

# パターン認識と機械学習

## 第5章 ニューラルネットワーク

橋本龍二

February 5, 2021

### 1 フィードフォワード

省略。

### 2 ネットワーク訓練

入力ベクトル  $\mathbf{x}$  に対し、一次元目標ベクトル  $t$  の従う分布を

$$p(t|\mathbf{x}, \mathbf{w}) = \mathcal{N}(t|y(\mathbf{x}, \mathbf{w}), \beta^{-1})$$

と置く。入力ベクトルの集合  $\{\mathbf{x}_n\} (n = 1, \dots, N)$  と対応する目標ベクトルの集合  $\{t_n\} (n = 1, \dots, N)$  が与えられたときの対数尤度関数は、

$$\begin{aligned} p(\mathbf{t}|X, \mathbf{w}, \beta) &= \prod_{n=1}^N p(t_n|\mathbf{x}_n, \mathbf{w}, \beta) \\ &= \prod_{n=1}^N \sqrt{\frac{\beta}{2\pi}} \exp\left\{-\frac{\beta(t_n - y(\mathbf{x}_n, \mathbf{w}))^2}{2}\right\} \\ &= \beta^{N/2} (2\pi)^{-N/2} \exp\left\{-\frac{\beta}{2} \sum_{n=1}^N (y(\mathbf{x}_n, \mathbf{w}) - t_n)^2\right\} \\ \therefore -\log p &= \frac{\beta}{2} \sum_{n=1}^N \{y(\mathbf{x}_n, \mathbf{w}) - t_n\}^2 - \frac{N}{2} \log \beta + \frac{N}{2} \log(2\pi) \quad (1) \end{aligned}$$

より、尤度関数を最大化する  $\mathbf{w}_{ML}$  は、誤差関数：

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(\mathbf{x}_n) - t_n\}^2$$

を最小化することによって得られる。 $\mathbf{w}_{ML}$  が求まれば、 $\beta_{ML}$  は(1)より

$$\frac{1}{\beta_{ML}} = \frac{1}{N} \sum_{n=1}^N \{y(\mathbf{x}_n, \mathbf{w}_{ML}) - t_n\}^2$$

と求まる。同様に、目標ベクトルを  $K$  次元としてその分布を

$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}) = \mathcal{N}(y(\mathbf{x}, \mathbf{w}), \beta^{-1})$$

と置くと、データ  $N$  個の対数尤度関数は

$$\begin{aligned} p(\mathbf{t}|X, \mathbf{w}, \beta) &= \prod_{n=1}^N p(\mathbf{t}_n|\mathbf{x}_n, \mathbf{w}, \beta) \\ &= \prod_{n=1}^N \frac{\beta^{K/2}}{(2\pi)^{K/2}} \exp\left\{-\frac{\beta}{2} \|y(\mathbf{x}_n, \mathbf{w}) - \mathbf{t}_n\|^2\right\} \\ &= \beta^{NK/2} (2\pi)^{-K/2} \exp\left\{-\frac{\beta}{2} \sum_{n=1}^N \|y(\mathbf{x}_n, \mathbf{w}) - \mathbf{t}_n\|^2\right\} \\ \therefore -\log p &= \frac{\beta}{2} \sum_{n=1}^N \|y(\mathbf{x}_n - \mathbf{w}) - \mathbf{t}_n\|^2 - \frac{NK}{2} \log \beta - \frac{K}{2} \log(2\pi) \quad (2) \end{aligned}$$

より、尤度関数を最大化する  $\mathbf{w}_{ML}$  は、誤差関数：

$$E(\mathbf{w}) = \frac{\beta}{2} \sum_{n=1}^N \|y(\mathbf{x}_n - \mathbf{w}) - \mathbf{t}_n\|^2$$

を最小化することによって得られる。 $\mathbf{w}_{ML}$  が求めれば、 $\beta_{ML}$  は(2)より

$$\frac{1}{\beta_{ML}} = \frac{1}{NK} \sum_{n=1}^N \|y(\mathbf{x}_n - \mathbf{w}) - \mathbf{t}_n\|^2$$

