

誤差逆伝播法

吉永 塁

2020 年 5 月 15 日

- L 層のニューラルネットワークを考える.
- 第 l 層のユニット数を $v(l)$ とおく.
- 各ユニットの入出力について,

$$z_j^{(l)} = f(u_j^{(l)}) = f\left(\sum_{i=0}^{v(l-1)} w_{ji}^{(l)} z_i^{(l-1)}\right)$$

- z は入力, u は出力を表す.
- 右肩括弧内の数字は第何層かを, 下添字は何番目のユニットかを示す.
- w は重み.
- バイアスは $w_{j0}^{(l)} = b_j^{(l)}$, $z_0^{(l)} \equiv 1$ として重みの一部とする.
- f は活性化関数.
- w : 全てのパラメータ (重み, バイアス) を成分に持つベクトル.
- $\mathbf{1}_n := \sum_{i=1}^n \mathbf{e}_i$: 全成分が 1 の n 次元ベクトル

定理 (連鎖律)

U, W がそれぞれ $\mathbb{R}^n, \mathbb{R}^m$ の開集合であり、二つの関数 $f: U \rightarrow \mathbb{R}^n$, $g: W \rightarrow \mathbb{R}^p$ が合成可能、つまり $f(U) \subset W$ であるとする。

いま、 f が $x \in U$ で微分可能で、 g が $y = f(x)$ で微分可能であるとすれば、合成関数 $\varphi = g \circ f$ は x で微分可能で、

$$\frac{\partial \varphi_r}{\partial x_j}(x) = \sum_{i=1}^m \frac{\partial g_r}{\partial y_i}(y) \frac{\partial f_i}{\partial x_j}(x)$$
$$(1 \leq r \leq p, 1 \leq j \leq n)$$

但し、添字 i は第 i 成分を表す。

証明略。

誤差逆伝播法

- 勾配降下法では重み更新において勾配 ∇E が必要。ここで誤差関数 $E(\mathbf{w})$ について,

$$\begin{cases} E(\mathbf{w}) = \sum_{n=1}^N E_n(\mathbf{w}) & (\text{バッチ学習}) \\ E(\mathbf{w}) = \frac{1}{|D_t|} \sum_{n \in D_t} E_n(\mathbf{w}) & (\text{ミニバッチ学習}) \end{cases}$$

より、何れも $\nabla E_n(\mathbf{w})$ を計算できれば良い。

- 勾配計算はそのまま行くと面倒。例えば、二乗誤差の場合,

$$\frac{\partial E_n}{\partial w_{ji}^{(l)}} = {}^t(\mathbf{y}(x_n) - \mathbf{d}_n) \frac{\partial \mathbf{y}}{\partial w_{ji}^{(l)}}$$

となり、 \mathbf{y} について、活性化関数のネストが深くなるので、連鎖律を繰り返し適用しなければいけない。

第 l 層の重み $w_{ji}^{(l)}$ について,

$$\frac{\partial E_n}{\partial w_{ji}^{(l)}} = \frac{\partial E_n}{\partial u_j^{(l)}} \frac{\partial u_j^{(l)}}{\partial w_{ji}^{(l)}} \quad (1)$$

(1) 式右辺第一項について,

$$\frac{\partial E_n}{\partial u_j^{(l)}} = \sum_{k=0}^{v^{(l+1)}} \frac{\partial E_n}{\partial u_k^{(l+1)}} \frac{\partial u_k^{(l+1)}}{\partial u_j^{(l)}} \quad (2)$$

ここで, 記号

$$\delta_j^{(l)} := \frac{\partial E_n}{\partial u_j^{(l)}} \quad (3)$$

を導入する. δ は誤差と呼ばれる.

(2) 式右辺第二項について, $u_k^{(l+1)} = \sum_i w_{ki}^{(l+1)} f(u_i^{(l)})$ に注意して,

$$\frac{\partial u_k^{(l+1)}}{\partial u_j^{(l)}} = w_{kj}^{(l+1)} f'(u_j^{(l)})$$

よって式 (2) は,

$$\delta_j^{(l)} = f'(u_j^{(l)}) \sum_{k=0}^{v(l+1)} \delta_k^{(l+1)} w_{kj}^{(l+1)} \quad (4)$$

と表される. この式は, $\delta_j^{(l)}$ が $\delta_k^{(l+1)}$ ($k = 0, \dots, v(l+1)$) により計算できることを意味する.

- 微分計算は出力層から入力層へと順に計算できる.
- 順伝播とは逆向きに伝播するので誤差逆伝播と呼ばれる.

また、逆伝播の最初の値

$$\delta_j^{(L)} = \frac{\partial E_n}{\partial u_j^{(L)}}$$

は陽に計算される.

(1) 式右辺第二項について, $u_j^{(l)} = \sum_k w_{jk}^{(l)} z_k^{(l-1)}$ に注意して,

$$\frac{\partial u_j^{(l)}}{\partial w_{ji}^{(l)}} = z_i^{(l-1)}$$

以上より, (1) 式は,

$$\frac{\partial E_n}{\partial w_{ji}^{(l)}} = \delta_j^{(l)} z_i^{(l-1)} \quad (5)$$

と表される.

