

Neural networks and Back Propagation

Hans Mathias Mamen Vege

December 3, 2018

1 Introduction

A short presentation on back propagation is given, in the derivation of the back propagation algorithm is given. This is done as an exercise as well as an future reference for the author. For a note on the notation in this short paper, consult the appendix A.

We begin with an short presentation of the layout of a neural network, then we will go through and present different cost functions and different output activations. A presentation of middle layer activations is also given. Then, we will begin deriving the back-propagation and the gradients which will be used to change the weights and biases in a neural network.

2 A neural network

A neural network, more specifically the *Multilayer Perceptron*, is a machine learning method which is takes some data as input and outputs some desired output based on a cost function and labeled training data. A typical would be some sort of classification, e.g. image classification.

A basic neural network consists of an input layer, and $L - 1$ hidden layers, where L is the total number of layers in addition to an output layer. For each layer we have a number of neurons, which will take as input the previous layer output, multiply it by some weights and add some biases. Then, it will be send through some activation function and taken as input for the next layer. For a simple neural network with just an input layer one hidden layer we will have two weights and two biases. One for the hidden layer and one for the output layer. To illustrate,

$$\mathbf{a}^1 = \sigma(\mathbf{z}^1) = \sigma(W^1 \mathbf{x} + \mathbf{b}^1) \quad (1)$$

$$\hat{\mathbf{y}} = \sigma_f(\mathbf{z}^2) = \sigma_f(W^2 \mathbf{a}^1 + \mathbf{b}^2) \quad (2)$$

The first equation (1) takes x which is the input data, and multiplies it with the weights matrix for the hidden layer W_1 , and adds biases b_1 . Then this is passed through some activation function σ . The output, $\hat{\mathbf{y}}$ will be the output of the neural network, and is equivalent to \mathbf{a}^L (or \mathbf{a}^2 for this case). From now on, \mathbf{a}^L will be used for the output layer. The z_1 is simply a collective term for the input to the activation function, and is generalized as

$$\mathbf{z} = W\mathbf{x} + \mathbf{b} \quad (3)$$

To generalize this notation to any layer l , the output of a layer will be given as

$$\mathbf{a}^l = \sigma(\mathbf{z}^l) = \sigma(W^l \mathbf{a}^{l-1} + \mathbf{b}^l) \quad (4)$$

This will be called the *forward pass* of a single layer.

If we extend this model a a full neural network, the final layer output will be,

$$\mathbf{a}^L = \sigma(W^L \mathbf{a}^{L-1} + \mathbf{b}^L) \quad (5)$$

Inserting the previous layer outputs, the full forward pass will be,

$$\mathbf{a}^L = \sigma(W^L \sigma(W^{L-1} \sigma(\dots \sigma(W^1 \mathbf{x} + \mathbf{b}^1) \dots) + \mathbf{b}^{L-1}) + \mathbf{b}^L) \quad (6)$$

In order to make the network *learn*, we will have to update these weights and biases somehow. The way that is done is by first having a measure of how good the network is performing. This is done using a *cost function*, \mathcal{C} . We will end up taking the gradient of this function with respect to either the bias \mathbf{b}^l or W^l ,

$$\nabla_{j,W}^l \mathcal{C} = \frac{\partial \mathcal{C}}{\partial W_j^l}$$

and

$$\nabla_{j,b}^l \mathcal{C} = \frac{\partial \mathcal{C}}{\partial b_j^l}$$

Note that we are taking the derivative component-wise. Let us start by investigating two cost functions which is widely used.

3 Cost functions

The cost function is a function which measures how good our network predicts an outcome based on some labeled outcome. We wish to minimize, and in doing so will use stochastic gradient descent and update our weights based on the gradient descent. We will be looking at two main cost functions, the *quadratic cost* (or MSE cost function), or the *cross entropy*.

