

3W-6

TIP-4Pにおける
ニューロ処理方式岩下正雄¹藤田善弘¹石黒尚夫²山崎雅生²¹日本電気㈱²日本電気技術情報システム開発㈱

1. はじめに

ニューラルネットワークは、学習能力による適応性と汎用性から各分野でその応用が期待され、盛んに研究が進められている。しかし、学習にかかると、試行錯誤的な実験が多く、処理時間が研究のネックとなっている。ところが層構造のネットワークでは同一層内の各ユニットは独立であり、さらにユニット内の積和演算も並列性を有している。この並列性を十分に引き出すことにより処理時間短縮が可能となる。ただし、このネットワーク構造における逆伝播学習は各ニューロンユニットの計算にネットワーク全体のデータが必要であり、複数のプロセッサ間でのデータ交換が避けられない。

筆者らは、データフロープロセッサImPPを64個実装した大規模画像処理システムTIP-4Pを既に開発した。今回これを用いてニューラル・ネットワークのうち逆伝播学習アルゴリズムの処理方法の検討を行った。

2. TIP-4Pの構成

TIP-4Pは、データフロープロセッサImPP8個のリングを1単位

とした処理部を更に8個実装して高並列にデータを処理する。各処理部は専用のローカルメモリ(LM)上でデータを処理する。処理部間のデータの受渡しは、メインイメージメモリ(IM)を介して行う。LMとIM間のデータ転送は2Kワード単位の高速度ブロック転送で行う。

3. TIP-4Pによる並列処理

逆伝播学習アルゴリズムを図1に示す。前述したように、このアルゴリズムは1つのニューロンユニットの計算にネットワーク全体のデータが必要である。従って、TIP-4Pの複数処理部で分割処理する場合、処理部間でデータの集配交換が必要である。この集配交換すべきデータ量をできるだけ小さくする処理方法が課題である。

3-1. 方式1

一つは中間層および出力層のニューロンユニットを複数のグループに分割し、各処理部に担当させる方法がある。即ち図1の各式が、複数処理部で分割されて並列に実行され、式(1)から式(6)までが順番に処理される。しかしこの方法では式(2)の $Y(h)$ 、式(4)の

$$Y(h) = f\left(\sum_i W_h(h, i) \times X(i)\right) \quad (1)$$

$$Z(o) = f\left(\sum_h W_o(o, h) \times Y(h)\right) \quad (2)$$

$$D_o(o) = (T(o) - Z(o)) \times f'\left(\sum_h W_o(o, h) \times Y(h)\right) \quad (3)$$

$$D_h(h) = \sum_o W_o(o, h) \times D_o(o) \times f'\left(\sum_i W_h(h, i) \times X(i)\right) \quad (4)$$

$$W_o(o, h) = W_o(o, h) + P_o \times D_o(o) \times Y(h) \quad (5)$$

$$W_h(h, i) = W_h(h, i) + P_h \times D_h(h) \times X(i) \quad (6)$$

$X(i)$: 入力パタン $Y(h)$: 中間層の出力 $Z(o)$: 出力層の出力

$W_h(h, i)$: 中間層の重み $W_o(o, h)$: 出力層の重み $T(o)$: 教師信号

$i=0 \sim l-1$ $h=0 \sim H-1$ $o=0 \sim O-1$ $f()$: シグモイド関数

l : 入力ユニット数 H : 中間ユニット数 O : 出力ユニット数

図1 逆伝播学習のアルゴリズム

A Neural Nets Processing Method with TIP-4P

Masao Iwashita¹ Yoshihiro Fujita Takao Ishiguro² Masao Yamazaki²

¹NEC corporation

²NEC Scientific Information System Development

W_o(o, h) 及び D_o(o) は各処理部でネットワーク全体のデータが必要である。このため1回の逆伝播学習中に各処理部の間で3種類のデータテーブルの集配交換処理が必要となる。特に W_o(o, h) は中間層数と出力層数の増加に伴い2乗のオーダーでデータ量が増加し、オーバーヘッドが爆発する。

