

# Note on Backpropagation in Neural Networks

Xin Jin  
felixxinjin@gmail.com

August 8, 2019

## Abstract

This note intends to facilitate low level implementation by providing an analytical perspective on neural networks. Different feedforward and recurrent neural networks are dissected through a derivation of the backpropagation update. We choose Multilayer Perceptron (MLP) which possesses the basic architecture of deep artificial neural network as a departure of introduction. Sigmoid Cross-Entropy loss is applied to MLP for an exemplification of multi-label classification. We then turn to introduce Convolutional Neural Network (CNN) — an intricate architecture which adopts filter and sub-sampling to realize a form of regularization. In the end, we illustrate Backpropagation Through Time (BPTT) to elicit Exploding / Vanishing Gradients problem and Long short-term memory (LSTM).

## 1 Fully Connected Neurons: Multilayer Perceptron

### 1.1 Structure of an elemental MLP

Multilayer Perceptron (MLP) has a structure of fully connected neurons. Neuron-like processing unit is the basic element of MLP:

$$a = \phi \left( \sum_i w_i x_i + b \right) \quad (1)$$

where  $x_i$  are input to the neuron, the  $w_i$  are the weights,  $b$  is the bias,  $\phi$  is the activation function, and  $a$  is the unit's activation.

For simplicity but without loss of generality, we consider an elemental MLP with only 3 layers:

The 1st Layer (Input Layer) is composed of  $M$  neurons with no connection to each other. The activation-input relation in this layer is the intact transit of input and bias:

$$\left[ \begin{array}{ccc} \vec{a}_1^{(in)} & \dots & \vec{a}_N^{(in)} \end{array} \right] = \left[ \begin{array}{ccc} 1 \dots 1 \\ a_{1,1}^{(in)} \dots a_{1,N}^{(in)} \\ \vdots \\ a_{M,1}^{(in)} \dots a_{M,N}^{(in)} \end{array} \right] = \left[ \begin{array}{ccc} 1 \dots 1 \\ x_{1,1}^{(in)} \dots x_{1,N}^{(in)} \\ \vdots \\ x_{M,1}^{(in)} \dots x_{M,N}^{(in)} \end{array} \right] \quad (2)$$

In practical use, the  $M$  neurons are playing a role of dispatching  $M$  features  $(x_{1,i}^{(in)} \dots x_{M,i}^{(in)})$  of the  $i$ th sample  $x_i$ ,  $i \in [1, \dots, N]$  respectively to the next layer.

The 2nd Layer (Hidden Layer) has  $K$  neurons each processes  $M$  activation of Input Layer as input:

$$\left[ \begin{array}{ccc} \vec{a}_1^{(h)} & \dots & \vec{a}_N^{(h)} \end{array} \right] = \left[ \begin{array}{ccc} 1 \dots 1 \\ \phi \left( \vec{w}_1^{(h)} \vec{a}_1^{(in)} \right) \dots \phi \left( \vec{w}_1^{(h)} \vec{a}_N^{(in)} \right) \\ \vdots \\ \phi \left( \vec{w}_K^{(h)} \vec{a}_1^{(in)} \right) \dots \phi \left( \vec{w}_K^{(h)} \vec{a}_N^{(in)} \right) \end{array} \right] \quad (3)$$

where we choose sigmoid function  $\phi(z) = \left( \frac{1}{1+e^{-z}} \right)$  as activation function in Hidden Layer.

The 3rd Layer (Output Layer) possesses  $C$  neurons, where  $C$  is the number of classes. The activation of the  $c$ th neuron represents the predicted conditional probability that the  $i$ th input sample belongs to class  $c$ :

$$\left[ \begin{array}{ccc} \vec{a}_1^{(out)} & \dots & \vec{a}_N^{(out)} \end{array} \right] = \left[ \begin{array}{ccc} \phi \left( \vec{w}_1^{(out)} \vec{a}_1^{(h)} \right) \dots \phi \left( \vec{w}_1^{(out)} \vec{a}_N^{(h)} \right) \\ \vdots \\ \phi \left( \vec{w}_c^{(out)} \vec{a}_1^{(h)} \right) \dots \phi \left( \vec{w}_c^{(out)} \vec{a}_N^{(h)} \right) \\ \vdots \\ \phi \left( \vec{w}_C^{(out)} \vec{a}_1^{(h)} \right) \dots \phi \left( \vec{w}_C^{(out)} \vec{a}_N^{(h)} \right) \end{array} \right] \quad (4)$$

where  $a_{c,i}^{(out)} = \phi \left( \vec{w}_c^{(out)} \vec{a}_i^{(h)} + b_c \right) = \hat{P}(y_i = c | x_i)$ ,  $i \in [1, \dots, N]$ ,  $c \in [1, \dots, C]$  and  $\phi(z) = \left( \frac{1}{1+e^{-z}} \right)$ .

