

エントロピー項の導入による バックプロパゲーション学習則の拡張と学習効率の向上

秋 山 泰†

本論文では、多層型ニューラルネットにおいて最も広く利用されているバックプロパゲーション学習則（以下、BP 則）の新しい改良法を提案する。提案する学習則は、既存の BP 則の自然かつ強力な拡張となっている。学習の目的関数に、教師信号との出力誤差の最小化だけでなく、エントロピーとよぶ量の最大化を加える。エントロピーは、各ニューロンの出力値が 0.5 に近い曖昧な値を取るほど大きくなる量なので、エントロピー項を導入すると各ニューロンは 0 や 1 に極端に近い出力値は取りにくくなり、学習途上で入力荷重の絶対値が過度に増大することが自然に防がれる。この性質は学習のローカルミニマムからの脱出に重要である。エントロピー項の導入は、学習空間全体の起伏を減らし、より簡単な学習問題へ変形する操作と考えられるが、一方では学習を通じて得られる結合荷重が当初の目的からずれる欠点がある。そこで、学習の序盤にはエントロピー項の寄与を大きくしておき、中盤以降から徐々にエントロピー項の寄与を減じながら学習を進める技法を導入する。一般に、隠れ層のニューロン数を減らすと、学習は困難になるが汎化能力は向上する。エントロピー項の動的制御により、このような困難な学習の効率改善が期待できる。学習のローカルミニマムが知られる単純な例題を用いて、本提案が収束効率を改善することを実験的に示す。また、誤差の基準として Kullback 情報量を使うことが提案されているが、エントロピー項の導入を組み合わせると、学習効率が特に改善される。

Improvement of the Back-propagation Learning Rule by Introducing an Entropy Term

YUTAKA AKIYAMA†

An extension of the error back-propagation (BP) learning rule which performs entropy maximization as well as error minimization is proposed in this paper. The learning rule is very simple and powerful as a very natural extension of the BP learning rule. The entropy term has an effect of keeping all weights in a moderate range of value so that the neural network can easily escape from a local minimum. We also propose another kind of extended back-propagation learning rule which employs Kullback's divergence for its error function. Finally we show that dynamic control of the entropy term performs better than using a fixed coefficient for the entropy term.

1. はじめに

本論文では、学習の目標を表す目的関数 ϕ において、各ニューロンの出力値に関するエントロピーとよぶ量の最大化を新たに加えた、バックプロパゲーション学習則 (BP 則)^{1),2)}の拡張法を提案する。

学習則にエントロピー項を導入すると、全てのニューロンはネットワーク出力と教師信号との誤差を減じるという本来の学習目標と同時に、ニューロンの出力値を 0.5 に近づけるといった比較的弱い目標も与えられる。

学習がなかなか進まない時に、個別のニューロンの結合荷重の変化を観察すると、初期に過大な正值（負値）を取った荷重が、後に負値（正值）に転向するのに多大な時間をかけているケースが多く見られる。言うまでもなく、これにはニューロンの出力関数の非線形性が関係しており、出力値がいったん 0 や 1 に近づけられた時には、入力値、ひいては結合荷重（の絶対値）が極めて大きな値を取らされているからである。エントロピー項を導入すると、各ニューロンが 0 や 1 に近い出力値を出すことを避けるようになり、結合荷重（の絶対値）が学習途上に極端に増大することを抑制できる。この性質はローカルミニマムからの脱出時間の短縮に重要であり、学習効率の改善が期待できる。

† 技術研究組合 新情報処理開発機構

Real World Computing Partnership

☆ 損失関数とも呼ばれる

表1 バックプロパゲーション学習則の改良に関する既存の提案例
Table 1 Previous ideas for improving back-propagation learning

収束の加速	慣性項 ¹⁾ , 直線探索など ^{3)~6)} 共役勾配法, 準ニュートン法 ^{7),8)}
パターン提示法	一括修正, 追記学習, 補習学習
ネットワークの簡素化	結合荷重の忘却 ⁹⁾ 隠れ層ニューロン数の増減 ^{10)~15)}
誤差項の変更	Kullback 情報量の利用 ^{16),17)}

出力層ニューロンについては, 教師信号を 0.1 や 0.9 のような値に設定して, 結合荷重の過大化を防ぐという ad hoc な工夫が広く知られているが, 隠れ層ニューロンには適用できなかった. 本研究は, これと同様の効果を, より自然な形で理論化し, 隠れ層ニューロンにも拡張したものである.

BP 学習則については, 1986 年の提案^{1),2)}以来, 世界中できわめて活発な研究が行なわれ, 多くの改良案が提案されている(表 1). 改良は多方面にわたっており, 学習空間内の探索技法としての収束の加速法や, パターン提示の順番に関する工夫, ネットワークのアーキテクチャ自体の動的変更, そして誤差項の変更などに大きく分類される. これらは組み合わせて適用できる.

本研究は, これらの分類の中では誤差項の変更にも関係が深い. すなわち, 誤差項の変更以外の既存の改善提案が, 学習の対象となる荷重-対-評価値の空間(学習空間)がまず与えられ, その複雑な空間内の最適点(最適な結合荷重)をいかに高速に見つけるかを検討するものであるのに対して, 誤差項の変更や, 本研究におけるエントロピー項の追加は, 学習空間の形態そのものを変更しようとしている点が特色である. 本提案は, 既存の多くの加速法とは独立であるため, これらと組み合わせて用いることができる.

本論文では, 2 章において, 提案するエントロピー項付き BP 則の定義を示し, 具体的な学習方程式等を導出する. またエントロピー項の動的調整のアイデアを述べる. 3 章では, 2 つの例題による計算機実験を通じて, エントロピー項付き BP 則が学習効率を実際に改善することを実験的に例証する. 4 章では, 誤差項として通常の自乗誤差ではなく Kullback 情報量を採用した場合のエントロピー項付き BP 則を導出し, 計算機実験の結果を示す. 5 章では, 本研究の意義に関する議論とまとめを行なう.

2. エントロピー項付きバックプロパゲーション学習則

本章では, エントロピー項付きバックプロパゲーション学習則を定義し, 具体的な学習方程式を導出する. エントロピー項付き BP 則の定義は, Rumelhart らによるオリジナルの BP 則^{1),2)}に単に 1 つの項を付け加えただけであり, きわめて単純である.