3.1 Quadratic cost

The quadratic cost, or mean square error (MSE) cost function is given by

$$\mathcal{C}_{\text{MSE}} = \frac{1}{2N} \sum_{i=0}^N \|\mathbf{a}_i - \mathbf{y}_i\|_2^2, \quad (7)$$

where we take Euclidean (or L^2 norm) between the final layer output \mathbf{a} minus the labeled training data \mathbf{y} . The sum i is over all of the training data, N . Note, that we will usually ignore the full cost function, since we are mostly dealing with one sample at a time. We then get

$$\mathcal{C}_{\text{MSE}} = \frac{1}{2} \|\mathbf{a}_i - \mathbf{y}_i\|_2^2, \quad (8)$$

with no notational differences being given.

3.2 Cross entropy

The cross entropy cost function is given by

$$\mathcal{C}_{\text{CE}} = -\frac{1}{N} \sum_{i=0}^N \sum_{k=1}^K y_{ik} \ln a_{ik} \quad (9)$$

This is the full cross entropy cost function, which is the average over all the samples N . The K -sum is a sum over all of the different labels. We will mostly simply use the expression without the explicit sum over all samples,

$$\mathcal{C}_{\text{CE}} = -\sum_{k=1}^K y_k \ln a_k \quad (10)$$

This cost function has a nice expression in the binary classification case. We have that the sum of the output vector \mathbf{y} has to be equal to one when summed up (since it is a probability, the maximum can be one),

$$\sum_{k=0}^K y_k = 1,$$

we will in the binary classification case with $K = 2$, have

$$\mathcal{C}_{\text{CE}}^2 = - \sum_{k=1}^{K=2} y_k \ln a_k.$$

Writing this out, we get

$$\mathcal{C}_{\text{CE}}^2 = -y_1 \ln a_1 - y_2 \ln a_2$$

We can now rewrite y_2 given that the total probability is 1, as $y_2 = 1 - y_1$. We can then relabel $y_1 \rightarrow y$. The same can be done for the output.

$$\mathcal{C}_{\text{CE}}^2 = -y \ln a - (1 - y) \ln(1 - a). \quad (11)$$

$$(12)$$

With these cost functions, we can begin at looking at the activation functions.

4 Activation functions

So far, we have not touched on activation functions. We will differ between the output activation function, which will in the classification case be forced to be between 0 and 1, and the hidden layer activation functions.

4.1 Output layer activation functions

4.1.1 Sigmoidal activation function

The **sigmoid** activation function is given as,

$$\sigma_{\text{sig}}(z) = \frac{1}{1 + e^{-z}} \quad (13)$$

The sigmoidal activation function can also be used as a hidden layer activation function. The derivative with respect to z of the sigmoidal activation function is easy,

$$\begin{aligned} \frac{\partial}{\partial z}(\sigma_{\text{sig}}(z)) &= \frac{\partial}{\partial z} \left(\frac{1}{1 + e^{-z}} \right) \\ &= (-e^{-z}) \frac{-1}{(1 + e^{-z})^2} \\ &= \sigma_{\text{sig}}(z) \left(\frac{1 - 1 + e^{-z}}{1 + e^{-z}} \right) \\ &= \sigma_{\text{sig}}(z) \left(1 - \frac{1}{1 + e^{-z}} \right) \\ &= \sigma_{\text{sig}}(z) (1 - \sigma_{\text{sig}}(z)) \end{aligned}$$

And this, summed up we have

$$\sigma'_{\text{sig}}(z) = \sigma_{\text{sig}}(z)(1 - \sigma_{\text{sig}}(z)) \quad (14)$$

4.1.2 Softmax activation function

The **softmax** activation function is defined component-wise, and is given as

$$\sigma_{\text{sm}}(z)_i = \frac{e^{z_i}}{\sum_{k=1}^K e^{z_k}} \quad (15)$$

The softmax activation function can only be used in the output layer¹. Let us now find the derivative of the softmax function with respect to z .