## 1.2 Multi-Label MLP and Sigmoid Cross-Entropy loss

Multi-Label MLP is able to make prediction for Multi-Label Classification, where the sample  $x_i$  can belong to more than one class. If  $x_i$  belonging to a certain class doesn't influence the decision for another class (i.e.  $P(y_i = c | x_i)$ ,  $c \in [1 \dots C]$  are independent), Multi-Label Classification problem can be split into  $C$  binary Classification problems. The solution of  $C$  binary Classification problems is to learn weights  $\mathbf{w}$  by maximizing the likelihood  $L(\mathbf{w})$  below:

$$L(\mathbf{w}) = \prod_{c=1}^C \prod_{i=1}^n \left( a_{c,i}^{(out)} \right)^{P(y_i=c|x_i)} \left( 1 - a_{c,i}^{(out)} \right)^{1-P(y_i=c|x_i)} \quad (5)$$

In practice, the equivalent approach of minimizing Sigmoid Cross-Entropy Loss  $J(\mathbf{w})$  makes calculation convenient:

$$J(\mathbf{w}) = - \sum_{c=1}^C \sum_{i=1}^n \left[ P(y_i = c | x_i) \log \left( a_{c,i}^{(out)} \right) + (1 - P(y_i = c | x_i)) \log \left( 1 - a_{c,i}^{(out)} \right) \right] \quad (6)$$

### 1.3 Training MLP via Backpropagation

In order to update the weights through Backpropagation, we first apply forward propagation to generate the activation of the output layer through (2),(3), and (4) with initial weight values. We then use gradient descent to update the weights in each layer:

$$\frac{\partial J(\mathbf{w})}{\partial w_{c,k}^{(out)}} = \left( a_{c,i}^{(out)} - P(y_i = c | x_i) \right) a_{k,i}^{(h)} \quad (7)$$

$$\frac{\partial J(\mathbf{w})}{\partial w_{k,m}^{(h)}} = \left( a_{c,i}^{(out)} - P(y_i = c | x_i) \right) w_{c,k}^{(out)} a_{k,i}^{(h)} \left( 1 - a_{k,i}^{(h)} \right) a_{m,i}^{(in)} \quad (8)$$

$$w_{c,k}^{(out)} = w_{c,k}^{(out)} - \eta \frac{\partial J(\mathbf{w})}{\partial w_{c,k}^{(out)}} \quad (9)$$

$$w_{k,m}^{(h)} = w_{k,m}^{(h)} - \eta \frac{\partial J(\mathbf{w})}{\partial w_{k,m}^{(h)}} \quad (10)$$

Where  $\eta$  is the learning rate. The derivation of Equation (7) and (8) is detailed in Appendix.

## 2 Convolutional Neural Networks

Convolutional Neural Networks (CNNs) are composed of multiple Convolutional Layer / Pooling Layer pair followed by Fully Connection Layer at the end [1, 2]. We exhibit below a basic CNN with a single Convolutional Layer / Pooling Layer pair with a Fully Connection Layer.

## 2.1 Convolution Layer

Convolution Layer can be imagined as a variant of Hidden Layer in which the inner product of the weight and the input is replaced by convolution:

$$A_{c_{out}} = \phi \left( X * W_{c_{out}}^{(conv)} + b_{c_{out}}^{(conv)} \right) \quad (11)$$

$$A_{c_{out}}(i, j) = \phi \left( \sum_{k_1=0}^{m_1-1} \sum_{k_2=0}^{m_2-1} X(i - k_1, j - k_2) W_{c_{out}}^{(conv)}(k_1, k_2) + b_{c_{out}}^{(conv)} \right) \quad (12)$$

where  $X \in R^{n_1 \times n_2}$  is the input matrix,  $W_{c_{out}}^{(conv)} \in R^{m_1 \times m_2}$ ,  $c_{out} \in \{1, \dots, C_{out}\}$  is the kernel matrix for the  $c_{out}$ th channel,  $b_{c_{out}}^{(conv)}$  is the bias for the  $c_{out}$ th channel, and  $\phi(z) = \left( \frac{1}{1+e^{-z}} \right)$ . We consider the zero-padded edges are not added to the original input matrix  $X$ , i.e. the kernel cannot exceed the boundary of input during the convolution which leads  $i \in \{m_1, \dots, n_1\}$ ,  $j \in \{m_2, \dots, n_2\}$ .

