

Chain Rule of Neural Network is Error Back Propagation

Masayuki Tanaka

2013/11/20

1 Introduction

When I studied the convolution network, I could not figure out how to impliment an error backpropagation for the Lp pooling. I had found there was a few information about that. Through a discussion with Hiroshi Kuwajima, I finally figure out the error backpropagation is a simple chain rule of the derivatives. This framework is very general and easy to understand, but again there was a few explanatoin of the error backpropagation as the simple chain rule.

2 Notations

Here, I define some notations.

The elementwise sigmoid function is defined by

$$\mathbf{y} = \boldsymbol{\sigma}(\mathbf{x}) = \begin{bmatrix} \sigma(x_1) \\ \sigma(x_2) \\ \vdots \\ \sigma(x_N) \end{bmatrix}, \quad (1)$$

where

$$\sigma(x) = \frac{1}{1 + e^{-x}}. \quad (2)$$

I use the symbol, \circ , for the compositional operation of functions. I also introduce the brackets, $\llbracket \cdot \rrbracket$, to specify the parameters differentiating from the arguments of the function.

$$\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}) = \mathbf{f}[\boldsymbol{\theta}](\mathbf{x}), \quad (3)$$

$$\mathbf{g}(\mathbf{f}(\mathbf{x})) = \mathbf{g} \circ \mathbf{f} \circ \mathbf{x}, \quad (4)$$

$$\mathbf{g}[\boldsymbol{\tau}](\mathbf{f}[\boldsymbol{\theta}](\mathbf{x})) = \mathbf{g}[\boldsymbol{\tau}] \circ \mathbf{f}[\boldsymbol{\theta}] \circ \mathbf{x}. \quad (5)$$

The one-layer neural network with a sigmoid activation function can be expressed by

$$\boldsymbol{\sigma}(\mathbf{W}\mathbf{x}) = \boldsymbol{\sigma} \circ \boldsymbol{\rho}[\mathbf{W}] \circ \mathbf{x}, \quad (6)$$

where

$$\boldsymbol{\rho}[\mathbf{W}](\mathbf{x}) = \boldsymbol{\rho}[\mathbf{W}] \circ \mathbf{x} = \mathbf{W}\mathbf{x}, \quad (7)$$

where $\boldsymbol{\rho}[\mathbf{W}]$ represents the weighted sum or the linear transformation, \mathbf{W} is the weight or the parameters, and \mathbf{x} is the input of the neural network. Here, I ommit to put the bias term to simplify the expression.

The derivatives of the vectors are defined as

$$\frac{\partial \mathbf{f}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \cdots & \frac{\partial f_1}{\partial x_M} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \cdots & \frac{\partial f_2}{\partial x_M} \\ \vdots & \vdots & & \vdots \\ \frac{\partial f_N}{\partial x_1} & \frac{\partial f_N}{\partial x_2} & \cdots & \frac{\partial f_N}{\partial x_M} \end{bmatrix} \quad (8)$$

3 Chain rule of the Derivatives

The chain rules of the derivatives can be expressed as

$$\frac{\partial}{\partial \mathbf{x}} (L \circ \mathbf{g}[\boldsymbol{\tau}] \circ \mathbf{f}[\boldsymbol{\theta}] \circ \mathbf{x}) = \frac{\partial L}{\partial \mathbf{g}} \frac{\partial \mathbf{g}}{\partial \mathbf{f}} \frac{\partial \mathbf{f}}{\partial \mathbf{x}}. \tag{9}$$

The chain rule with respect to the parameter can be derived very similarly as

$$\frac{\partial}{\partial \boldsymbol{\theta}} (L \circ \mathbf{g}[\boldsymbol{\tau}] \circ \mathbf{f}[\boldsymbol{\theta}] \circ \mathbf{x}) = \frac{\partial L}{\partial \mathbf{g}} \frac{\partial \mathbf{g}}{\partial \mathbf{f}} \frac{\partial \mathbf{f}}{\partial \boldsymbol{\theta}}. \tag{10}$$

4 Chain Rule of Neural Network (Error Back Propagation)

4.1 Sigmoid Activation

The one-layer neural network can be expressed by the compositional form of the sigmoid activation and the linear transformation. Namely,

$$\mathbf{y} = \boldsymbol{\sigma}(\mathbf{W}\mathbf{x}) = \boldsymbol{\sigma}(\rho[\mathbf{W}](\mathbf{x})) = \boldsymbol{\sigma} \circ \rho[\mathbf{W}] \circ \mathbf{x}, \tag{11}$$

where

$$\rho[\mathbf{W}] = \mathbf{W}\mathbf{x}. \tag{12}$$

Here, I omit the bias term for the simplification.

