

ニューラルネットワークについて

前川眞一

19980328, 051102,0320, 20151013,20210624, 20210719,20dnc,21,27dnc,2021.

12. 22

1 教師信号ありの順伝搬型ニューラルネットワーク

Q 個の説明変数 x を第 1 層に、 m 個の基準変数 y を予測する最終層を第 n 層とするネットワークを考える。各層のニューロンの数は $n_k, k = 1, 2, \dots, n$ とする。(ただし、後述する常に 1 を出力するニューロンは数えない。) 従って、 $n_1 = Q, n_n = m$ である。また、個体 a の第 k 層の第 j 番目のニューロンへの入力を i_{aj}^k 、同じく出力を o_{aj}^k と書く。なお、本節では、上付きの添え字 $k, k-1, k+1, n$ 等はすべて添え字であり冪乗を示すものではないことに注意されたい。各個体を表す添え字は $a, a = 1, 2, \dots, N$ である。各ニューロンはそれぞれ上位のニューロンに向かって同一の出力を発する。また、各ニューロンは下位のニューロンの発する出力を重み付けして加えた形でそのニューロンへの入力を形成する。すなわち、個体 a の第 $k-1$ 層の第 i 番目のニューロンの出力 o_{ai}^{k-1} から個体 a の第 k 層の第 j 番目のニューロンへの入力 i_{aj}^k を作る際の重みを w_{ij}^k とすれば、

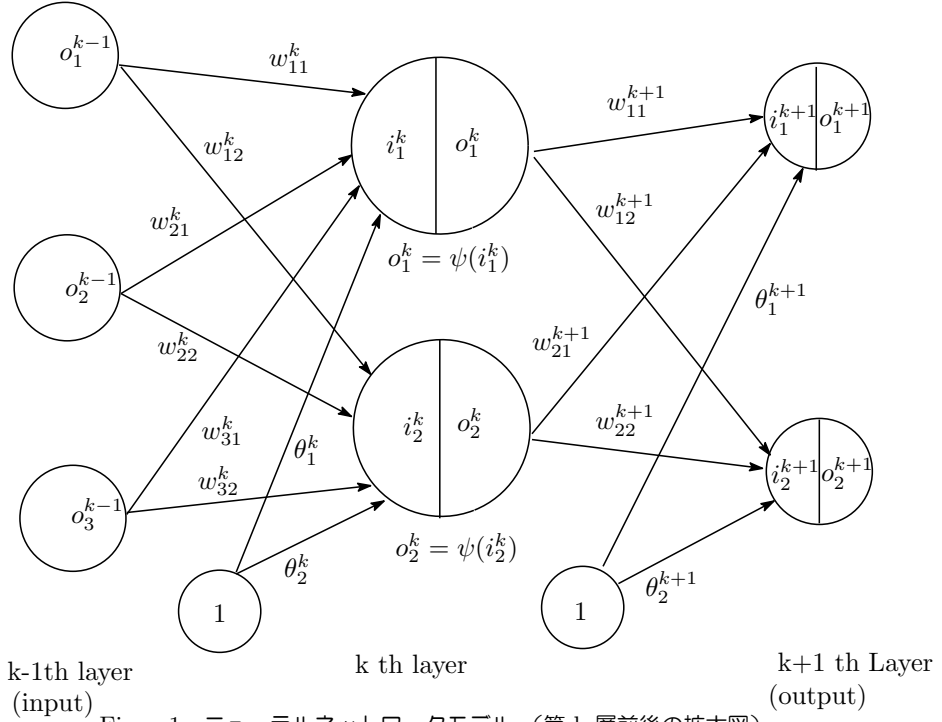
$$i_{aj}^k = \sum_{i=1}^{n_{k-1}} o_{ai}^{k-1} w_{ij}^k + \theta_j^k \quad (1)$$

である。ただし、 θ_k はバイアスと呼ばれる定数項(閾値)である。また、第 k 層の入力はシグモイド関数(たとえばロジスティック関数) ψ により、

$$o_{aj}^k = \psi(i_{aj}^k) = \frac{1}{1 + \exp(-i_{aj}^k)} \quad (2)$$

という変換を受けた形で出力される。ただし、最終層 o_{aj}^n は y_{aj} を近似するためのものであるため、シグモイド関数による変換を施さない場合もあるが、その場合は、 $o_{aj}^n = i_{aj}^n$ である。この過程を、第 k 層の前後に関して図示したものが Figure 1 である。図中、第 $k-1$ 層のニューロンの数は 3、第 k 層のニューロンの数は 2、第 $k+1$ 層のニューロンの数は 2 であるが、定数項 θ は、常に 1 を出力するニューロンへの重みとして記述してある。また、個体を表す添字 a は省略してある。