## 2.2 Pooling Layer

Subsampling is performed in pooling layer, where the average value (mean-pooling) or the max value (max-pooling) is calculated in each sub-region. Here we use nonoverlapping mean-pooling as an example:

$$S_{c_{out}}(i, j) = \frac{1}{p_1 p_2} \sum_{k_1=1}^{p_1} \sum_{k_2=1}^{p_2} A_{c_{out}}(p_1 i' - k_1, p_2 j' - k_2) \quad (13)$$

where  $p_1, p_2$  are pooling size in each dimension,  $i' = m_1 + i \in \left\{ m_1 + 1, \dots, m_1 + \frac{n_1}{p_1} \right\}$  and  $j' = m_2 + j \in \left\{ m_2 + 1, \dots, m_2 + \frac{n_2}{p_2} \right\}$ .

## 2.3 Fully Connection Layer

Fully Connection Layer is similar to the output layer which contains neuron-like processing units. Additionally, Fully Connection Layer vectorizes and concatenates the output of pooling layer.

$$z = Z \left( \{S_{c_{out}}\}_{c_{out}=1, \dots, C_{out}} \right) \quad (14)$$

$$\{S_{c_{out}}\}_{c_{out}=1, \dots, C_{out}} = Z^{-1}(z) \quad (15)$$

Equation (14) denotes the process of vectorizing the output of each channel by column scan and concatenates them to form a whole string. Equation (15) represents the reverse process, where  $z \in R^{C_{out}(n_1/p_1)(n_2/p_2) \times 1}$ . The whole string is then input into the activation function for predicting class labels.

$$\vec{\hat{y}} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_c \\ \vdots \\ \hat{y}_C \end{bmatrix} = \phi \left( W^{(f)} z + \vec{b} \right) \quad (16)$$

where  $\hat{y}_c = \hat{P}(y = c | X)$ ,  $C$  is the number of classes,  $W^{(f)} \in R^{C \times C_{out}(n_1/p_1)(n_2/p_2)}$ ,  $\vec{b} \in R^{C \times 1}$ , and  $\phi(z) = \left( \frac{1}{1+e^{-z}} \right)$ .

## 2.4 Backpropagation in CNNs

We choose Cross-Entropy as the Loss Function in this note, since it outperforms Squared Error Loss for some circumstances [3]. The Cross-Entropy Loss  $J$  is specified as below:

$$J = - \sum_{c=1}^C P(y = c | X) \log(\hat{y}_c) \quad (17)$$

We still begin with performing forward propagation to generate  $\vec{\hat{y}}$  with initial value of weight and bias in Convolution Layer and Fully Connection Layer. Then we resort to gradient descent to update the weights and bias from Fully Connection Layer to Convolution Layer:

$$\frac{\partial J}{\partial W_{c,j}^{(f)}} = -(1 - \hat{y}_c) z_j \quad (18)$$

$$\frac{\partial J}{\partial b_c^{(f)}} = -(1 - \hat{y}_c) \quad (19)$$

$$\frac{\partial J}{\partial W_{c_{out}}^{(conv)}(k_1, k_2)} = \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} \alpha * X_{rot180} \quad (20)$$

$$\frac{\partial J}{\partial b_{c_{out}}^{(conv)}} = \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} \alpha(i, j) \quad (21)$$

$$W_{c,j}^{(f)} = W_{c,j}^{(f)} - \eta \frac{\partial J}{\partial W_{c,j}^{(f)}} \quad (22)$$

$$b_c^{(f)} = b_c^{(f)} - \eta \frac{\partial J}{\partial b_c^{(f)}} \quad (23)$$

$$W_{c_{out}}^{(conv)}(k_1, k_2) = W_{c_{out}}^{(conv)}(k_1, k_2) - \eta \frac{\partial J}{\partial W_{c_{out}}^{(conv)}(k_1, k_2)} \quad (24)$$

$$b_{c_{out}}^{(conv)} = b_{c_{out}}^{(conv)} - \eta \frac{\partial J}{\partial b_{c_{out}}^{(conv)}} \quad (25)$$

The derivation of Equation (18),(19),(20), and (21) can be refer to Appendix for the details.

### 3 Recurrent Neural Network

Different from a feedforward network (MLP as an example) dealing with independent and identically distributed data, Recurrent Neural Network (RNN) processes input data in a certain order and dependent of each other. Base on such purpose, RNN is designed to has input layer/hidden layers/output layer unit at each time instance with connection between hidden layers in the same level at two adjacent time instances [4, 5].