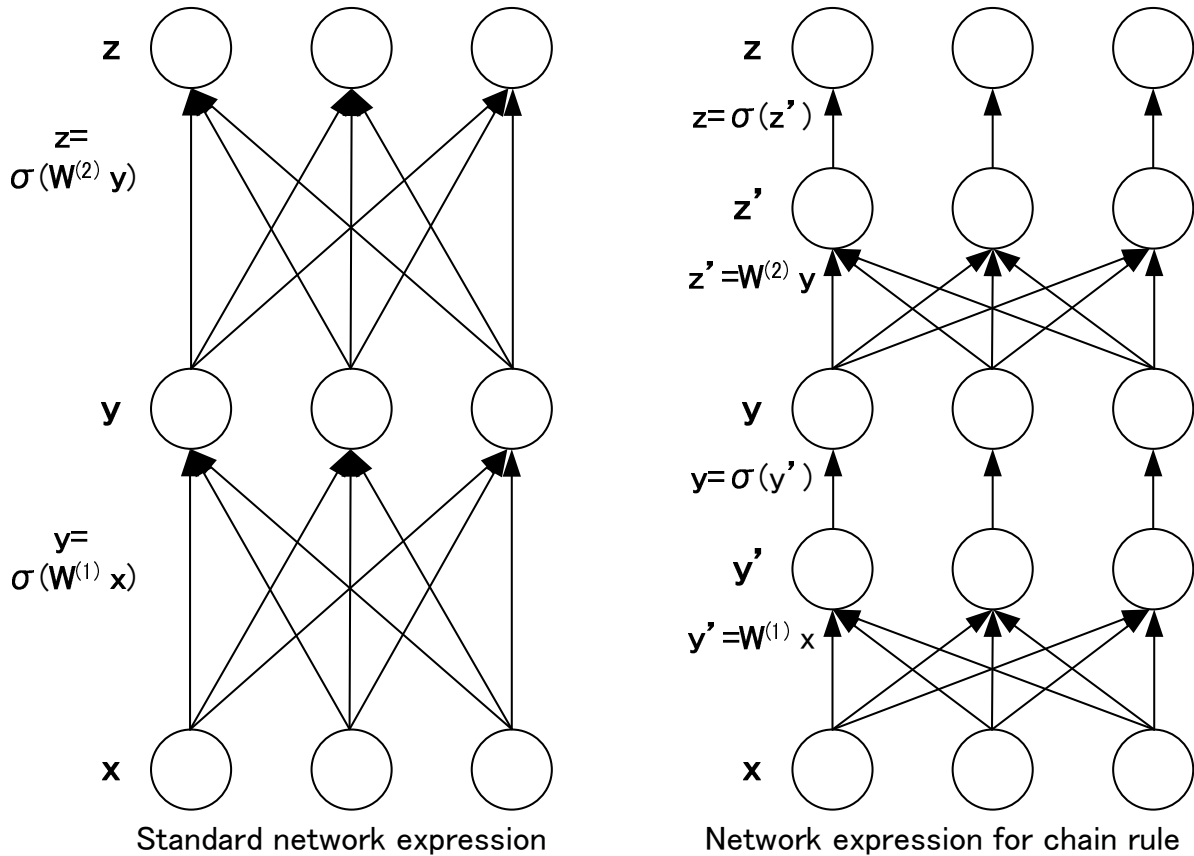


Figure 1: The two expressions of the two-layer neural network.

Using these notation, the two-layer neural network as shown in Fig. 1 can be expressed by

$$\mathbf{z} = \boldsymbol{\sigma} \circ \rho[\mathbf{W}^{(2)}] \circ \boldsymbol{\sigma} \circ \rho[\mathbf{W}^{(1)}] \circ \mathbf{x}, \tag{13}$$

where hidden variables are

$$z = \sigma \circ z', \quad (14)$$

$$z' = \rho[\mathbf{W}^{(2)}] \circ \mathbf{y}, \quad (15)$$

$$\mathbf{y} = \sigma \circ \mathbf{y}', \quad (16)$$

$$\mathbf{y}' = \rho[\mathbf{W}^{(1)}] \circ \mathbf{x}. \quad (17)$$

The cost function is defined like the square error. Namely,

$$L = \sum_k L_k \text{ where, } L_k = \|\mathbf{t}_k - \mathbf{z}_k(\mathbf{x}_k)\|_2^2 \quad (18)$$

where \mathbf{t}_k is the teach data. The training of the network is performed by minimizing the cost function with respect to parameters, $\{\mathbf{W}^{(1)}, \mathbf{W}^{(2)}\}$. The chain rule is independent to the cost function. One can consider any kind of the cost function such as a cross entropy.

