

# The Derivatives in the Fully Connected Artificial Neural Network

Yaoyu Hu, Ph.D.  
 Shanghai Jiao Tong University, Shanghai, China.  
 huyaoyu@sjtu.edu.cn  
 Jun. 23<sup>th</sup>, 2017

Two layers of neural networks are illustrated in Fig. 1.  $b_n^j$  is not listed in Fig. 1.

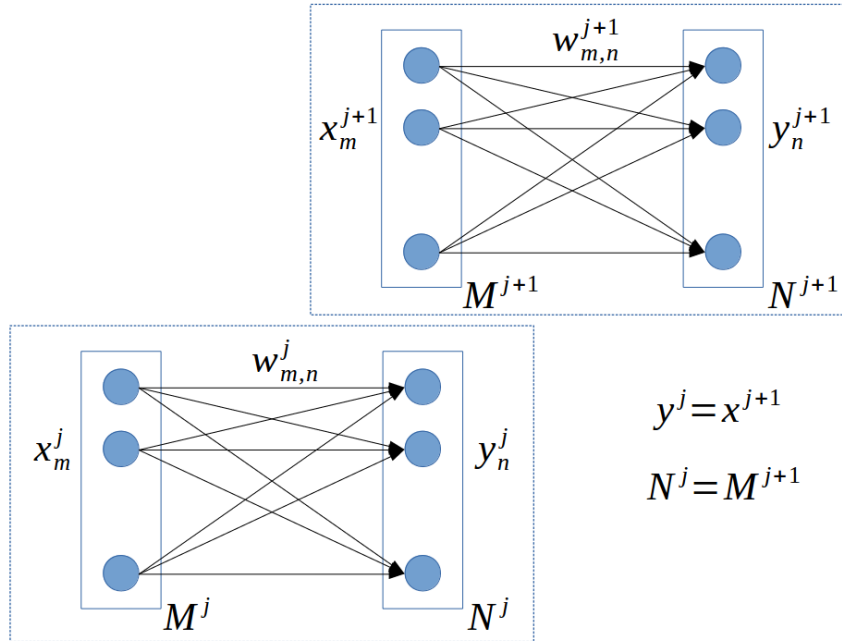


Fig. 1: Two layers of neural network.

The following relations are defined.

$$y_n^j = f\left(\sum_m^{M^j} w_{m,n}^j x_m^j + b_n^j\right) \quad (1)$$

where  $f()$  is the activation function.

$$f(x) = \begin{cases} x & x \geq 0 \\ 0 & x < 0 \end{cases} \quad (2)$$

If we define

$$g_n^j = \sum_m^{M^j} w_{m,n}^j x_m^j + b_n^j \quad (3)$$

Then we have

$$y_n^j = f(g_n^j) \quad (4)$$

Let  $L$  be the loss function,  $J$  be the last index of neural network layer and  $\mathbf{Y}$  be the training value of the neural network. Taking the sum-of-squares loss function as an example.

$$L = \frac{1}{2} \sum_n^{N^j} (y_n^j - Y_n)^2 \quad (5)$$

With  $j = J$

$$\frac{\partial L}{\partial y_n^J} = (y_n^J - Y_n) \quad (6)$$

$$\frac{\partial L}{\partial x_m^J} = \sum_n^{N^j} \frac{\partial L}{\partial y_n^J} \frac{\partial y_n^J}{\partial g_n^J} \frac{\partial g_n^J}{\partial x_m^J} = \sum_n^{N^j} \frac{\partial L}{\partial y_n^J} \frac{\partial y_n^J}{\partial g_n^J} w_{m,n}^J \quad (7)$$

$$\frac{\partial L}{\partial w_{m,n}^J} = \frac{\partial L}{\partial y_n^J} \frac{\partial y_n^J}{\partial g_n^J} \frac{\partial g_n^J}{\partial w_{m,n}^J} = \frac{\partial L}{\partial y_n^J} \frac{\partial y_n^J}{\partial g_n^J} x_{m,n}^J \quad (8)$$

$$\frac{\partial L}{\partial b_n^J} = \frac{\partial L}{\partial y_n^J} \frac{\partial y_n^J}{\partial g_n^J} \frac{\partial g_n^J}{\partial b_n^J} = \frac{\partial L}{\partial y_n^J} \frac{\partial y_n^J}{\partial g_n^J} \quad (9)$$

