に示す。予想したように  $\gamma$  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、 $\gamma$  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、 $\gamma$  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  $\gamma$  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、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

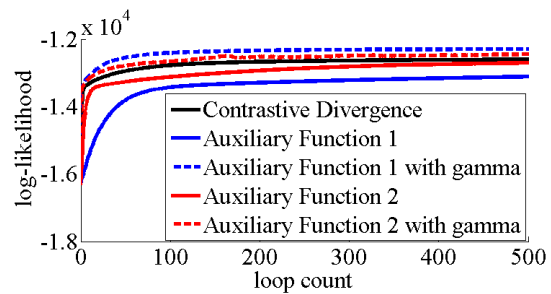


Fig. 2 各学習アルゴリズムにおける反復毎の対数尤度。黒い実線はCD法を表し、青、赤の実線と破線はそれぞれ提案手法の  $\gamma$  による学習率の近似をしなかったものとしたものを表す。

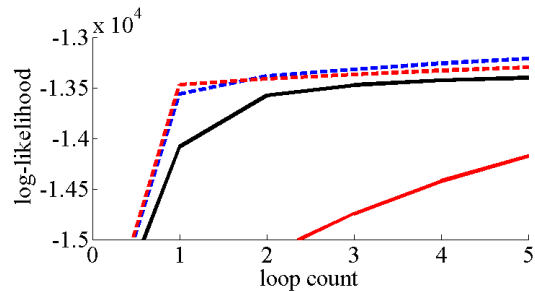


Fig. 3 Fig. 2の一部拡大

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(4)とは異なるものから出発した学習アルゴリズムである2つ目の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  $\gamma$  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。

#### 5 まとめ

本稿では、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。

#### 参考文献

- [1] Hinton, G. E., et al. "A fast learning algorithm for deep belief nets," *Neural Computation*, 2006, 18.7, 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: 1771-1800.
- [5] H. Kameoka, et al. "Complex NMF: A new sparse representation for acoustic signals," In *Proc of ICASSP*, 2009, p. 3437-3440.