Now, it is ready to derive the derivatives of L_k with respect to the parameters by using the chain rule. Once one can derive the derivatives of L_k , one can obtain the derivatives of L by simply summing up.

$$\frac{\partial L_k}{\partial \mathbf{W}^{(2)}} = \frac{\partial L_k}{\partial z} \frac{\partial z}{\partial z'} \frac{\partial z'}{\partial \mathbf{W}^{(2)}} = \delta_z \frac{\partial z'}{\partial \mathbf{W}^{(2)}} \quad (19)$$

$$\frac{\partial L_k}{\partial \mathbf{W}^{(1)}} = \frac{\partial L_k}{\partial z} \frac{\partial z}{\partial z'} \frac{\partial z'}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{y}'} \frac{\partial \mathbf{y}'}{\partial \mathbf{W}^{(1)}} = \delta_{\mathbf{y}'} \frac{\partial \mathbf{y}'}{\partial \mathbf{W}^{(1)}} \quad (20)$$

$$\delta_{z'} = \frac{\partial L_k}{\partial z} \frac{\partial z}{\partial z'} \quad (21)$$

$$\delta_{\mathbf{y}'} = \frac{\partial L_k}{\partial z} \frac{\partial z}{\partial z'} \frac{\partial z'}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{y}'} \quad (22)$$

The "error" of $\delta_{z'}$ and $\delta_{\mathbf{y}'}$ can be propagated layer-by-layer backwardly.

$$\delta_z = \frac{\partial L_k}{\partial z} \quad (23)$$

$$\delta_{z'} = \frac{\partial L_k}{\partial z} = \delta_z \frac{\partial z}{\partial z'} = \frac{\partial L_k}{\partial z} \frac{\partial z}{\partial z'} \quad (24)$$

$$\delta_{\mathbf{y}} = \frac{\partial L_k}{\partial \mathbf{y}} = \delta_{z'} \frac{\partial z'}{\partial \mathbf{y}} = \frac{\partial L_k}{\partial z} \frac{\partial z}{\partial z'} \frac{\partial z'}{\partial \mathbf{y}} \quad (25)$$

$$\delta_{\mathbf{y}'} = \frac{\partial L_k}{\partial \mathbf{y}'} = \delta_{\mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{y}'} = \frac{\partial L_k}{\partial z} \frac{\partial z}{\partial z'} \frac{\partial z'}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{y}'} \quad (26)$$

This chain of calculations are called an error back propagation.

Each derivatives are calculated as follows:

$$\frac{\partial z}{\partial z'} = \mathbf{z}' \otimes (\mathbf{1} - \mathbf{z}') \quad (27)$$

$$\frac{\partial z'}{\partial \mathbf{y}} = \mathbf{W}^{(2)} \quad (28)$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{y}'} = \mathbf{y}' \otimes (\mathbf{1} - \mathbf{y}') \quad (29)$$

$$\frac{\partial z'}{\partial \mathbf{W}^{(2)}} = \mathbf{y}^T \quad (30)$$

$$\frac{\partial \mathbf{y}'}{\partial \mathbf{W}^{(1)}} = \mathbf{x}^T \quad (31)$$

where \otimes is the elementwise multiplication, $\mathbf{1}$ is the column vector whose element is one, and \mathbf{x}^T represents the transpose vector of the vector \mathbf{x} .

4.2 Lp Pooling

A lp pooling is a technique of subsampling which is usually used in a convolutional network. The lp pooling and its derivatives can be expressed by

$$z = \eta_p(\mathbf{y}) = \left(\sum_i |y_i|^p \right)^{1/p}, \quad (32)$$

$$\frac{\partial z}{\partial y_i} = \left(\sum_i |y_i|^p \right)^{1/p-1} |y_i|^{p-1}. \quad (33)$$

Here, I derive the error back propagation, or the chain rule, of the neural network which includes the lp pooling as shown in Fig. 2.

