

# 行列分解を用いた深層残差 $H_\infty$ 学習の計算量の削減

## Reduction of a Computational Complexity of the Deep Residual $H_\infty$ -Learning Using Matrix Factorization

西山 清

Kiyoshi NISHIYAMA

岩手大学理工学部システム創成工学科

Faculty of Science and Engineering, Iwate University

### 1 はじめに

多層ニューラルネットワーク [1] において層を深くして行くと、ある深さから学習精度が著しく劣化することが知られている。この対策としてスキップ接続 (ショートカット接続) が導入された残差ネットワーク (Residual Network) が考案された [2]。残差ネットワークでは層を深くすればするほど、学習時の精度が向上することが報告されている。我々の先行研究では、モデル集合の包含関係と誤差曲面の高次鞍点化の観点 [3] から深層残差ネットワークの  $H_\infty$  学習が優れた学習特性をもつことを明らかにした。しかし、 $H_\infty$  学習を大規模なネットワークに適用するためには計算量の問題を解決する必要があった。

本研究では、 $H_\infty$  学習の計算量を削減するため、まず  $N_L$  出力のネットワークの出力層から入力層に向けて各ニューロンに対する結合重みとしきい値 (部分重みベクトル) を順次更新する。次に、各ニューロンの部分観測行列が  $N_L \times 1$  行列と  $1 \times N_n$  行列の積に分解できることに着目し、部分重みベクトルの更新に必要な  $N_L \times N_L$  行列の逆行列の計算を回避する。最後に、これらの削減法を導入した二つの深層残差  $H_\infty$  学習アルゴリズムを提案する。

### 2 $H_\infty$ 学習

$H_\infty$  学習問題とは、ある値  $\gamma_f$  が与えられたとき、

$$\sup_{w_0, \{v_p\}} \frac{\sum_{p=0}^k \|e_{f,p}\|_{(\sigma_v^2 \mathbf{I})^{-1}}^2}{\|w - \hat{w}_0\|_{(\epsilon_0 \mathbf{I})^{-1}}^2 + \sum_{p=0}^k \|v_p\|_{(\sigma_v^2 \mathbf{I})^{-1}}^2} < \gamma_f^2 \quad (1)$$

を満たす  $H_\infty$  準最適な学習アルゴリズム  $\mathcal{F}_f$  を求める問題である。ここで、 $e_{f,p}$  は出力誤差、 $w$  は重みベクトル、 $v_p$  は線形化誤差を含む観測雑音である。

この  $H_\infty$  学習問題の解は、ニューラルネットワークを線形化した状態空間モデルに  $H_\infty$  フィルタ (EHF) [4] を適用して得られる。この学習アルゴリズムは 1 次微分のみを用いて重み空間全体を探索して  $H_\infty$  準最適な解を求めることから、 $g$ -EHF 学習アルゴリズムと呼ばれる [5]。次に、出力層のニューロン数が一つの場合を示す。

$$\hat{w}_k = \hat{w}_{k-1} + \mathbf{K}_{s,k}(y_k - h_k(\hat{w}_{k-1})) \quad (2)$$

$$\begin{aligned} \mathbf{K}_{s,k} &= \hat{\mathbf{P}}_{k|k-1} \mathbf{H}_k^T (\mathbf{H}_k \hat{\mathbf{P}}_{k|k-1} \mathbf{H}_k^T + 1)^{-1} \\ \hat{\mathbf{P}}_{k+1|k} &= \hat{\mathbf{P}}_{k|k-1} - \hat{\mathbf{P}}_{k|k-1} \\ &\quad \times \begin{bmatrix} \mathbf{H}_k^T & \mathbf{H}_k^T \end{bmatrix} \mathbf{R}_{e,k}^{-1} \begin{bmatrix} \mathbf{H}_k \\ \mathbf{H}_k \end{bmatrix} \hat{\mathbf{P}}_{k|k-1} \end{aligned} \quad (3)$$

ただし、

$$\begin{aligned} \mathbf{H}_k &= \left. \frac{\partial h_k(w)}{\partial w} \right|_{w=\hat{w}_{k-1}}, \hat{\mathbf{P}}_{0|-1} = \frac{\epsilon_0}{\sigma_v^2} \mathbf{I}, \epsilon_0 > 0 \quad (4) \\ \mathbf{R}_{e,k} &= \mathbf{R} + \begin{bmatrix} \mathbf{H}_k \\ \mathbf{H}_k \end{bmatrix} \hat{\mathbf{P}}_{k|k-1} \begin{bmatrix} \mathbf{H}_k^T & \mathbf{H}_k^T \end{bmatrix} \\ \mathbf{R} &= \begin{bmatrix} 1 & 0 \\ 0 & -\gamma_f^2 \end{bmatrix}, \gamma_f > 1 \end{aligned} \quad (5)$$

### 3 残差ネットワーク (ResNet)

ResNet は、ある 2 つの層間の出力をスキップ接続 (skip connection) で結合した構造を含んだニューラルネットワーク (NN) である [2]。スキップ接続とは、ある NN における  $n$  層の出力  $x \in \mathcal{R}^M$  と  $n+m$  層の出力  $y(x) \in \mathcal{R}^M$  を加算することである。その和を  $z(x) = y(x) + x$  とする。なお、 $\mathcal{R}$  は実数全体の集合、 $n, m, M$  は自然数である。このスキップ接続により、 $z(x)$  を学習する問題は残差  $y(x) = z(x) - x$  を学習する問題に帰着できる。このことから、このスキップ接続を含む NN は残差ネットワーク (ResNet) と呼ばれる。