と求まる。

次に、二クラス分類を考える。目標変数  $t$  は  $t = 1$  でクラス  $\mathcal{C}_1$ 、 $t = 0$  でクラス  $\mathcal{C}_2$  を表す。活性  $a$  に対しネットワークの出力層は

$$y = \sigma(a) \equiv \frac{1}{1 + \exp(-a)}$$

なる活性化関数を持つ。今、入力  $\mathbf{x}$  が与えられたときの条件付き分布を

$$\begin{aligned} p(t|\mathbf{x}, \mathbf{w}) &= Ber(y) \\ &= y(\mathbf{x}, \mathbf{w})^t \{1 - y(\mathbf{x}, \mathbf{w})\}^{1-t} \end{aligned}$$

とすると、データ  $N$  個の対数尤度関数は

$$\begin{aligned}
 p(\mathbf{t}|\mathbf{x}, \mathbf{w}) &= \prod_{n=1}^N p(\mathbf{t}_n|\mathbf{x}_n, \mathbf{w}) \\
 &= \prod_{n=1}^N (y_n)^{t_n} (1 - y_n)^{1-t_n} \\
 &= (y_n)^{\sum_{n=1}^N t_n} (1 - y_n)^{\sum_{n=1}^N 1-t_n} \\
 \therefore -\log p &= -\sum_{n=1}^N \{t_n \log y_n + (1 - t_n) \log(1 - y_n)\}
 \end{aligned}$$

ただし、 $y_n = y(\mathbf{x}_n, \mathbf{w})$  と置いた。したがって尤度関数を最大化する  $\mathbf{w}_{ML}$  は、交差エントロピー誤差：

$$E(\mathbf{w}) = -\sum_{n=1}^N \{t_n \log y_n + (1 - t_n) \log(1 - y_n)\}$$

によって与えられる。このように、解くべき問題に応じてネットワークの出力層の活性化関数および誤差関数は自然に選択される。

## 2.1 パラメータ最適化

重み空間内で重みベクトル  $\mathbf{w}$  を

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} + \Delta \mathbf{w}^{(\tau)}$$

の形で移動させていくことで、停留点：

$$\nabla E(\mathbf{w}) = 0$$

なる点を見つける。ただし、一般に誤差関数の重みパラメータへの依存性には高い非線形性があるため、重み空間内には停留点が数多く存在する。

## 2.2 局所二次近似

$E(\mathbf{w})$  を重み空間内のある点  $\tilde{\mathbf{w}}$  周りでテイラー展開する。

$$E(\mathbf{w}) \simeq E(\tilde{\mathbf{w}}) + {}^T(\mathbf{w} - \tilde{\mathbf{w}}) \frac{\partial E(\tilde{\mathbf{w}})}{\partial \mathbf{w}} + \frac{1}{2} {}^T(\mathbf{w} - \tilde{\mathbf{w}}) H(\mathbf{w} - \tilde{\mathbf{w}})$$

ただし、

$$(H)_{ij} = \frac{\partial E(\tilde{\mathbf{w}})}{\partial w_i \partial w_j}$$

なるヘッセ行列。ゆえに、勾配の局所近似は

$$\nabla E = \frac{\partial E}{\partial \mathbf{w}} \simeq \frac{\partial E(\tilde{\mathbf{w}})}{\partial \mathbf{w}} + H(\mathbf{w} - \tilde{\mathbf{w}})$$

誤差関数の極小点  $\mathbf{w}^*$  周りでの局所近似を考える。

$$\frac{\partial E(\mathbf{w}^*)}{\partial \mathbf{w}} = \mathbf{0}$$

より、

$$E(\mathbf{w}) \simeq E(\mathbf{w}^*) + \frac{1}{2} {}^T(\mathbf{w} - \mathbf{w}^*)H(\mathbf{w} - \mathbf{w}^*)$$