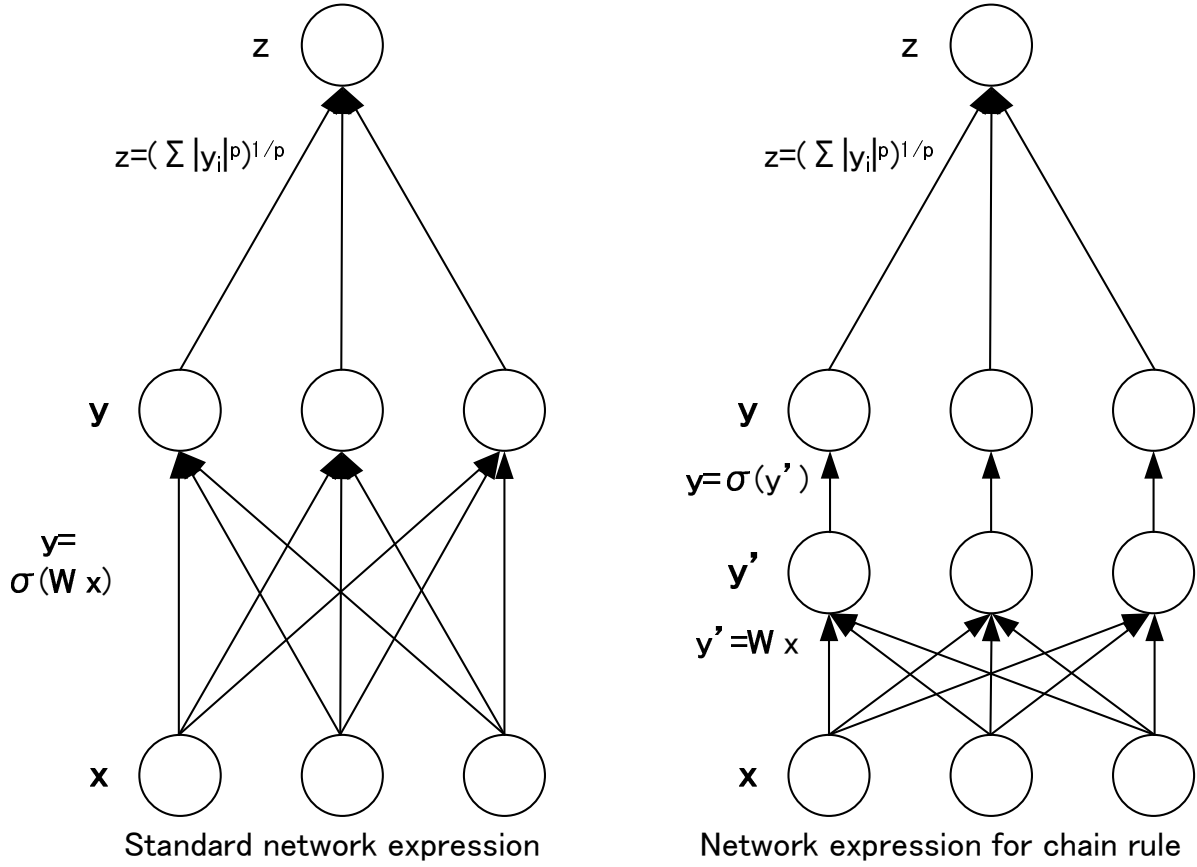


Figure 2: The two expressions of the neural network which includes the lp pooling.

The neural network can be expressed by

$$z = \eta_p \circ \sigma \circ \rho[\mathbf{W}] \circ x \quad (34)$$

The hidden variables are

$$z = \eta_p \circ \mathbf{y}, \quad (35)$$

$$\mathbf{y} = \sigma \circ \mathbf{y}', \quad (36)$$

$$\mathbf{y}' = \rho[\mathbf{W}] \circ x. \quad (37)$$

The error back propagation described in Section 4.1 is general. One can apply to the lp pooling as well.

$$\delta_z = \frac{\partial L_k}{\partial z} \quad (38)$$

$$\delta_{\mathbf{y}} = \frac{\partial L_k}{\partial \mathbf{y}} = \delta_z \frac{\partial z}{\partial \mathbf{y}} = \frac{\partial L_k}{\partial z} \frac{\partial z}{\partial \mathbf{y}} \quad (39)$$

$$\delta_{\mathbf{y}'} = \frac{\partial L_k}{\partial \mathbf{y}'} = \delta_{\mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{y}'} = \frac{\partial L_k}{\partial z} \frac{\partial z}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{y}'} \quad (40)$$

Then, one can calculate the derivative of the cost function with respect to the parameters \mathbf{W} as

$$\frac{\partial L_k}{\partial \mathbf{W}} = \delta_{\mathbf{y}'} \frac{\partial \mathbf{y}'}{\partial \mathbf{W}} = \delta_{\mathbf{y}'} \mathbf{x}^T. \quad (41)$$

Acknowledgement

Thank you for good discussion with Hiroshi Kuwajima.

References

- [1] UFLDL Tutorial, http://ufldl.stanford.edu/wiki/index.php/UFLDL_Tutorial