本研究では、図 1 のような  $L$  層を持ち、各隠れ層のニューロン数が同じである、ResNet について考える (多出力の場合も扱う)。この ResNet の層数  $L$  はスキップ接続の数  $l$  で決定される ( $L = 2l + 3, l = 0, 1, \dots$ )。図中の四角のニューロンは、応答関数が恒等関数であり、しきい値を持たないことを表す。一方、丸のニューロンは応答関数がシグモイド関数  $f(x) = 1/(1 + \exp(-\eta_0 x))$  であり、しきい値をもつことを表す ( $\eta_0 > 0$  はシグモイド関数の傾き)。図中の弧線はスキップ接続を表す。

この  $L$  層 ResNet に  $p$  番目の入力  $z_p^1 = [z_{1,p}^1, \dots, z_{N_1,p}^1]^T \in \mathcal{R}^{N_1 \times 1}$  が与えられたとき、 $n$  層の出力  $z_p^n = [z_{1,p}^n, \dots, z_{N_n,p}^n]^T \in \mathcal{R}^{N_n \times 1}$  を以下のよう

$$\begin{aligned} z_p^n &= f(s_p^{n-1}), s_p^{n-1} = \mathbf{W}^{n-1} z_p^{n-1} + \mathbf{b}^{n-1} \\ &\quad (n = 2 \text{ または } n = 3, 5, \dots, L) \end{aligned} \quad (6)$$

$$z_p^n = s_p^{n-1}, \quad s_p^{n-1} = \mathbf{W}^{n-1} z_p^{n-1} + z_p^{n-2} \quad (n = 4, 6, \dots, L-1) \quad (7)$$

ここで、 $N_n$  は  $n$  層のニューロン数、 $s_p^{n-1} = [s_{1,p}^{n-1}, \dots, s_{N_n,p}^{n-1}]^T \in \mathcal{R}^{N_n \times 1}$  は  $n$  層の膜電位、

$$\mathbf{f}(s_p^{n-1}) = [f(s_{1,p}^{n-1}), \dots, f(s_{N_n,p}^{n-1})]^T \quad (8)$$

は  $n$  層の膜電位  $s_p^{n-1}$  の各成分に対するニューロンの応答から成るベクトル値関数である。また、

$$\mathbf{W}^n = \begin{bmatrix} w_{1,1}^n & \cdots & w_{1,N_n}^n \\ \vdots & \ddots & \vdots \\ w_{N_{n+1},1}^n & \cdots & w_{N_{n+1},N_n}^n \end{bmatrix} \in \mathcal{R}^{N_{n+1} \times N_n} \quad (9)$$

は  $n$  層から  $n+1$  層への重み行列、 $\mathbf{b}^n = [w_{1,0}^n, \dots, w_{N_{n+1},0}^n]^T$  は  $n+1$  層のしきい値ベクトルである。

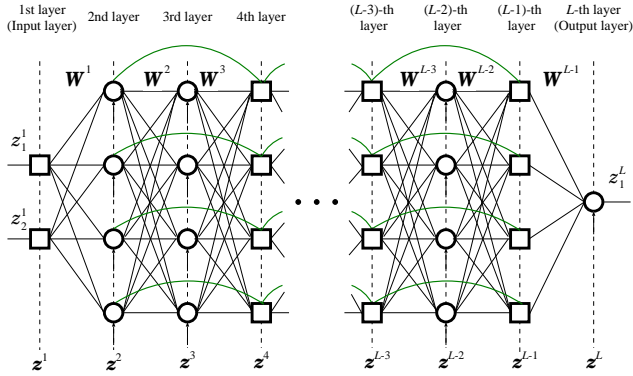


図 1  $L$  層 ResNet ;  $L$  は層数であり ( $L = 2l + 3$ ,  $l = 0, 1, \dots$ )、隠れ層のニューロン数  $N_n$ ,  $n = 2, \dots, L-1$  は等しい。 $z^n$  は  $n$  層の出力、 $\mathbf{W}^n$  は  $n$  層から  $n+1$  層への重み行列である。四角のニューロンは、応答関数が恒等写像であり、しきい値を持たないことを表す。一方、丸のニューロンは応答関数がシグモイド関数であり、しきい値を持つことを表す。 $n$  を 0 ではない偶数とすると、 $n$  層と  $n+2$  層はスキップ接続されている。

#### 4 スキップ接続を考慮した $H_\infty$ 学習

$L$  層 ResNet における  $H_\infty$  学習と文献 [5] で述べた 3 層 NN の  $H_\infty$  学習の違いは、NN の線形状態空間モデル

$$\mathbf{w}_k = \mathbf{w}_{k-1}, \quad \mathbf{m}_k = \mathbf{H}_k \mathbf{w}_k + \mathbf{v}_k \quad (10)$$

における観測行列

$$\mathbf{H}_k = \left. \frac{\partial \mathbf{h}_k(\mathbf{w})}{\partial \mathbf{w}} \right|_{\mathbf{w}=\hat{\mathbf{w}}_{k-1}} \in \mathcal{R}^{N_L \times N_w} \quad (11)$$

の計算方法だけである。ここで、

$$\mathbf{m}_k = \mathbf{y}_k - \mathbf{h}_k(\hat{\mathbf{w}}_{k-1}) + \mathbf{H}_k \hat{\mathbf{w}}_{k-1} \quad (12)$$

$$\mathbf{w} = [w_{1,0}^1, w_{1,1}^1, \dots, w_{1,N_1}^1, \dots, w_{N_L,0}^{L-1}, w_{N_L,1}^{L-1}, \dots, w_{N_L,N_{L-1}}^{L-1}]^T \quad (13)$$