$$\begin{aligned} \frac{\partial}{\partial z_j} (\sigma_{\text{sm}}(z)_i) &= \frac{\partial}{\partial z_j} \left(\frac{e^{z_i}}{\sum_{k=1}^K e^{z_k}} \right) \\ &= \begin{cases} \frac{\partial}{\partial z_j} \left(\frac{e^{z_i}}{\sum_{k=1}^K e^{z_k}} \right) & \text{if } j \neq i \\ \frac{\partial}{\partial z_j} \left(\frac{e^{z_i}}{\sum_{k=1}^K e^{z_k}} \right) & \text{if } j = i \end{cases} \\ &= \begin{cases} \frac{-e^{z_j} e^{z_i}}{\left(\sum_{k=1}^K e^{z_k}\right)^2} & \text{if } j \neq i \\ \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}} - \frac{e^{z_j} e^{z_i}}{\left(\sum_{k=1}^K e^{z_k}\right)^2} & \text{if } j = i \end{cases} \\ &= \begin{cases} -\sigma_{\text{sm}}(z)_j \sigma_{\text{sm}}(z)_i & \text{if } j \neq i \\ \sigma_{\text{sm}}(z)_j - \sigma_{\text{sm}}(z)_j \sigma_{\text{sm}}(z)_i & \text{if } j = i \end{cases} \\ &= \sigma_{\text{sm}}(z)_j (\delta_{ji} - \sigma_{\text{sm}}(z)_i) \end{aligned}$$

Which is, summing

$$\frac{\partial}{\partial z_j} (\sigma_{\text{sm}}(z)_i) = \sigma_{\text{sm}}(z)_j (\delta_{ji} - \sigma_{\text{sm}}(z)_i) \quad (16)$$

4.2 Hidden layer activation functions

For hidden layer activation functions, we have several options. We have already looked at the sigmoidal function as an output layer activation function, and it can be used as a hidden layer activation function as well.

4.2.1 Hyperbolic tangens activation function

The **hyperbolic tangens** activation function is given as

$$\sigma_{\text{tanh}}(z) = \tanh(z) \quad (17)$$

with its derivative

$$\sigma'_{\text{tanh}}(z) = 1 - \tanh^2(z) \quad (18)$$

4.2.2 Relu activation function

The **relu** or rectifier activation is given as,

$$\sigma_{\text{relu}}(z) = \begin{cases} z & \text{if } z \geq 0 \\ 0 & \text{if } z < 0 \end{cases} \quad (19)$$

with its derivative

$$\sigma'_{\text{relu}}(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{if } z < 0 \end{cases} \quad (20)$$

¹See this page for why softmax is commonly not used.

4.2.3 Heaviside activation function

The **Heaviside** activation function is given as

$$\sigma_{\text{Heaviside}}(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{if } z < 0 \end{cases} \quad (21)$$

with its derivative

$$\sigma'_{\text{Heaviside}}(z) = 0 \quad (22)$$

5 Training a network

As mentioned, we will train the network using *Stochastic Gradient Descent* (SGD) and *mini-batches*. The idea is to shuffle the data set, and randomly select a subset of elements in the dataset. We will then perform a forward-pass of the network, and using the output \mathbf{a}^L , perform back propagation, and update each weight and bias in the network. The main algorithm for training a network can be seen below,

Algorithm 1 Training a MLP

- 1: Initialize learning rate η .
 - 2: Initialize with input data $\{\mathbf{x}_i\}$ and labels $\{\mathbf{y}_i\}$ of size N .
 - 3: **for** $i_{\text{epoch}} = 1 : N_{\text{epochs}}$ **do**
 - 4: Shuffle data.
 - 5: Split data into K_{mb} mini batches, $\{\mathbf{x}_{i_{\text{mb}}}\}$ and $\{\mathbf{y}_{i_{\text{mb}}}\}$.
 - 6: **for** $i_{\text{mb}} = 1 : N/K_{\text{mb}}$ **do**
 - 7: Perform a forward pass.
 - 8: Perform a back propagation, retrieving $\nabla W_{i_{\text{mb}}}^l$ and $\nabla b_{i_{\text{mb}}}^l$ for each i_{mb} .
 - 9: **end for**
 - 10: Average the $\nabla W_{i_{\text{mb}}}^l$ and $\nabla b_{i_{\text{mb}}}^l$, and update weights W^l and biases b^l with these averages.
 - 11: Update learning rate η if needed.
 - 12: **end for**
-