後に示すように, 対数関数に基づくエントロピー項は, 指数関数から組み立てられるニューロンの出力関数と数学的に深く結び付いているため, 最終的に導出される学習方程式は, 驚くほど単純なものとなる.

2.1 ニューラルネットの構造と入出力動作

通常の BP 学習則と同様, 多層型のニューラルネットを対象とする. ニューロン i は前層に属するニューロンから出力値を受け取り, 次式に従って入力 u_i と, 出力値 v_i を計算する. 出力関数 $f(x)$ には, BP 学習則において広く用いられるロジスティック関数を用いる. 出力値は $0 < v_i < 1$ の範囲をとる.

$$u_i = \sum_{j \in C(i)} w_{ij} v_j \quad (1)$$

$$v_i = f(u_i) = \frac{1}{1 + \exp(-u_i)} \quad (2)$$

ここで $C(i)$ はニューロン i から見て前層に属するニューロンの番号の集合, w_{ij} はニューロン j からニューロン i への結合荷重を表す. 単純化のため, しきい値については, 常に 1 を出力する仮想ニューロンとの結合と考え, 結合荷重に含めて扱うこととするが, 議論の一般性は失われない.

2.2 目的関数の定義

目的関数 L および学習のポテンシャル関数 F を以下のように定義する. 学習のポテンシャル関数 F は, 入出力パターン p が固定されたときに最小化すべき量を表わすのに対して, 目的関数 L は環境内でのパターンの出現頻度までを考慮して, 学習全体を通じて最小化すべき量を表現する.

$$L(w) \stackrel{\text{def}}{=} \sum_p P(p) \cdot F_p(w) \quad (3)$$

$$F_p(w) \stackrel{\text{def}}{=} E_p(w) - \lambda \cdot S_p(w) \quad (4)$$

ここで w は系内の全ての結合荷重を表わすベクトル, $P(p)$ は入出力パターン p の出現確率, E_p は入出力パターン p に対する出力の誤差項, S_p が新たに提案するエントロピー項である. 非負の係数 λ によってエントロピー項の影響度が調整できる.

誤差項については通常どおりに, 以下の自乗誤差を

用いる.

$$E_p(\mathbf{w}) \stackrel{\text{def}}{=} \frac{1}{2} \sum_{i \in O} (t_{pi} - v_{pi})^2 \quad (5)$$

v_{pi} はパターン p に対するニューロン i の出力値, t_{pi} は教師信号, O は出力層に属するニューロンの番号の集合を表す.

新たに導入するエントロピー項は,

$$S_p(\mathbf{w}) \stackrel{\text{def}}{=} \sum_{i \in O \cup H} s(v_{pi}) \quad (6)$$

$$s(x) \stackrel{\text{def}}{=} x \cdot \log_e \frac{1}{x} + (1-x) \cdot \log_e \frac{1}{1-x} \quad (7)$$

ただし H は隠れ層に属するニューロンの番号の集合を表す.

エントロピー関数 $s(x)$ は, $x = 0.5$ の時最大値 ($\log_e 2$) をとる. エントロピー項の導入により, 出力層および隠れ層の全てのニューロンは, その出力値を 0.5 に近づけることを求められる.

$x = 0.5$ で最大値をとるような関数は他にも考えられるが, $s(x)$ の微分は, 出力関数 $f(x)$ の逆関数との間に次の密接な関係があり, 後に導出する学習方程式が, 美しく単純な形になるのでこれを採用した.

$$\frac{d}{dx} s(x) = -\log_e \frac{x}{1-x} = -f^{-1}(x) \quad (8)$$

2.3 学習方程式

学習ポテンシャル F_p の, 結合荷重 w_{ij} による偏微分を求め,

$$\Delta_p w_{ij} = -\eta \frac{\partial F_p}{\partial w_{ij}} \quad (9)$$

に従って結合荷重を修正することにより, 正の係数 η (学習率) が十分に小さいとき, 結合荷重 \mathbf{w} は目的関数 L に関して最適なものに近づいていく^{(18),(19)}.

以下では偏微分 $\partial F_p / \partial w_{ij}$ を計算することにより, 実際の学習則を導く. 導出の流れは, オリジナルの BP 則における導出^{(1),(2)} と本質的には同じである. エントロピー項のために修正される部分に興味がある.

$$\begin{aligned} \frac{\partial F_p}{\partial w_{ij}} &= \frac{\partial E_p}{\partial w_{ij}} - \lambda \frac{\partial S_p}{\partial w_{ij}} \\ &= \left(\frac{\partial E_p}{\partial u_{pi}} - \lambda \frac{\partial S_p}{\partial u_{pi}} \right) \cdot \frac{\partial u_{pi}}{\partial w_{ij}} \\ &= \left(\frac{\partial E_p}{\partial u_{pi}} - \lambda \frac{\partial S_p}{\partial u_{pi}} \right) \cdot v_{pj} \end{aligned} \quad (10)$$

ここで, 誤差信号 δ_{pi} を以下のように定義する.

$$\delta_{pi} \stackrel{\text{def}}{=} -\left(\frac{\partial E_p}{\partial u_{pi}} - \lambda \frac{\partial S_p}{\partial u_{pi}} \right)$$

学習方程式は以下のようにまとめられる.

$$-\frac{\partial F_p}{\partial w_{ij}} = \delta_{pi} \cdot v_{pj}$$

誤差信号 δ_{pi} の具体的な求め方については, 出力層ニューロンと隠れ層ニューロンの場合に分けて, 次節で導出する.