であり、 $\mathbf{w} = [\text{vec}([\mathbf{b}^1, \mathbf{W}^1]^T)^T, \dots]^T$  はすべてのしきい値と結合重みからなる重みベクトル、 $\mathbf{h}_k(\mathbf{w}) = [h_{1,k}(\mathbf{w}), \dots, h_{N_L,k}(\mathbf{w})]^T$  は更新ステップ  $k$  の  $L$  層 ResNet の出力  $z_k^L$  を表す  $\mathbf{w}$  の関数 ( $\mathbf{y}_k = z_k^L$ )、 $\hat{\mathbf{w}}_k$  は更新ステップ  $k$  における重みベクトル  $\mathbf{w}$  の推定値である。

観測行列  $\mathbf{H}_k$  中の  $\frac{\partial \mathbf{h}_k}{\partial w_{j,i}^n}$  は次式により計算される [6]。

$$\frac{\partial \mathbf{h}_k}{\partial w_{j,i}^1} = \Phi_k^2 \cdot \frac{\partial s_k^1}{\partial w_{j,i}^1} \quad (14)$$

$$\frac{\partial \mathbf{h}_k}{\partial w_{j,i}^2} = \Phi_k^3 \cdot \frac{\partial s_k^2}{\partial w_{j,i}^2}, \quad \dots, \quad \frac{\partial \mathbf{h}_k}{\partial w_{j,i}^n} = \Phi_k^{n+1} \cdot \frac{\partial s_k^n}{\partial w_{j,i}^n} \quad (15)$$

ここで、行列  $\Phi_k^n$  は次の逆方向の再帰式で計算できる。

$$\Phi_k^n = \Phi_k^{n+1} \frac{\partial s_k^n}{\partial z_k^n} \frac{\partial z_k^n}{\partial s_k^{n-1}} \in \mathcal{R}^{N_L \times N_n} \quad (16)$$

$$= \Phi_k^{n+1} \mathbf{W}^n \begin{bmatrix} \frac{\partial z_{1,k}^n}{\partial s_{1,k}^{n-1}} & & & \mathbf{O} \\ & \ddots & & \\ & & \ddots & \\ \mathbf{O} & & & \frac{\partial z_{N_n,k}^n}{\partial s_{N_n,k}^{n-1}} \end{bmatrix} \quad (17)$$

( $n = L-1$  または  $n = L-2, \dots, 5, 3$ )

$$\Phi_k^n = (\Phi_k^{n+1} \frac{\partial s_k^n}{\partial z_k^n} + \Phi_k^{n+2}) \frac{\partial z_k^n}{\partial s_k^{n-1}} \quad (18)$$

$$= (\Phi_k^{n+1} \mathbf{W}^n + \Phi_k^{n+2}) \begin{bmatrix} \frac{\partial z_{1,k}^n}{\partial s_{1,k}^{n-1}} & & & \mathbf{O} \\ & \ddots & & \\ & & \ddots & \\ \mathbf{O} & & & \frac{\partial z_{N_n,k}^n}{\partial s_{N_n,k}^{n-1}} \end{bmatrix} \quad (19)$$

$$(n = L-3, L-5, \dots, 4, 2)$$

ただし、

$$\Phi_k^L = \frac{\partial z_k^L}{\partial s_k^{L-1}} \in \mathcal{R}^{N_L \times N_L} \quad (20)$$

$$\frac{\partial z_k^n}{\partial s_k^{n-1}} = \begin{bmatrix} \frac{\partial z_{1,k}^n}{\partial s_{1,k}^{n-1}} & & & \mathbf{O} \\ & \ddots & & \\ & & \ddots & \\ \mathbf{O} & & & \frac{\partial z_{N_n,k}^n}{\partial s_{N_n,k}^{n-1}} \end{bmatrix} \quad (21)$$

ここで、 $\frac{\partial z_k^n}{\partial s_k^{n-1}}$  は対角行列であり、 $n = L-1, L-3, \dots, 6, 4$  のとき単位行列となる。また、

$$\frac{\partial s_k^n}{\partial z_k^n} = \mathbf{W}^n \in \mathcal{R}^{N_{n+1} \times N_n} \quad (22)$$

である。

5 深層残差  $H_\infty$  学習の計算量の削減

先行研究 [6] では、1 出力の場合に限定した  $H_\infty$  学習に関して述べた。本章では、多出力 NN に対する深層残差  $H_\infty$  学習の計算量の削減方法を導出する [5],[7]。

まず、第  $n+1$  層の  $i$  ニューロンに着目した場合の  $H_\infty$  学習のリカッチ方程式を考える。

$$\begin{aligned} \hat{\Sigma}_{i,k+1|k}^n &= \hat{\Sigma}_{i,k|k-1}^n \\ &- \hat{\Sigma}_{i,k|k-1}^n \begin{bmatrix} \mathbf{H}_{i,k}^n \\ \mathbf{H}_{i,k}^n \end{bmatrix}^T \mathbf{R}_{e,k}^{n,i-1} \begin{bmatrix} \mathbf{H}_{i,k}^n \\ \mathbf{H}_{i,k}^n \end{bmatrix} \hat{\Sigma}_{i,k|k-1}^n \end{aligned} \quad (23)$$

ただし、

$$\mathbf{R}_{e,k}^{n,i} = \mathbf{R}_k^n + \begin{bmatrix} \mathbf{H}_{i,k}^n \\ \mathbf{H}_{i,k}^n \end{bmatrix} \hat{\Sigma}_{i,k|k-1}^n \begin{bmatrix} \mathbf{H}_{i,k}^n \\ \mathbf{H}_{i,k}^n \end{bmatrix}^T \quad (24)$$

$$\mathbf{R}_k^n = \begin{bmatrix} \hat{\sigma}_{v_k}^2 \mathbf{I} & 0 \\ 0 & -\gamma_f^2 \mathbf{I} \end{bmatrix} \quad (25)$$

$$\mathbf{H}_k = [\mathbf{H}_{1,k}^1, \dots, \mathbf{H}_{N_n,k}^{n+1}, \dots, \mathbf{H}_{N_L,k}^{L-1}] \quad (26)$$

