

# Neural Network Mathematics

by Fraser Sabine

version 1.0

## Table of contents

|                   |                                 |    |
|-------------------|---------------------------------|----|
| <b>1</b>          | <b>Notation</b>                 | 2  |
| <b>2</b>          | <b>Cost/Loss Function</b>       | 3  |
| 2.1               | Squared error                   | 3  |
| 2.2               | Cross entropy                   | 3  |
| <b>3</b>          | <b>Activation Function</b>      | 4  |
| 3.1               | Sigmoid                         | 4  |
| 3.2               | Softmax                         | 4  |
| <b>4</b>          | <b>Matrix weights</b>           | 5  |
| <b>5</b>          | <b>Backpropagation</b>          | 6  |
| 5.1               | Sigmoid and squared error       | 6  |
| 5.2               | Softmax and cross entropy       | 7  |
| <b>Appendix A</b> | <b>Cross entropy derivation</b> | 9  |
| <b>Appendix B</b> | <b>Softmax derivation</b>       | 11 |

# 1 Notation

Where  $n$  is the number of nodes in the neural network

Where  $L$  denotes the final layer in the neural network and  $L - \dots$  denotes previous layers

Where  $\mathbf{z}^L$  is the calculated output  $\mathbf{z}^L = \mathbf{a}^{L-1} W^L + B^L$

Where  $\mathbf{a}^L$  is the activation function applied to  $\mathbf{z}^L$ . Note  $\mathbf{a}^{L-H}$  where  $H$  is the number of layers, is the input layer of the neural network.

Where  $W^L$  is the weights matrix

Where  $\mathbf{b}^L$  is the biases matrix

Where  $\mathbf{t}$  is the target output given as a training example

Where  $c$  is the loss/cost function

Note  $\oslash$  is the Hadamard division or the element-wise division of a vector or matrix

Note  $\circ$  is the Hadamard product or the element-wise multiplication of a vector or matrix

Note  $\cdot$  is the dot product of a vector or matrix

$$W^L = \begin{pmatrix} w_{00}^L & w_{01}^L & \cdots & w_{0n}^L \\ w_{10}^L & w_{11}^L & \cdots & w_{1n}^L \\ \vdots & \vdots & \ddots & \vdots \\ w_{n0}^L & w_{n1}^L & \cdots & w_{nn}^L \end{pmatrix} \quad \mathbf{b}^L = \begin{pmatrix} b_0^L \\ b_1^L \\ \vdots \\ b_n^L \end{pmatrix}$$

$$\mathbf{z}^L = W^L \cdot \mathbf{a}^{L-1} + \mathbf{b}^L = \begin{pmatrix} w_{00}^L a_0^{L-1} + w_{01}^L a_1^{L-1} + \cdots + w_{0n}^L a_n^{L-1} + b_0^L \\ w_{10}^L a_0^{L-1} + w_{11}^L a_1^{L-1} + \cdots + w_{1n}^L a_n^{L-1} + b_1^L \\ \vdots \\ w_{n0}^L a_0^{L-1} + w_{n1}^L a_1^{L-1} + \cdots + w_{nn}^L a_n^{L-1} + b_n^L \end{pmatrix} = \begin{pmatrix} z_0 \\ z_1 \\ \vdots \\ z_n \end{pmatrix}$$

$$z_j^L = \sum_{k=0}^n w_{jk}^L a_k^{L-1} + b_j^L$$

where  $f()$  is the activation function e.g. sigmoid, tanh, relu etc

$$\mathbf{a}^L = \begin{pmatrix} f(z_0) \\ f(z_1) \\ \vdots \\ f(z_n) \end{pmatrix} = \begin{pmatrix} a_0^L \\ a_1^L \\ \vdots \\ a_n^L \end{pmatrix}$$

## 2 Loss/Cost Function

The loss function is generally for a single training example whereas the cost function is usually the mean of the loss function over a batch (or mini-batch) of training examples. Doing gradient decent with batches (or mini-batches) is more computationally efficient as backpropagation is calculated only as many times as there are batches. The loss/cost function needs to produce a scalar value to make gradient decent possible.