離散時間モデルの場合, 例えば以下のような学習方程式を実際には用いることができる. η は学習率, α は慣性率¹⁾ と呼ばれる係数である.

$$\begin{aligned} \Delta_p w_{ij}(t+1) &= -\eta \cdot \frac{\partial F_p}{\partial w_{ij}} + \alpha \cdot \Delta_p w_{ij}(t) \\ &= -\eta \cdot \delta_{pi} \cdot v_{pj} \\ &\quad + \alpha \cdot \Delta_p w_{ij}(t) \end{aligned} \quad (11)$$

2.4 誤差信号

1) ニューロン i が出力層に属する場合

$$\begin{aligned} \delta_{pi} &= -\left(\frac{\partial E_p}{\partial v_{pi}} - \lambda \frac{\partial S_p}{\partial v_{pi}} \right) \cdot \frac{\partial v_{pi}}{\partial u_{pi}} \\ &= (t_{pi} - v_{pi} - \lambda f^{-1}(v_{pi})) \cdot f'(u_{pi}) \\ &= (t_{pi} - v_{pi} - \lambda u_{pi}) \cdot f'(u_{pi}) \end{aligned} \quad (12)$$

2) ニューロン i が隠れ層に属する場合

$$\delta_{pi} = -\left(\frac{\partial E_p}{\partial v_{pi}} - \lambda \frac{\partial S_p}{\partial v_{pi}} \right) \cdot \frac{\partial v_{pi}}{\partial u_{pi}} \quad (13)$$

出力誤差への影響度 $\partial E_p / \partial v_{pi}$ は, 以下のように再帰的に計算できる¹⁾.

$$\begin{aligned} \frac{\partial E_p}{\partial v_{pi}} &= \sum_k \frac{\partial E_p}{\partial u_{pk}} \cdot \frac{\partial u_{pk}}{\partial v_{pi}} \\ &= \sum_k \frac{\partial E_p}{\partial u_{pk}} \cdot \frac{\partial}{\partial v_{pi}} \sum_l w_{kl} \cdot v_{pl} \\ &= \sum_k \frac{\partial E_p}{\partial u_{pk}} \cdot w_{ki} \end{aligned} \quad (14)$$

エントロピーへの影響度 $\partial S_p / \partial v_{pi}$ は, 後ろの層に関する再帰的計算と, 自らの出力値を介した直接の影響との和となり,

$$\begin{aligned} \frac{\partial S_p}{\partial v_{pi}} &= \sum_k \frac{\partial S_p}{\partial u_{pk}} \cdot \frac{\partial u_{pk}}{\partial v_{pi}} + \frac{\partial s(v_{pi})}{\partial v_{pi}} \\ &= \left(\sum_k \frac{\partial S_p}{\partial u_{pk}} \cdot w_{ki} \right) - u_{pi} \end{aligned} \quad (15)$$

したがって,

$$\begin{aligned} \delta_{pi} &= \left(-\sum_k \frac{\partial E_p - \lambda S_p}{\partial u_{pk}} \cdot w_{ki} - \lambda u_{pi} \right) \cdot \frac{\partial v_{pi}}{\partial u_{pi}} \\ &= \left(\sum_k \delta_{pk} \cdot w_{ki} - \lambda u_{pi} \right) \cdot f'(u_{pi}) \end{aligned} \quad (16)$$

以上の結果をまとめると,

$$\delta_{pi} =$$

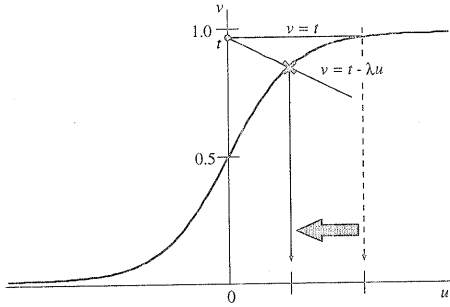


図1 エントロピー項による学習目標値の変化
Fig. 1 Learning target modification caused by the entropy term

$$\begin{cases} (t_{pi} - v_{pi} - \lambda u_{pi}) \cdot f'(u_{pi}) & i \in O \\ (\sum_k \delta_{pk} \cdot w_{ki} - \lambda u_{pi}) \cdot f'(u_{pi}) & i \in H \end{cases} \quad (17)$$

既存のバックプロパゲーション学習則との違いは、下線で示した “ $-\lambda u_{pi}$ ” の部分のみである。 $\lambda = 0$ とすれば既存の BP 則と一致する。

2.5 エントロピー付き BP 則の解釈

次に、エントロピー項を導入した BP 則は、既存の BP 則と比べて具体的にどのように異なるのかを、単一ニューロンのレベルでミクロに説明する。

(17) 式に示されたように、学習における誤差信号 δ_{pi} は、既存の学習則に比べて “ $-\lambda u_{pi}$ ” だけシフトしている。このとき、

- 出力値が $v_{pi} > 0.5$ ならば、入力 u_{pi} は正で、 $-\lambda u_{pi}$ は負値、
- 出力値が $v_{pi} < 0.5$ ならば、入力 u_{pi} は負で、 $-\lambda u_{pi}$ は正值、

であるから、誤差信号のシフトは必ずニューロンの出力値 v_{pi} が 0.5 に向かう方向に起きることがわかる。

具体的に学習の目標がどのように変わるかを、出力層ニューロンの場合を例にして考える。既存の BP 学習則 ($\lambda = 0$) の場合、学習の目標は出力値 v_{pi} が教師信号 t_{pi} と正確に一致することである。この時入力 u_{pi} の目標値は、当然ながら曲線 $v = f(u)$ と直線 $v = t$ の交点の u 座標となる (図 1)。

t が 0 または 1 に近い場合には、 u の目標値 (の絶対値) は極めて大きくなってしまふ。このような場合には、学習中に特定の結合荷重が過大となり、ローカルミニマムからの脱出が遅れるなどの問題が起きる。

いっぽう $\lambda > 0$ の場合、 u の目標は、曲線 $v = f(u)$ と直線 “ $v = t - \lambda u$ ” の交点で表される値となり、入力 u の目標値は原点に近い範囲となるので、結合荷

重が過大になることを防げる (図 1)。

出力層のニューロンについてだけならば、前述したように、与える教師信号を調整する方法 (1.0 の代わりに 0.9 を与える等) も可能であるが、本研究で提案するエントロピー項の場合には、隠れ層のニューロンについても、同様に誤差信号を調整できる点が特徴となる。