ここで、 $\mathbf{R}_{e,k}^{n,i-1}$  を  $2 \times 2$  ブロック行列

$$\mathbf{R}_{e,k}^{n,i-1} = \begin{bmatrix} \mathbf{L}_{i,k}^n & \mathbf{M}_{i,k}^n \\ \mathbf{N}_{i,k}^n & \mathbf{O}_{i,k}^n \end{bmatrix} \quad (27)$$

で表し、

$$\mathbf{S}_{i,k}^n = \mathbf{L}_{i,k}^n + \mathbf{M}_{i,k}^n + \mathbf{N}_{i,k}^n + \mathbf{O}_{i,k}^n \quad (28)$$

を定義すれば、式 (23) を次のように書き換えることができる。

$$\hat{\Sigma}_{i,k+1|k}^n = \hat{\Sigma}_{i,k|k-1}^n - \hat{\Sigma}_{i,k|k-1}^n \mathbf{H}_{i,k}^n \mathbf{S}_{i,k}^n \mathbf{H}_{i,k}^n \hat{\Sigma}_{i,k|k-1}^n \quad (29)$$

一方、 $\mathbf{R}_{e,k}^{n,i}$  の逆行列は次のように分解・整理できる。

$$\begin{aligned} \mathbf{R}_{e,k}^{n,i-1} &= \begin{bmatrix} \mathbf{I} & \mathbf{A}_{i,k}^n \\ 0 & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{B}_{i,k}^n & 0 \\ 0 & \mathbf{C}_{i,k}^n \end{bmatrix} \begin{bmatrix} \mathbf{I} & 0 \\ \mathbf{A}_{i,k}^n & \mathbf{I} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{B}_{i,k}^n + \mathbf{A}_{i,k}^n \mathbf{C}_{i,k}^n \mathbf{A}_{i,k}^n & \mathbf{A}_{i,k}^n \mathbf{C}_{i,k}^n \\ \mathbf{C}_{i,k}^n \mathbf{A}_{i,k}^n & \mathbf{C}_{i,k}^n \end{bmatrix} \end{aligned} \quad (30)$$

また、 $n+1$  層のニューロン  $i$  に対する部分観測行列は  $N_L \times 1$  行列  $\mathbf{u}_{i,k}^n$  を用いて次のように分解できる。

$$\mathbf{H}_{i,k}^n = \frac{\partial z_k^L}{\partial s_{i,k}^n} \frac{\partial s_{i,k}^n}{\partial \mathbf{w}_{i,k}^n} = \mathbf{u}_{i,k}^n \hat{\mathbf{z}}_k^n \quad (31)$$

さらに、ブロック行列  $\mathbf{R}_{e,k}^{n,i}$  の逆行列の公式と行列分解  $\mathbf{H}_{i,k}^n = \mathbf{u}_{i,k}^n \hat{\mathbf{z}}_k^n$  を用いれば、行列  $\mathbf{A}_{i,k}^n$ 、 $\mathbf{B}_{i,k}^n$  および  $\mathbf{C}_{i,k}^n$  は次のように整理できる。

$$\mathbf{A}_{i,k}^n = -\frac{\alpha_{i,k}^n}{\hat{\sigma}_{v_k}^2 + \alpha_{i,k}^n \beta_{i,k}^n} \mathbf{u}_{i,k}^n \mathbf{u}_{i,k}^n \quad (32)$$

$$\begin{aligned} \mathbf{B}_{i,k}^n &= (\hat{\sigma}_{v_k}^2 \mathbf{I} + \mathbf{H}_{i,k}^n \hat{\Sigma}_{i,k|k-1}^n \mathbf{H}_{i,k}^n)^{-1} \\ &= (\hat{\sigma}_{v_k}^2 \mathbf{I} + \mathbf{u}_{i,k}^n \alpha_{i,k}^n \mathbf{u}_{i,k}^n)^{-1} \end{aligned} \quad (33)$$

$$= \frac{1}{\hat{\sigma}_{v_k}^2} \left[ \mathbf{I} - \frac{\alpha_{i,k}^n}{\hat{\sigma}_{v_k}^2 + \alpha_{i,k}^n \beta_{i,k}^n} \mathbf{u}_{i,k}^n \mathbf{u}_{i,k}^n \right] \quad (34)$$

$$= \frac{1}{\hat{\sigma}_{v_k}^2} [\mathbf{I} + \mathbf{A}_{i,k}^n] \quad (35)$$

$$\begin{aligned} \mathbf{C}_{i,k}^n &= \frac{1}{\gamma_f^2} \left[ \frac{\hat{\sigma}_{v_k}^2 \alpha_{i,k}^n}{\hat{\sigma}_{v_k}^2 \alpha_{i,k}^n \beta_{i,k}^n - \gamma_f^2 (\hat{\sigma}_{v_k}^2 + \alpha_{i,k}^n \beta_{i,k}^n)} \right. \\ &\quad \left. \mathbf{u}_{i,k}^n \mathbf{u}_{i,k}^n - \mathbf{I} \right] \end{aligned} \quad (36)$$