### 2.1 Squared error

Note 1: this is very similar to the mean squared error function however the  $\frac{1}{n}$  is omitted as it is quicker to compute without it and doesn't meaningfully affect the gradient descent.

Note 2:  $\frac{1}{2}$  is multiplied as a constant to the general squared error function to remove the constant from the derivative.

Element form

$$c = \frac{1}{2} \sum_{j=0}^n (t_j - a_j^L)^2 = \frac{1}{2} ((t_0 - a_0^L)^2 + (t_1 - a_1^L)^2 + \dots + (t_n - a_n^L)^2)$$

$$\frac{dc}{da_j^L} = t_j - a_j^L$$

Matrix Form

$$c = \frac{1}{2} ((\mathbf{t} - \mathbf{a}^L) \circ (\mathbf{t} - \mathbf{a}^L))$$

$$\frac{dc}{d\mathbf{a}^L} = \mathbf{t} - \mathbf{a}^L = \begin{pmatrix} t_0 - a_0^L \\ t_1 - a_1^L \\ \vdots \\ t_n - a_n^L \end{pmatrix}$$

### 2.2 Cross entropy

Element form

$$c = - \sum_{j=0}^n t_j \log_b(a_j^L) = -(t_0 \log_b(a_0^L) + t_1 \log_b(a_1^L) + \dots + t_n \log_b(a_n^L))$$

$$\frac{dc}{da_j^L} = -\frac{1}{\ln(b)} \frac{t_j}{a_j^L}$$

Matrix form

$$c = -\mathbf{t}^T \cdot \log_b(\mathbf{a}^L)$$

$$\frac{dc}{d\mathbf{a}^L} = -\frac{1}{\ln(b)} (\mathbf{t} \oslash \mathbf{a}^L) = -\frac{1}{\ln(b)} \begin{pmatrix} \frac{t_0}{a_0^L} \\ \frac{t_1}{a_1^L} \\ \vdots \\ \frac{t_n}{a_n^L} \end{pmatrix}$$

Please see appendix A for full cross entropy derivation.

### 3 Activation Function

Note that the backpropagation is meaningfully different if the activation function is a vector or scalar function.

#### 3.1 Sigmoid

Element form

$$a_j^L = \sigma(z_j^L) = \frac{1}{1 + e^{-z_j^L}}$$

$$\frac{da_j^L}{dz_j^L} = \frac{d\sigma(z_j^L)}{dz_j^L} = \sigma(z_j^L) (1 - \sigma(z_j^L)) = a_j^L (1 - a_j^L)$$

Matrix Form

$$\mathbf{a}^L = \sigma(\mathbf{z}^L) = \frac{1}{1 + e^{-\mathbf{z}^L}}$$

$$\frac{d\mathbf{a}^L}{d\mathbf{z}^L} = \frac{d\sigma(\mathbf{z}^L)}{d\mathbf{z}^L} = \sigma(\mathbf{z}^L) \circ (1 - \sigma(\mathbf{z}^L)) = \mathbf{a}^L \circ (1 - \mathbf{a}^L)$$

#### 3.2 Softmax

Element form

$$a_j^L = \frac{e^{z_j^L}}{\sum_{i=0}^n e^{z_i^L}}$$

where  $\delta$  is the Kronecker delta:  $\delta_{jk} \begin{cases} 1 & \text{if } j = k \\ 0 & \text{if } j \neq k \end{cases}$

$$\frac{da_j^L}{dz_k^L} = a_j^L (\delta_{jk} - a_k^L)$$

Matrix form

$$\mathbf{a}^L = S(\mathbf{z}^L) = \frac{e^{\mathbf{z}^L}}{\sum_{i=0}^n e^{z_i^L}} = \frac{e^{\mathbf{z}^L}}{e^{z_0^L} + e^{z_1^L} + \dots + e^{z_n^L}}$$

$$\frac{d\mathbf{a}^L}{d\mathbf{z}^L} = \frac{dS(\mathbf{z}^L)}{d\mathbf{z}^L} = \begin{pmatrix} a_0^L (1 - a_0^L) & -a_0^L a_1^L & \dots & -a_0^L a_n^L \\ -a_1^L a_0^L & a_1^L (1 - a_1^L) & \dots & -a_1^L a_n^L \\ \vdots & \vdots & \ddots & \vdots \\ -a_n^L a_0^L & -a_n^L a_1^L & \dots & a_n^L (1 - a_n^L) \end{pmatrix}$$

Please see appendix B for full softmax derivation.