With arbitrary  $j \leq J - 1$

$$\frac{\partial L}{\partial y_n^j} = \frac{\partial L}{\partial x_n^{j+1}} \quad (10)$$

$$\frac{\partial L}{\partial x_m^j} = \sum_n^{N^j} \frac{\partial L}{\partial y_n^j} \frac{\partial y_n^j}{\partial g_n^j} \frac{\partial g_n^j}{\partial x_m^j} = \sum_n^{N^j} \frac{\partial L}{\partial y_n^j} \frac{\partial y_n^j}{\partial g_n^j} w_{m,n}^j \quad (11)$$

$$\frac{\partial L}{\partial w_{m,n}^j} = \frac{\partial L}{\partial y_n^j} \frac{\partial y_n^j}{\partial g_n^j} \frac{\partial g_n^j}{\partial w_{m,n}^j} = \frac{\partial L}{\partial y_n^j} \frac{\partial y_n^j}{\partial g_n^j} x_{m,n}^j \quad (12)$$

$$\frac{\partial L}{\partial b_n^j} = \frac{\partial L}{\partial y_n^j} \frac{\partial y_n^j}{\partial g_n^j} \frac{\partial g_n^j}{\partial b_n^j} = \frac{\partial L}{\partial y_n^j} \frac{\partial y_n^j}{\partial g_n^j} \quad (13)$$

Let the partial derivative of  $L$  respect to  $y^j$  be in the matrix form of

$$\left\{ \frac{\partial L}{\partial y_n} \right\}_{N^j \times 1}^j \quad (14)$$

where  $N^j \times 1$  is the dimension specification. Similarly, we could write a matrix expression related to  $\partial y_n / \partial g_n$

$$\left[ \frac{\partial y_n}{\partial g_n} \right]_{N \times N}^j = \begin{bmatrix} \frac{\partial y_1}{\partial g_1} & 0 & \dots & 0 \\ 0 & \frac{\partial y_2}{\partial g_2} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \frac{\partial y_N}{\partial g_N} \end{bmatrix}_{N^j \times N^j}^j \quad (15)$$

We define four parameter matrices

$$\{y_n\}_{N^j \times 1}^j \quad (16)$$

$$\{Y_n\}_{N^j \times 1}^J \quad (17)$$

$$\{x_m\}_{M^j \times 1}^j \quad (18)$$

$$\{w_{m,n}\}_{M^j \times N^j}^j \quad (19)$$

Then we have

$$\left\{ \frac{\partial L}{\partial x_m} \right\}_{M^j \times 1}^j = \left\{ \sum_n^{N^j} \frac{\partial L}{\partial y_n^j} \frac{\partial y_n^j}{\partial g_n^j} w_{m,n}^j \right\}_{M^j \times 1}^j = \{w_{m,n}\}_{M^j \times N^j}^j \left[ \frac{\partial y_n}{\partial g_n} \right]_{N^j \times N^j}^j \left\{ \frac{\partial L}{\partial y_n} \right\}_{N^j \times 1}^j \quad (20)$$

$$\left[ \frac{\partial L}{\partial w_{m,n}} \right]_{M^j \times N^j}^j = \left[ \frac{\partial L}{\partial y_n^j} \frac{\partial y_n^j}{\partial g_n^j} x_m^j \right]_{M^j \times N^j}^j = \{x_m\}_{M^j \times 1}^j \left[ \left[ \frac{\partial y_n}{\partial g_n} \right]_{N^j \times N^j}^j \left\{ \frac{\partial L}{\partial y_n} \right\}_{N^j \times 1}^j \right]^T \quad (21)$$

$$\left\{ \frac{\partial L}{\partial b_n} \right\}_{N^j \times 1}^j = \left\{ \frac{\partial L}{\partial y_n} \frac{\partial y_n}{\partial g_n} \right\}_{N \times 1}^j = \left[ \frac{\partial y_n}{\partial g_n} \right]_{N^j \times N^j}^j \left\{ \frac{\partial L}{\partial y_n} \right\}_{N^j \times 1}^j \quad (22)$$

