$$\frac{da}{dw} = \underset{④}{\frac{\partial a}{\partial z}} \; \underset{③}{\frac{\partial z}{\partial s}} \; \underset{②}{\frac{\partial s}{\partial H}} \; \underset{①}{\frac{\partial H}{\partial w}}$$

$$S \to \text{Sum} \; (\underset{\underset{H}{\curvearrowright}}{w \odot x})$$

By Jacobian differentiation,

① $\frac{\partial H}{\partial w} = \begin{bmatrix} x_1 & 0 & \cdots 0 & 0 \\ 0 & x_2 & & \vdots \\ \vdots & & x_s & \vdots \\ 0 & \cdots & 0 & x_n \end{bmatrix}$  $n$ inputs

(with $n$ across the top)

② By Jacobian & derivative of a sum

$\overset{H_1, H_2 \dots H_n}{\underset{i=1}{\overset{n}{\sum}}} w_i x_i = \begin{bmatrix} - & S & - \end{bmatrix}$

$\frac{\partial s}{\partial H} = \underbrace{\begin{bmatrix} 1 & 1 & 1 & 1 & \cdots & 1 \end{bmatrix}}_{n \; times}$   Each derivative has only 1 component that returns a 1.

③ $Z = S - b$

$\frac{\partial z}{\partial s} = 1$ (Scalar)

④ 0 or 1, depending on if $z \geq 0$, or $z \leq 0$ (scalar)

End result:

$\frac{da}{dw} = (1) \begin{bmatrix} 1 \end{bmatrix} \overset{n}{\begin{bmatrix} & & \\ & & \end{bmatrix}}^n \overset{n}{\begin{bmatrix} & & \\ & & \end{bmatrix}}$

$= 1 \begin{bmatrix} x_1, x_2, x_3, \dots, x_n \end{bmatrix}$

For $\frac{\partial a}{\partial b}$,   $\frac{\partial a}{\partial b} = \frac{\partial a}{\partial z} \frac{\partial z}{\partial b} 1$

$= 0$ or $1$ (Scalar)

$$\text{MSE}: \quad \frac{1}{2m}\sum_{i=1}^{m}(\overset{\overset{\text{desired}}{\downarrow}}{y_i}-\overset{\overset{\text{outputted}}{\downarrow}}{\hat{y}_i})^2$$

$$\frac{\partial c}{\partial w} = \frac{\partial c}{\partial a}\frac{\partial a}{\partial w}$$

let $v = y - a^L$

$$\frac{1}{2m}\sum_{i=1}^{m}(v)^2$$

$$\frac{\partial c}{\partial w} = \frac{\partial c}{\partial v}\overset{①}{}\frac{\partial v}{\partial w}\overset{②}{}$$

$$\frac{\partial v}{\partial w} = \frac{\partial}{\partial w}(y-a^L) \qquad \text{where } v = y - a^L$$

$$\frac{\partial v}{\partial w} = -(1)\left(\frac{\partial a}{\partial w}\right) \qquad (\text{chain rule again})$$

② $\quad \dfrac{\partial c}{\partial v} = \dfrac{1}{2m}\sum_{i=1}^{m}(2v)$

$$= \frac{1}{m}\sum_{i=1}^{m}v$$

$$\frac{\partial c}{\partial v} = \frac{1}{m}\sum_{i=1}^{m}(v)$$

$$\frac{\partial c}{\partial w} = \frac{\partial c}{\partial v}\frac{\partial v}{\partial w}$$

$$= \frac{1}{m}\sum_{i=1}^{m}v(-1)\left(\frac{\partial a}{\partial w}\right)$$

$$= \frac{1}{m}\sum_{i=1}^{m}\begin{cases} -(0)^T \\ -v\,\frac{\partial z}{\partial w} \end{cases}$$

$$\frac{\partial c}{\partial w} = \frac{1}{m}\sum_{i=1}^{m}\begin{cases} \overset{T}{0} \\ -(y_i-(w^T x + b))(x_1, x_2, \dots x_n) \end{cases}$$

if $w^T x + b \geq 0$, there is no need for a second max function.

$$\frac{\partial c}{\partial w} = \frac{1}{m}\sum_{i=1}^{m}\underbrace{(w^T x + b - y_i)}_{\text{scalar term}}\,[x_1, x_2 \dots x_n]$$

what does this looks like?

if we define $(w^t x + b - y_i)$ as $e_i$

$$\frac{1}{m} \begin{bmatrix} e_1 x_1 + e_2 x_2 + \ldots + e_n x_1 \\ e_2 x_2 + e_2 x_2 + \ldots + e_n x_2 \\ \vdots \\ e_n x_n \ e_2 x_n + \ldots + e_n x_n \end{bmatrix}$$

$\frac{\partial c}{\partial v}$ wrt all $v$s, averaged over all training examples.

