

誤差逆伝播法

迫田真太郎

平成 32 年 2 月 28 日

1 パーセプトロンの学習則とデルタ則

ニューラルネットの起源は人工ニューロンであった。これは活性化関数が階段関数であり、入出力値が 0,1 の 2 値であるニューラルネットの一例であるとみなせる。信号が離散値であるため勾配が取れない。従ってパーセプトロンでは専用の学習則がいくつか提案されている。

中でも代表的なものは 1958 年にローゼンブラットが提唱したものである。今、中間層なしの古典パーセプトロンで出力ユニットが 1 つであるものを考えることとする。

一つの訓練サンプル $\{x, y\}$ について、入力 x に対してパーセプトロンの出力が $\hat{y}(x)$ であったとき、出力が間違っているパターンは次の 2 通りになる。

1. 目標出力が $y = 1$ であったのに実際の出力が $\hat{y}(x) = 0$ であった場合
出力層にやってくるシグナルの量が足りなかったということなので、発火したユニットとの結合パラメータを大きくすればよい。したがってパラメータの更新式は η を学習率として

$$w_i \leftarrow w_i + \eta x_i \quad (1)$$

となる。 $x_i = 0$ となるような i については w_i は不変である。

2. 目標出力が $y = 0$ であったのに実際の出力が $\hat{y}(x) = 1$ であった場合
出力層にやってくるシグナルが強すぎたということなので、発火したユニットとの結合パラメータを小さくすればよい。したがって

$$w_i \leftarrow w_i - \eta x_i \quad (2)$$

これらの更新式は次のようにまとめることができる。

$$w_i \leftarrow w_i - \eta(\hat{y}(x) - y)x_i \quad (3)$$

これをパーセプトロンの学習則といい、線形分離可能な問題に対しては訓練データを十分に与えること

でパラメータが有限回で最適値に収束していくことが証明されている。

パラメータが有限回の更新で最適値に収束することの証明

考えているパーセプトロンは $\hat{y}(x) = \text{sgn}(w^t x)$ として与えられる。 $w^t x$ の符号を見ているだけなので w をすべて正の定数倍しても出力は変わらない。線形分離可能な訓練データ $\{x_1, y_1\}, \{x_2, y_2\}, \dots, \{x_N, y_N\}$ について、簡単のため $y_i < 0$ であるものはすべて $x_i = -x_i, y_i = -y_i$ と変換しておくこととする。訓練データを繰り返し学習していく中で、誤識別された学習パターンを重複を許して順番に並べて $x_1, x_2, \dots, x_k, \dots$ とする。重みパラメータの初期値を 0 とした場合、正の定数倍に対しての不変性から学習率はどのようにとっても変わらないため 1 とし、 k 回目の誤認識をもとにパラメータを $w^k \rightarrow w^{k+1}$ へと更新する式は

$$w^{k+1} = w^k + x^k \quad (4)$$

である。ここで解となる重みベクトルの一つを \hat{w} とすると

$$\hat{w}^t x_i > 0 \quad (i = 1, 2, \dots, N) \quad (5)$$

であり、式 (??) の両辺からこれを引くと

$$(w^{k+1} - \hat{w}) = (w^k - \hat{w}) + x^k \quad (6)$$

であり、両辺を自乗して

$$(w^{k+1} - \hat{w})^2 = ((w^k - \hat{w}) + x^k)^2 \quad (7)$$

$$= (w^k - \hat{w})^2$$

$$+ 2(w^k - \hat{w})^t x^k + (x^k)^2 \quad (8)$$

$$= (w^k - \hat{w})^2 + 2(w^k)^t x^k$$

$$- 2\hat{w}^t x^k + (x^k)^2 \quad (9)$$

となる。 x^k は誤認識されたパターンなので $(w^k)^t x^k \leq 0$ が成り立ち、従って式 (??) は

$$(w^{k+1} - \hat{w})^2 \leq (w^k - \hat{w})^2 - 2\hat{w}^t x^k + (x^k)^2 \quad (10)$$

ここで \hat{w} の正の定数倍は同じく解であることを利用すると正の定数 α を用いて $\hat{w} \rightarrow \alpha \hat{w}$ と書き換えられ

$$(\mathbf{w}^{k+1} - \alpha \hat{\mathbf{w}})^2 \leq (\mathbf{w}^k - \alpha \hat{\mathbf{w}})^2 - 2\alpha \hat{\mathbf{w}}^k \mathbf{x}^k + (\mathbf{x}^k)^2 \quad (11)$$

となる. 不等式の条件を変化させないよう定数 β, γ を

$$\beta \stackrel{\text{def}}{=} \max_{p=1, \dots, n} (\mathbf{x}_p)^2 \quad (12)$$

$$\gamma \stackrel{\text{def}}{=} \min_{p=1, \dots, n} \hat{\mathbf{w}}^t \mathbf{x}_p > 0 \quad (13)$$

と定義する. すると

$$(\mathbf{w}^{k+1} - \alpha \hat{\mathbf{w}})^2 \leq (\mathbf{w}^k - \alpha \hat{\mathbf{w}})^2 - 2\alpha\gamma + \beta \quad (14)$$

であり, $\alpha = \frac{\beta}{\gamma}$ とし添え字 k 順番にずらしていくと

$$(\mathbf{w}^{k+1} - \alpha \hat{\mathbf{w}})^2 \leq (\mathbf{w}^k - \alpha \hat{\mathbf{w}})^2 - \beta \quad (15)$$

$$\leq (\mathbf{w}^{k-1} - \alpha \hat{\mathbf{w}})^2 - 2\beta \quad (16)$$

...

$$\leq (\mathbf{w}^1 - \alpha \hat{\mathbf{w}})^2 - k\beta \quad (17)$$

という不等式が成り立つ. この左辺は 0 以上なので

$$0 \leq (\mathbf{w}^1 - \alpha \hat{\mathbf{w}})^2 - k\beta \quad (18)$$

よって

$$k \leq \frac{(\mathbf{w}^1 - \alpha \hat{\mathbf{w}})^2}{\beta} \quad (19)$$

k が定数で上から抑えられるため有限回で収束することが示された.

2 誤差逆伝播法とは

3 デルタの意味

4 誤差逆伝播法の利点

5 勾配消失, 勾配爆発問題

5.1 勾配消失, 勾配爆発問題とは

5.2 勾配消失への対策

5.2.1 事前学習

5.2.2 活性化関数の工夫

参考文献

[1] 瀧雅人. これならわかる深層学習入門. p114-131

[2] 鍋谷昂一. 単純パーセプトロンの収束定理と限界.

http://starpentagon.net/analytics/simple_perceptron/