2.6 エントロピー項の動的調整

学習の序盤においては有効に働いたエントロピー項も、学習が終盤に差しかかってくると、逆に細部の収束を妨げる原因となる。誤差信号 δ_{pi} 中で本来の出力誤差由来の成分が次第に小さくなりエントロピー項の影響ばかりが残るからである。もっとも極端な場合には、各ニューロンの出力値が 0.5 付近に留まり学習はまったく進まない。

このためエントロピー項を固定的に導入する場合には、係数 λ を、学習を壊さない程度に小さく、かつ、導入の効果は得られる程度に大きく設定する必要がある。このままでは、 λ の値の最適範囲が狭く、例題にも依存して変化するので、実際の運用は難しくなる。

λ に関する最適範囲をあらかじめ正確に知ることは容易ではないので、運用上の方策として、係数 λ を動的に調整する。動的調整の方法としては、

- 1) あらかじめ定めたスケジュールに従い λ を単調に減衰させる。
- 2) 学習の進行をモニターしながら λ を上下する。

の 2 通りが考えられる。2) の方法は、シミュレーテッドアニーリングにおいて、比熱 (温度変化あたりのエネルギー変化) を観測しながら、温度低下スケジュールを細かく制御する技法との類似性がある。この方法は興味深いが、観測と実施の間の遅れの扱いなどが複雑で、運用上はパラメータ調整の労苦が増える可能性がある。本論文の目的は、エントロピー項の効果を簡潔に示すことにあるので、1) の固定スケジュールのうちでも最も簡単な方法を採用する。

変数 λ の動的制御において、本論文では、次のような線形の減衰スケジュールを採用する。

$$\lambda(t) = \lambda_0 \times \left(1 - \frac{t}{t_{max}}\right) \quad (18)$$

λ_0 はパラメータ λ の初期値、 t は現在の学習サイクル、 t_{max} は学習の打ち切りサイクルである。学習の “サイクル” とは、全ての入出力パターンの提示と学習が一巡することをいう。上記では簡単化のため学習の終了時にエントロピー項の影響度が 0 になるようにしているが、本来は 0 まで落す必要はない。

エントロピー項を時間的に変化させるということは、

学習の対象となる目的関数を学習中に常に動かしてしまうことになる。この変更があまり急速であれば、学習の収束への悪影響も懸念される。しかし経験上は、上記の線形スケジュールでやや乱暴に落すことを複数回繰り返して良い結果を選ぶほうが、慎重な学習を数少なく行なうよりは好ましい場合が多かった。

3. 計算機実験

本章では2種類の簡単な例題に対する計算機実験を通じて、エントロピー項の導入による学習効率向上の効果を示す。例題の選択に当たっては、次のような事を考慮した。多層型ニューラルネットの学習の困難さは、学習すべき入出力パターンの数と、ネットワーク内のパラメータ数(結合荷重の数)のバランスにより変化する。一般に、多くの隠れ層ニューロンを用いると、学習は比較的安易になるが、訓練に用いたパターンに対して過学習気味になり、未知データに対する汎化の能力が弱くなる。逆に、少ない隠れ層ニューロンを用いると、学習は困難になるが、汎化の面では良い結果を得られることが多いと言われる。ここでは後者のような少ない隠れ層ニューロンを用いることとし、学習のローカルミニマの発生が知られている有名な例題を選択した。

3.1 例題 1: XOR 問題

XOR(排他的論理和)問題¹⁾にエントロピー付きBP学習則を適用する。3層のネットワークを用い、入力層・隠れ層・出力層のニューロン数はそれぞれ2, 2, 1とする。XORの入出力関係は、1枚の超平面では分離不可能なことが良く知られており、隠れ層を2個とした場合が、入出力関係を学習するための理論的な最小構成である。最小構成での学習では、学習空間にローカルミニマムが生ずることが知られている。隠れ層を3個以上とすると学習は容易に進むが、ここでは前述の目的で実験を行なうため、意図的に2個に制限する。

4種類の入出力パターン(00→0, 01→1, 10→1, 11→0)を決められた順番に提示し一巡することが、例題1における学習の1サイクルとなる。サイクル内において、各パターンを提示するごとに学習を行なう方法をとった。

学習の終了は、全てのパターンに対して出力ニューロンと教師信号の誤差が0.2未満になった時とした。学習は1000サイクルで打ち切った(学習率 η が低い時には1000サイクルでは足りない時があるが、あまり大きい値に設定すると実験の計算時間が膨大になるため1000サイクルに定めた)。

結合荷重および入力バイアスの初期値は-0.3と+0.3の間の一様乱数により設定した。乱数の初期値によって学習の様子は大きく変わり、運の良い初期値から始めたものだけ学習が終了し、そうでない場合はローカルミニマムに捕捉されて学習が終了しない。そこで100通りの乱数の初期値を用いて、どれだけの割合で学習が正しく収束するかを、既存のBP則とエントロピー付きBP則とで、比較する。

ここで、実験結果を正しく解釈するために、若干の工夫が必要であった。ある一組の学習率 η と慣性率 α を用いて実験を行なっただけでは、複数の手法間の公平な比較はできない。エントロピー項の導入により、学習にはやや「ブレーキ」がかかるため、最適な学習率や慣性率の値が異なってしまうためである。また、各手法ごとに最も都合のよい学習率と慣性率を1組だけ用いて比較する方法では、各手法が学習率や慣性率に対して敏感に左右されるか、ロバストかという点が示されないという欠点がある。

そこで以下では、学習率と慣性率を適切と思われる範囲で様々に変化させ、その全体としての振る舞いを、手法間で比較する。学習率については $\eta = 0.1$ から2.0まで0.1刻み、慣性率については $\alpha = 0.1$ から0.9まで0.1刻みとして、合計180通りのパラメータを用いて収束の性能を測定した。それぞれのパラメータの組み合わせに対して、前述のように乱数を変えて100試行ずつの実験を行なった。

1) エントロピー項を用いない場合

図2(a)はエントロピー項を用いない場合($\lambda = 0$)の、収束率を表している。2種類の横軸は、学習率 η (Learning rate)と、慣性率 α (Inertia)の値を表わし、縦軸(Convergence ratio)は、それらのパラメータを使った場合に、100試行中どれだけの割合で、1000サイクルの打ち切り以内に学習が正しく収束したかを示している。