ヘッセ行列は対称より、その  $i$  番目の固有値を  $\lambda_i$ 、固有ベクトルを  $\mathbf{u}_i$  と置くと  $\mathbf{u}_i, \mathbf{u}_j$  ( $i \neq j$ ) は互いに直行する。 $\mathbf{u}_i$  らは重み空間上の正規直交基底より、

$$\mathbf{w} - \mathbf{w}^* = \sum_i \alpha_i \mathbf{u}_i = U\boldsymbol{\alpha}$$

の形でかける。ゆえ、

$$\begin{aligned} {}^T(\mathbf{w} - \mathbf{w}^*)H(\mathbf{w} - \mathbf{w}^*) &= {}^T\boldsymbol{\alpha} {}^T U H U \boldsymbol{\alpha} \\ &= {}^T\boldsymbol{\alpha} \Lambda \boldsymbol{\alpha} = \sum_i \lambda_i \alpha_i^2 \end{aligned}$$

なので、誤差関数は

$$E(\mathbf{w}) = E(\mathbf{w}^*) + \frac{1}{2} \sum_i \lambda_i \alpha_i^2 \quad (1)$$

という形で書ける。新しい座標系  $\{\mathbf{u}_i\}$  における  $\alpha_i$  の二次関数→誤差関数の等高線は各軸がヘッセ行列の固有ベクトル  $\mathbf{u}_i$  にそろった楕円になる。

正方行列  $A$  について、全ての  $\mathbf{v} \neq \mathbf{0}$  に対し  ${}^T\mathbf{v} A \mathbf{v} > 0$  を満たす時、 $A$  は正定値であるという。 $A$  の固有ベクトルの集合  $\{\mathbf{u}_i\}$  は完全基底をなすので、

$$\mathbf{v} = \sum_i c_i \mathbf{u}_i = U\boldsymbol{c}$$

という形で書ける。したがって、

$$\begin{aligned} {}^T\mathbf{v} A \mathbf{v} &= {}^T\boldsymbol{c} {}^T U A U \boldsymbol{c} \\ &= {}^T\boldsymbol{c} \Lambda \boldsymbol{c} \\ &= \sum_i \lambda_i c_i^2 \end{aligned}$$

が得られる。したがって、全ての固有値  $\lambda_i$  が正の時、またその時に限り行列  $A$  は正定値になる。

(1) より、点  $\mathbf{w}^*$  が極小点である条件は

$$\sum_i \lambda_i \alpha_i^2 > 0$$

すなわち  $\mathbf{w}^*$  で評価したヘッセ行列  $H$  が正定値であることである。(逆なら極大)

### 2.3 勾配情報の利用

誤差曲面は  $\frac{\partial E(\tilde{\mathbf{w}})}{\partial \mathbf{w}}$  と  $H$  で特定される。 $\mathbf{w}$  の次元を  $W$  と置くと、パラメータの数は

$$W + W + (W - 1) + \dots + 1 = \frac{W(W + 3)}{2}$$

より、二次近似の極小点の位置は  $O(W^3)$  個のパラメータに依存する。勾配情報 ( $W$  次元) を利用すれば各勾配の評価は  $O(W)$  ステップで実行可能。各表はで  $O(W)$  使うので、極小点は  $O(W^2)$  で見つかる。

## 3 誤差逆伝播

### 3.1 誤差関数微分の評価

$N$  個のデータについての誤差関数は、各データに対応する誤差項の和：

$$E(\mathbf{w}) = \sum_{n=1}^N E_n(\mathbf{w})$$

として表される。一つの項に対する勾配  $\nabla E_n(\mathbf{w})$  を考えよう。ネットワークの出力  $y_k$  が入力変数  $x_i$  の線形和

$$y_k = \sum_i w_{ki} x_i$$

で、入力パターン  $n$  に対する誤差関数が

$$\begin{aligned} E_n &= \frac{1}{2} \sum_k (y_{nk} - t_{nk})^2 \\ &= \frac{1}{2} \sum_k \{y_k(\mathbf{x}_n, \mathbf{w}) - t_{nk}\}^2 \\ &= \frac{1}{2} \sum_k (\sum_i w_{ki} x_i - t_{nk})^2 \end{aligned}$$

