

[Get started](#)[Open in app](#)

To make Medium work, we log user data.  
By using Medium, you agree to our  
Privacy Policy, including cookie policy.



## Chetan Patil

12 Followers

[About](#)[Follow](#)

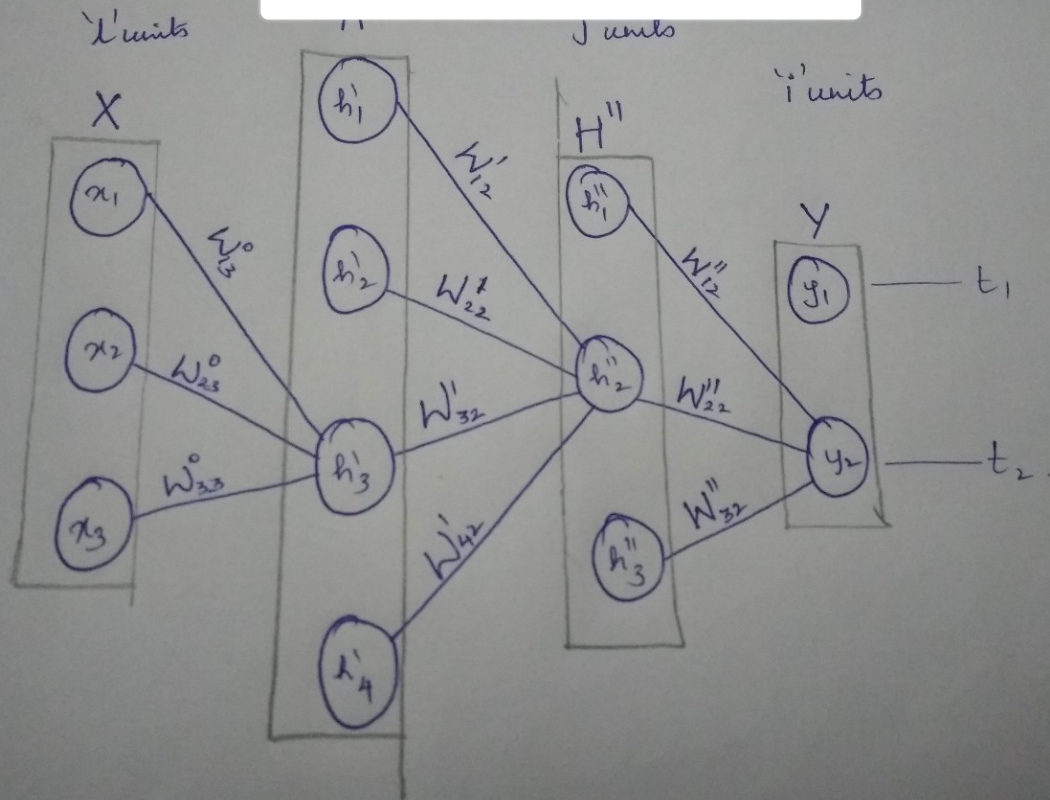
# Derivation of Back Propagation with Cross Entropy



Chetan Patil Dec 29, 2018 · 2 min read

In order to understand the Back Propagation algorithm, we first need to understand some basic concepts such as Partial Derivatives, chain rule, Cross Entropy loss, Sigmoid function and Softmax function.

To make Medium work, we log user data.  
By using Medium, you agree to our  
Privacy Policy, including cookie policy.



A two layered Neural Network with sigmoid activations

Assuming we have already forward-passed the inputs to get some outputs at the last layer Y, we will have to calculate the loss function E and propagate the loss to all the preceding layers by changing the weights associated with each of the layers.

$$E = - \sum_i [t_i \cdot \log(y_i) + (1-t_i) \log(1-y_i)] \quad , \quad y_i = \frac{e^{s_i''}}{\sum_c e^{s_c''}}$$

$$\begin{aligned} \text{Consider } \frac{dy_1}{ds_1''} &= \frac{d}{ds_1''} \left[ \frac{e^{s_1''}}{e^{s_1''} + e^{s_2''}} \right] \\ &= \frac{(e^{s_1''} + e^{s_2''}) \cdot e^{s_1''} - e^{s_1''} \cdot e^{s_1''}}{(e^{s_1''} + e^{s_2''})^2} = \frac{e^{s_1''} [e^{s_2''} - e^{s_1''}]}{(e^{s_1''} + e^{s_2''})^2} \end{aligned}$$

$$\frac{dy_1}{ds_1''} = \frac{e^{s_1''}}{e^{s_1''} + e^{s_2''}} \cdot \left[ 1 - \frac{e^{s_1''}}{e^{s_1''} + e^{s_2''}} \right] = y_1 (1 - y_1)$$

To make Medium work, we log user data.  
By using Medium, you agree to our  
Privacy Policy, including cookie policy.

