

Homework

神经网络与应用

Name: Li Shibo Student ID: 119033910046 Email: ShiboLi@sjtu.edu.cn

Problem 1 Let $X = [w^T, b]^T$, $Z_q = [P_q^T, 1]^T$, where P_q is the input. Then $n = w^T p + b = X^T Z$. Let η denote learning rate. The perceptron learning rule can now be written $X^{new} = X^{old} + e\eta Z$. We will assume that a weight vector exists that can be correctly categorize all Q input vectors and this solution will be denoted X^* . For the solution we will assume if $t_q = 1, X^{*T} Z_Q > \delta > 0$ (1) and if $t_q = 0, X^{*T} Z_Q < \delta < 0$ (2).

If we count only those iterations for which the weight vector is changed, the learning rule becomes $X(k) = X(k-1) + \eta Z'(k-1)$ (3), where $Z'(k-1)$ is the appropriate member of the set $\{Z_1, Z_2, \dots, Z_Q, -Z_1, -Z_2, \dots, -Z_Q\}$.

Assume without generality that the learning algorithm is initialized with $X(0) = 0$ (4). We can get $X(k) = \eta(Z'(0) + Z'(1) + \dots + Z'(k-1))$ (5) from (3). From (1) and (2), we can get $X^{*T} Z'(i) > \delta$ (6). From (5) and (6), we can get $X^{*T} X(k) > k\eta\delta$ (7).

Using Cauchy-Schwartz inequality on (7), $(X^{*T} X(k))^2 \leq \|X^*\|^2 \|X(k)\|^2$ (8) where $\|X\|^2 = X^T X$. With (7) and (8), $\|X(k)\|^2 \geq (X^{*T} X(k))^2 / \|X^*\|^2 > (k\eta\delta)^2 / \|X^*\|^2$ (9).

We can easily get $X(k-1)^T Z'(k-1) \leq 0$ (10) since the weights would not be updated unless the previous input vector had been misclassified. We also know $\|X(k)\|^2 = X(k)^T X(k) = [X(k-1) + \eta Z'(k-1)]^T [X(k-1) + \eta Z'(k-1)] = X(k-1)^T X(k-1) + 2\eta Z'(k-1)^T X(k-1) + \eta^2 Z'(k-1)^T Z'(k-1)$ (11). With (10), (11) can be simplified to $\|X(k)\|^2 \geq X(k-1)^T X(k-1) + \eta^2 Z'(k-1)^T Z'(k-1)$ (12). When we repeat this process, we will get $\|X(k)\|^2 \leq \eta^2 (\|Z'(0)\|^2 + \|Z'(1)\|^2 + \dots + \|Z'(k-1)\|^2)$ (13). Let $M = \max \|Z'(i)\|^2$, Then $\|X(k)\|^2 \leq k\eta M$. (14)

With (9) and (14), we can get $\eta^2 k M \geq \|X(k)\|^2 > (k\eta\delta)^2 / \|X^*\|^2$ (15). By simplify, $k < M \|X^*\|^2 / \delta^2$. Obviously, the maximum number of iterations has nothing to do with the learning rate.

Problem 2 Let C denote cost function, $z^{(l)}$ denote the input vector of l -th layer, $a^{(l)}$ denote the active vector of l -th layer.

According to the network, we can get $z^{(l)} = (a^{(l-1)})^T U^{(l)} a^{(l-1)} + V^{(l)} a^{(l-1)} + b^{(l)}$ (1) and $a^{(l)} = f(z^{(l)})$ (2).

Define cost function as $C_{(i)} = 1/m \sum_{r=1}^m \|a_r^{(L)} - y_r\|$ and $C = 1/N^{(L)} \sum_{i=1}^{N^{(L)}} C_{(i)}$, where y_r is the label of r -th sample and $N^{(L)}$ is the number of layers in the output layer.

Case 1: Neuron j is an output Node.