## 4 Matrix weights

Element form

$$\begin{aligned}
 z_j^L &= a_k^{L-1} w_{jk}^L + b_j^L \\
 \frac{dz_j^L}{dw_{jk}^L} &= a_k^{L-1} \\
 \frac{dz_j^L}{db_j^L} &= 1 \\
 \frac{dz_j^L}{da_k^{L-1}} &= w_{jk}^L
 \end{aligned}$$

Matrix form

$$\begin{aligned}
 \mathbf{z}^L &= W^L \cdot \mathbf{a}^{L-1} + \mathbf{b}^L \\
 &= \begin{pmatrix} w_{00} & w_{01} & \cdots & w_{0n} \\ w_{10} & w_{11} & \cdots & w_{1n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{n0} & w_{n1} & \cdots & w_{nn} \end{pmatrix} \begin{pmatrix} a_0^{L-1} \\ a_1^{L-1} \\ \vdots \\ a_n^{L-1} \end{pmatrix} + \begin{pmatrix} b_0^L \\ b_1^L \\ \vdots \\ b_n^L \end{pmatrix} \\
 \frac{d\mathbf{z}^L}{dW^L} &= \mathbf{a}^{L-1} \\
 \frac{d\mathbf{z}^L}{d\mathbf{b}^L} &= 1 \\
 \frac{d\mathbf{z}^L}{d\mathbf{a}^{L-1}} &= W^L
 \end{aligned}$$

## 5 Backpropagation

### 5.1 Sigmoid and squared error

#### 5.1.1 Element form

Layer  $L$

$$\begin{aligned}
 \frac{dc}{dw_{jk}^L} &= \frac{dc}{da_j^L} \frac{da_j^L}{dz_j^L} \frac{dz_j^L}{dw_{jk}^L} \\
 \frac{dc}{da_j^L} &= t_j - a_j^L \\
 \frac{da_j^L}{dz_j^L} &= a_j^L (1 - a_j^L) \\
 \frac{dc}{dz_j^L} &= \frac{dc}{da_j^L} \frac{da_j^L}{dz_j^L} = (t_j - a_j^L) a_j^L (1 - a_j^L) = (t_j - a_j^L) (1 - a_j^L) a_j^L \\
 \frac{dz_j^L}{dw_{jk}^L} &= a_k^{L-1} \\
 \frac{dc}{dw_{jk}^L} &= \frac{dc}{dz_j^L} \frac{dz_j^L}{dw_{jk}^L} = (t_j - a_j^L) (1 - a_j^L) a_j^L a_k^{L-1}
 \end{aligned}$$

Layer  $L - 1$

$$\begin{aligned}
 \frac{dc}{dw_{kq}^{L-1}} &= \left( \sum_{j=0}^n \frac{dc}{dz_j^L} \frac{dz_j^L}{da_k^{L-1}} \right) \frac{da_k^{L-1}}{dz_k^{L-1}} \frac{dz_k^{L-1}}{dw_{kq}^{L-1}} \\
 \frac{dz_j^L}{da_k^{L-1}} &= w_{jk}^L \\
 \frac{dc}{da_k^{L-1}} &= \sum_{j=0}^n \frac{dc}{dz_j^L} \frac{dz_j^L}{da_k^{L-1}} = \sum_{j=0}^n (t_j - a_j^L) (1 - a_j^L) a_j^L w_{jk}^L \\
 \frac{da_k^{L-1}}{dz_k^{L-1}} &= a_k^{L-1} (1 - a_k^{L-1}) \\
 \frac{dz_k^{L-1}}{dw_{kq}^{L-1}} &= a_q^{L-1} \\
 \frac{dc}{dw_{kq}^{L-1}} &= \frac{dc}{da_k^{L-1}} \frac{da_k^{L-1}}{dz_k^{L-1}} \frac{dz_k^{L-1}}{dw_{kq}^{L-1}} = \left( \sum_{j=0}^n (t_j - a_j^L) (1 - a_j^L) a_j^L w_{jk}^L \right) (a_k^{L-1} (1 - a_k^{L-1})) (a_q^{L-1})
 \end{aligned}$$