今回測定した学習率・慣性率の範囲内では、収束率が100%となるようなパラメータの組は存在しない。結合荷重の変化の軌跡を調べたところ、およそ10%程度の試行がローカルミニマム状態から脱出できずにいる(例えば5000サイクルまで待っても出られない)ことがわかった。なお学習率と慣性率が共に低い領域で収束率が極端に落ち込んでいるのは、ローカルミニマムの問題ではなく、学習速度がきわめて遅いために、打ち切りサイクルに間に合わなかったものである。

2) エントロピー項を導入した場合

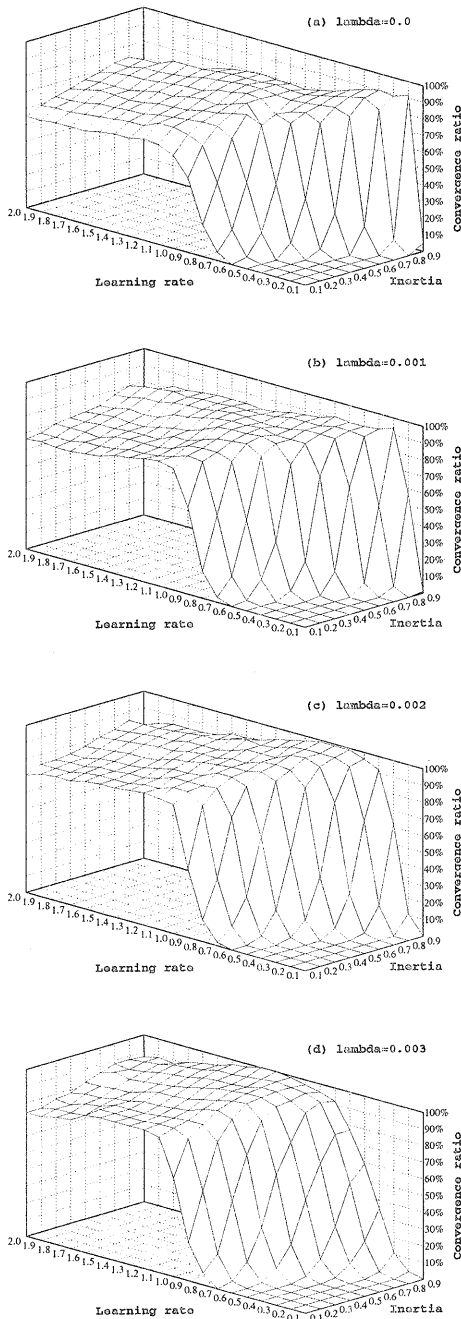


図2 エントロピー項による学習効率の変化 (例題1)
Fig. 2 Improvements of convergence ratio (Problem 1)

図2(b)~(d)に、 λ をそれぞれ0.001, 0.002, 0.003に設定した場合の収束率を示す。これらの実験結果から、以下の事項が観察される。

- エントロピーを加えていくと、収束率のピークが徐々に向上する。

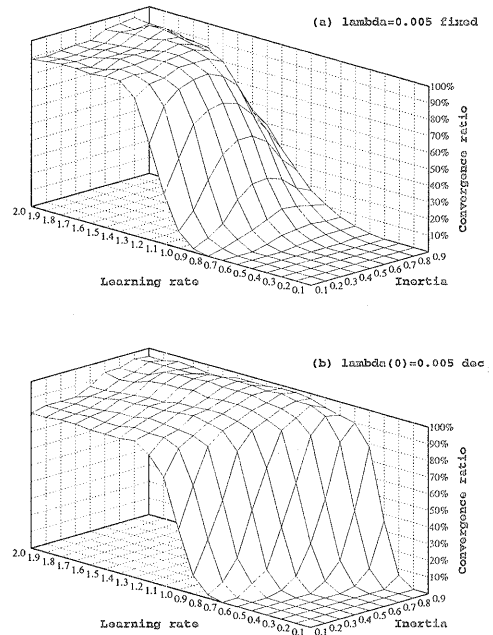


図3 エントロピー項の動的な調整の効果 (例題1)
Fig. 3 Effect of dynamical control of the entropy term (Problem 1)

- 最も収束が良いのは、学習率を低く、慣性率を高く設定した場合だった。
- エントロピーを加えていくと、学習を打ち切り時間内に終わらせるための限界線が移動する。エントロピー項の導入は学習を遅くする方向に働くため、学習率は通常の場合よりもわずかに高めに設定するのが好ましい。

3) エントロピー項の動的な調整

学習の打ち切りサイクルにおいて0となるような以下のスケジュールで、エントロピー項の重みを動的に減衰させてみる。

$$\lambda(t) = \lambda_0 \times \left(1 - \frac{t}{1000}\right) \quad (19)$$

結果例として、 $\lambda = 0.005$ で固定した場合を図3(a)に、 $\lambda_0 = 0.005$ から(19)式にしたがい動的に減衰させた場合を図3(b)に示す。

実験結果から、 λ を固定していた場合に比べてピーク性能を得るためのパラメータ範囲が広がることがわかる。 λ を固定する場合では、なかなか得られない高い収束率が、広いパラメータ範囲にわたって得られた。

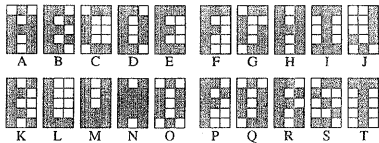


図4 例題2:パターン識別問題
Fig. 4 Problem 2: Pattern Recognition

3.2 例題2:パターン識別問題

XOR問題よりは、少しだけ複雑な例として、 3×5 のパターン^{12),15)}(図4)からなる20種の文字の弁別を行なう。入力層・隠れ層・出力層のニューロン数はそれぞれ15, 4, 20とする。全パターンに対して、各出力ニューロンの誤差が全て0.5未満になったとき^{12),15)}、学習が収束したと定めた。学習は1000サイクルで打ち切った。各パターン提示ごとに学習する方式をとった。結合荷重の初期値は-0.3と+0.3の間の一様乱数により設定した。学習率は0.1から1.0まで0.1刻み、慣性率は0.1から0.9まで0.1刻みとし、合計90通りのパラメータの組み合わせについて、それぞれ50試行ずつ実験を行なった。