という形を取る場合を考える。ただし、入力ベクトル  $\mathbf{x}_n$  の  $i$  番目の要素  $(\mathbf{x}_n)_i = x_i$  と置いた。この勾配は

$$\frac{\partial E_n}{\partial w_{ji}} = (y_{nj} - t_{nj})x_i$$

である。これをより多層なニューラルネットワークへと拡張する。ユニット  $j$  の活性  $a_j$  は

$$a_j = \sum_i w_{ji} z_i$$

で書ける。ただし、 $z_i$  はユニット  $j$  と結合がある前層のユニット  $i$  からの出力で、 $w_{ji}$  はユニット  $i$  からユニット  $j$  への重みを表す。

ユニット  $j$  の出力  $z_j$  は非線形活性化関数  $h$  を用いて

$$z_j = h(a_j)$$

で表せる。 $E_n$  を  $w_{ji}$  で微分する。

$$\frac{\partial E_n}{\partial w_{ji}} = \frac{\partial E_n}{\partial a_j} \frac{\partial a_j}{\partial w_{ji}} = \delta_j z_i \quad (1)$$

ただし、

$$\delta_j = \frac{\partial E_n}{\partial a_j}$$

と置いた。 $\delta_j$  をユニット  $j$  における誤差という。(1) より、 $E_n$  の  $w_{ji}$  での微分は「ユニット  $j$  における誤差」×「ユニット  $i$  からの出力」の形で書けることがわかった。

$\delta$  について考える。多くの出力ユニット  $k$  で、

$$\delta_k = y_k - t_k$$

と書ける。例えば、出力層の活性化関数を恒等関数とする場合(回帰)、 $y_k = a_k$  であるから

$$\begin{aligned} E_n &= \frac{1}{2}(y_k - t_k)^2 \\ \therefore \frac{\partial E_n}{\partial a_k} &= y_k - t_k \end{aligned}$$

またロジスティックシグモイド関数を選ぶ場合(二クラス分類)、

$$y_k = \frac{1}{1 + \exp(-a)}$$

であるから

$$\begin{aligned} E_n &= -\{t_k \log y_k + (1-t_k) \log(1-y_k)\} \\ \therefore \frac{\partial E_n}{\partial a_k} &= -\left(\frac{t_k}{y_k} - \frac{1-t_k}{1-y_k}\right)y_k(1-y_k) = y_k - t_k \end{aligned}$$

一層後の  $\delta_k$  が分かっている時、隠れユニット  $j$  の  $\delta_j$  は

$$\begin{aligned} \delta_j &\equiv \frac{\partial E_n}{\partial a_j} = \sum_k \frac{\partial E_n}{\partial a_k} \frac{\partial a_k}{\partial a_j} \\ &= \sum_k \frac{\partial E_n}{\partial a_k} \frac{\partial a_k}{\partial a_j} \\ &= \sum_k \delta_k w_{kj} \frac{\partial z_j}{\partial a_j} \\ &= h'(a_j) \sum_k \delta_k w_{kj} \end{aligned}$$

### まとめ

- 入力データ  $n$  が得られた時、 $E_n$  の重み  $w_{ji}$  による微分は

$$\frac{\partial E_n}{\partial w_{ji}} = \frac{\partial E_n}{\partial a_j} \frac{\partial a_j}{\partial w_{ji}} = \delta_j z_j$$

と書ける。

- 出力ユニット  $k$  では

$$\delta_k = y_k - t_k$$

が成り立つ。

- 一層前の  $\delta_k$  が分かっていれば、 $\delta_j$  は

$$\delta_j = h'(a_j) \sum_k \delta_k w_{kj}$$

で得られる。したがって、全ユニットの  $w_{ji}$  による  $E_n$  の重み  $w_{ji}$  が求まるため、それを並べたベクトルが勾配  $\nabla E_n$  である。