

人間からニューラルネットワークへの知識転移に関する基礎研究

鈴木 一誠†

ピトヨ・ハルトノ‡

中京大学 工学研究科†

中京大学 工学部‡

1. はじめに

現在、Deep Neural Network (DNN) の実世界問題への応用が盛んに行われる。DNN の学習には構造化されたデータが必要であるため、データ化の難しい人間の常識、主観や経験などをニューラルネットワーク (NN) に反映させることは難しい。そのため、NN から他の NN に知識の転移に関する研究[1]は多くあるが、人から NN に知識を転移する手段はまだ少ない。

そこで、本研究では 2 次元の位相的な中間層を持つ NN、Restricted Radial Basis Function Network (rRBF) [2] を用いて、人間の常識や経験を NN に転移させる手法を提案する。提案手法では、人間が常識、経験に基づいて NN の 2 次元の中間層を手動で組織化することで、NN を初期化することができる。人間の初期化によって、その人間の特性が後の NN の学習に反映されることを期待する。本研究では、人間の特性転移によって初期化された NN にエントロピーによる解析を行い、その結果を報告する。

2. 人間の特性転移が可能な NN

本研究で用いる rRBF は、2 次元の中間層を持つ階層型 NN であり、教師付き学習を実行することができる。rRBF の中間層の組織化は Self-Organizing Maps (SOM) [3] とは異なり、入力のラベル (context) を反映するためこの中間層を context-relevant self-organizing map (CRSOM) という。CRSOM は 2 次元であり、マップとして可視化することができるため、rRBF を次元圧縮アルゴリズムや可視化による入出力関係を直観的に理解できる NN として用いる [4, 5]。

本研究では、2 次元の CRSOM の位相的な性質を用いて人間から NN へ特性転移を可能とする手法の提案をする。初めに、人間が経験、好み、知識を基に手動で中間層の組織化をする。具体的には、人間は類似する入力同士を 2 次元の中間層上で近くに配置し、類似しない入力同士を中間層上で遠くに配置する。人間による組織化が行われた後に、rRBF の教師付き学習を実行する。

組織化の仕方が個人の常識、経験や主観によって異なるため、rRBF にそれを組織化する人間の特性を転移でき、学習後の NN にもその特性が反映されると考える。本研究で用いる rRBF の概要を図 1 に示す。ここでは、中間層と出力層の間に畳み込み層を置き、組織化の概要を示す。

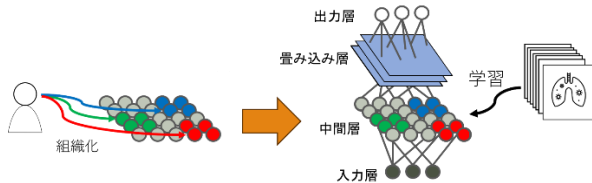


図 1 rRBF の概要

Research on the knowledge transfer from humans to neural network

†Issei Suzuki

Graduate School of Engineering Chukyo University

‡Pitoyo Hartono

School of Engineering Chukyo University

2.1 rRBF の学習

人間による組織化後、rRBF は以下のように学習する。

rRBF に対し、 $X(t) \in R^d$ を時間 t における入力とし、それに対する Best Matching Unit (BMU)、 win を式(1)に示す。

$$win = \operatorname{argmin} \|X(t) - W_{ij}\|^2 \quad (1)$$

式(1)の $W_{ij} \in R^d$ は中間層上の (i, j) に位置するニューロンに対応する参照ベクトルである。BMU の決定後、 (i, j) に位置する中間ニューロンの出力 O_{ij}^h を式(2)に示す。

$$O_{ij}^h(t) = e^{-I_{ij}^h(t)} \sigma(win, (i, j), t) \quad (2)$$

$$I_{ij}^h(t) = \|X(t) - W_{ij}\|^2$$

次に近傍関数 $\sigma(win, (i, j), t)$ を式(3)に示す。

$$\sigma(win, (i, j), t) = e^{-\frac{\operatorname{dist}(win, (i, j))}{s(t)}} \quad (3)$$

$$s(t) = s_0 \left(\frac{S_{\operatorname{end}}}{s_0} \right)^{\frac{t}{t_{\operatorname{end}}}}$$