$$\frac{\partial E}{\partial W_{ji}''} = \frac{\partial E}{\partial y_i} \frac{\partial y_i}{\partial \delta_i''} \frac{\partial \delta_i''}{\partial W_{ji}''}$$

$$= \left( -\frac{t_i}{y_i} + \frac{1-t_i}{1-y_i} \right) y_i (1-y_i) \cdot h_{ji}''$$

$$= \frac{-t_i(1-y_i) + y_i(1-t_i)}{y_i(1-y_i)} \cdot y_i(1-y_i) h_{ji}''$$

$$\frac{\partial E}{\partial W_{ji}''} = h_{ji}'' \cdot (y_i - t_i)$$

example

$$\frac{\partial E}{\partial W_{32}''} = h_{32}'' \cdot (y_2 - t_2)$$

Weight gradient of a weight connecting the third unit in second layer and second unit in the output layer using softmax activation.

Knowing the cross entropy loss  $E$  and the softmax activation " $y_i$ ", we can calculate the change in loss with respect to any weight connecting the output layer using the chain rule of partial derivatives. Intuitively, we can even find the weight gradients for the whole layer using the matrix notation shown below.

$$\frac{\partial E}{\partial W''} = H'' \cdot (Y - t)^T$$

Matrix notation of weights connecting the second hidden layer and output layer.

Now, finding weight gradients connecting



$H'$  and

WKT,  $h_j'' = \frac{1}{1 + e^{-s_j'}}$  where  $s_j' = \sum_k h_k' \cdot w_{kj}'$ .

example.

$$h_2'' = \text{sigmoid}(\underbrace{h_1' w_{12}' + h_2' w_{22}' + h_3' w_{32}' + h_4' w_{42}'}_{s_2'})$$

Consider

$$\begin{aligned} \frac{\partial E}{\partial w_{kj}'} &= \sum_{i=1}^{n_{out}} \frac{\partial E}{\partial s_i''} \cdot \frac{\partial s_i''}{\partial h_j''} \cdot \frac{\partial h_j''}{\partial s_j'} \cdot \frac{\partial s_j'}{\partial w_{kj}'} \\ &= \sum_{i=1}^{n_{out}} (y_i - t_i) \cdot w_{ji}'' \cdot h_j'' (1 - h_j'') \cdot h_k' \end{aligned}$$

example

$$\frac{\partial E}{\partial w_{41}'} = \left[ (y_1 - t_1) w_{11}'' + (y_2 - t_2) w_{12}'' \right] \cdot h_1'' (1 - h_1'') \cdot h_4'$$

Weight gradient of a weight connecting the fourth unit in first hidden layer to the first unit in the second hidden layer using sigmoid activation.

In matrix notation

$$\frac{\partial E}{\partial H'} = H' \cdot \left\{ (W'' \cdot [Y(1-t)]) * H^2 * (1 - H^2) \right\}^T$$

$dW'$ 

To make Medium work, we log user data.

By using Medium, you agree to our

Matrix

Privacy Policy, including cookie policy.

ayers

Now finding the weight gradients connecting  $X$  and  $H'$

WKT,  $h'_k = \frac{1}{1 + e^{-\delta_k^o}}$  where  $\delta_k^o = \sum_l x_l \cdot w_{lk}^o$ .

example

$$h'_3 = \text{sigmoid}(x_1 \cdot w_{13}^o + x_2 \cdot w_{23}^o + x_3 \cdot w_{33}^o)$$

Consider

$$\frac{\partial E}{\partial w_{lk}^o} = \sum_i \frac{\partial E}{\partial y_i} \cdot \frac{\partial y_i}{\partial s_i''} \cdot \frac{\partial s_i''}{\partial h_j''} \cdot \frac{\partial h_j''}{\partial s_j'} \cdot \frac{\partial s_j'}{\partial h'_k} \cdot \frac{\partial h'_k}{\partial \delta_k^o} \cdot \frac{\partial \delta_k^o}{\partial w_{lk}^o}$$

$$= \sum_i y_i (1 - y_i) \cdot w_{ji}'' \cdot \sum_j h_j'' (1 - h_j'') \cdot w_{kj}' \cdot h'_k (1 - h'_k) \cdot x_l$$

On rearranging,

$$= \sum_j \left[ h_j'' (1 - h_j'') \cdot \left( \sum_i y_i (1 - y_i) \cdot w_{ji}'' \right) \right] \cdot w_{kj}' \cdot h'_k (1 - h'_k) \cdot x_l$$

Weight gradient of a weight connecting a unit  $L$  in input layer to the unit  $K$  in the first hidden layer using sigmoid activation.

$$\partial W^o = X \cdot \left( W' \cdot \left( H'' (1 - H'') \right)^* (W'' \cdot Y^* (1 - Y)) \right)^* \left( H' (1 - H') \right)^T$$

Matrix notation of weights connecting the input layer and first hidden layer.

To make Medium work, we log user data.

By using Medium, you agree to our

Privacy Policy, including cookie policy.

[About](#) [Write](#) [Help](#) [Legal](#)

Get the Medium app

