

# Tutorial of Back Propagation

Xingyu Jin

2016 年 2 月 21 日

## 目录

<b>1</b>	<b>Back propagation of MLP</b>	<b>1</b>
1.1	Symbols . . . . .	1
1.2	Forward Calculation . . . . .	2
1.3	Loss function . . . . .	3
1.4	Backward calculation . . . . .	3
1.5	Multiple Input . . . . .	4
1.6	Softmax Classification . . . . .	5
<b>2</b>	<b>Back propagation of CNN</b>	<b>6</b>

## 1 Back propagation of MLP

In this section, we will introduce the theoretical formulas and practical implementation of MLP.

### 1.1 Symbols

First of all, we define some variable symbols. Bold name denotes vectors.  $\mathbf{y}^i$  represents the output vector of layer  $i$  and  $y_k^i$  represents the  $k$ th output node of layer  $i$ .  $n^i$  means  $i$ th layer has  $n^i$  nodes in total except bias node.  $b^i$  is the bias node in layer  $i$ .  $\mathbf{W}^i$  represents weight parameter matrix between layer  $i$  and layer  $i+1$  and  $w_{kj}^i$  means the weight parameter between the  $k$ th output node in layer  $i$  and the  $j$ th input node in layer  $i+1$ .  $\mathbf{u}^i$

represents the input vector of layer  $i$  and  $u_k^i$  represents the  $k$ th input node of layer  $i$ .  $\phi(\cdot)$  is element-wise nonlinear function to convert input node data to output node data in one layer. The notations can be seen in the schematic diagram as shown in [1](#)

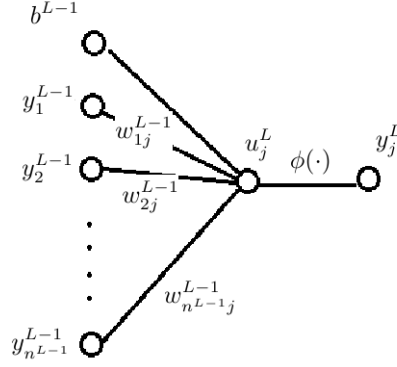


图 1: Schematic diagram for symbols

## 1.2 Forward Calculation

Suppose the network has  $L$  layers and  $\mathbf{y}^L$  is exactly the binary coding output layer.  $\mathbf{y}^{k-1}$  is a hidden layer. We can easily get:

$$\mathbf{u}^k = \mathbf{W}^{k-1} \mathbf{y}^{k-1} + \mathbf{b}^{k-1}$$

$$\mathbf{y}^k = \phi(\mathbf{u}^k)$$

1. If  $\phi(\cdot)$  is sigmoid, i.e.  $\phi(x) = \frac{1}{1+e^{-x}}$ , we have,

$$\phi'(\mathbf{u}^k) = \mathbf{y}^k \cdot (1 - \mathbf{y}^k)$$

where,  $\cdot$  represents element wise multiplication.

2. if  $\phi(\cdot)$  is relu, i.e.  $\phi(x) = \max(0, x)$ , we have,

$$\phi'(\mathbf{u}^k) = \mathbf{y}^k > 0$$

where  $>$  is element wise in MATLAB codes.

### 1.3 Loss function

We utilize the mean squared error as the loss function of the neural network.

$$E = \frac{1}{n} \sum_{i=1}^n (t_i - y_i^L)^2 \quad (1)$$

$E$  is the loss function of total training data.  $n$  is the number of training data.  $y_i^L$  represent the  $i$ th prediction when the input is the  $i$ th training data.  $t_i$  is the corresponding label.

### 1.4 Backward calculation

The error slope in the  $j$ th input node in layer  $k$  is  $\delta_j^k$

$$\delta_j^k = \frac{\partial E}{\partial u_j^k}$$

From loss function [Equation 1](#), we can deduce that:

$$\delta_j^L = \frac{\partial E}{\partial y_j^L} \frac{\partial y_j^L}{\partial u_j^L} \quad (2)$$

$$= -\frac{2}{n} (t_j - y_j^L) \phi'(u_j^L) \quad (3)$$

Then we have:

$$\frac{\partial E}{\partial w_{ij}^{L-1}} = \frac{\partial E}{\partial u_i^L} \frac{\partial u_i^L}{\partial w_{ij}^{L-1}} \quad (4)$$

$$= \delta_i^L y_j^{L-1} \quad (5)$$

We can use back propagation algorithm between layer  $k$  and layer  $k+1$  ( $k < L-1$ ):

$$\delta_i^k = \frac{\partial E}{\partial u_i^k} \quad (6)$$

$$= \sum_{j=1}^{n^{k+1}} \frac{\partial E}{\partial u_j^{k+1}} \frac{\partial u_j^{k+1}}{\partial y_i^k} \frac{\partial y_i^k}{\partial u_i^k} \quad (7)$$

$$= \sum_{j=1}^{n^{k+1}} \delta_j^{k+1} w_{ji}^k \phi'(u_i^k) \quad (8)$$

i.e.