式(3)の $\operatorname{dist}(win, (i, j))$ は、中間層上での BMU と (i, j) に位置するニューロンのデカルト距離を示す。 $s_0, s_{\operatorname{end}}$ は、それぞれ学習の始まりと終わりの近傍半径を示す定数である。また、 t_{end} は学習回数である。このような近傍関数を用いることで、学習回数により近傍半径は変化する。次の畳み込み層では、フィルタ $F(n, c, r, s)$ 、バイアス b とした場合、 (i, j) に位置するニューロンの出力 $O_{ij}^{\operatorname{conv}}(t)$ を以下に示す。

$$O_{ij}^{\operatorname{conv}}(t) = \sum_{c=1}^C \sum_{r=0}^{R-1} \sum_{s=0}^{S-1} O_{(i+s)(j+r)}^h(t) \cdot F(n, c, r, s) + b \quad (4)$$

式(4)の c は入力のチャンネル数、 n は出力のチャンネル数、 r はフィルタの高さ、 s はフィルタの幅を表す。 k 番目の出力ニューロン O_k を式(5)で計算する。

$$O_k(t) = f \left(\sum_j v_{(ij)k}(t) O_{ij}^{\operatorname{conv}}(t) - \theta_k(t) \right) \quad (5)$$

$$f(x) = \frac{1}{1 + e^{-x}} \quad (6)$$

式(5)の $v_{(ij)k}$ は (i, j) に位置する畳み込み層のニューロンと k 番目の出力ニューロンとの間の重み、 θ_k は k 番目のニューロンのバイアスである。 $f(x)$ は式(6)で示す Sigmoid 関数である。

ここで、Loss 関数を以下のように定義し、 T_k は教師信号の k 目の要素を示す。

$$Loss = \frac{1}{K} \sum_{k=1}^K (T_k(t) - O_k(t))^2 \quad (7)$$

学習では、重み、参照ベクトル、フィルタパラメータの更新を行う。実装では Python の自動微分機能を用いる。

3. 実験

実験データは、「0」～「9」の画像データ、MNIST[6]を用いる。MNISTは、手書によって作成されたため、個人の特性が反映するからである

実験では、「0」～「9」のをそれぞれ10枚、合計100枚を人間が個人の主観で中間層に配置することで組織化を行う。人間が行った初期化の例を図2(a)に示す。ここでは、ラベルに関わらず、類似する入力同士を2次元の中間層上で近くに配置し、類似しない入力同士を中間層上で遠くに配置する。図2(b)は学習後のCRSOMを示す。また、色は教師ラベルの違いを表す。

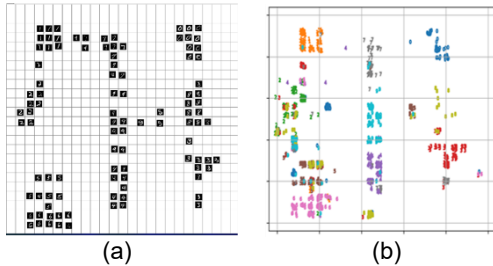


図2 人の初期化によるrRBF学習

図2(a)と図2(b)より学習後のCRSOMが人間による初期化を反映することがわかる。これは、人間の知識や主観がrRBFに反映することを意味する。

また、各ラベルに対するランダム初期化のエラーバーを図3(a)、人の初期化のエラーバーを図3(b)に示す。このグラフの横軸は手文字画像のground truthであり、縦軸は他の文字への誤認識率を示す。ここでは、 i 番目のrRBFが学習後に、文字 j の誤認識率を $p_i(j)$ 、文字 j を文字 k として誤認識率を $pf_i(j, k)$ とする。 $N_i(j, k)$ を文字 j を文字 k と誤認識した数とする。

$$p_i(j) = \frac{\sum_{k \neq j=0}^9 N_i(j, k)}{\sum_{j=0}^9 \sum_{k \neq j=0}^9 N_i(j, k)}, \quad pf_i(j, k) = \frac{N_i(j, k)}{\sum_{k \neq j=0}^9 N_i(j, k)} \quad (8)$$