#### 3.1 Structure of a Single Layer RNN

We consider a single layer RNN in which only one hidden layer locates at each time instance. The hidden layer activation at time instance  $\tau$  is the function of  $\vec{x}_\tau$  the input data at the time instance  $\tau$  and  $\vec{h}_{\tau-1}$  the hidden layer activation at time instance  $\tau - 1$ :

$$\vec{h}_\tau = \phi_h \left( W_{hh}^T \vec{h}_{\tau-1} + W_{xh}^T \vec{x}_\tau + \vec{b}_h \right) \quad (26)$$

where  $\phi_h(\cdot)$  is the activation function of the hidden layer and  $\vec{b}_h$  is the bias vector.

Output layer pre-activation at instance  $\tau$  is calculated from weighting matrix from hidden layer to output layer  $W_{hy}$ , hidden layer activation  $\vec{h}_\tau$  (X), and the bias at output layer  $\vec{b}_y$  as

$$\vec{y}_\tau = \phi_y \left( W_{hy}^T \vec{h}_\tau + \vec{b}_y \right) \quad (27)$$

where  $\phi_y(\cdot)$  is the activation function of the output layer which is the softmax function used in this note, and  $\vec{y}_\tau = [\hat{y}_1^\tau, \dots, \hat{y}_c^\tau, \dots, \hat{y}_C^\tau]$  is the predicted probability distribution at instance  $\tau$ .

#### 3.2 Backpropagation Through Time (BPTT) and Exploding / Vanishing Gradients problem

The Loss function  $J$  for RNN is the sum of all the loss functions at each time instance. Here we choose Cross-Entropy as the Loss Function:

$$J = \sum_{\tau=1}^T J^{(\tau)} = \sum_{\tau=1}^T \left( - \sum_{c=1}^C y_{c,\tau} \log(\hat{y}_{c,\tau}) \right) \quad (28)$$

As  $W_{hh}$  is include in the hidden layer activation at each time instance, we need to apply Backpropagation Through Time (BPTT) to RNN which will expand the partial derivative in the time demension when we calculate gradient descent to update the weights:

$$\begin{aligned}
\frac{\partial J^{(\tau)}}{\partial W_{hh}} &= \frac{\partial J^{(\tau)}}{\partial \vec{y}_\tau} \frac{\partial \vec{y}_\tau}{\partial \vec{h}_\tau} \left( \sum_{t=1}^{\tau} \frac{\partial \vec{h}_\tau}{\partial \vec{h}_t} \frac{\partial \vec{h}_t}{\partial W_{hh}} \right) \\
&= \frac{\partial J^{(\tau)}}{\partial \vec{y}_\tau} \frac{\partial \vec{y}_\tau}{\partial \vec{h}_\tau} \left( \sum_{t=1}^{\tau} \prod_{i=t+1}^{\tau} \frac{\partial \vec{h}_i}{\partial \vec{h}_{i-1}} \frac{\partial \vec{h}_t}{\partial W_{hh}} \right) \\
&= \frac{\partial J^{(\tau)}}{\partial \vec{y}_\tau} \frac{\partial \vec{y}_\tau}{\partial \vec{h}_\tau} \left( \sum_{t=1}^{\tau} \prod_{i=t+1}^{\tau} W_{hh} \frac{\partial \vec{h}_t}{\partial W_{hh}} \right)
\end{aligned} \tag{29}$$

$\prod_{i=t+1}^{\tau} W_{hh}$  will lead  $\frac{\partial J^{(\tau)}}{\partial W_{hh}}$  to trend towards positive infinity when  $|W_{hh}| > 1$  and  $\tau - t$  is large. Conversely,  $\frac{\partial J^{(\tau)}}{\partial W_{hh}}$  will trend towards zero when  $|W_{hh}| < 1$  and  $\tau - t$  is large. The former is called Exploding Gradients problem and the other is called Vanishing Gradients problem.

### 3.3 Long Short-Term Memory (LSTM)

Long Short-Term Memory (LSTM) [6] is a RNN architecture designed to overcome Exploding / Vanishing Gradients problem. LSTM is introduce 3 gates (Input gate  $i_\tau$ , Forget gate  $f_\tau$ , Output gate  $o_\tau$ ) to suppress  $\vec{a}_\tau$  the input activation at the time instance  $\tau$ ,  $\vec{C}_{\tau-1}$  the cell state at the time instance  $\tau - 1$ ,  $\vec{C}_\tau$  the cell state at the time instance  $\tau$  respectively:

The calculation of gates are specified below, where  $\sigma$  denotes sigmoid function:

$$\vec{i}_\tau = \sigma(\vec{z}_i^\tau) = \sigma\left(W_{xi} \vec{x}_\tau + W_{hi} \vec{h}_{\tau-1} + \vec{b}_i\right) \tag{30}$$

$$\vec{f}_\tau = \sigma(\vec{z}_f^\tau) = \sigma\left(W_{xf} \vec{x}_\tau + W_{hf} \vec{h}_{\tau-1} + \vec{b}_f\right) \tag{31}$$

$$\vec{o}_\tau = \sigma(\vec{z}_o^\tau) = \sigma\left(W_{xo} \vec{x}_\tau + W_{ho} \vec{h}_{\tau-1} + \vec{b}_o\right) \tag{32}$$