5.1 Learning rate

When updating the weights and biases with SGD, we did so by a learning rate parameter η . There are several way to define η , with the simplest one having $\eta = \text{constant}$. Another option is one that is inversely decreasing as a function of the epochs. That is, for a given epoch i_{epoch} out a total N_{epochs} , we set the learning rate as

$$\eta(i_{\text{epoch}}) = \eta_0 \left(1 - \frac{i_{\text{epoch}}}{1 + N_{\text{epochs}}} \right) \quad (23)$$

This will force the step size to decrease toward 0 as we close in on the maximum number of epochs N_{epochs} .

5.2 Weight initialization

When initializing weights and biases, we will look at two ways of how this can be done. The first is through a Gaussian distribution, $(0, 1)$ which we will call *large*, as the biases will have large, spread-out distribution.

Then, we will use a Gaussian distribution but divided by the number of training samples, $(0, 1/N_{\text{train}})$, dubbing that one to be called *default*, as this is the one we will by default use in our neural network.

The effect of these two is essentially shrinking in the initial search space, and we should expect them to converge at large epoch times.

5.3 Measuring the performance

The performance of a neural network (or any classifier), can in its simplest form be measured by the accuracy, which is defined as

$$\text{Accuracy} = \frac{\sum_{i=1}^n I(t_i = y_i)}{n}, \quad (24)$$

where n is the number of samples we are testing against, I is the indicator function, which returns 1 if the prediction t_i equals the true values y_i .

6 Deriving the back propagation

We have now collected a bunch of tools, and are now ready to apply them to the back-propagation. We begin by imagining us to update the weights based on the gradient of the cost function. We begin with updating the output weights, W^L . Labeling the vectors and matrices, we get

$$\frac{\partial \mathcal{C}}{\partial W_{jk}^L} = \frac{\partial \mathcal{C}}{\partial a_j^L} \frac{\partial a_j^L}{\partial z_j^L} \frac{\partial z_j^L}{\partial W_{jk}^L} \quad (25)$$

We see that two of the products are component-wise, such that we actually are multiplying by the Hadamard product.

$$\frac{\partial \mathcal{C}}{\partial W^L} = \left(\frac{\partial \mathcal{C}}{\partial \mathbf{a}^L} \odot \frac{\partial \mathbf{a}^L}{\partial \mathbf{z}^L} \right) \frac{\partial \mathbf{z}^L}{\partial W^L} \quad (26)$$

We see that we need three derivatives, one for the cost function, one for the output activation layer and one for the z^L .

To sum up, we have to find three derivatives, the first being the of the cost function

$$\frac{\partial \mathcal{C}}{\partial \mathbf{a}^L}.$$

We then need to find the output activation function derivative.

$$\frac{\partial \mathbf{a}^L}{\partial \mathbf{z}^L} = \frac{\partial}{\partial \mathbf{z}^L} (\sigma(\mathbf{z}^L)) \quad (27)$$

The final derivative, is given as

$$\frac{\partial \mathbf{z}^L}{\partial W^L} = \mathbf{a}^{L-1} \quad (28)$$

Let us now find the exact derivatives of the cost functions.