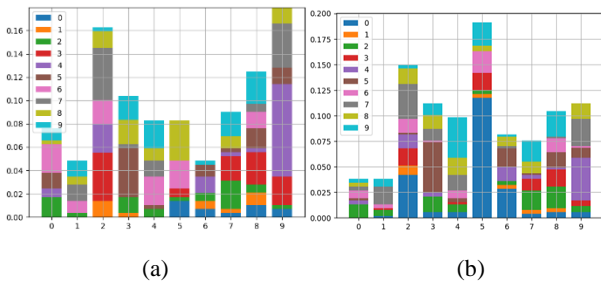


図3 初期化に対するrRBFの誤認識率

人間の知識転移の実験では、5人の被験者がそれぞれ異なるrRBFを初期化した後に、それらのrRBFの教師付学習を実装した。比較として、ランダムに初期化した5つのrRBFの学習も行った。ランダムに初期化したrRBFをrandom1~5、5人の被験者をperson1~5とし、それらのrRBFの差異を評価する。 h 番目のrRBFと i 番目のrRBFの差異の評価方法は、エラー全体の分布に注目する $coarse(h, i)$ (式9)とエラー詳細な内容に注目する $detail(h, i)$ (式10)である。

$$coarse(h, i) = \frac{1}{2} \sum_{j=0}^9 \{-pf_i(j) \log pf_h(j) - pf_h(j) \log pf_i(j)\} \quad (9)$$

$$detail(h, i) = \frac{1}{20} \sum_{j=0}^9 \sum_{k \neq j=0}^9 \{-p_i(j, k) \log p_h(j, k) - p_h(j, k) \log p_i(j, k)\} \quad (10)$$

式(9)、(10)に基づいて、異なるrRBFの差異を計算した。randomの平均とpersonの平均を表1に示す。

表1 rRBFの特性の差異

	randomの平均	personの平均
coarse	4.236	3.143
detail	4.250	3.484
coarse-detail	0.014	0.341

coarseでは、personの平均値がrandomの平均値よりも小さい。これは、数字の類似性に関する人の常識には共通点が多く、常識転移を行わない複数のニューラルネットワークの方に大きな多様性が生じる。一方、(coarse-detail)の値においてpersonが大きい。これは、人間の常識によって初期化された各ニューラルネットワークが異なるエラーの発生の仕方をしたことを意味する。つまり、人間が初期化したrRBFがその人間の特性を反映すると考えられる。

4. まとめ

本研究では人間の知識をNNに事前知識として転移できる手法を提案した。ここでは、手書き文字を用いて基礎実験を行った。実験の結果からは人間がNNに転移した常識が、その後のNNの特性に影響を与えることが分かった。このことは、データからだけでなく、人間から学ぶNNを学習させることができ、人間とAIの新しい関係を確立できると考える。

今後は、データ化することが困難な問題にこの学習方法を適応することを検討する。例えば、職人的な技術や経験や主観的を必要とするスキルのNNの学習を試みる。

謝辞

本研究はローム株式会社との共同研究による成果である。

参考文献

- [1] 神島敏弘, “転移学習” 人工知能学会誌, Vol.25, No.4 (2010).
- [2] P. Hartono, et al, “Learning-Regulated Context Relevant Topographical Map,” in IEEE Transactions on Neural Networks and Learning Systems, Vol.26, No.10 (2015).
- [3] T. Kohonen, “Self-organized formation of topologically correct feature maps.” *Biological cybernetics* 43.1 (1982).
- [4] P. Hartono, “Classification and dimensional reduction using restricted radial basis function networks”, *Neural Computing & Applications*, Vol.30, No.3 (2018).
- [5] P. Hartono, “Mixing autoencoder with classifier: conceptual data visualization”, *IEEE Access*, Vol.8 (2020).
- [6] Y. Lecun, L. Bottou, Y. Bengio and P. Haffner, “Gradient-based learning applied to document recognition,” in *Proceedings of the IEEE*, Vol.86, No.11 (1998).