And it is obvious that

$$\left\{ \frac{\partial L}{\partial y_n} \right\}_{N^j \times 1}^J = \{y_n - Y_n\}_{N^j \times 1}^J \quad (23)$$

Note that in Eq. (23)  $j = J$ . We could conveniently define the following column vector

$$\left\{ \frac{\partial L}{\partial \mathbf{g}_n} \right\}_{N^j \times 1}^j = \left[ \frac{\partial y_n}{\partial \mathbf{g}_n} \right]_{N^j \times N^j}^j \left\{ \frac{\partial L}{\partial y_n} \right\}_{N^j \times 1}^j = \left\{ \frac{\partial L}{\partial y_n} \frac{\partial y_n}{\partial \mathbf{g}_n} \right\}_{N^j \times 1}^j \quad (24)$$

Then

$$\left\{ \frac{\partial L}{\partial \mathbf{x}_m} \right\}_{M^j \times 1}^j = [\mathbf{w}_{m,n}]_{M^j \times N^j}^j \left\{ \frac{\partial L}{\partial \mathbf{g}_n} \right\}_{N^j \times 1}^j \quad (25)$$

$$\left[ \frac{\partial L}{\partial \mathbf{w}_{m,n}} \right]_{M^j \times N^j}^j = [\mathbf{x}_m]_{M^j \times 1}^j \left[ \left\{ \frac{\partial L}{\partial \mathbf{g}_n} \right\}_{N^j \times 1}^j \right]^T \quad (26)$$

$$\left\{ \frac{\partial L}{\partial \mathbf{b}_n} \right\}_{N^j \times 1}^j = \left\{ \frac{\partial L}{\partial \mathbf{g}_n} \right\}_{N^j \times 1}^j \quad (27)$$

It turns out that if we take the terms marked by \* in Eq. (12) to (14) as whole column vectors, then the results of Eq. (25) to (26) should fall out naturally.

Pseudo-code of Backscatter Propagation

*Table 1. Pseudo-code for calculating the gradients.*

<p>Calculate Eq. (23) with <math>j = J</math>.  Initialize <math>\partial L / \partial y_n^j</math>  Start from <math>j = J</math> loop until <math>j &lt; 0</math>:      Evaluate Eq. (24) based on current <math>\partial L / \partial y_n^j</math>.      Use Eq. (26) and Eq. (27) to obtain the gradients.      if <math>j \neq 0</math>:          Use Eq. (25) to get <math>\partial L / \partial y_n^{j-1}</math> for the previous pair of neural networks.      <math>j = j - 1</math></p>
---

For cross entropy loss function with Softmax function, the Softmax function is defined as

$$z_n = \frac{e^{y_n^j}}{\sum_n e^{y_n^j}} \quad (28)$$

Then the cross entropy loss function is expressed as

$$L = - \sum_n Y_n^J \ln(z_n) \quad (29)$$

In order to evaluate  $\partial L / \partial y_n^J$  the derivative of the Softmax function should be obtained.

$$\frac{\partial z_n}{\partial y_i^J} = \begin{cases} z_i - z_i^2 & n=i \\ -z_n z_i & n \neq i \end{cases} \quad (30)$$

The partial derivative of  $L$  respect to  $y_n^J$  is

$$\begin{aligned} \frac{\partial L}{\partial y_n^J} &= -Y_n \frac{1}{z_n} (z_n - z_n^2) + \sum_{i \neq n} Y_i z_n \\ &= \sum_i Y_i z_n - Y_n \\ &= z_n - Y_n \end{aligned} \quad (31)$$

with the fact that

$$\sum_i Y_i = 1 \quad (32)$$

Because in the context of classification, the vector  $\mathbf{Y}$  is represent in such a way that only one of its components is 1 and all other components are 0.