これを全ての $i_{aj}^k, o_{aj}^k, a = 1, 2, \dots, N$ に関してまとめて書くと、 \mathbf{O}^k を $N \times n_k$ の行列、 \mathbf{W}^k



を $n_{k-1} \times n_k$ の行列、 θ^k を $1 \times n_k$ のベクトルとして、

$$\mathbf{O}^k = \Psi(\mathbf{I}^k) \quad (3)$$

$$= \Psi(\mathbf{O}^{k-1} \mathbf{W}^k + \mathbf{1}\theta^k) \quad (4)$$

$$= \Psi(\Psi(\mathbf{I}^{k-1}) \mathbf{W}^k + \mathbf{1}\theta^k) \quad (5)$$

$$= \Psi(\Psi(\mathbf{O}^{k-2} \mathbf{W}^{k-1} + \mathbf{1}\theta^{k-1}) \mathbf{W}^k + \mathbf{1}\theta^k) \quad (6)$$

$$\vdots \quad (7)$$

となるが、これは

$$\mathbf{O}^k = \Psi(\mathbf{O}^{k-1} \mathbf{W}^k + \mathbf{1}\theta^k) = \Gamma_k(\mathbf{O}^{k-1}) \quad (8)$$

と定義をすれば、最終層の出力が

$$\mathbf{O}^n = \Gamma_n(\mathbf{O}^{n-1}) \quad (9)$$

$$= \Gamma_n(\Gamma_{n-1}(\mathbf{O}^{n-2})) \quad (10)$$

$$\vdots \quad (11)$$

$$= \Gamma_n \circ \Gamma_{n-1} \circ \cdots \circ \Gamma_2(\mathbf{O}^1) \quad (12)$$

$$= \Gamma(\mathbf{O}^1) = \Gamma(\mathbf{X}) \quad (13)$$

と書けることを示している。ただし、

$$\Gamma(\cdot) = \Gamma_n \circ \Gamma_{n-1} \circ \cdots \circ \Gamma_2(\cdot) \quad (14)$$

である。従って、ニューラルネットワークのモデルは、基準変数 Y を説明変数 X の非線形関数 $\hat{Y} = O^n = \Gamma(X)$ で近似するという意味で

$$Y = \Gamma(X) + \mathcal{E} \quad \text{or} \quad \hat{Y} = \Gamma(X) \quad (15)$$

と書くことができる。

なお、以下の記述では、各層の第 n_k 番目のニューロンに加えて、常に 1 を出力するニューロンがもう一つ存在し、その重みを θ^k としている。すなわち、 θ^k は W^k 行列の最初の行、すなわち、第 0 行として取り扱われるが、これは O^k 行列を $N \times (1 + n_k)$ の行列と定義し、その第 0 列がすべて 1 であるとするに相当する。

2 最小2乗法によるパラメタの推定に必要な1次微分

データ Y とモデルの値 \hat{Y} の差の2乗和

$$\text{RSS} = \sum_{s=1}^Q \sum_{a=1}^N (y_{as} - \hat{y}_{as})^2 \quad (16)$$

を最小にするようにモデルのパラメタ

$$w_{ij}^k, \quad k = 2, 3, \dots, n, \quad i = 0, 1, 2, \dots, n_{k-1}, \quad j = 1, 2, \dots, n_k$$

を求めるためには以下に示す \hat{y}_{as} を各パラメタで偏微分したもの、すなわちヤコビアン行列が必要となる。ただし、 w_{0j}^k は前節まで θ_j^k と記していたものであり、

$$\hat{y}_{as} = o_{as}^n \quad (17)$$

である。なお、以下の記述では、 o と i に付く個体を表す添字 a は省略してある。

$$\frac{\partial o_s^n}{\partial w_{ij}^k} = \frac{\partial o_s^n}{\partial i_j^k} \times \frac{\partial i_j^k}{\partial w_{ij}^k} \quad (18)$$

$$\frac{\partial o_s^n}{\partial i_j^k} = \frac{\partial o_s^n}{\partial o_j^k} \times \frac{\partial o_j^k}{\partial i_j^k} \quad (19)$$

$$\frac{\partial o_s^n}{\partial o_j^k} = \sum_{q=1}^{n_{k+1}} \frac{\partial o_s^n}{\partial i_q^{k+1}} \times \frac{\partial i_q^{k+1}}{\partial o_j^k} \quad (20)$$