We define the local gradient $\delta_j(n)$ as:

$$\delta_j^{(L)} = \frac{\partial C}{\partial z_j^{(L)}} = \frac{\partial C}{\partial a_j^{(L)}} \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} = \frac{\partial C}{\partial a_j^{(L)}} \frac{\partial f(z_j^{(L)})}{\partial z_j^{(L)}} = \frac{\partial C}{\partial a_j^{(L)}} f'(z_j^{(L)}) \quad (3)$$

Case 2: Neuron j is a Hidden Node.

$$\begin{aligned} \delta_j^{(l)} &= \frac{\partial C}{\partial z_j^{(l)}} \\ &= \sum_{k=1}^{N^{(l+1)}} \frac{\partial C}{\partial z_k^{(l+1)}} \frac{\partial z_k^{(l+1)}}{\partial a_j^{(l)}} \frac{\partial a_j^{(l)}}{\partial z_j^{(l)}} \\ &= \frac{\partial a_j^{(l)}}{\partial z_j^{(l)}} \sum_{k=1}^{N^{(l+1)}} \frac{\partial C}{\partial z_k^{(l+1)}} \frac{\partial z_k^{(l+1)}}{\partial a_j^{(l)}} \\ &= f'(z_j^{(l)}) \sum_{k=1}^{N^{(l+1)}} \delta_k^{(l+1)} \frac{\partial \sum_{s=1}^{N^{(l)}} (u_{ks}^{(l+1)} (a_s^{(l)})^2 + v_{ks}^{(l+1)} a_s^{(l)}) + b_k^{(l+1)}}{\partial a_j^{(l)}} \\ &= f'(z_j^{(l)}) \sum_{k=1}^{N^{(l+1)}} \delta_k^{(l+1)} (2u_{kj}^{(l+1)} a_j^{(l)} + v_{kj}^{(l+1)}) \end{aligned} \quad (4)$$

The partial derivative for cost function:

$$\begin{aligned} \frac{\partial C}{\partial u_{jk}^{(l)}} &= \frac{\partial C}{\partial z_j^{(l)}} \frac{\partial z_j^{(l)}}{\partial u_{jk}^{(l)}} \\ &= \delta_j^{(l)} \frac{\partial \sum_{s=1}^{N^{(l)}} (u_{js}^{(l)} (a_s^{(l-1)})^2 + v_{js}^{(l)} a_s^{(l-1)}) + b_j^{(l)}}{\partial u_{jk}^{(l)}} \\ &= \delta_j^{(l)} (a_k^{(l-1)})^2 \end{aligned} \quad (5)$$

$$\begin{aligned} \frac{\partial C}{\partial v_{jk}^{(l)}} &= \frac{\partial C}{\partial z_j^{(l)}} \frac{\partial z_j^{(l)}}{\partial v_{jk}^{(l)}} \\ &= \delta_j^{(l)} \frac{\partial \sum_{s=1}^{N^{(l)}} (u_{js}^{(l)} (a_s^{(l-1)})^2 + v_{js}^{(l)} a_s^{(l-1)}) + b_j^{(l)}}{\partial v_{jk}^{(l)}} \\ &= \delta_j^{(l)} a_k^{(l-1)} \end{aligned} \quad (6)$$

$$\begin{aligned}
\frac{\partial C}{\partial b_j^{(l)}} &= \frac{\partial C}{\partial z_j^{(l)}} \frac{\partial z_j^{(l)}}{\partial b_j^{(l)}} \\
&= \delta_j^{(l)} \frac{\partial \sum_{s=1}^{N^{(l)}} (u_{js}^{(l)} (a^{(l-1)})_s^2 + v_{js}^{(l)} a_s^{(l-1)}) + b_j^{(l)}}{\partial b_j^{(l)}} \\
&= \delta_j^{(l)}
\end{aligned} \tag{7}$$

Then we can update the parameter

$$u_{jk}^{(l)} = u_{jk}^{(l)} - \eta \frac{\partial C}{\partial u_{jk}^{(l)}} \tag{8}$$

$$v_{jk}^{(l)} = v_{jk}^{(l)} - \eta \frac{\partial C}{\partial v_{jk}^{(l)}} \tag{9}$$

$$b_j^{(l)} = b_j^{(l)} - \eta \frac{\partial C}{\partial b_j^{(l)}} \tag{10}$$

Then the learning algorithm is:

The forward pass: compute $z^{(l)}, a^{(l)}$ in each layer with equation (1) and (2).