$$\delta^k = (W^{kT} \delta^{k+1}) * \phi'(u^k)$$

where,  $*$  represents element wise multiplication.

and :

$$\frac{\partial E}{\partial w_{ij}^k} = \frac{\partial E}{\partial u_i^{k+1}} \frac{\partial u_i^{k+1}}{\partial w_{ij}^k} \quad (9)$$

$$= \delta_i^{k+1} y_j^k \quad (10)$$

$$\frac{\partial E}{\partial b_i^k} = \sum_{j=1}^{n^{k+1}} \frac{\partial E}{\partial u_j^{k+1}} \frac{\partial u_j^{k+1}}{\partial b_i^k} \quad (11)$$

$$= \delta_i^{k+1} \quad (12)$$

When we get  $\frac{\partial E}{\partial w_{ij}^k}$ , we can update weight parameters by iterations like:

$$w_{ij}^k \leq w_{ij}^k - \frac{\partial E}{\partial w_{ij}^k} \Delta w$$

Here  $\leq$  is used as value assignment in order to avoid confusion with mathematical equal.

i.e.

$$\mathbf{W}^k \leq \mathbf{W}^k - \delta^{k+1} \mathbf{y}^{kT} \Delta w$$

$$\mathbf{b}^k \leq \mathbf{b}^k - \delta^{k+1} \Delta b$$

## 1.5 Multiple Input

If the input is  $X = [x^{(1)} x^{(2)} \dots x^{(M)}]$ , a  $N \times M$  matrix, then we can compute pararely, however, the parameters should be shared, and a slight change should be made.  $x^{(m)}$  is the  $m$  th input vector.

1. Forward Part:

$$\mathbf{u}^{(m)k} = \mathbf{W}^{k-1} \mathbf{y}^{(m)k-1} + \mathbf{b}^{(m)k-1}$$

$$\mathbf{y}^{(m)k} = \phi(\mathbf{u}^{(m)k})$$

2. Loss Function:

$$E = \sum_{m=1}^M \sum_{i=1}^n (t_i^{(m)} - y_i^{(m)L})^2$$

Let  $T = [t^{(1)} t^{(2)} \dots t^{(M)}]$ ,  $Y^k = [y^{(1)k} y^{(2)k} \dots y^{(M)k}]$ , then we have,

$$E = [1 \quad 1 \dots 1](T - Y^L)^T(T - Y^L)$$

3. Backward: Let  $\Delta^k = [\delta^{(1)k} \delta^{(2)k} \dots \delta^{(M)k}]$

$$\delta_j^{(i)L} = \frac{\partial E}{\partial y_j^{(i)L}} \frac{y_j^{(i)L}}{u_j^{(i)L}}$$

Thus,

$$\delta^{(i)L} = -2(t^{(i)} - y^{(i)L}) \cdot \phi'(u^{(i)L})$$

Finally,

$$\Delta^L = -2(T - Y) \cdot \phi'(U^L)$$

$$\frac{\partial E}{\partial w_{ij}^k} = \sum_{m=1}^M \frac{\partial E}{\partial u_i^{(m)k+1}} \frac{\partial u_i^{(m)k+1}}{\partial w_{ij}^k} \quad (13)$$

$$= \sum_{m=1}^M \delta_i^{(m)k+1} y_j^{(m)k} \quad (14)$$

Thus,

$$\frac{\partial E}{\partial \mathbf{W}} = \Delta^k Y^{kT}$$

## 1.6 Softmax Classification

A softmax layer output  $\mathbf{y}$  is a  $K \times 1$  vector and  $\mathbf{x}$  is input vector.

$$y_k = P(t_k = 1) = \frac{\exp(x_k)}{\sum_{i=1}^K \exp(x_i)}$$

Cross entropy error function is defined as follows.

$$E = -\ln(p) = -\sum_{i=1}^K t_i \ln(y_i)$$

If the input of the neural network is multiple, i.e.  $\mathbf{X} = [x^{(1)}, x^{(2)}, \dots, x^{(n)}]$ , then the loss function  $E$  satisfies,

$$E = \sum_{i=1}^n E^{(i)}$$

$$E^{(i)} = - \sum_{i=1}^K t_i^{(i)} \ln(y_i^{(i)})$$

Backward propagation:

$$\frac{\partial E}{\partial x_i^{(m)}} = \frac{\partial E^{(m)}}{\partial x_i^{(m)}} \quad (15)$$

$$= \sum_{j=1}^K \frac{\partial E^{(m)}}{\partial y_j^{(m)}} \frac{\partial y_j^{(m)}}{\partial x_i^{(m)}} \quad (16)$$

$$= -\frac{t_i^{(m)}}{y_i^{(m)}} (y_i^{(m)} - y_i^{(m)^2}) + \sum_{j=i}^K \frac{t_j^{(m)}}{y_j^{(m)}} y_j^{(m)} y_i^{(m)} \quad (17)$$

$$= -t_i^{(m)} (1 - y_i^{(m)}) + \sum_{j \neq i}^K t_j^{(m)} y_i^{(m)} \quad (18)$$

$$= -t_i^{(m)} + \sum_{j=1}^K t_j^{(m)} y_i^{(m)} \quad (19)$$

$$= -t_i^{(m)} + y_i^{(m)} \quad (20)$$

## 2 Back propagation of CNN

Convolutional Neural Network is different from MLP mentioned before. The most important difference is that CNN layer contain 3 dimentions, 2 of spatial dimension, like height and width, and 1 dimension describing depth or the number of filters. The parameters in a filter is 3 dimensional, local connected on height and width, but fully connected in previous depth. The input size is  $W_1 \times H_1 \times n_{in}$ ,  $W_1$  is width,  $H_1$  is height,  $n_{in}$  is depth, usually three channel of a colorful image as input or the number of filters of previous Conv layer, and the hyperparameters of current Conv layer:

1. Number of filters  $n_{out}$

2. their spatial extent  $K$
3. the stride  $S$
4. the amount of zero padding  $P$

Each filter is  $K \times K \times n_{in}$ . And the Conv layer produces a volume of size  $W_2 \times H_2 \times D_2$ , where

1.  $W_2 = \frac{W_1 - K + 2P}{S} + 1$
2.  $H_2 = \frac{H_1 - K + 2P}{S} + 1$
3.  $D_2 = n_{out}$

With parameter sharing, the Conv layer has a total of  $k \times k \times n_{in} \times n_{out}$  weights and  $n_{out}$  biases. Let  $w_{ijd}^k$  be the  $i$ th height,  $j$ th width,  $d$ th depth,  $k$ th layer weight entry.

$$u^k = (W^{k-1} * y^{k-1})_{valid} + b^{k-1}$$

where  $*$  represents convolutional operation. i.e.

$$u_{ijm}^k = \sum_{p=1}^K \sum_{q=1}^K \sum_{d=1}^{D^{k-1}} w_{pqdm}^{k-1} y_{(p+i-1)(q+j-1)d}^{k-1} + b^{k-1}$$

$$1 \leq i \leq H^k, 1 \leq j \leq W^k, 1 \leq m \leq D^k$$

$$\delta_{pqw}^k = \frac{\partial E}{\partial u_{pqw}^k} \quad (21)$$

$$= \sum_{i=1}^{H^{k+1}} \sum_{j=1}^{W^{k+1}} \sum_{d=1}^{D^{k+1}} \frac{\partial E}{\partial u_{ijd}^{k+1}} \frac{\partial u_{ijd}^{k+1}}{\partial y_{pqw}^k} \frac{\partial y_{pqw}^k}{\partial u_{pqw}^k} \quad (22)$$

$$= \sum_{i=1}^{H^{k+1}} \sum_{j=1}^{W^{k+1}} \sum_{d=1}^{D^{k+1}} \delta_{ijd}^{k+1} w_{(p-i+1)(q-j+1)d}^k \phi'(u_{pqw}^k) \quad (23)$$

$$= \sum_{i=1}^K \sum_{j=1}^K \sum_{d=1}^{D^{k+1}} w_{ijwd}^k \delta_{(p-i+1)(q-j+1)d}^{k+1} \phi'(u_{pqw}^k) \quad (24)$$

$$(25)$$

i.e.

$$\delta^k = (\delta^{k+1} * W^k)_{full} * \phi'(u^k)$$

$$\frac{\partial E}{\partial w_{pqwn}^k} = \sum_{i=1}^{H^{k+1}} \sum_{j=1}^{W^{k+1}} \sum_{m=1}^{D^{k+1}} \frac{\partial E}{\partial u_{ijm}^{k+1}} \frac{\partial u_{ijm}^{k+1}}{\partial w_{pqwn}^k} \quad (26)$$

$$= \sum_{i=1}^{H^{k+1}} \sum_{j=1}^{W^{k+1}} \delta_{ijn}^{k+1} y_{(p+i-1)(q+j-1)w}^k \quad (27)$$

$$\frac{\partial E}{\partial b_n^k} = \sum_{i=1}^{H^{k+1}} \sum_{j=1}^{W^{k+1}} \sum_{m=1}^{D^{k+1}} \frac{\partial E}{\partial u_{ijm}^{k+1}} \frac{\partial u_{ijm}^{k+1}}{\partial b_n^k} \quad (28)$$

$$= \sum_{i=1}^{H^{k+1}} \sum_{j=1}^{W^{k+1}} \delta_{ijn}^{k+1} \quad (29)$$

im2col function can be applied like this.

im2col(Y(:, :, m, n), [K, K], 'sliding') will produce a  $(K \cdot K) \times [(H - K + 1) \cdot (W - K + 1)]$  matrix, use a for loop across all depth m and batch n could build a  $K \cdot K \cdot n_{in} \times (H - K + 1) \cdot (W - K + 1) \cdot batch\_size$  matrix A.

$$U^k = W^{k-1} * A^{k-1} + b^{k-1}$$

where  $U^k$  is  $n_{out} \times ((H - K + 1) \cdot (W - K + 1) \cdot batch\_size)$  matrix which can be reformed to  $(H - K + 1) \times (W - K + 1) \times n_{out} \times batch\_size$  output tensor