

# RNNLIB: Softmax Layer

Wantee Wang

2015-02-05 21:08:17 +0800

## Contents

<b>1</b>	<b>Fundamentals</b>	<b>1</b>
1.1	List of Symbols . . . . .	1
1.2	Formulas . . . . .	2
1.3	Layers in RNNLIB . . . . .	2
<b>2</b>	<b>Forward Pass</b>	<b>2</b>
<b>3</b>	<b>Backpropagating</b>	<b>3</b>

I used to think that, in order to get the proper gradient, we have to take derivative of log of softmax with respect to weights. However, the RNNLIB shows that we can actually factorize the network into single layers. In this post, we look into the Softmax Layer.

## 1 Fundamentals

### 1.1 List of Symbols

Symbol	Meaning
$J$	cost function
$y_k$	activation of a neon
$u_k$	input of a neon
$S_i(\mathbf{u})$	softmax function, $i$ th value for a vector $\mathbf{u}$

## 1.2 Formulas

Softmax of a vector  $\mathbf{u}$  is defined as,

$$S_i(\mathbf{u}) = \frac{e^{u_i}}{\sum_k e^{u_k}} = y_i \quad (1)$$

the derivative of softmax is,

$$\frac{\partial S_i(\mathbf{u})}{\partial u_j} = \frac{\partial y_i}{\partial u_j} = \begin{cases} y_i(1 - y_i) & i = j \\ -y_i y_j & i \neq j \end{cases} \quad (2)$$

## 1.3 Layers in RNNLIB

Every layer in RNNLIB consists of input and output sides, both sides contain activations and errors. Their relations with terms in math are shown in following table,

Variable	Term
<i>inputActivations</i>	$u_k$
<i>outputActivations</i>	$y_k$
<i>inputErrors</i>	$\frac{\partial J}{\partial u_k}$
<i>outputErrors</i>	$\frac{\partial J}{\partial y_k}$

## 2 Forward Pass

Forward pass computes  $y_k$  from  $u_k$  using equation (1). There is a trick in the code, we can call it the *safe* softmax.

To understand it, consider dividing both numerator and denominator by  $e^c$  in equation (1),

$$S_i(\mathbf{u}) = \frac{\frac{e^{u_i}}{e^c}}{\frac{\sum_k e^{u_k}}{e^c}} = \frac{e^{u_i - c}}{\sum_k e^{u_k - c}} = S_i(\hat{\mathbf{u}}) \quad (3)$$

thus, in order to avoid overflow when calculating exponentials<sup>1</sup>, we can replace  $u_k$  with  $\hat{u}_k = u_k - c$ . Typically,  $c$  is set to  $u_{max}$ .

<sup>1</sup>Strictly speaking, this converts overflow into underflow. Underflow is no problem, because that rounds off to zero, which is a well-behaved floating point number. otherwise, it will be Infinity or NaN. see [this article](#) for details.

In RNNLIB,

$$c = \frac{u_{max} + u_{min}}{2}$$

### 3 Backpropagating

Backpropagation computes  $\frac{\partial J}{\partial u_k}$  from  $\frac{\partial J}{\partial y_k}$ .

In RNNLIB, the result is

$$\frac{\partial J}{\partial u_j} = y_j \left( \frac{\partial J}{\partial y_j} - \langle \mathbf{y}, \frac{\partial J}{\partial \mathbf{y}} \rangle \right) \quad (4)$$

where,  $\langle \cdot, \cdot \rangle$  denotes inner product.

To get the above equation, we first notice that variations in  $u_j$  give rise to variations in the error function  $J$  through variations in all  $y_k$ s. Thus, according to the [Multivariable Chain Rules](#), we can write,

$$\frac{\partial J}{\partial u_j} = \sum_k \frac{\partial J}{\partial y_k} \frac{\partial y_k}{\partial u_j} \quad (5)$$

Using equation (2) to replace  $\frac{\partial y_k}{\partial u_j}$ , we get,

$$\begin{aligned} \frac{\partial J}{\partial u_j} &= y_j(1 - y_j) \frac{\partial J}{\partial y_j} + \sum_{k:k \neq j} -y_k y_j \frac{\partial J}{\partial y_k} \\ &= y_j \left( \frac{\partial J}{\partial y_j} - y_j \frac{\partial J}{\partial y_j} + \sum_{k:k \neq j} -y_k \frac{\partial J}{\partial y_k} \right) \\ &= y_j \left( \frac{\partial J}{\partial y_j} - \sum_k y_k \frac{\partial J}{\partial y_k} \right) \\ &= y_j \left( \frac{\partial J}{\partial y_j} - \langle \mathbf{y}, \frac{\partial J}{\partial \mathbf{y}} \rangle \right) \end{aligned} \quad (6)$$

Finally, we reach equation (4).

In this way, softmax operation can be implemented to be a standalone layer.