ただし、

$$\alpha_{i,k}^n = \hat{\mathbf{z}}_k^n \hat{\Sigma}_{i,k|k-1}^n \hat{\mathbf{z}}_k^n, \quad \beta_{i,k}^n = \mathbf{u}_{i,k}^n \mathbf{u}_{i,k}^n \quad (37)$$

ここで、式 (27) と式 (30) を比較すれば式 (28) は次のように整理できる ( $\hat{\sigma}_{v_k}^2 + \alpha_{i,k}^n \beta_{i,k}^n$  による約分に注意)。

$$\begin{aligned} \mathbf{S}_{i,k}^n &= (\mathbf{B}_{i,k}^n + \mathbf{A}_{i,k}^n \mathbf{C}_{i,k}^n \mathbf{A}_{i,k}^n) \\ &\quad + \mathbf{A}_{i,k}^n \mathbf{C}_{i,k}^n + \mathbf{C}_{i,k}^n \mathbf{A}_{i,k}^n + \mathbf{C}_{i,k}^n \\ &= (\mathbf{A}_{i,k}^n + \mathbf{I}) \left\{ \frac{1}{\hat{\sigma}_{v_k}^2} \mathbf{I} + \mathbf{C}_{i,k}^n (\mathbf{A}_{i,k}^n + \mathbf{I}) \right\} \\ &= \frac{(1 - \gamma_f^{-2} \hat{\sigma}_{v_k}^2)}{\hat{\sigma}_{v_k}^2} \\ &\quad \times \left[ \mathbf{I} - \frac{(1 - \gamma_f^{-2} \hat{\sigma}_{v_k}^2) \alpha_{i,k}^n}{\hat{\sigma}_{v_k}^2 + (1 - \gamma_f^{-2} \hat{\sigma}_{v_k}^2) \alpha_{i,k}^n \beta_{i,k}^n} \mathbf{u}_{i,k}^n \mathbf{u}_{i,k}^n \right] \end{aligned} \quad (38)$$

これより、次式を得る。

$$\mathbf{u}_{i,k}^n \mathbf{S}_{i,k}^n \mathbf{u}_{i,k}^n = \frac{(1 - \gamma_f^{-2} \hat{\sigma}_{v_k}^2) \beta_{i,k}^n}{\hat{\sigma}_{v_k}^2 + (1 - \gamma_f^{-2} \hat{\sigma}_{v_k}^2) \alpha_{i,k}^n \beta_{i,k}^n} \quad (39)$$

これを次のように整理された式 (29) の右辺第 2 項

$$\begin{aligned} &\hat{\Sigma}_{i,k|k-1}^n \mathbf{H}_{i,k}^n \mathbf{S}_{i,k}^n \mathbf{H}_{i,k}^n \hat{\Sigma}_{i,k|k-1}^n \\ &= \hat{\Sigma}_{i,k|k-1}^n \hat{\mathbf{z}}_k^n \left( \mathbf{u}_{i,k}^n \mathbf{S}_{i,k}^n \mathbf{u}_{i,k}^n \right) \hat{\mathbf{z}}_k^n \hat{\Sigma}_{i,k|k-1}^n \end{aligned}$$

に代入すれば逆行列を含まないリカッチ方程式が得られる。

$$\begin{aligned} \hat{\Sigma}_{i,k+1|k}^n &= \hat{\Sigma}_{i,k|k-1}^n \\ &\quad - \frac{(1 - \gamma_f^{-2} \hat{\sigma}_{v_k}^2) \beta_{i,k}^n}{\hat{\sigma}_{v_k}^2 + (1 - \gamma_f^{-2} \hat{\sigma}_{v_k}^2) \alpha_{i,k}^n \beta_{i,k}^n} \psi_{i,k}^n \psi_{i,k}^n \end{aligned} \quad (40)$$

ただし、

$$\psi_{i,k}^n = \hat{\Sigma}_{i,k|k-1}^n \hat{\mathbf{z}}_k^n \quad (41)$$

同様に、フィルタゲインも式 (34) を用いれば次のように整理できる ( $\mathbf{H}_{i,k}^n \hat{\Sigma}_{i,k|k-1}^n \mathbf{H}_{i,k}^n = \mathbf{u}_{i,k}^n \alpha_{i,k}^n \mathbf{u}_{i,k}^n$ )。

$$\begin{aligned} \mathbf{K}_{i,k}^n &= \hat{\Sigma}_{i,k|k-1}^n \mathbf{H}_{i,k}^n \mathbf{T} (\hat{\sigma}_{v_k}^2 \mathbf{I} + \mathbf{H}_{i,k}^n \hat{\Sigma}_{i,k|k-1}^n \mathbf{H}_{i,k}^n)^{-1} \\ &= \frac{1}{\hat{\sigma}_{v_k}^2 + \alpha_{i,k}^n \beta_{i,k}^n} \psi_{i,k}^n \mathbf{u}_{i,k}^n \end{aligned} \quad (42)$$

以上より, 図 4 の  $l$ -EHF アルゴリズムが得られる。このアルゴリズムでは出力の観測雑音の分散  $\hat{\sigma}_{v_k}^2$  が推定されている。一方、ハイパー  $H_\infty$  フィルタ [8] を用いて同様に求めた図 5 の  $l$ -EHHF アルゴリズムでは忘却係数  $\rho$  が効果的に利用されている。

## 6 シミュレーション

本章では英文字パターン認識問題を用いて深層残差  $H_\infty$  学習を評価する。ここで、 $A \sim Z$  までの合計 26 パターンを用いた ( $N_p = 26$ )。これらは 0 または 1 の値をとる  $7 \times 7$  ドットで表されている。図 2 はスキップ接続がないプレーンな 49-12-26 ネットワークの目的とする入出力関係を示す。これより、あるパターンが入力された場合に、出力層内の対応する 1 つのニューロンのみが値 1 をとり、その他のニューロンは値 0 をとる。