図5(a)~(c)は、 λ をそれぞれ0.0, 0.005, 0.010とした場合の収束率を示す。エントロピー項を導入することにより、収束性能が大きく向上している。

このパターン識別問題においても、隠れ層のニューロン数を5以上とした場合には学習は容易になり、既存のBP学習則でも適切な学習率と慣性率を用いれば、1000サイクル以内にほぼ100%の収束をさせることができる。本論文では結果は割愛するが、8-3-8エンコーダ問題¹⁾なども収束が容易であることがわかった。これらの場合にはエントロピー項を導入しても性能の向上はわずかであった。深いローカルミニマムがあまり存在せず、既存のBPでも100%近い収束率が見込めるような素直な例題では、エントロピー項の導入効果を期待できないという点には注意が必要である。

汎化能力が問われず、どんなに過学習をしても良い状況では、本提案のような工夫はあまり意味がない。本論文で用いたXOR問題やパターン識別問題などの例題では、利用する隠れ層のニューロン数をわざと絞りこんだために、小規模ながら難しい学習の例となっていた。本論文では、簡単な例題だけを取り扱ったが、汎化能力を問われるような実問題での学習効率を評価していくことも今後は重要と考えられる。

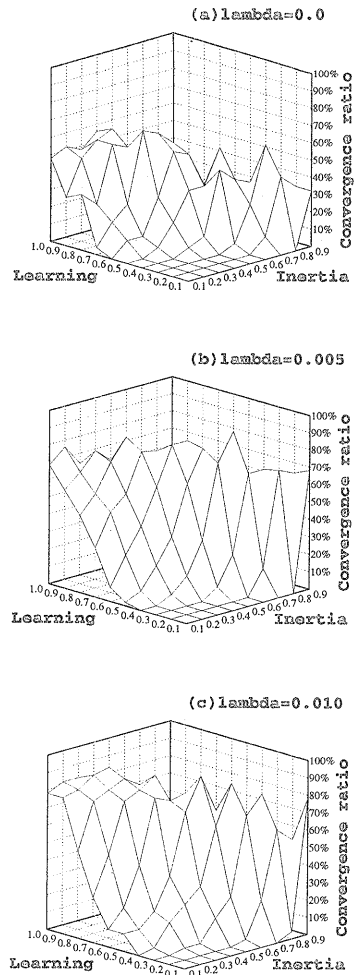


図5 エントロピー項による学習効率の変化(例題2)
Fig. 5 Improvements of convergence ratio (Problem 2)

4. Kullback 情報量を誤差項とした場合

Kullback 情報量^{*}は、2つの確率密度関数や確率分布の間の違いを表す尺度として1951年にS. KullbackとR. Leiblerにより提案された。この量を自乗誤差の代わりに用いて、バックプロパゲーション学習則を構成することが根岸¹⁶⁾および丹¹⁷⁾によって提案されている。

Kullback 情報量をBP学習則に用いる場合には、誤差項 E_p は(5)式に代わり以下のようになる¹⁶⁾。

* Kullback-Leibler 情報量とも呼ばれる

$$E_p(\mathbf{w}) \stackrel{\text{def}}{=} \sum_{i \in O} \left\{ t_{pi} \log_e \frac{t_{pi}}{v_{pi}} + (1 - t_{pi}) \log_e \frac{1 - t_{pi}}{1 - v_{pi}} \right\} \quad (20)$$

v_{pi} はパターン p に対するニューロン i の出力値, t_{pi} は教師信号である.

v_{pi} での偏微分は, 自乗誤差の場合よりも簡単な形となり,

$$\frac{\partial E_p(\mathbf{w})}{\partial v_{pi}} = -\frac{t_{pi} - v_{pi}}{v_{pi}(1 - v_{pi})} = -\frac{t_{pi} - v_{pi}}{f'(u_{pi})} \quad (21)$$

Kullback 情報量を誤差評価に用いた場合には, 2章で導いた学習則において, 出力層における誤差信号 δ_{pi} だけが若干変更される. (12), (13) 式の計算において, (21) 式を考慮すれば, 誤差信号として以下が導ける.

$$\delta_{pi} = \begin{cases} (t_{pi} - v_{pi}) - \lambda u_{pi} \cdot f'(u_{pi}) & i \in O \\ (\sum_k \delta_{pk} \cdot w_{ki} - \lambda u_{pi}) \cdot f'(u_{pi}) & i \in H \end{cases} \quad (22)$$

出力層の誤差信号 δ_{pi} は, 誤差 $(t_{pi} - v_{pi})$ に関して出力関数の微分係数 $f'(u_{pi})$ が乗じられず, 単純な誤り訂正学習の形になっている. 微分係数を乗じた場合には, v_{pi} が 0 や 1 に近づいた時に学習が減速されるのであるが, Kullback 情報量を用いた場合には, v_{pi} が 0 や 1 に近づいた後でも, 比較的速い速度で学習が進んでいくことがわかる.

学習率や慣性率の設定の仕方にもよるが, ごく一般的に言えば, Kullback 情報量を用いた場合には学習の収束が速くなる傾向がある. 学習空間にローカルミニマムがない場合には, Kullback 情報量を誤差項に用いる方法には収束速度の面で魅力がある.

ただし学習空間にローカルミニマムが存在する場合 (たとえば XOR 問題の例) では, あまりにも急速に学習が進むためか, 自乗誤差を用いた既存の BP 学習則に比べてローカルミニマムへ陥る確率が高い.

そこで学習則にエントロピー項を導入すると, 出力値が 0 や 1 に近づくことが防止されるため, Kullback 情報量を用いた場合にも特定の結合荷重が急速に増大する現象が抑えられ, 穏やかに学習が進行することが期待できる.

以下に, 例題 1 の XOR 問題を学習させた例を示す.

図 6(a) はエントロピー項を用いない場合 ($\lambda = 0$) の収束の様子を表している. 図 2(a) と比べて収束率が著しく下がっており, 自乗誤差を用いる既存の学習則と比べて性能が劣ることがわかる. 結合荷重の成長の様子を子細に調べてみると, ローカルミニマムに陥っていることがわかった.