誤差逆伝播法

以上より、訓練サンプル $(\mathbf{x}_n, \mathbf{d}_n)$ が与えられたとき、この訓練サンプルについての誤差関数 E_n の勾配は次の手順で得られる。

誤差逆伝播法による勾配計算

1. $\mathbf{z}^{(1)} = \mathbf{x}_n$ として各層の入出力 $\mathbf{u}^{(l)}, \mathbf{z}^{(l)} (l = 2, \dots, L)$ を計算。
このとき計算は入力層から順に行われる。[順伝播]
2. 出力層 (第 L 層) の誤差 $\delta_j^{(L)}$ を計算。
3. 中間層について、(4) 式より誤差 $\delta_j^{(l)} (l = 2, \dots, L - 1)$ を計算。
このとき計算は出力層から順に行われる。[逆伝播]
4. (5) 式より各層の重み $w_{ji}^{(l)}$ に関する微分を計算。

複数個の訓練サンプルについての誤差の総和は

$$\frac{\partial E}{\partial w_{ji}^{(l)}} = \sum_n \frac{\partial E_n}{\partial w_{ji}^{(l)}}$$

より、単に加算すればよい。

誤差逆伝播法

計算量について考える．ネットワークの重みの総数を W とおく．

このとき，順伝播+逆伝播により勾配 ∇E を求めるのには， $\mathcal{O}(W)$ の計算が必要となる．

∴ 一般的に，(重み数) \gg (ユニット数) であることに注意する．

順伝播計算の大部分は， $u_j^{(l)} = \sum_i w_{ji}^{(l)} z_i^{(l-1)}$ であるので，計算量は W による．

逆伝播計算についても， $\delta_j^{(l)} = f'(u_j^{(l)}) \sum_k \delta_k^{(l+1)} w_{kj}^{(l+1)}$ より同様．

このことより，数値微分

$$\frac{\partial E_n}{\partial w_{ji}^{(l)}} = \frac{1}{\varepsilon} (E_n(w_{ji}^{(l)} + \varepsilon) - E_n(w_{ji}^{(l)})) + \mathcal{O}(\varepsilon)$$

は，全ての重みについて $E_n(w_{ji}^{(l)} + \varepsilon)$ を計算する必要があるため，全体では $\mathcal{O}(W^2)$ となる．

数値微分は計算効率は良くないが，逆伝播の実装の確認として有用．

誤差逆伝播法

ここまでの計算を行列形式で表す。但し、バイアスと重みを区別して扱う。

- $\mathbf{X} = [\mathbf{x}_1 \cdots \mathbf{x}_N]$: \mathbf{x}_n は n 番目の訓練サンプル.
- $\mathbf{U}^{(l)} = [\mathbf{u}_1^{(l)} \cdots \mathbf{u}_N^{(l)}]$: $\mathbf{u}_n^{(l)}$ は \mathbf{x}_n についての ${}^t[u_1^{(l)} \cdots u_{v(l)}^{(l)}]$
- $\mathbf{Z}^{(l)} = [\mathbf{z}_1^{(l)} \cdots \mathbf{z}_N^{(l)}]$: $\mathbf{z}_n^{(l)}$ は $\mathbf{u}_n^{(l)}$ と同様.
- $\mathbf{W}^{(l)} = [w_{ji}^{(l)}]_{\substack{1 \leq j \leq v(l) \\ 1 \leq i \leq v(l-1)}}$: ユニット間に接続がなければ 0.
- $\mathbf{b}^{(l)} = {}^t[b_1 \cdots b_{v(l)}]$

以上より、順伝播計算は $\mathbf{U}^{(1)} = \mathbf{X}$ として,

$$\mathbf{U}^{(l)} = \mathbf{W}^{(l)} \mathbf{Z}^{(l-1)} + \mathbf{b}^{(l)} \mathbf{1}_N \quad (6)$$

$$\mathbf{Z}^{(l)} = f^{(l)}(\mathbf{U}^{(l)}) \quad (7)$$

と表される ($l = 2, \dots, L$). 但し、活性化関数 $f^{(l)}$ は各成分に適応されることとする.

逆伝播について,

- $Y = [y_1 \cdots y_N]$: y_n は x_n についての出力.
- $D = [d_1 \cdots d_N]$: d_n は x_n についての出力目標.
- $\Delta^{(l)} = [\delta_1^{(l)} \cdots \delta_N^{(l)}]$: $\delta_n^{(l)}$ は x_n についての $^t[\delta_1^{(l)} \cdots \delta_{v(l)}^{(l)}]$

とおく. 以上より, 逆伝播計算は $\Delta^{(L)} = D - Y$ として,

$$\Delta^{(l)} = f^{(l)'}(U^{(l)}) \odot ({}^tW^{(l+1)} \Delta^{(l+1)}) \quad (8)$$

と表される ($l = L - 1, \dots, 2$). 但し, \odot はアダマール積.

誤差逆伝播法

誤差関数の勾配計算について,

$$\bullet \partial \mathbf{W}^{(l)} = \left[\frac{\partial E}{\partial w_{ji}^{(l)}} \right]_{\substack{1 \leq j \leq v(l) \\ 1 \leq i \leq v(l-1)}}, \quad \partial \mathbf{b}^{(l)} = t \left[\frac{\partial E}{\partial b_1^{(l)}} \cdots \frac{\partial E}{\partial b_{v(l)}^{(l)}} \right]$$

とおく. 以上より,

$$\begin{aligned} \partial \mathbf{W}^{(l)} &= \frac{1}{N} \mathbf{\Delta}^{(l)} {}^t \mathbf{Z}^{(l-1)} \\ \partial \mathbf{b}^{(l)} &= \frac{1}{N} \mathbf{\Delta}^{(l)} \mathbf{1}_N \end{aligned}$$

となり, パラメータ更新は,

$$\mathbf{W}^{(l)} \leftarrow \mathbf{W}^{(l)} - \varepsilon \partial \mathbf{W}^{(l)} \quad (9)$$

$$\mathbf{b}^{(l)} \leftarrow \mathbf{b}^{(l)} - \varepsilon \partial \mathbf{b}^{(l)} \quad (10)$$

と表される. 重みの更新について, 重み減衰やモメンタムを用いる場合には, (9)(10) 式右辺第二項を適当に修正する.