$$\frac{\partial o_j^k}{\partial i_j^k} = \begin{cases} o_j^k(1 - o_j^k) & \text{最終層も } \Psi \text{ で変換する場合} \\ 1 & \text{最終層は変換しない場合} \end{cases} \quad (21)$$

$$\frac{\partial i_q^{k+1}}{\partial o_j^k} = w_{jq}^{k+1} \quad (22)$$

$$\frac{\partial i_j^k}{\partial w_{ij}^k} = \begin{cases} o_i^{k-1} & \text{if } i \geq 1 \quad (w) \\ 1 & \text{if } i = 0 \quad (\theta) \end{cases} \quad (23)$$

なお、Eq.(20) に於いて i_0^{k+1} は定数 1 であり o_j^k に依存しないため、総和記号の和は $q = 1$ からとなる。従って、 $\delta_j^{k,s} = \frac{\partial o_s^n}{\partial i_j^k}$ と定義すると、この値は一つ上の第 $k + 1$ 層から最終層である

第 n 層までの値を使って、 $2 \leq k < n$ に関して

$$\delta_j^{k,s} = \frac{\partial o_s^n}{\partial i_j^k} \quad (24)$$

$$= \sum_{q=1}^{n_{k+1}} \frac{\partial o_s^n}{\partial i_q^{k+1}} w_{jq}^{k+1} \times \frac{\partial o_j^k}{\partial i_j^k} = \sum_{q=1}^{n_{k+1}} \delta_q^{k+1,s} w_{jq}^{k+1} \times \frac{\partial o_j^k}{\partial i_j^k} \quad (25)$$

$$= \sum_{q=1}^{n_{k+1}} \delta_q^{k+1,s} w_{jq}^{k+1} \times o_j^k (1 - o_j^k) \quad (26)$$

と再帰的に書くことが出来るが、これを誤差逆伝搬 error back propagation と呼ぶ。なお、最終層である第 n 層は

$$\delta_j^{n,s} = \frac{\partial o_s^n}{\partial i_j^n} = \begin{cases} \begin{cases} o_s^n (1 - o_s^n) & \text{最終層も } \Psi \text{ で変換する場合} \\ 1 & \text{最終層は変換しない場合} \end{cases} & \text{if } j = s \\ 0 & \text{otherwise} \end{cases} \quad (27)$$

である。

まとめると

$$\frac{\partial o_s^n}{\partial w_{ij}^k} = \frac{\partial o_s^n}{\partial i_j^k} \times \frac{\partial i_j^k}{\partial w_{ij}^k} = \delta_j^{n,s} \times \begin{cases} o_i^{k-1} & \text{if } i \geq 1 \quad (w) \\ 1 & \text{if } i = 0 \quad (\theta) \end{cases} \quad (28)$$

$$= \sum_{q=1}^{n_{k+1}} \delta_q^{k+1,s} w_{jq}^{k+1} \times o_j^k (1 - o_j^k) \times \begin{cases} o_i^{k-1} & \text{if } i \geq 1 \quad (w) \\ 1 & \text{if } i = 0 \quad (\theta) \end{cases} \quad (29)$$

となる。また、RSS の w_{ij}^k に関する微分 (gradient) は

$$\frac{\partial \text{RSS}}{\partial w_{ij}^k} = -2 \sum_{s=1}^m \sum_{a=1}^N (y_{as} - o_{as}^n) \times \frac{\partial o_{as}^n}{\partial w_{ij}^k} \quad (30)$$

となる。

2.1 行列表記

いま、 \mathbf{O}^k , $k = 1, 2, \dots, n$ を $N \times n_k$ の第 k 層の出力を要素として持つ行列とする。ただし、 $\mathbf{O}^1 = \mathbf{X}$ であり、最終層の第 s 列を $o_{(s)}^n$ という $N \times 1$ のベクトルで表す。また、 \mathbf{W}^k , $k = 2, 3, \dots, n$ を第 $k-1$ 層の出力から第 k 層の入力を

$$\mathbf{I}^k = [\mathbf{1}_N, \mathbf{O}^{k-1}] \mathbf{W}^k \quad (31)$$

という形で作るための $(n_{k-1} + 1) \times n_k$ の重み行列とする。なお、 \mathbf{O}^{k-1} 行列の最初の列を $\mathbf{1}$ と定義せずに、 $N \times n_{k-1}$ の行列 \mathbf{O}^{k-1} の第 1 列の左側に $\mathbf{1}$ を付け加えたものに \mathbf{W}^k を掛けて \mathbf{I}^k を定義していることに注意されたい。したがって、 $(n_{k-1} + 1) \times n_k$ の行列 \mathbf{W}^k を行単位に

$$\mathbf{W}^k = (\mathbf{w}_0^k, \mathbf{w}_1^k, \mathbf{w}_2^k, \dots, \mathbf{w}_{n_{k-1}}^k)' \quad (32)$$