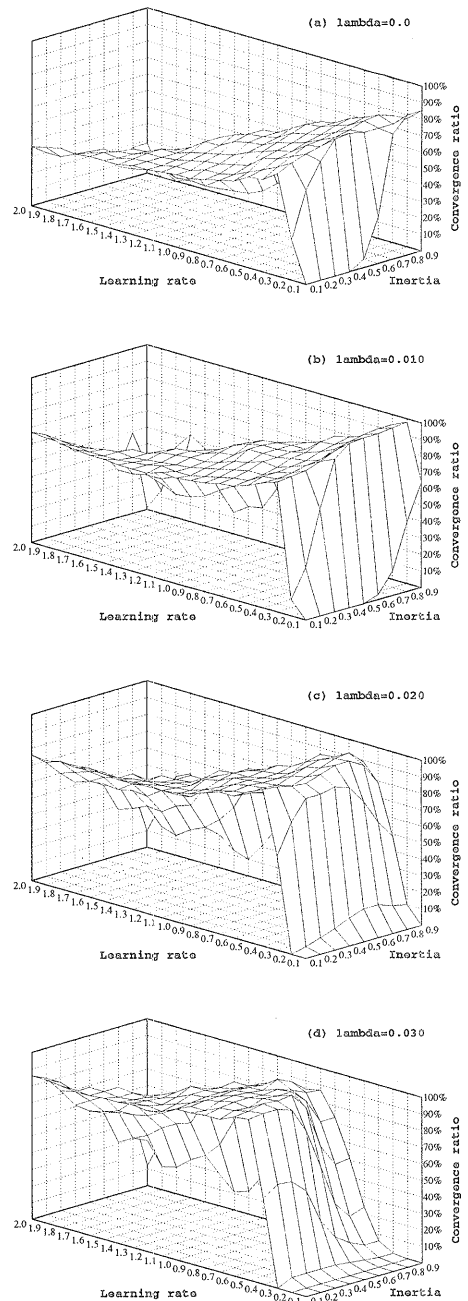


図 6 エントロピー項による学習効率の変化 (例題 1, Kullback)
Fig. 6 Improvements of convergence ratio (Problem 1, Kullback)

図 6(b)~(d) は, それぞれ λ を 0.010, 0.020, 0.030 に設定した場合の収束性能を示している.

これらの実験結果から, 以下の事項が観察される.

- エントロピー項を加えていくと, 収束率が全体的に著しく向上する. Kullback 情報量を用いた場

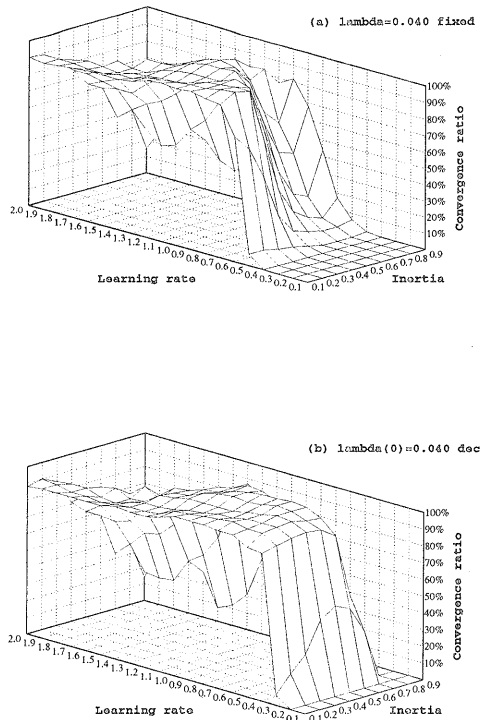


図7 エントロピー項の動的な調整の効果 (例題1, Kullback)
 Fig. 7 Effect of dynamical control of the entropy term
 (Problem 1, Kullback)

合の方が、エントロピー項の導入の効果が顕著である。

- 自乗誤差を用いた場合と異なり、学習率と慣性率を共に高くした時の性能が悪い。
- 学習率は、エントロピー項が大きくなってあまり高くする必要がない。慣性率を0付近に設定しても性能が良く、慣性率を利用しなくてもよい。

出力誤差の評価に Kullback 情報量を用いた場合にも、エントロピー項の動的な調整はやはり有効である。

打ち切り時に $\lambda(t) = 0$ となる前述のスケジュールを採用して $\lambda_0 = 0.040$ から減衰させた場合と、0.040 で固定した場合の比較を図7に示す。動的調整を行なった場合、学習率 η が低い範囲でも性能が高く保たれることがわかる。

今回の実験の結果から、慣性率を用いる必要がないこと、学習率に関する最適な設定範囲が広いこと、一般的に収束サイクルが短くなること等、Kullback 情報量を誤差評価に用いる利点を見いだすことができた。

Kullback 情報量を用いる場合、エントロピー項の導入により大幅に収束効率が改善される場合があることが示唆された。Kullback 情報量を誤差項に用いることは、自乗誤差よりも理論的には納得の行くものであるが、エントロピー項の導入と組にして検討をしないとその真価が測れないのではないかと考えられる。

5. むすび

本論文では、学習の目的関数に各ニューロンの出力値に関するエントロピー項を付け加えた新しいバックプロパゲーション学習則を提案し、その具体的な学習則を導いた。

エントロピー項を導入すると、各ニューロンは出力値として0や1に極端に近い値を出すことを避けるため、結果として、入力側の結合荷重が極端に増大することを抑えられる。この性質がローカルミニマムからの脱出を助け、収束率を高める効果をもつことを計算機実験を通じて示した。ただしエントロピー項を導入する際には、単に固定的に加えるのではなく、その影響度を動的に減衰させるなどの工夫が必要であり、その工夫によってエントロピー項の導入がより有効なものとなる。

さらに誤差の評価に Kullback の情報量を用いた別のエントロピー項付き学習則を導いた。Kullback 情報量を用いた学習則は、収束の高速化などの利点があるが、エントロピー項の導入を組み合わせることにより、性能を著しく向上できる場合があることを明らかにした。