### 5.1.2 Matrix form

Layer  $L$

$$\begin{aligned} \frac{dc}{dW^L} &= \left( \frac{dc}{d\mathbf{a}^L} \circ \frac{d\mathbf{a}^L}{dz^L} \right) \cdot \frac{dz^L}{dW^L}{}^T \\ &= \left( \begin{pmatrix} t_0 - a_0^L \\ t_1 - a_1^L \\ \vdots \\ t_n - a_n^L \end{pmatrix} \circ \begin{pmatrix} a_0^L (1 - a_0^L) \\ a_1^L (1 - a_1^L) \\ \vdots \\ a_n^L (1 - a_n^L) \end{pmatrix} \right) \cdot \begin{pmatrix} a_0^{L-1} \\ a_1^{L-1} \\ \vdots \\ a_n^{L-1} \end{pmatrix}{}^T \\ &= \begin{pmatrix} (t_0 - a_0^L) (1 - a_0^L) a_0^L a_0^{L-1} & (t_0 - a_0^L) (1 - a_0^L) a_0^L a_1^{L-1} & \cdots & (t_0 - a_0^L) (1 - a_0^L) a_0^L a_n^{L-1} \\ (t_1 - a_1^L) (1 - a_1^L) a_1^L a_0^{L-1} & (t_1 - a_1^L) (1 - a_1^L) a_1^L a_1^{L-1} & \cdots & (t_1 - a_1^L) (1 - a_1^L) a_1^L a_n^{L-1} \\ \vdots & \vdots & \ddots & \vdots \\ (t_n - a_n^L) (1 - a_n^L) a_n^L a_0^{L-1} & (t_n - a_n^L) (1 - a_n^L) a_n^L a_1^{L-1} & \cdots & (t_n - a_n^L) (1 - a_n^L) a_n^L a_n^{L-1} \end{pmatrix} \end{aligned}$$

Layer  $L - 1$

$$\begin{aligned} \frac{dc}{dW^{L-1}} &= \left( \left( \frac{dc}{dz^L} \cdot \frac{dz^L}{da^{L-1}} \right) \circ \frac{da^{L-1}}{dz^{L-1}} \right) \cdot \frac{dz^{L-1}}{dw^{L-1}} \\ \frac{dz^L}{da^{L-1}} &= W^L \\ \frac{da^{L-1}}{dz^{L-1}} &= \begin{pmatrix} a_0^{L-1} (1 - a_0^{L-1}) \\ a_1^{L-1} (1 - a_1^{L-1}) \\ \vdots \\ a_n^{L-1} (1 - a_n^{L-1}) \end{pmatrix} \\ \frac{dz^{L-1}}{dw^{L-1}} &= a^{L-2} \end{aligned}$$

## 5.2 Softmax and cross entropy

### 5.2.1 Element form

On the left  $\mathbf{a}^L$  is left in vector form and the dot product is taken. Instead on the right the summation of the elements is calculated but they are equivalent.

$$\begin{aligned} \frac{dc}{dw_{jk}^L} &= \left( \frac{dc}{d\mathbf{a}^L}^T \cdot \frac{d\mathbf{a}^L}{dz_j^L} \right) \frac{dz_j^L}{dw_{jk}^L} & \frac{dc}{dw_{jk}^L} &= \left( \sum_{i=0}^n \frac{dc}{da_i^L} \frac{da_i^L}{dz_j^L} \right) \frac{dz_j^L}{dw_{jk}^L} \\ \frac{dc}{d\mathbf{a}^L} &= -\frac{1}{\ln(b)} \begin{pmatrix} \frac{t_0}{a_0^L} \\ \frac{t_1}{a_1^L} \\ \vdots \\ \frac{t_n}{a_n^L} \end{pmatrix} & \frac{dc}{da_i^L} &= -\frac{1}{\ln(b)} \frac{t_i}{a_i^L} \\ \frac{d\mathbf{a}^L}{dz_j^L} &= \begin{pmatrix} a_0^L (\delta_{0j} - a_j^L) \\ a_1^L (\delta_{1j} - a_j^L) \\ \vdots \\ a_n^L (\delta_{nj} - a_j^L) \end{pmatrix} & \frac{da_i^L}{dz_j^L} &= a_i^L (\delta_{ij} - a_j^L) \\ \frac{dc}{dz_j^L} &= \frac{dc}{d\mathbf{a}^L}^T \cdot \frac{d\mathbf{a}^L}{dz_j^L} & \frac{dc}{dz_j^L} &= \sum_{i=0}^n \frac{dc}{da_i^L} \frac{da_i^L}{dz_j^L} \end{aligned}$$