と定義すれば、第 0 行 $\mathbf{w}_0^{k'}$ は定数項 1 に掛かる係数、第 1 行 $\mathbf{w}_1^{k'}$ が \mathbf{O}^{k-1} の第 1 列、 $o_{(1)}^{k-1}$ 、に掛かる係数、第 2 行 $\mathbf{w}_2^{k'}$ が $o_{(2)}^{k-1}$ に掛かる係数、という具合になる。また、 Δ_s^k を Eq.(26)

で定義される $\delta_j^{k,s}$ 、すなわち $\delta_{aj}^{k,s} = \frac{\partial o_{as}^n}{\partial i_{aj}^k}$ をその a, j 要素として持つ $N \times n_k$ の行列と定義すると、Eq.(26) と Eq.(27) は

$$\Delta_s^k = \begin{cases} \begin{cases} \mathbf{o}_{(s)}^n \odot (\mathbf{1}_N - \mathbf{o}_{(s)}^n) \mathbf{e}'_s & \text{最終層も } \Psi \text{ で変換する場合} \\ \mathbf{1}_N & \text{最終層は変換しない場合} \end{cases} & \text{if } k = n \\ \mathbf{O}^k \odot (\mathbf{1}_N \mathbf{1}'_{n_k} - \mathbf{O}^k) \odot (\Delta_s^{k+1} \mathbf{W}^{k+1\dagger}) & \text{if } k = n-1, \dots, 2 \end{cases} \\ k = 2, 3, \dots, n-1, \quad s = 1, 2, \dots, n_n \quad (33)$$

と書くことが出来る。ただし、 \odot は行列のアダマール積を表し、 \mathbf{e}_s はその第 s 行のみが 1 となっている $n_n \times 1$ のベクトルである。また、 $\mathbf{W}^{k+1\dagger}$ は \mathbf{W}^{k+1} 行列の第 0 行目を取り除いた $n_k \times n_{k+1}$ 行列を転置して得られる $n_{k+1} \times n_k$ の行列である。(この \mathbf{W} の第 0 行目を取り除くことは Eq.(20) の和が $q = 1$ から始まっていることに対応している。)

したがって、Eq.(28) に示す $\frac{\partial o_{as}^n}{\partial w_{ij}^k}$ の w_i^k の要素 w_{ij}^k , $j = 1, 2, \dots, n_k$ に関する部分をその aj 要素として持つ $N \times n_k$ のヤコビアン行列は

$$\frac{\partial \mathbf{o}_{(s)}^n}{\partial \mathbf{w}_i^k} = \begin{cases} \left(\mathbf{o}_{(i)}^{k-1} \mathbf{1}'_{n_k} \right) \odot \Delta_s^k & \text{if } i \geq 1 \\ \Delta_s^k & \text{if } i = 0 \end{cases} \quad (34)$$

であり、第 k 層の重み行列 \mathbf{W}^k のすべての要素に関するヤコビアン行列は Eq.(34) を $i = 0, 1, 2, \dots, n_k$ と変えて横に並べた $N \times (n_{k-1} + 1)n_k$ の行列

$$\frac{\partial \mathbf{o}_{(s)}^n}{\partial \text{vec}(\mathbf{W}^{k'})} = \left[\frac{\partial \mathbf{o}_{(s)}^n}{\partial \mathbf{w}_0^k}, \frac{\partial \mathbf{o}_{(s)}^n}{\partial \mathbf{w}_1^k}, \frac{\partial \mathbf{o}_{(s)}^n}{\partial \mathbf{w}_2^k}, \dots, \frac{\partial \mathbf{o}_{(s)}^n}{\partial \mathbf{w}_{n_k}^k} \right] \quad (35)$$

$$= \left[\Delta_s^k, \left(\mathbf{o}_{(1)}^{k-1} \mathbf{1}'_{n_k} \right) \odot \Delta_s^k, \left(\mathbf{o}_{(2)}^{k-1} \mathbf{1}'_{n_k} \right) \odot \Delta_s^k, \dots, \left(\mathbf{o}_{(n_{k-1})}^{k-1} \mathbf{1}'_{n_k} \right) \odot \Delta_s^k \right] \quad (36)$$