The Feed-Forward operation of LSTM:

$$\vec{a}_\tau = \sigma(\vec{z}_a^\tau) = \tanh\left(W_{xa} \vec{x}_\tau + W_{ha} \vec{h}_{\tau-1} + \vec{b}_a\right) \tag{33}$$

$$\vec{C}_\tau = \left(\vec{a}_\tau \odot \vec{i}_\tau\right) \oplus \left(\vec{C}_{\tau-1} \odot \vec{f}_\tau\right) \tag{34}$$

$$\vec{h}_\tau = \tanh\left(\vec{C}_\tau\right) \odot \vec{o}_\tau \tag{35}$$

After the feed-forward operation is performed with initial value of weight and bias, the backpropagation is conducted to update weight and bias by minimizing the loss which we use the sum of squared of the errors (SSE) loss as an example:

$$J = \sum_{\tau=1}^T J^{(\tau)} = \sum_{\tau=1}^T \left( \frac{1}{2} (\vec{y}_\tau - \vec{h}_\tau)^2 \right) \quad (36)$$

The calculation of the partial derivatives of loss function regarding  $W_{xa}, W_{xi}, W_{xf}, W_{xo}, \vec{b}_a$  are detailed below. The partial derivatives of loss function regarding  $W_{hf}, W_{xo}, W_{ho}, W_{hi}, \vec{b}_i, \vec{b}_f, \vec{b}_o$  have the very similar structure and are not specified here.

$$\frac{\partial J}{\partial W_{xa}} = \begin{cases} \text{When } \tau = T : \\ = \frac{\partial J^{(T)}}{\partial \vec{C}_T} \frac{\partial \vec{C}_T}{\partial \vec{a}_T} \frac{\partial \vec{a}_T}{\partial W_{xa}} \\ = (\vec{h}_T - \vec{y}_T) \odot (1 - \tanh^2(\vec{C}_T)) \odot \vec{\sigma}_T \odot \vec{i}_T \odot (1 - \tanh^2(\vec{z}_{aT})) \vec{x}_T \\ \text{When } \tau < T : \\ = \frac{\partial J^{(T)}}{\partial \vec{C}_\tau} \frac{\partial \vec{C}_\tau}{\partial \vec{a}_\tau} \frac{\partial \vec{a}_\tau}{\partial W_{xa}} + \dots + \frac{\partial J^{(\tau)}}{\partial \vec{C}_\tau} \frac{\partial \vec{C}_\tau}{\partial \vec{a}_\tau} \frac{\partial \vec{a}_\tau}{\partial W_{xa}} \\ = (\vec{h}_T - \vec{y}_T) \odot (1 - \tanh^2(\vec{C}_T)) \odot \vec{\sigma}_T \odot \vec{f}_T \dots \vec{f}_{\tau+1} \odot \vec{i}_\tau \odot (1 - \tanh^2(\vec{z}_{a\tau})) \vec{x}_\tau \\ + \dots + (\vec{h}_\tau - \vec{y}_\tau) \odot (1 - \tanh^2(\vec{C}_\tau)) \odot \vec{\sigma}_\tau \odot \vec{i}_\tau \odot (1 - \tanh^2(\vec{z}_{a\tau})) \vec{x}_\tau \end{cases} \quad (37)$$

$$\frac{\partial J}{\partial W_{xi}} = \begin{cases} \text{When } \tau = T : \\ = \frac{\partial J^{(T)}}{\partial \vec{C}_T} \frac{\partial \vec{C}_T}{\partial \vec{i}_T} \frac{\partial \vec{i}_T}{\partial W_{xi}} \\ = (\vec{h}_T - \vec{y}_T) \odot (1 - \tanh^2(\vec{C}_T)) \odot \vec{\sigma}_T \odot \vec{a}_T \odot (1 - \tanh^2(\vec{z}_{iT})) \vec{x}_T \\ \text{When } \tau < T : \\ = \frac{\partial J^{(T)}}{\partial \vec{C}_\tau} \frac{\partial \vec{C}_\tau}{\partial \vec{i}_\tau} \frac{\partial \vec{i}_\tau}{\partial W_{xi}} + \dots + \frac{\partial J^{(\tau)}}{\partial \vec{C}_\tau} \frac{\partial \vec{C}_\tau}{\partial \vec{i}_\tau} \frac{\partial \vec{i}_\tau}{\partial W_{xi}} \\ = (\vec{h}_T - \vec{y}_T) \odot (1 - \tanh^2(\vec{C}_T)) \odot \vec{\sigma}_T \odot \vec{f}_T \dots \vec{f}_{\tau+1} \odot \vec{\sigma}_\tau \odot \vec{a}_\tau \odot (1 - \tanh^2(\vec{z}_{i\tau})) \vec{x}_\tau \\ + \dots + (\vec{h}_\tau - \vec{y}_\tau) \odot (1 - \tanh^2(\vec{C}_\tau)) \odot \vec{\sigma}_\tau \odot \vec{a}_\tau \odot (1 - \tanh^2(\vec{z}_{i\tau})) \vec{x}_\tau \end{cases} \quad (38)$$