The multiplication of both forms is too long to write side by side

$$\begin{aligned} \frac{dc}{dz_j^L} &= \sum_{i=0}^n \frac{dc}{da_i^L} \frac{da_i^L}{dz_j^L} = \sum_{i=0}^n \left( -\frac{1}{\ln(b)} \frac{t_i}{a_i^L} a_i^L (\delta_{ij} - a_j^L) \right) = -\frac{1}{\ln(b)} \sum_{i=0}^n \left( \frac{t_i}{a_i^L} a_i^L (\delta_{ij} - a_j^L) \right) = \\ \frac{dc}{dz_j^L} &= \frac{dc}{d\mathbf{a}^L}^T \cdot \frac{d\mathbf{a}^L}{dz_j^L} = -\frac{1}{\ln(b)} \begin{pmatrix} \frac{t_0}{a_0^L} \\ \frac{t_1}{a_1^L} \\ \vdots \\ \frac{t_n}{a_n^L} \end{pmatrix}^T \cdot \begin{pmatrix} a_0^L (\delta_{0j} - a_j^L) \\ a_1^L (\delta_{1j} - a_j^L) \\ \vdots \\ a_n^L (\delta_{nj} - a_j^L) \end{pmatrix} = -\frac{1}{\ln(b)} \sum_{i=0}^n \left( \frac{t_i}{a_i^L} a_i^L (\delta_{ij} - a_j^L) \right) = \end{aligned}$$

The two forms now converge to the same equation. The Kronecker delta can be removed from the equation by taking the single case  $i = j$  out of the equation and summing across all other cases ( $i \neq j$ ).

$$-\frac{1}{\ln(b)} \left( t_j (1 - a_j^L) - \sum_{i \neq j}^n t_i a_j^L \right) = -\frac{1}{\ln(b)} \left( t_j (1 - a_j^L) - a_j^L \sum_{i \neq j}^n t_i \right)$$

As all the outputs of the softmax function sum to 1 ( $\sum_i^n t_i = 1$ ). We get the following that allows us to remove the summation:  $\sum_{i \neq j}^n t_i = \sum_i^n t_i - t_j = 1 - t_j$

$$\begin{aligned} &= -\frac{1}{\ln(b)} (t_j (1 - a_j^L) - a_j^L (1 - t_j)) = \\ &= -\frac{1}{\ln(b)} (t_j - a_j^L t_j - a_j^L + a_j^L t_j) = -\frac{1}{\ln(b)} (t_j - a_j^L) \\ &= \frac{dc}{dz_j^L} = \frac{1}{\ln(b)} (a_j^L - t_j) \end{aligned}$$



Now to use  $\frac{dc}{dz_j^L}$  in the final calculation

$$\frac{dz_j^L}{dw_{jk}^L} = a_k^{L-1}$$

$$\frac{dc}{dw_{jk}^L} = \left( \frac{dc}{d\mathbf{a}^L}^T \cdot \frac{d\mathbf{a}^L}{dz_j^L} \right) \frac{dz_j^L}{dw_{jk}^L} = \frac{dc}{dz_j^L} \frac{dz_j^L}{dw_{jk}^L} = \frac{1}{\ln(b)} (a_j^L - t_j) a_k^{L-1}$$