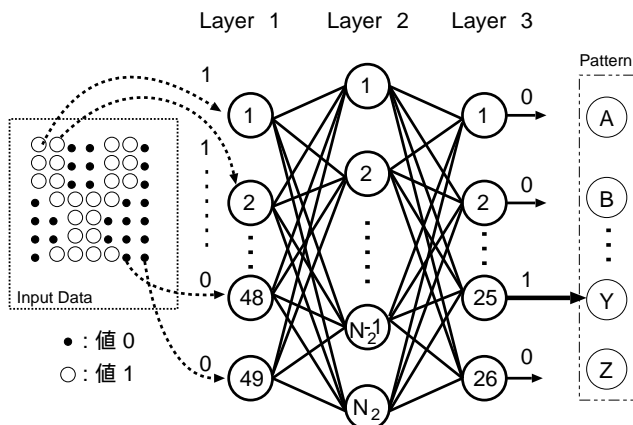


図 2 ニューラルネットによるパターン'Y'の認識

図 3 には、層数  $L$  が 5, 21, 71 の ResNet に対する  $l$ -EHHF の 100 試行の学習曲線、表 1 には学習回数 (epoch 数 = 更新回数 / パターン数) の平均、最大、最小を示す。一方、 $l$ -EHF ( $\gamma_f = 1.7$ ) は 21 層で一つの初期値で更新回数が上限に達した。ただし、重みベクトル  $w$  の学習は各試行ごとに  $(-0.2 \sim 0.2)$  の範囲の乱数により初期化され、 $10^2$  試行を行った。学習は出力 2 乗誤差  $J$  が 0.05 以下となったとき終了し、更新回数は  $10^4$  回を上限とした。また、シグモイド関数の傾き  $\eta_0$  は 2.5 とし、 $l$ -EHF、 $l$ -EHHF において  $\epsilon_0 = 0.1$  に固定した。

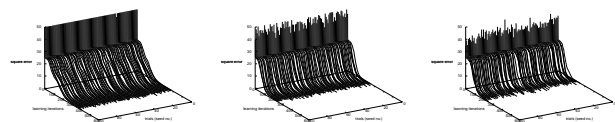
## 7 まとめ

本研究では、各ニューロンに対する部分重みベクトルに着目してリカッチ方程式の共分散行列をブロック対角化した上で、多出力系に起因する逆行列の計算を部分観測行列の行列分解を用いて回避することにより深層残差  $H_\infty$  学習の計算量を大幅に削減することに成功した。

今後は、大規模な畳み込みネットワーク (CNN) 等の  $H_\infty$  学習に挑戦したい。

表 1  $L$  層 ResNet における学習終了時の学習回数 (epoch 数) (区間  $[-0.20, 0.20]$  の一様分布により重みを初期化した 100 試行); 学習法は  $l$ -EHHF 法 ( $\gamma_f = 32.0, \rho = 1 - \gamma_f^{-2}$ )、入力層のニューロン数は  $N_1 = 49$ 、隠れ層のニューロン数は  $N_n = 5$ 、出力層のニューロン数は  $N_L = 26$ 、シグモイド関数の傾きは  $\eta_0 = 2.5$ 、打ち切り誤差は  $5.0 \times 10^{-2}$ 。

層数 $L$	学習回数の平均	学習回数の最大/最小
5	15.0	16.3/14.0
11	11.7	13.2/10.8
21	9.89	11.2/9.04
51	8.48	10.8/6.81
71	8.46	10.2/6.54



(a)  $L = 5$     (b)  $L = 21$     (c)  $L = 71$

図 3  $L$  層 ResNet における 100 試行の学習曲線 (重みは区間  $[-0.20, 0.20]$  の一様分布により初期化); 学習法は  $l$ -EHHF 法 ( $\gamma_f = 32.0, \rho = 1 - \gamma_f^{-2}$ )、入力層のニューロン数は  $N_1 = 49$ 、隠れ層のニューロン数は  $N_n = 5$ 、出力層のニューロン数は  $N_L = 26$ 、シグモイド関数の傾きは  $\eta_0 = 2.5$ 、打ち切り誤差は  $5.0 \times 10^{-2}$ 。

## 参考文献

- [1] M. A. Nielsen, Neural networks and deep learning, Determination Press, 2015.
- [2] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp.770-778, 2016.
- [3] 西山 清, "深層ニューラルネットワークにおける鞍点問題を攻略するための深層残差  $H_\infty$  学習," 第 20 回情報科学技術フォーラム, F-017, pp.1-4, 2021.
- [4] 西山 清, 最適フィルタリング, 培風館, 2001.
- [5] K. Nishiyama and K. Suzuki, " $H_\infty$ -learning of layered neural networks," IEEE Trans. Neural Networks, 12, 6, pp.1265-1277, 2001.
- [6] 西山 清, 菅原 康滉, "深層残差ネットワークの  $H_\infty$  学習," 第 19 回情報科学技術フォーラム, G-007, pp.1-3, 2020.
- [7] 西山 清, 鈴木清彦, " $H_\infty$  学習 - 局所的最適化アプローチ," 電子情報通信学会 技術研究報告, NC2001-224, pp.223-230, 2002.
- [8] K. Nishiyama, "A Unified View of Adaptive Algorithms for Finite Impulse Response Filters using the  $H_\infty$  Framework," Signal Processing (Elsevier), 97, pp.55-64, April 2014.