Derivative of the cost wrt the bias:

$$\frac{\partial c}{\partial b} = \frac{\partial v}{\partial a} \frac{\partial a}{\partial b}$$

$$= (-1)(v) \begin{cases} 0 \\ 1 \end{cases}$$

$$= \begin{cases} 0 \\ -(y - a^L) \end{cases}$$

$$= \frac{1}{m} \sum_{i=1}^{m} \underbrace{(w^t x + b - y)}_{e_i} = \frac{1}{m} \sum_{i=1}^{m} e_i \quad (\text{scalar})$$

Four equations of back propagation:

$$\delta_j^l = \frac{\partial C}{\partial a_j^l}\, \sigma'(z_j^l)$$

① $\quad \delta^L = \nabla_a C \odot \sigma'(z^L)$

② $\quad \delta^l = ((w^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l)$

③ $\quad \dfrac{\partial C}{\partial b} = \delta$

④ $\quad \dfrac{\partial C}{\partial w} = (a^{l-1})^T \delta^l$

Imagine we want our neural network to be $784 \times 30 \times 10$

self. num_layers = len(sizes) = 3

passing in self.sizes = $[784, 30, 10]$

Looping through sizes 30, 10

Need a bias vector of $30 \times 1$, and $10 \times 1$            list of

Biases. append (np.random.randn(y,1)) => creates a V numpy arrays

$(30,1), (10 \times 1)$ numpy arrays

weights need to be $30 \times 784$, $10 \times 30$ for dot products

using the zip function,

zip (sizes[:-1], sizes[1:]) -> x from first, up excluding last

-> y from second, excluding last

x: 784, 30,

y: 30, 10

(y,x) = $30 \times 784$, $10 \times 30$


starting from the input, feedforward runs through a loop

For a training example $784 \times 1$



784

$784 \begin{bmatrix} \\ \\ \end{bmatrix} \begin{bmatrix} \\ \\ \end{bmatrix} 30 + \begin{bmatrix} \\ \\ \end{bmatrix} 30$ , then apply sigmoid element wise to all elements

$30 \begin{bmatrix} \\ \\ \end{bmatrix} \begin{bmatrix} \\ \\ \end{bmatrix} 10 + \begin{bmatrix} \\ \\ \end{bmatrix} 10$, apply sigmoid element wise

.30


length of test data => n-test = 10,000

n = length of training data = 50,000


Looping through epochs,

Mini-batches of size mini_batch_size => # of columns = batch size

# of epochs x # of batches = # of times gradient descent is performed. e.g. batches = 500, epochs = 10 => 5000 times

create weight and bias matrices, matching the shapes already initialized

we first initialize the matching matrix with zeros. Running back propagation, we set th lut of matrices that match the shapes, of the gradient matrices.

Then, we subtract the gradients, after multiplying by eta / # of training examples

In the back propagation algorithm, we are passing in a batch of training examples.
Initialize the activation $x$, and activations as a list. Create a list of the weighted sums, and activations

① $\delta^L = \nabla_a C \odot \sigma'(z^L)$

② $\delta^\ell = ((w^{\ell+1})^T \delta^{\ell+1}) \odot \sigma'(z^\ell)$

③ $\frac{\partial C}{\partial b} = \delta$

④ $\frac{\partial C}{\partial w} = (a^{\ell-1})^T \delta^\ell$

From neural networks to implementation of back propogation:

eg. 2 inputs of 784 pixels

$$a_1 = {}_j^{15}\!\left[\quad {}^{784\ K}\quad\right]\left[{}^2\right]784 + \left(\quad\right)^2 15 \implies \left(\quad\right)^2 15$$

kth neuron, into the jth neuron

$${}^{10}\!\left[{}^{75}\quad\right]\left[\quad\right]^2 15 + \left(\quad\right)^2 10 = \left(\quad\right)^2 10 = a_2$$

$${}^{784}\!\left[{}^{4788}\quad\right]\left({}^{784}\quad\right)^{14} + \left(\quad\right)^{10}$$

$$\left(\left[{}^{10}\;{}^{W}\right]^{110}\right)^{10}\left[{}^{X\ 1\cdot0}\quad\right] = {}^{10}\left({}^{120}\quad\right){}^{10\cdot}\left({}^{120}\quad\right)$$