### 5.2.2 Matrix form

$$\frac{dc}{dW^L} = \left( \frac{dc}{d\mathbf{a}^L}^T \cdot \frac{d\mathbf{a}^L}{dz^L} \right) \cdot \frac{dz^L}{dW^L}$$

$$\frac{dc}{d\mathbf{a}^L} = -\frac{1}{\ln(b)} \begin{pmatrix} \frac{t_0}{a_0^L} \\ \frac{t_1}{a_1^L} \\ \vdots \\ \frac{t_n}{a_n^L} \end{pmatrix}$$

$$\frac{d\mathbf{a}^L}{dz^L} = \begin{pmatrix} a_0^L(1-a_0^L) & -a_0^L a_1^L & \cdots & -a_0^L a_n^L \\ -a_1^L a_0^L & a_1^L(1-a_1^L) & \cdots & -a_1^L a_n^L \\ \vdots & \vdots & \ddots & \vdots \\ -a_n^L a_0^L & -a_n^L a_1^L & \cdots & a_n^L(1-a_n^L) \end{pmatrix}$$

$$\frac{dc}{dz^L} = \frac{dc}{d\mathbf{a}^L}^T \cdot \frac{d\mathbf{a}^L}{dz^L} = -\frac{1}{\ln(b)} \begin{pmatrix} t_0(1-a_0^L) - t_1 a_0^L + \cdots - t_n a_0^L \\ -t_0 a_1^L + t_1(1-a_1^L) + \cdots - t_n a_1^L \\ \vdots \\ -t_0 a_n^L - t_1 a_n^L + \cdots + t_n(1-a_n^L) \end{pmatrix}$$

Please see element form above for the simplification of the summation

$$= -\frac{1}{\ln(b)} \begin{pmatrix} t_0 - a_0^L \\ t_1 - a_1^L \\ \vdots \\ t_n - a_n^L \end{pmatrix} = -\frac{1}{\ln(b)} (\mathbf{t} - \mathbf{a}^L) = \frac{1}{\ln(b)} (\mathbf{a}^L - \mathbf{t})$$

$$\frac{dz^L}{dW^L} = \begin{pmatrix} a_0^{L-1} \\ a_1^{L-1} \\ \vdots \\ a_n^{L-1} \end{pmatrix}$$

$$\frac{dc}{dW^L} = \frac{dc}{dz^L} \cdot \frac{dz^L}{dW^L} = \frac{1}{\ln(b)} \begin{pmatrix} (a_0^L - t_0) a_0^{L-1} & (a_0^L - t_0) a_1^{L-1} & \cdots & (a_0^L - t_0) a_n^{L-1} \\ (a_1^L - t_1) a_0^{L-1} & (a_1^L - t_1) a_1^{L-1} & \cdots & (a_1^L - t_1) a_n^{L-1} \\ \vdots & \vdots & \ddots & \vdots \\ (a_n^L - t_n) a_0^{L-1} & (a_n^L - t_n) a_1^{L-1} & \cdots & (a_n^L - t_n) a_n^{L-1} \end{pmatrix}$$

$$= \frac{1}{\ln(b)} ((\mathbf{a}^L - \mathbf{t}) \cdot \mathbf{a}^{L-1^T})$$

## A Cross entropy derivation

To get the derivative of cross entropy, the following rules are required:

Derivative of a scalar function with respect to a vector is a vector

$$\frac{dy}{d\mathbf{x}} = \begin{pmatrix} \frac{dy}{dx_0} \\ \frac{dy}{dx_1} \\ \vdots \\ \frac{dy}{dx_n} \end{pmatrix}$$

Derivative of a logarithm

$$f(x) = \log_b(x) \quad f'(x) = \frac{1}{x \ln(b)}$$

Chain rule

$$c(z) = f(z) + g(z) \quad c'(z) = f'(z) + g'(z)$$

The cross entropy function takes a vector as an input and produces a scalar as an output so it is a scalar function.

$$f: \mathbb{R}^N \rightarrow \mathbb{R}$$

The cross entropy function is:

$$c = - \sum_{i=0}^n t_i (\log_b(a_i)) = -(t_0 \log_b(a_0^L) + t_1 \log_b(a_1^L) + \dots + t_n \log_b(a_n^L))$$

The derivative of the a scalar function with with respect to a vector will be a vector of all the partial derivatives.