$$\begin{aligned}
 & \text{for}(k = 1, \dots, N_p, 1 + N_p, \dots, 2N_p, 1 + 2N_p, \dots) \\
 & \{ \\
 & \quad \hat{\mathbf{z}}_k^n, \quad n = 1, 2, \dots, L \quad (\text{forward propagation with } \hat{\mathbf{w}}_{k-1}) \\
 & \quad \hat{\sigma}_{v_k}^2 = \hat{\sigma}_{v_{k-1}}^2 + \mu_k \left[ \frac{(\mathbf{y}_k - \hat{\mathbf{z}}_k^L)^T (\mathbf{y}_k - \hat{\mathbf{z}}_k^L)}{N_L} - \hat{\sigma}_{v_{k-1}}^2 \right] \\
 & \quad \hat{\mathbf{h}}_{1,k}^{L-1} = \hat{\mathbf{z}}_k^L \\
 & \text{for}(n = L - 1; n \geq 1; n \leftarrow n - 1) \\
 & \{ \\
 & \quad \text{for}(i = 1; i \leq N_{n+1}; i \leftarrow i + 1) \\
 & \quad \{ \\
 & \quad \quad \mathbf{u}_{i,k}^n = \begin{cases} \eta_0 \hat{z}_{i,k}^{n+1} (1 - \hat{z}_{i,k}^{n+1}) [0, \dots, 0, \overset{i}{1}, 0, \dots, 0]^T, & n = L - 1 \\ \sum_{j=1}^{N_{n+2}} \hat{w}_{j,i,k}^{n+1} \mathbf{u}_{j,k}^{n+1}, & n = L - 2 \\ \sum_{j=1}^{N_{n+2}} \hat{w}_{j,i,k}^{n+1} \mathbf{u}_{j,k}^{n+1} + \mathbf{u}_{i,k}^{n+2}, & n \% 2 \neq 0 \\ \eta_0 \hat{z}_{i,k}^{n+1} (1 - \hat{z}_{i,k}^{n+1}) \sum_{j=1}^{N_{n+2}} \hat{w}_{j,i,k}^{n+1} \mathbf{u}_{j,k}^{n+1}, & n \% 2 = 0 \\ \eta_0 \hat{z}_{i,k}^{n+1} (1 - \hat{z}_{i,k}^{n+1}) \left( \sum_{j=1}^{N_{n+2}} \hat{w}_{j,i,k}^{n+1} \mathbf{u}_{j,k}^{n+1} + \mathbf{u}_{i,k}^{n+2} \right), & n = 1 \end{cases} \\
 & \quad \quad \psi_{i,k}^n = \hat{\Sigma}_{i,k|k-1}^n \hat{\mathbf{z}}_k^n \\
 & \quad \quad \alpha_{i,k}^n = \hat{\mathbf{z}}_k^{nT} \psi_{i,k}^n, \quad \beta_{i,k}^n = \mathbf{u}_{i,k}^{nT} \mathbf{u}_{i,k}^n \\
 & \quad \quad \hat{\mathbf{w}}_{i,k}^n = \hat{\mathbf{w}}_{i,k-1}^n + \frac{\mathbf{u}_{i,k}^{nT} [\mathbf{y}_k - \hat{\mathbf{h}}_{i,k}^n]}{\hat{\sigma}_{v_k}^2 + \alpha_{i,k}^n \beta_{i,k}^n} \psi_{i,k}^n \\
 & \quad \quad \hat{\Sigma}_{i,k+1|k}^n = \hat{\Sigma}_{i,k|k-1}^n - \frac{(1 - \gamma_f^{-2} \hat{\sigma}_{v_k}^2) \beta_{i,k}^n}{\hat{\sigma}_{v_k}^2 + (1 - \gamma_f^{-2} \hat{\sigma}_{v_k}^2) \alpha_{i,k}^n \beta_{i,k}^n} \psi_{i,k}^n \psi_{i,k}^{nT} \\
 & \quad \quad \hat{\mathbf{h}}_{i+1,k}^n = \hat{\mathbf{h}}_{i,k}^n + \mathbf{u}_{i,k}^n \hat{\mathbf{z}}_k^{nT} (\hat{\mathbf{w}}_{i,k}^n - \hat{\mathbf{w}}_{i,k-1}^n) \\
 & \quad \quad \} \\
 & \quad \quad \hat{\mathbf{h}}_{1,k}^{n-1} = \hat{\mathbf{h}}_{N_{n+1}+1,k}^n \\
 & \quad \quad \} \\
 & \quad \quad \hat{\mathbf{w}}_k = \left[ \hat{\mathbf{w}}_{1,k}^1 T, \dots, \hat{\mathbf{w}}_{N_2,k}^1 T, \dots, \hat{\mathbf{w}}_{1,k}^{L-1} T, \dots, \hat{\mathbf{w}}_{N_L,k}^{L-1} T \right]^T \\
 & \quad \quad \} \\
 & \}
 \end{aligned}$$

 図 4 The  $l$ -EHF learning algorithm with skip connections ;

$$\begin{aligned}
 & \hat{\mathbf{w}}_{i,k}^n = [\hat{\theta}_{i,k}^n, \hat{w}_{i,1,k}^n, \dots, \hat{w}_{i,N_n,k}^n]^T, \quad \hat{\mathbf{z}}_k^n = [1, \hat{z}_{1,k}^n, \dots, \hat{z}_{N_n,k}^n]^T, \quad \hat{z}_{i,k}^1 = z_{i,k}^1, \\
 & \Phi_k^{n+1} = [\mathbf{u}_{1,k}^{n+1}, \dots, \mathbf{u}_{N_{n+1},k}^{n+1}], \quad \hat{\Sigma}_{i,1|0}^n = \epsilon_0 \mathbf{I}, \quad \mu_k = \begin{cases} \frac{1}{k}, & k \leq T_{max} \\ \frac{1}{T_{max}}, & \text{otherwise} \end{cases}, \quad T_{max} = 10N_p
 \end{aligned}$$