The backward pass: compute the local gradient in output layer with (3); compute the local gradient in each hidden layer with (4) and update the weight u, v and the bias b with (8),(9),(10).

In the on-line learning, the weight and bias will be updated with each sample; But the batch learning will update the weight and bias with a number of samples; in this way, the equation will be altered as:

$$u_{jk}^{(l)} = u_{jk}^{(l)} - \frac{\eta}{N_s} \sum_{i=1}^{N_s} \frac{\partial C_i}{\partial u_{ijk}^{(l)}} \tag{11}$$

$$v_{jk}^{(l)} = v_{jk}^{(l)} - \frac{\eta}{N_s} \sum_{i=1}^{N_s} \frac{\partial C_i}{\partial v_{ijk}^{(l)}} \tag{12}$$

$$b_j^{(l)} = b_j^{(l)} - \frac{\eta}{N_s} \sum_{i=1}^{N_s} \frac{\partial C_i}{\partial b_{ij}^{(l)}} \tag{13}$$

where N_s is the number of samples in each epoch.

Problem 3 I built a 3 layer perceptron network with the sigmoid active function both in hidden layer and output layer. But I found that when I setted different value of the learning rate and the number of hiddern layer the accurate rate of the result would be different. Even when I trained several times on the same number of hidden layers and learning rate, I would get different result too. I thought the reason is that the result did not converge to a global optimum; Maybe some time it

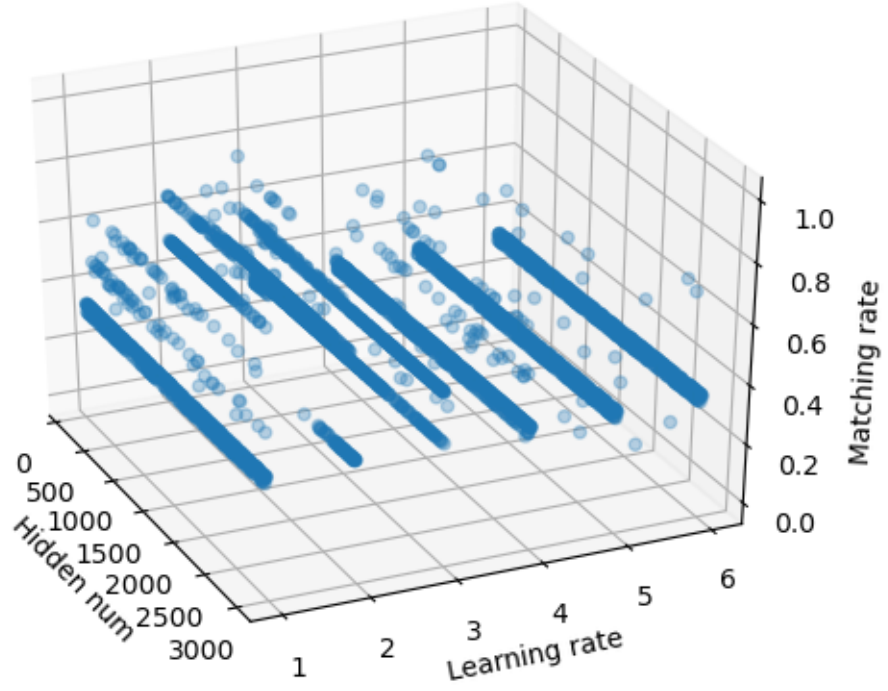


图 9: The result with different number of hidden layers is different learning rate. We represent 10^{-i} with i in the learning rate.

Obviously, with different number of hidden layers and different learning rate, we will get different result. Then I find the different result maybe caused by small number of iterations of the training process; So I set the number of iterations to 500000 rather than 500, set the number of hidden layers 100 and set the learning rate with a small number 0.00001 because the large number of iterations. Then I got a better result.

[illegible]

图 10: The number of iterations is 500000;The number of hidden layers is 100; The learning rate is 0.00001.