$$\frac{dc}{d\mathbf{a}^L} = \begin{pmatrix} \frac{dc}{da_0^L} \\ \frac{dc}{da_1^L} \\ \vdots \\ \frac{dc}{da_n^L} \end{pmatrix}$$

So to find the derivative for an element in  $\frac{dc}{da^L}$  where  $j$  is the element

$$\frac{dc}{da_j^L} = - \frac{d}{da_j^L} \sum_{i=0}^n t_i \log_b(a_i^L) = - \frac{d}{da_j^L} (t_0 \log_b(a_0^L) + t_1 \log_b(a_1^L) + \dots + t_n \log_b(a_n^L))$$

Every term where  $i \neq j$  will be 0 so we can ignore all terms except where  $i = j$ :

$$\frac{dc}{da_j^L} = -\frac{d}{da_j^L} t_j(\log_b(a_j)) = -t_j \frac{d}{da_j^L} (\log_b(a_j^L)) = -\frac{t_j}{a_j^L \ln(b)} = -\frac{1}{\ln(b)} \frac{t_j}{a_j^L}$$

So the derivative is:

$$\frac{dc}{d\mathbf{a}^L} = \begin{pmatrix} \frac{dc}{da_0^L} \\ \frac{dc}{da_1^L} \\ \vdots \\ \frac{dc}{da_n^L} \end{pmatrix} = -\frac{1}{\ln(b)} \begin{pmatrix} \frac{t_0}{a_0^L} \\ \frac{t_1}{a_1^L} \\ \vdots \\ \frac{t_n}{a_n^L} \end{pmatrix} = -\frac{1}{\ln(b)} \frac{\mathbf{t}}{\mathbf{a}^L}$$

## B Softmax derivation

To get the derivative of softmax, the following rules are required:

Derivative of a vector function with respect to a vector is a jacobian matrix

$$\frac{d\mathbf{y}}{d\mathbf{x}} = \begin{pmatrix} \frac{dy_0}{dx_0} & \frac{dy_1}{dx_0} & \cdots & \frac{dy_n}{dx_0} \\ \frac{dy_0}{dx_1} & \frac{dy_1}{dx_1} & \cdots & \frac{dy_n}{dx_1} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{dy_0}{dx_n} & \frac{dy_1}{dx_n} & \cdots & \frac{dy_n}{dx_n} \end{pmatrix}$$

total derivative is the sum of all partial derivatives

$$f'(x, y) = f'(x) + f'(y)$$

law of independence

$$\frac{de^z}{de^x} = 0$$

derivative of exponentials

$$f(x) = e^x \quad f'(x) = e^x$$

chain rule

$$c(z) = f(z) + g(z) \quad c'(z) = f'(z) + g'(z)$$

quotient rule

$$q(z) = \frac{f(z)}{g(z)} \quad q'(z) = \frac{g(z) f'(z) - f(z) g'(z)}{g(z)^2}$$

The softmax function takes a vector as an input and produces a vector as an output so it is a vector function.

$$f: \mathbb{R}^N \rightarrow \mathbb{R}^N$$

The softmax function is:

$$a^L = \frac{e^{z^L}}{\sum_{i=0}^n e^{z_i^L}} = \begin{pmatrix} \frac{e^{z_0^L}}{\sum_{i=0}^n e^{z_i^L}} \\ \frac{e^{z_1^L}}{\sum_{i=0}^n e^{z_i^L}} \\ \vdots \\ \frac{e^{z_n^L}}{\sum_{i=0}^n e^{z_i^L}} \end{pmatrix}$$

The derivative of the a vector function with respect to a vector will be the jacobian matrix of all the partial derivatives.

$$\frac{d\mathbf{a}}{d\mathbf{z}} = \begin{pmatrix} \frac{da_0}{dz_0} + \frac{da_0}{dz_1} + \dots + \frac{da_0}{dz_n} \\ \frac{da_1}{dz_0} + \frac{da_1}{dz_1} + \dots + \frac{da_1}{dz_n} \\ \vdots \\ \frac{da_n}{dz_0} + \frac{da_n}{dz_1} + \dots + \frac{da_n}{dz_n} \end{pmatrix}$$