3-2. 方式2

次に中間層の分割は前述と同じだが、出力層の重み行列を転置して分割する方法を考える。即ち出力層の各ニューロンユニット内の重み(H個)を分割して複数処理部に担当させる。各処理部は、自分が担当する中間層の出力に対する重みだけを全出力層ユニットにわたって保持する。この方法によると、r個の各処理部では全出力層ユニットの積和の部分ベクトル Z_n(o) を出力する (n = 0 ~ r-1)。これをいったんIMにブロック転送したのち、一つの処理部がこれらを集めてr個のベクトルの総和

$$Z(o) = \sum_n Z_n(o)$$

を求め(3)式まで計算する。ここで求められた D_o(o) を全処理部にコピーすれば、あとは(6)式までデータの集配交換なしに各処理部が独立に計算を進められる。しかも Z(o) は、出力層ユニットの数に比例して増加するだけである。方法1よりも効率がよいことは明かである。

4. オーバーヘッド

ブロック転送は2kワード単位のため、転送時間は転送ワード数に対する階段状の関数(g())となる。2kワードのブロック転送時間を T_b、LMのアクセス時間を T_m とし、処理はLMの稼働率80%で行うとする。集配交換処理によるオーバーヘッドは

$$(3r+1)g(o) + o(r+2)T_m \quad (7)$$

となる。

また処理時間そのものは

$$((5i+7o+5)H/r + 3o)T_m \quad (8)$$

となる。

6. 台数と処理時間

オーバーヘッドを含んだ全体の処理

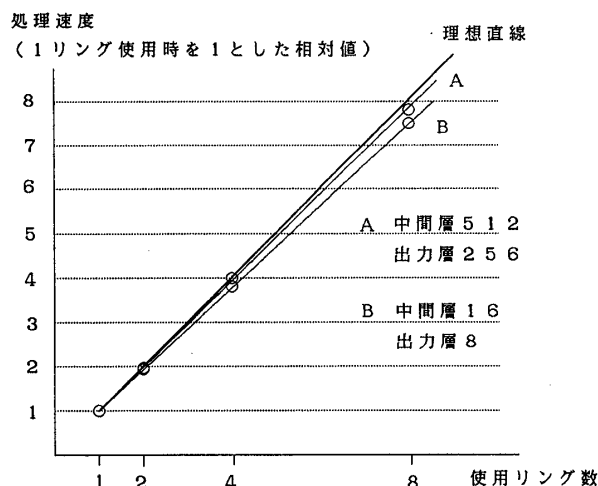


図2 使用リング数と処理速度

時間は式(7)+式(8)である。また T I P - 4 P の場合 $T_b \approx 25 T_m$ であることから、使用リング数と処理速度を図2に示した。グラフは理想直線と比較しても十分な台数効果を示しており、逆伝播学習を高速に処理できることがわかる。

7. おわりに

T I P - 4 P のニューラルネットワークの逆伝播学習への適応性について検討し効果的な処理方を述べた。本方式と高速ブロック転送機能により、逆伝播学習の処理における分散したローカルメモリ間のデータ転送によるオーバーヘッドを十分小さくすることができ、高い性能が得られることが確認された。

最後に、今回の検討にあたりご協力いただいた研究部の方々に深謝する。

[参考文献]

- [1] 岩下他、"データ駆動画像処理プロセッサ T I P - 4 P"、第37回情報全大、1988
- [2] 岩下他、"画像処理プロセッサ (I m P P) とニューラルネットワークへの応用"、理研シンポジウム、1988
- [3] Rumelhart, D. E., Hinton, G. E. and Williams, R. J., "Learning Representations by Back-Propagating Error," Nature, vol. 323, pp. 533-536, Oct. 9, 1986.