6.1 Quadratic cost derivative

We start with finding the derivative of the quadratic cost,

$$\begin{aligned} \frac{\partial \mathcal{C}_{\text{MSE}}}{\partial \mathbf{a}^L} &= \frac{\partial}{\partial \mathbf{a}^L} \left(\frac{1}{2} (\mathbf{a}^L - \mathbf{y})^2 \right) \\ &= \mathbf{a}^L - \mathbf{y}, \end{aligned}$$

which is summed up as,

$$\frac{\partial \mathcal{C}_{\text{MSE}}}{\partial \mathbf{a}^L} = \mathbf{a}^L - \mathbf{y}. \quad (29)$$

6.2 Cross entropy derivative

The derivative of the cross entropy is given as,

$$\begin{aligned}\frac{\partial \mathcal{C}_{\text{CE}}}{\partial \mathbf{a}^L} &= \frac{\partial}{\partial \mathbf{a}^L} (-\mathbf{y} \ln \mathbf{a}^L) \\ &= -\frac{\mathbf{y}}{\mathbf{a}^L}.\end{aligned}$$

Which, summed up is,

$$\frac{\partial \mathcal{C}_{\text{CE}}}{\partial \mathbf{a}^L} = -\frac{\mathbf{y}}{\mathbf{a}^L}. \quad (30)$$

If we have a binary classification ($K = 2$ classes), we can rewrite the expression as,

$$\begin{aligned}\frac{\partial \mathcal{C}_{\text{CE}}^2}{\partial a^L} &= -\frac{\partial}{\partial a^L} (y \ln a^L - (1 - y) \ln(1 - a^L)) \\ &= \frac{1 - y}{1 - a^L} - \frac{y}{a^L} \\ &= \frac{a^L - y}{a^L(1 - a^L)}\end{aligned}$$

6.3 Calculating the full back propagation

We have four different output (in reality only three²), which can be seen in 1.

Table 1: Different cost function activations. MSE is the same as quadratic cost.

Cost function	Output activation
MSE	Sigmoid
MSE	Softmax
CE	Softmax

For the derivations, we will mostly ignore that we are dealing with vectors and matrices, and instead bring it back in at the very end.

6.3.1 MSE and sigmoidal output activation

Let us then begin with looking at MSE and sigmoidal output activation. The change in the output layer L is given as,

$$\begin{aligned}\frac{\partial \mathcal{C}_{\text{MSE}}}{\partial W^L} &= \left(\frac{\partial \mathcal{C}_{\text{MSE}}}{\partial \mathbf{a}_{\text{sig}}^L} \odot \frac{\partial \mathbf{a}_{\text{sig}}^L}{\partial \mathbf{z}^L} \right) \frac{\partial \mathbf{z}^L}{\partial W^L} \\ &= \left(\underbrace{(\mathbf{a}^L - \mathbf{y})}_{\text{MSE}} \odot \underbrace{\mathbf{a}^L(1 - \mathbf{a}^L)}_{\text{Sigmoid}} \right) \mathbf{a}^{L-1}.\end{aligned}$$

For the change in bias b^L , we get

$$\begin{aligned}\frac{\partial \mathcal{C}_{\text{MSE}}}{\partial \mathbf{b}^L} &= \left(\frac{\partial \mathcal{C}_{\text{MSE}}}{\partial \mathbf{a}_{\text{sig}}^L} \odot \frac{\partial \mathbf{a}_{\text{sig}}^L}{\partial \mathbf{z}^L} \right) \frac{\partial \mathbf{z}^L}{\partial \mathbf{b}^L} \\ &= \left(\underbrace{(\mathbf{a}^L - \mathbf{y})}_{\text{MSE}} \odot \underbrace{\mathbf{a}^L(1 - \mathbf{a}^L)}_{\text{Sigmoid}} \right) \times 1.\end{aligned}$$

²For multiclass cross entropy output, sigmoidal activation as final layer is simply as special case of the softmax function. Thus, we will only use softmax activation for final layer output for the cross entropy.

We observe that the only change in output between the derivatives is \mathbf{a}^{L-1} .