$$\frac{\partial J}{\partial W_{xf}} = \begin{cases} \text{When } \tau = T : \\ = \frac{\partial J^{(T)}}{\partial \vec{C}_T} \frac{\partial \vec{C}_T}{\partial \vec{f}_T} \frac{\partial \vec{f}_T}{\partial W_{xf}} \\ = \left( \vec{h}_T - \vec{y}_T \right) \odot \left( 1 - \tanh^2 \left( \vec{C}_T \right) \right) \odot \vec{\sigma}_T \odot \vec{C}_{T-1} \odot \left( 1 - \tanh^2 \left( \vec{z}_{fT} \right) \right) \vec{x}_T \\ \text{When } \tau < T : \\ = \frac{\partial J^{(T)}}{\partial \vec{C}_\tau} \frac{\partial \vec{C}_\tau}{\partial \vec{f}_\tau} \frac{\partial \vec{f}_\tau}{\partial W_{xf}} + \dots + \frac{\partial J^{(\tau)}}{\partial \vec{C}_\tau} \frac{\partial \vec{C}_\tau}{\partial \vec{f}_\tau} \frac{\partial \vec{f}_\tau}{\partial W_{xf}} \\ = \left( \vec{h}_T - \vec{y}^T \right) \odot \left( 1 - \tanh^2 \left( \vec{C}_T \right) \right) \odot \vec{\sigma}_T \odot \vec{f}_T \dots \vec{f}_{\tau+1} \odot \vec{\sigma}_\tau \odot \vec{C}_{\tau-1} \odot \left( 1 - \tanh^2 \left( \vec{z}_{f\tau} \right) \right) \vec{x}_\tau \\ + \dots + \left( \vec{h}_\tau - \vec{y}_\tau \right) \odot \left( 1 - \tanh^2 \left( \vec{C}_\tau \right) \right) \odot \vec{\sigma}_\tau \odot \vec{C}_{\tau-1} \odot \left( 1 - \tanh^2 \left( \vec{z}_{f\tau} \right) \right) \vec{x}_\tau \end{cases} \quad (39)$$

$$\begin{aligned} \frac{\partial J}{\partial W_{xo}} &= \frac{\partial J^{(\tau)}}{\partial \vec{h}_\tau} \frac{\partial \vec{h}_\tau}{\partial \vec{\sigma}_\tau} \frac{\partial \vec{\sigma}_\tau}{\partial W_{xo}} \\ &= \left( \vec{h}_T - \vec{y}^T \right) \odot \tanh \left( \vec{C}_\tau \right) \odot \left( 1 - \tanh^2 \left( \vec{z}_{o\tau} \right) \right) \vec{x}_\tau \end{aligned} \quad (40)$$

$$\frac{\partial J}{\partial b_a} = \begin{cases} \text{When } \tau = T : \\ = \frac{\partial J^{(T)}}{\partial \vec{C}_T} \frac{\partial \vec{C}_T}{\partial \vec{a}_T} \frac{\partial \vec{a}_T}{\partial b_a} \\ = \left( \vec{h}_T - \vec{y}_T \right) \odot \left( 1 - \tanh^2 \left( \vec{C}_T \right) \right) \odot \vec{\sigma}_T \odot \vec{i}_T \odot \left( 1 - \tanh^2 \left( \vec{z}_{aT} \right) \right) \\ \text{When } \tau < T : \\ = \frac{\partial J^{(T)}}{\partial \vec{C}_\tau} \frac{\partial \vec{C}_\tau}{\partial \vec{a}_\tau} \frac{\partial \vec{a}_\tau}{\partial W_{xa}} + \dots + \frac{\partial J^{(\tau)}}{\partial \vec{C}_\tau} \frac{\partial \vec{C}_\tau}{\partial \vec{a}_\tau} \frac{\partial \vec{a}_\tau}{\partial W_{xa}} \\ = \left( \vec{h}_T - \vec{y}_T \right) \odot \left( 1 - \tanh^2 \left( \vec{C}_T \right) \right) \odot \vec{\sigma}_T \odot \vec{f}_T \dots \vec{f}_{\tau+1} \odot \vec{i}_\tau \odot \left( 1 - \tanh^2 \left( \vec{z}_{a\tau} \right) \right) \\ + \dots + \left( \vec{h}_\tau - \vec{y}_\tau \right) \odot \left( 1 - \tanh^2 \left( \vec{C}_\tau \right) \right) \odot \vec{\sigma}_\tau \odot \vec{i}_\tau \odot \left( 1 - \tanh^2 \left( \vec{z}_{a\tau} \right) \right) \end{cases} \quad (41)$$