となる。ただし、 $\text{vec}(\mathbf{W}^{k'})$ は \mathbf{W}^k 行列の各行を横に並べて作られた横ベクトルの転置である。また、

$$\text{RSS} = \sum_{s=1}^m (\mathbf{y}_{(s)} - \mathbf{o}_{(s)}^n)' (\mathbf{y}_{(s)} - \mathbf{o}_{(s)}^n) \quad (37)$$

$$= (\text{vec}(\mathbf{Y}) - \text{vec}(\mathbf{O}^n))' (\text{vec}(\mathbf{Y}) - \text{vec}(\mathbf{O}^n)) \quad (38)$$

であるから、全ての n_n 個の変数に関する $(N \times n_n) \times (n_{k-1} + 1)n_k$ のヤコビアン行列は

$$\frac{\partial \text{vec}(\mathbf{O}^n)}{\partial \text{vec}(\mathbf{W}^{k'})} = \left[\left(\frac{\partial \mathbf{o}_{(1)}^n}{\partial \text{vec}(\mathbf{W}^{k'})} \right)', \left(\frac{\partial \mathbf{o}_{(2)}^n}{\partial \text{vec}(\mathbf{W}^{k'})} \right)', \dots, \left(\frac{\partial \mathbf{o}_{(n_n)}^n}{\partial \text{vec}(\mathbf{W}^{k'})} \right)' \right]' \quad (39)$$

となり、これらを層の添字 k を変えて横に並べたものが全ての重み行列に関するヤコビアン行列となる。

なお、 $\frac{\partial \mathbf{o}_{(s)}^n}{\partial \text{vec}(\mathbf{W}^{k'})}$ 行列の列を commutation matrix $\mathbf{K}(n_k, n_{k-1} + 1)$ を用いて並べ替えると

$$\frac{\partial \mathbf{o}_{(s)}^n}{\partial \text{vec}(\mathbf{W}^k)} = \frac{\partial \mathbf{o}_{(s)}^n}{\partial \text{vec}(\mathbf{W}^{k'})} \mathbf{K}(n_k, n_{k-1} + 1) \quad (40)$$

が得られる。ただし、 $e_{a,b}$ をその第 b 要素の値のみが 1 でそれ以外の要素は 0 である a 次元の縦ベクトルとして

$$\mathbf{K}(r, c) = \sum_{i=1}^r \sum_{j=1}^c (\mathbf{e}_{r,i} \mathbf{e}'_{c,j}) \otimes (\mathbf{e}_{c,j} \mathbf{e}'_{r,i}) \quad (41)$$

と書ける。(commutation matrix に関しては Harville, 1997, や wikipedia を参照のこと。)したがって、 $\frac{\partial \text{vec}(\mathbf{O}^n)}{\partial \text{vec}(\mathbf{W}^k)}$ は上記の $\frac{\partial \mathbf{o}_{(s)}^n}{\partial \text{vec}(\mathbf{W}^k)}$ を縦に並べて得ることが出来る。

最後に、RSS の \mathbf{W}^k に関する微分 (gradient) は以下の $(n_{k-1} + 1)n_k \times 1$ のベクトル

$$\frac{\partial \text{RSS}}{\partial \text{vec}(\mathbf{W}^{k'})} = -2 \left[\frac{\partial \text{vec}(\mathbf{O}^n)}{\partial \text{vec}(\mathbf{W}^k)} \right]' (\text{vec}(\mathbf{Y}) - \text{vec}(\mathbf{O}^n)) \quad (42)$$

となる。また、

$$\frac{\partial \text{RSS}}{\partial \text{vec}(\mathbf{W}^k)} = \mathbf{K}(n_k, n_{k-1} + 1) \frac{\partial \text{RSS}}{\partial \text{vec}(\mathbf{W}^{k'})} \quad (43)$$

である。

References

Mayekawa, Shin-ichi. (1995) An efficient algorithm for feed-forward neural network regression analysis. *Behaviormetrika*. Vol.22, No. 1, 67-90

Commutation matrix. (n.d.). In *Wikipedia* Retrieved from

https://en.wikipedia.org/w/index.php?title=Commutation_matrix&oldid=1027627556

Harville,D.A. (1997) *Matrix Algebra from a Statistician's Perspective*. Springer-Verlag New York, Inc.