If we now wish to find the change in the $L - 1$ layer with respect to W^{L-1} , we get

$$\begin{aligned}\frac{\partial \mathcal{C}_{\text{MSE}}}{\partial W^{L-1}} &= \frac{\partial \mathcal{C}_{\text{MSE}}}{\partial \mathbf{a}_{\text{sig}}^L} \frac{\partial \mathbf{a}_{\text{sig}}^L}{\partial \mathbf{z}^L} \frac{\partial \mathbf{z}^L}{\partial W^{L-1}} \\ &= \frac{\partial \mathcal{C}_{\text{MSE}}}{\partial \mathbf{a}_{\text{sig}}^L} \frac{\partial \mathbf{a}_{\text{sig}}^L}{\partial \mathbf{z}^L} \frac{\partial \mathbf{z}^L}{\partial \mathbf{a}^{L-1}} \frac{\partial \mathbf{a}^{L-1}}{\partial \mathbf{z}^{L-1}} \frac{\partial \mathbf{z}^{L-1}}{\partial W^{L-1}} \\ &= (\mathbf{a}^L - \mathbf{y}) \mathbf{a}^L (1 - \mathbf{a}^L) W^L \sigma'_{\text{sig}}(\mathbf{z}^{L-1}) \mathbf{a}^{L-2}\end{aligned}$$

Here, we recognize the first second as what we had for the first layer, so we can call that,

$$\delta^L = (\mathbf{a}^L - \mathbf{y}) \mathbf{a}^L (1 - \mathbf{a}^L) \quad (31)$$

Using this, for the $L - 1$ layer we get

$$\begin{aligned}\frac{\partial \mathcal{C}_{\text{MSE}}}{\partial W^L} &= \delta^L \mathbf{a}^{L-1} \\ \frac{\partial \mathcal{C}_{\text{MSE}}}{\partial \mathbf{b}^L} &= \delta^L\end{aligned}$$

and for the $L - 2$ layer, we get

$$\begin{aligned}\frac{\partial \mathcal{C}_{\text{MSE}}}{\partial W^{L-1}} &= \delta^L W^L \sigma'_{\text{sig}}(\mathbf{z}^{L-1}) \mathbf{a}^{L-2} \\ \frac{\partial \mathcal{C}_{\text{MSE}}}{\partial \mathbf{b}^{L-1}} &= \delta^L W^L \sigma'_{\text{sig}}(\mathbf{z}^{L-1}).\end{aligned}$$

Looking at the derivative for layer $L - 2$ with respect to W^{L-2} we get,

$$\begin{aligned}\frac{\partial \mathcal{C}_{\text{MSE}}}{\partial W^{L-2}} &= \frac{\partial \mathcal{C}_{\text{MSE}}}{\partial \mathbf{a}_{\text{sig}}^L} \frac{\partial \mathbf{a}_{\text{sig}}^L}{\partial \mathbf{z}^L} \frac{\partial \mathbf{z}^L}{\partial \mathbf{a}^{L-1}} \frac{\partial \mathbf{a}^{L-1}}{\partial \mathbf{z}^{L-1}} \frac{\partial \mathbf{z}^{L-1}}{\partial \mathbf{a}^{L-2}} \frac{\partial \mathbf{a}^{L-2}}{\partial \mathbf{z}^{L-2}} \frac{\partial \mathbf{z}^{L-2}}{\partial W^{L-2}} \\ &= \delta^L W^L \sigma'_{\text{sig}}(\mathbf{z}^{L-1}) W^{L-1} \sigma'_{\text{sig}}(\mathbf{z}^{L-2}) \mathbf{a}^{L-3}\end{aligned}$$

Updating δ^L for the subsequent layers gives us

$$\begin{aligned}\delta^{L-1} &= \delta^L W^L \sigma'_{\text{sig}}(\mathbf{z}^{L-1}) \\ \delta^{L-2} &= \delta^{L-1} W^{L-1} \sigma'_{\text{sig}}(\mathbf{z}^{L-2}).\end{aligned}$$