エントロピー項の影響を時間とともに徐々に減少させるというアイデアは、学習空間の構造を動的に変更しながら学習を進めることを意味する。これはシミュレーテッドアニーリングにおいてランダムネスを徐々に弱めていく徐冷法ともアナロジーがあるが、より直接的には、連続値型の Hopfield ニューラルネット^{20),21)}におけるエネルギー最小化の研究と深く関連する。Hopfield 型ニューラルネットにおいて、出力関数の勾配を極限まで急峻にして、二値的になる場合には、二次形式で表現される系のエネルギーが最小化される。しかし、出力関数の勾配が有限であり中間的な出力値を取り得る場合には、この二次形式に「積分項」を加えた量が最小化の対象となる²⁰⁾。この積分項は実はエントロピー関数の形式に直すことができ²²⁾、勾配の値がエントロピー項の重みを制御する係数となる。Hopfield 型ネットの場合も、この係数を時間とともに減衰させていく手法²²⁾がきわめて有効である。ただし、Hopfield 型ネットワークが、出力値だ

けを変化させて安定状態を探索するのに対し、BP 則の場合には系内のニューロン間の結合荷重が学習されていくという点は大きく異なっている。

シミュレーテッドアニーリングにおける温度パラメータや、Hopfield 型ネットワークの最適化における勾配の係数などと同様に、バックプロパゲーション学習においてもエントロピー項の動的な制御により、解きやすい簡単な問題で求めた解（学習された結合荷重）を、徐々に本当に求めたい解に近づけていくという技法の存在が重要になっていくのではないかと考える。

今後はエントロピー項の係数 λ の最適な設定法、各層ごとに独自の λ を使う方法などについて研究を進めることが課題である。

(なお本論文は、1991年に国内研究会において発表したアイデア²³⁾を整理したものである。)

参 考 文 献

- 1) Rumelhart, D., Hinton, G. and Williams, R.: Learning Internal Representations by Error Propagation, in (ed. Rumelhart) Parallel Distributed Processing, Vol. I, MIT Press (1986).
- 2) Rumelhart, D., Hinton, G. and Williams, R.: Learning representations by back-propagating errors, *Nature*, Vol. 323, pp. 533-536 (1986).
- 3) Tollenaere, T.: SuperSAB: fast adaptive back-propagation with good scaling properties, *Neural Networks*, Vol. 3, pp. 561-574 (1990).
- 4) 猪飼, 山崎, 小迫: 逆伝播学習法における動的学習率の適応的決定, 信学技報 NC90-67 (1990).
- 5) 鎌田: ニューロンの非線形性を適応的に制御する誤差逆伝播学習法, 信学技報 NC90-31 (1990).
- 6) 岩田, 郷原, 内川: 学習曲面の大域的構造を考慮した“谷学習法”, 信学技報 NC90-65 (1990).
- 7) Watrous, R.: Learning algorithms for connectionist networks: applied gradient methods of nonlinear optimization, *Proc. ICNN-87*, Vol. II, pp. 619-628 (1987).
- 8) Owens, A. and Filkin, D.: Efficient training of the back propagation network by solving a system of stiff ordinary differential equations, *IJCNN-89*, Vol. II, pp. 381-386 (1987).
- 9) 石川真澄: コネクションリストモデルの忘却を用いた適応的構造学習, 信学技報 NC90-118 (1991).
- 10) Kung, S. and Hwang, J.: An algebraic projection analysis for optimal hidden units size and learning rates in back-propagation learning, *Proc. ICNN-88*, Vol. I, pp. 363-370 (1988).
- 11) 栗田多喜夫: 情報量基準による3層ニューラルネットワークの隠れ層ユニット数の決定法, 信学論 J73-D-II, pp. 1872-1878 (1990).
- 12) 萩原将文: 淘汰機能を有するバックプロパゲーション, 信学技報 NC89-104 (1990).
- 13) 松岡, 榎原: 素子数が増加するニューラルネットワークを用いた高速な学習, 信学技報 NC90-14 (1990).
- 14) 大熊, 玄地: ニューラルネットワークを用いた処理における汎化能力向上学習法, 信学技報 NC90-32 (1990).
- 15) 旭, 村上, 相原: ネットワーク構造の最適化と学習の高速化を目指す BP アルゴリズム, 信学技報 NC90-64 (1990).
- 16) 根岸, 高橋, 富田: Kullback 誤差関数による誤差伝播と性能評価, 信学技報 MBE88-175 (1989).
- 17) 丹, 加藤, 江島: 誤差評価関数による PDP モデルの高速化, 信学論 J73-D-II, pp. 2022-2028 (1990).
- 18) 甘利俊一: 学習識別の理論, 信学誌, Vol. 50, pp. 1272-1279 (1967).
- 19) 甘利俊一: 神経回路網の数理, 産業図書 (1978).
- 20) Hopfield, J.: Neurons with graded response have collective computational properties like those of two-state neurons, *Proc. Natl. Acad. Sci. USA*, Vol. 81, pp. 3088-3092 (1984).
- 21) Hopfield, J. and Tank, D.: ‘Neural’ computation of decisions in optimization problems, *Biol. Cybern.*, Vol. 52, pp. 141-152 (1985).
- 22) 秋山泰: ホップフィールド型ニューラルネットワークにおけるエネルギー最小状態への収束性を向上させる3つの技法, 信学技報 NC90-40 (1990).
- 23) 秋山泰, 古谷立美: 損失関数にエントロピー項を導入したバックプロパゲーション学習則, 信学技報 NC91-6 (1991).

(平成 11 年 1 月 21 日受付)

(平成 11 年 2 月 18 日採録)



秋山 泰 (正会員)

昭和 36 年生。平成 2 年慶應義塾大学大学院理工学研究科電気工学専攻博士課程修了。工学博士。同年通産省電子技術総合研究所研究官。平成 4 年京都大学化学研究所助教授。平成 8 年技術研究組合新情報処理開発機構並列応用つくば研究室室長。現在に至る。並列計算機を用いたタンパク質立体構造および遺伝子配列情報解析等の研究に従事。電子情報通信学会, 日本生物物理学会, 分子生物学会, 神経回路学会, IEEE 各会員。