## 4 Appendix

### 4.1 Derivation of Equations

Derivation of Equation (7): Using the chain rule, the partial derivative of loss with respect to the weight in output layer can be calculated as:

$$\begin{aligned}
\frac{\partial J(\mathbf{w})}{\partial w_{c,k}^{(out)}} &= \frac{\partial J(\mathbf{w})}{\partial a_{c,i}^{(out)}} \frac{\partial a_{c,i}^{(out)}}{\partial w_{c,k}^{(out)}} & (42) \\
&= - \left( \frac{P(y_i = c | x_i)}{a_{c,i}^{(out)}} - \frac{(1 - P(y_i = c | x_i))}{1 - a_{c,i}^{(out)}} \right) \frac{\partial \phi(z_{c,i}^{(out)})}{\partial z_{c,i}^{(out)}} \frac{\partial z_{c,i}^{(out)}}{\partial w_{c,k}^{(out)}} \\
&= - \left( \frac{P(y_i = c | x_i)}{a_{c,i}^{(out)}} - \frac{(1 - P(y_i = c | x_i))}{1 - a_{c,i}^{(out)}} \right) a_{c,i}^{(out)} (1 - a_{c,i}^{(out)}) a_{k,i}^{(h)} \\
&= (a_{c,i}^{(out)} - P(y_i = c | x_i)) a_{k,i}^{(h)}
\end{aligned}$$

where  $z_{c,i}^{(out)} = a_{1,i}^{(h)} w_{c,1}^{(out)} + \dots + a_{k,i}^{(h)} w_{c,k}^{(out)} + \dots + a_{K,i}^{(h)} w_{c,K}^{(out)}$ .

Derivation of Equation (8):

$$\begin{aligned}
\frac{\partial J(\mathbf{w})}{\partial w_{k,m}^{(h)}} &= \frac{\partial J(\mathbf{w})}{\partial a_{c,i}^{(out)}} \frac{\partial a_{c,i}^{(out)}}{\partial a_{k,i}^{(h)}} \frac{\partial a_{k,i}^{(h)}}{\partial w_{k,m}^{(h)}} & (43) \\
&= - \left( \frac{P(y_i = c | x_i)}{a_{c,i}^{(out)}} - \frac{(1 - P(y_i = c | x_i))}{1 - a_{c,i}^{(out)}} \right) a_{c,i}^{(out)} (1 - a_{c,i}^{(out)}) w_{c,k}^{(out)} a_{k,i}^{(h)} (1 - a_{k,i}^{(h)}) a_{m,i}^{(in)} \\
&= (a_{c,i}^{(out)} - P(y_i = c | x_i)) w_{c,k}^{(out)} a_{k,i}^{(h)} (1 - a_{k,i}^{(h)}) a_{m,i}^{(in)}
\end{aligned}$$

Derivation of Equation (18):

$$\begin{aligned}
\frac{\partial J}{\partial W_{c,j}^{(f)}} &= \frac{\partial J}{\partial \hat{y}_c} \frac{\partial \hat{y}_c}{\partial W_{c,j}^{(f)}} & (44) \\
&= - \frac{1}{\hat{y}_c} \frac{\partial \phi\left(\sum_{j=0}^{C_{out}(n_1/p_1)(n_2/p_2)-1} W_{c,j}^{(f)} z + b_c^{(f)}\right)}{\partial W_{c,j}^{(f)}} \\
&= - \frac{1}{\hat{y}_c} \hat{y}_c (1 - \hat{y}_c) z_j \\
&= -(1 - \hat{y}_c) z_j
\end{aligned}$$