Using these expressions, we can organize our previous expressions to

$$\begin{aligned}\frac{\partial \mathcal{C}_{\text{MSE}}}{\partial W^{L-1}} &= \delta^{L-1} \mathbf{a}^{L-2} \\ \frac{\partial \mathcal{C}_{\text{MSE}}}{\partial W^{L-2}} &= \delta^{L-2} \mathbf{a}^{L-3}.\end{aligned}$$

We see a pattern, and can summarize our equations. Note that so far, we have mostly assumed that much of the matrix and vector multiplications are component wise, but will now include the full expression. Without specifying the sigmoidal function as the layer activation and writing the derivative of the cost function with respect to \mathbf{a} as $\nabla_{\mathbf{a}} \mathcal{C}$, we get

$$\delta^L = \nabla_{\mathbf{a}}^L \mathcal{C} \odot \sigma'(\mathbf{a}^L) \quad (32)$$

$$\frac{\partial \mathcal{C}_{\text{MSE}}}{\partial W^L} = \delta^L (\mathbf{a}^{L-1})^T \quad (33)$$

$$\frac{\partial \mathcal{C}_{\text{MSE}}}{\partial W^l} = \delta^L \quad (34)$$

For the subsequent layers, we get

$$\delta^{l-1} = \left((W^l)^T \delta^l \right) \odot \sigma'(\mathbf{z}^{l-1}) \quad (35)$$

$$\frac{\partial \mathcal{C}_{\text{MSE}}}{\partial W^l} = \delta^l \mathbf{a}^{l-1} \quad (36)$$

$$\frac{\partial \mathcal{C}_{\text{MSE}}}{\partial W^l} = \delta^l \quad (37)$$

6.3.2 MSE and softmax output activation

For MSE and softmax as the output activation, we need only change the MSE the first two derivatives in equation (25). We use what we found in equation

$$\frac{\partial \mathcal{C}}{\partial a_j^L} \frac{\partial a_j^L}{\partial z_j^L} = (a_j^L - y_j) \sigma_{\text{sm}}(z)_j (\delta_{ji} - \sigma_{\text{sm}}(z)_i)$$

To simplify, this becomes

$$\delta_j^L = (a_j^L - y_j) \sigma_{\text{sm}}(z)_j (\delta_{ji} - \sigma_{\text{sm}}(z)_i) \quad (38)$$

6.3.3 Cross entropy and softmax

We start by taking the derivative with respect to component j^3 ,

$$\frac{\partial \mathcal{C}_{\text{CE}}}{\partial a_{j,\text{SM}}^L} \frac{\partial a_{j,\text{SM}}^L}{\partial z_j^L} = \sum_k \left(-\frac{y_k}{a_k^L} \right) a_i^L (\delta_{ik} - a_k^L)$$

Splitting the delta function into the case of $i \neq k$ and $i = k$,

$$\begin{aligned} \sum_k \left(-\frac{y_k}{a_k^L} \right) a_i^L (\delta_{ik} - a_k^L) &= -\frac{y_i}{a_i^L} a_i^L (1 - a_i^L) - \sum_{k \neq i} \frac{y_k}{a_k^L} (-a_i^L a_k^L) \\ &= -y_i + y_i a_i^L - \sum_{k \neq i} y_k (-a_i^L) = -y_i + \left(\sum_k y_k \right) a_i^L \\ &= a_i^L - y_i \end{aligned}$$

Summing up, we have

$$\frac{\partial \mathcal{C}_{\text{CE}}}{\partial a_{j,\text{SM}}^L} \frac{\partial a_{j,\text{SM}}^L}{\partial z_j^L} = a_i^L - y_i \quad (39)$$

Appendices

A Notation

For notation, we will use bold italic as vectors, e.g. \mathbf{x} , \mathbf{y} , \mathbf{z} . For matrices, we will use capital letters, W , B , X .

³An alternative derivation of the cross entropy with softmax as output layer can be found here