Let's find the partial derivative  $\frac{da_i}{dz_j}$ . To do this we will break up the numerator and denominator into separate functions to make applying the quotient rule easier. Note that  $\frac{df(z_i)}{dz}$  has two separate cases for when  $j = k$  and  $j \neq k$ .

$$a_j = S(z_j) = \frac{e^{z_j}}{\sum_{i=0}^n e^{z_i}} = \frac{f(z_j)}{g(z_j)}$$

$$\text{where } j = k \quad \begin{matrix} f(z_j) = e^{z_j} \\ \frac{df(z_j)}{dz_k} = \frac{df(z_j)}{dz_j} = e^{z_j} \end{matrix}$$

$$\text{where } j \neq k \quad \frac{df(z_j)}{dz_k} = 0$$

$$g(z_j) = \sum_{i=0}^n e^{z_i} = e^{z_0} + e^{z_1} + \dots + e^{z_n} \quad \frac{dg(z_j)}{dz_k} = e^{z_k}$$

when  $j = k$

$$\begin{aligned} \frac{da_j}{dz_k} &= \frac{da_j}{dz_j} = \frac{dS(z_j)}{dz_j} = \frac{g(z_j) e^{z_j} - e^{z_j} e^{z_j}}{g(z_j)^2} = \frac{e^{z_j} (g(z_j) - e^{z_j})}{g(z_j)^2} = \frac{e^{z_j}}{g(z_j)} \frac{g(z_j) - e^{z_j}}{g(z_j)} \\ &= \frac{e^{z_j}}{g(z_j)} \left( \frac{g(z_j)}{g(z_j)} - \frac{e^{z_j}}{g(z_j)} \right) = S(z_j) (1 - S(z_j)) = a_j (1 - a_j) \end{aligned}$$

when  $j \neq k$

$$\frac{da_j}{dz_k} = \frac{dS(z_j)}{dz_k} = \frac{g(z_j) 0 - e^{z_j} e^{z_k}}{g(z_j)^2} = \frac{-e^{z_j} e^{z_k}}{g(z_j)^2} = \frac{-e^{z_j}}{g(z_j)} \frac{e^{z_k}}{g(z_j)} = -S(z_j) S(z_k) = -a_j a_k$$

We can create the jacobian matrix which contains all the partial derivatives for all elements. Note  $I_n$  is an identity matrix of size  $n \times n$ .

$$\begin{aligned}
\frac{d\mathbf{a}^L}{dz^L} &= \begin{pmatrix} a_0^L (1 - a_0^L) & -a_0^L a_1^L & \cdots & -a_0^L a_n^L \\ -a_1^L a_0^L & a_1^L (1 - a_1^L) & \cdots & -a_1^L a_n^L \\ \vdots & \vdots & \ddots & \vdots \\ -a_n^L a_0^L & -a_n^L a_1^L & \cdots & a_n^L (1 - a_n^L) \end{pmatrix} \\
&= \begin{pmatrix} a_0^L & a_1^L & \cdots & a_n^L \end{pmatrix} \begin{pmatrix} (1 - a_0^L) & -a_1^L & \cdots & -a_n^L \\ -a_0^L & (1 - a_1^L) & \cdots & -a_n^L \\ \vdots & \vdots & \ddots & \vdots \\ -a_0^L & -a_1^L & \cdots & (1 - a_n^L) \end{pmatrix} \\
&= \begin{pmatrix} a_0^L & a_1^L & \cdots & a_n^L \end{pmatrix} \left( I_n - \begin{pmatrix} a_0^L & a_1^L & \cdots & a_n^L \\ a_0^L & a_1^L & \cdots & a_n^L \\ \vdots & \vdots & \ddots & \vdots \\ a_0^L & a_1^L & \cdots & a_n^L \end{pmatrix} \right)
\end{aligned}$$

This can also be written with the Kronecker delta where  $\delta_{jk} \begin{cases} 1 & \text{if } j = k \\ 0 & \text{if } j \neq k \end{cases}$

$$\frac{da_j^L}{dz_k^L} = a_j^L (\delta_{jk} - a_k^L)$$