$$\begin{aligned}
 & \text{for}(k = 1, \dots, N_p, 1 + N_p, \dots, 2N_p, 1 + 2N_p, \dots) \\
 & \{ \\
 & \quad \hat{\mathbf{z}}_k^n, \quad n = 1, 2, \dots, L \quad (\text{forward propagation with } \hat{\mathbf{w}}_{k-1}) \\
 & \quad \hat{\mathbf{h}}_{1,k}^{L-1} = \hat{\mathbf{z}}_k^L \\
 & \quad \text{for}(n = L - 1; n \geq 1; n \leftarrow n - 1) \\
 & \quad \{ \\
 & \quad \quad \text{for}(i = 1; i \leq N_{n+1}; i \leftarrow i + 1) \\
 & \quad \quad \{ \\
 & \quad \quad \quad \mathbf{u}_{i,k}^n = \begin{cases} \eta_0 \hat{z}_{i,k}^{n+1} (1 - \hat{z}_{i,k}^{n+1}) [0, \dots, 0, 1, 0, \dots, 0]^T, & n = L - 1 \\ \sum_{j=1}^{N_{n+2}} \hat{w}_{j,i,k}^{n+1} \mathbf{u}_{j,k}^{n+1}, & n = L - 2 \\ \sum_{j=1}^{N_{n+2}} \hat{w}_{j,i,k}^{n+1} \mathbf{u}_{j,k}^{n+1} + \mathbf{u}_{i,k}^{n+2}, & n \% 2 \neq 0 \\ \eta_0 \hat{z}_{i,k}^{n+1} (1 - \hat{z}_{i,k}^{n+1}) \sum_{j=1}^{N_{n+2}} \hat{w}_{j,i,k}^{n+1} \mathbf{u}_{j,k}^{n+1}, & n \% 2 = 0 \\ \eta_0 \hat{z}_{i,k}^{n+1} (1 - \hat{z}_{i,k}^{n+1}) \left( \sum_{j=1}^{N_{n+2}} \hat{w}_{j,i,k}^{n+1} \mathbf{u}_{j,k}^{n+1} + \mathbf{u}_{i,k}^{n+2} \right), & n = 1 \end{cases} \\
 & \quad \quad \quad \psi_{i,k}^n = \hat{\Sigma}_{i,k|k-1}^n \hat{\mathbf{z}}_k^n \\
 & \quad \quad \quad \alpha_{i,k}^n = \hat{\mathbf{z}}_k^n T \psi_{i,k}^n, \quad \beta_{i,k}^n = \mathbf{u}_{i,k}^n T \mathbf{u}_{i,k}^n \\
 & \quad \quad \quad \hat{\mathbf{w}}_{i,k}^n = \hat{\mathbf{w}}_{i,k-1}^n + \frac{\mathbf{u}_{i,k}^n T [\mathbf{y}_k - \hat{\mathbf{h}}_{i,k}^n]}{\rho + \alpha_{i,k}^n \beta_{i,k}^n} \psi_{i,k}^n \\
 & \quad \quad \quad \hat{\Sigma}_{i,k+1|k}^n = \left( \hat{\Sigma}_{i,k|k-1}^n - \frac{(1 - \gamma_f^{-2}) \beta_{i,k}^n}{\rho + (1 - \gamma_f^{-2}) \alpha_{i,k}^n \beta_{i,k}^n} \psi_{i,k}^n \psi_{i,k}^n T \right) / \rho \\
 & \quad \quad \quad \hat{\mathbf{h}}_{i+1,k}^n = \hat{\mathbf{h}}_{i,k}^n + \mathbf{u}_{i,k}^n \hat{\mathbf{z}}_k^n T (\hat{\mathbf{w}}_{i,k}^n - \hat{\mathbf{w}}_{i,k-1}^n) \\
 & \quad \quad \quad \hat{\mathbf{h}}_{1,k}^{n-1} = \hat{\mathbf{h}}_{N_{n+1}+1,k}^n \\
 & \quad \quad \} \\
 & \quad \quad \mathbf{\hat{w}}_k = \left[ \hat{\mathbf{w}}_{1,k}^1 T, \dots, \hat{\mathbf{w}}_{N_2,k}^1 T, \dots, \hat{\mathbf{w}}_{1,k}^{L-1 T}, \dots, \hat{\mathbf{w}}_{N_L,k}^{L-1 T} \right]^T \\
 & \quad \} \\
 & \}
 \end{aligned}$$

図 5 The  $l$ -EHHF learning algorithm with skip connections ;  
 $\hat{\mathbf{w}}_{i,k}^n = [\hat{\theta}_{i,k}^n, \hat{w}_{i,1,k}^n, \dots, \hat{w}_{i,N_n,k}^n]^T$ ,  $\hat{\mathbf{z}}_k^n = [1, \hat{z}_{1,k}^n, \dots, \hat{z}_{N_n,k}^n]^T$ ,  $\hat{z}_{i,k}^1 = z_{i,k}^1$ ,  
 $\Phi_k^{n+1} = [\mathbf{u}_{1,k}^{n+1}, \dots, \mathbf{u}_{N_{n+1},k}^{n+1}]$ ,  $\mathbf{R}_k^n = \begin{bmatrix} \rho \mathbf{I} & 0 \\ 0 & -\gamma_f^2 \rho \mathbf{I} \end{bmatrix}$ ,  $\hat{\Sigma}_{i,1|0}^n = \epsilon_0 \mathbf{I}$