Derivation of Equation (19):

$$\begin{aligned}
\frac{\partial J}{\partial b_c^{(f)}} &= \frac{\partial J}{\partial \hat{y}_c} \frac{\partial \hat{y}_c}{\partial b_c^{(f)}} & (45) \\
&= - \frac{1}{\hat{y}_c} \frac{\partial \phi\left(\sum_{j=0}^{C_{out}(n_1/p_1)(n_2/p_2)-1} W_{c,j}^{(f)} z + b_c^{(f)}\right)}{\partial b_c^{(f)}} \\
&= -(1 - \hat{y}_c)
\end{aligned}$$

Derivation of Equation (18):

$$\begin{aligned}
\frac{\partial J}{\partial W_{c_{out}}^{(conv)}(k_1, k_2)} &= \frac{\partial J}{\partial A_{c_{out}}(i, j)} \frac{\partial A_{c_{out}}(i, j)}{\partial W_{c_{out}}^{(conv)}(k_1, k_2)} & (46) \\
&= \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} \frac{\partial J}{\partial A_{c_{out}}(i, j)} \frac{\partial \left( \phi \left( \sum_{k_1=0}^{m_1-1} \sum_{k_2=0}^{m_2-1} X(i-k_1, j-k_2) W_{c_{out}}^{(conv)}(k_1, k_2) + b_{c_{out}}^{(conv)} \right) \right)}{\partial W_{c_{out}}^{(conv)}(k_1, k_2)} \\
&= \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} \frac{\partial J}{\partial A_{c_{out}}(i, j)} A_{c_{out}}(i, j) (1 - A_{c_{out}}(i, j)) X(i-k_1, j-k_2) \\
&= \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} \alpha(i, j) X_{rot180}(k_1-i, k_2-j) \\
&= \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} \alpha * X_{rot180}
\end{aligned}$$

where  $X_{rot180}$  denotes rotating  $X$  180 degrees.  $\alpha = \frac{\partial J}{\partial A_{c_{out}}(i, j)} A_{c_{out}}(i, j)$  and  $\frac{\partial J}{\partial A_{c_{out}}(i, j)}$  can be obtained by devectorizing and upsampling  $\frac{\partial J}{\partial z}$ . Equation (47), (48), (49) provides the solving process to get analytical solution of  $\frac{\partial J}{\partial A_{c_{out}}(i, j)}$ .

$$\begin{aligned}
\frac{\partial J}{\partial z} &= \frac{\partial J}{\partial \hat{y}_c} \frac{\partial \hat{y}_c}{\partial z} & (47) \\
&= -\frac{1}{\hat{y}_c} W_{c, j}
\end{aligned}$$

$$\frac{\partial J}{\partial S_{c_{out}}} = Z^{-1} \left( \frac{\partial J}{\partial z} \right) \quad (48)$$

$$\frac{\partial J}{\partial A_{c_{out}}(i, j)} = p_1 p_2 \frac{\partial J}{\partial S_{c_{out}}} ([i/p_1] [j/p_2]) \quad (49)$$

Derivation of Equation (21):

$$\begin{aligned}
\frac{\partial J}{\partial b_{c_{out}}^{(conv)}} &= \frac{\partial J}{\partial A_{c_{out}}(i, j)} \frac{\partial A_{c_{out}}(i, j)}{\partial b_{c_{out}}^{(conv)}} & (50) \\
&= \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} \alpha(i, j) \\
&= \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} p_1 p_2 \frac{\partial J}{\partial S_{c_{out}}} ([i/p_1] [j/p_2]) A_{c_{out}}(i, j)
\end{aligned}$$

## References

- [1] Zhifei Zhang. 2016. Derivation of Backpropagation in Convolutional Neural Network (CNN). Technical Report. University of Tennessee, Knoxville, TN.
- [2] Bouvrie, J. 2006. Notes on convolutional neural networks. Chatfield, K.; Simonyan, K.; Vedaldi, A.; and Zisserman, A. 2014. Return of the devil in the details: Delving deep into convolutional nets. arXiv preprint arXiv:1405.3531.
- [3] F.J. Huang and Y. LeCun. “Large-scale Learning with SVM and Convolutional for Generic Object Categorization”, In: Proc. 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, vol. 1, pp. 284-291, 2006.
- [4] Pascanu, Razvan, Mikolov, Tomas, and Bengio, Yoshua. On the difficulty of training recurrent neural networks. In Proceedings of the 30th International Conference on Machine Learning, ICML 2013, Atlanta, GA, USA, 16- 21 June 2013, pp. 1310–1318, 2013.
- [5] Hojjat Salehinejad, Julianne Baarbe, Sharan Sankar, Joseph Barfett, Errol Colak, and Shahrokh Valaee, “Recent advances in recurrent neural networks,” arXiv preprint arXiv:1801.01078, 2017.
- [6] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling. arXiv:1412.3555 [cs], December 2014. URL <http://arxiv.org/abs/1412.3555>.