The complete Mathematics for Machine Learning and Deep Learning
By: Adam Dhalla

## 1.3 Notation

$m$ = Size of training Set

$n$ = Amount of input Vars (ie. $x_1, x_2, ..., x_n$)
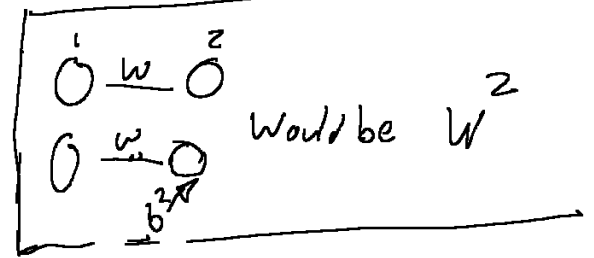
$L$ = Amount of layers in the Neural Network

$\ell$ = Specific layer (could be super/sub script)

$W^{\ell}$ = Weights     $\ell$ = Layer going into

$b^{\ell}$ = bias for laye $\ell$

$x_i$ = Single input Variable



$\bigcirc \xrightarrow{W} \bigcirc$
$\bigcirc \xrightarrow{w} \bigcirc$     Would be $W^2$
$b^2$

From 2.2 → Weights notation $W^{\ell}_{jk}$     where     $j$ = node in layer $\ell$
See page 7 for example          $k$ = node in layer $\ell-1$

# 1.5 Matrix Calculus Review

## 1.5.1 Gradients $(\mathbb{R}^n \to \mathbb{R}^1)$

For a function $F(x, y, ...)$ the gradient would be:

$$\begin{bmatrix} \partial f / \partial x \\ \partial f / \partial y \\ \vdots \end{bmatrix}$$

e.g. for a function $F(x, y) = x^2 + \cos(y)$, the gradient would be

$$\nabla F(x, y) = \begin{bmatrix} \partial f / \partial x \\ \partial f / \partial y \end{bmatrix} = \begin{bmatrix} 2x \\ -\sin(y) \end{bmatrix}$$

✭ A nabla $(\nabla)$ is the sign for gradient.

## 1.5.2 Jacobians ($\mathbb{R}^n \to \mathbb{R}^m$)

Lets say we have a function that is $\mathbb{R}^2 \to \mathbb{R}^2$

$$F(x,y) = \begin{bmatrix} 2x + y^3 \\ e^y - 13x \end{bmatrix} = \begin{bmatrix} f_1 \\ f_2 \end{bmatrix}$$

where:

$$f_1 = 2x + y^3$$
$$f_2 = e^y - 13x$$

Then

$$\frac{\partial f_1}{\partial x} = 2 \qquad\qquad \frac{\partial f_1}{\partial y} = 3y^2$$

$$\frac{\partial f_2}{\partial x} = -13 \qquad\qquad \frac{\partial f_2}{\partial y} = e^y$$

$$J_F(x,y) = \begin{bmatrix} \nabla f_1^T \\ \nabla f_2^T \end{bmatrix} = \begin{bmatrix} \dfrac{\partial f_1}{\partial x} & \dfrac{\partial f_1}{\partial y} \\ \dfrac{\partial f_2}{\partial x} & \dfrac{\partial f_2}{\partial y} \end{bmatrix} = \begin{bmatrix} 2 & 3y^2 \\ -13 & e^y \end{bmatrix}$$

## 1.5.3 New way of seeing the scalar chain rule

Split a function into 2 functions

e.g. to find the derivative of $\sin(x^2)$ we can substitute $g$ for $x^2$.

$$f = \sin(g) \qquad\qquad g = x^2$$

$$\frac{df}{dx} = \frac{\partial f}{\partial g}\frac{\partial g}{\partial x} = \cos(g) \cdot 2x = 2x\cos(x^2)$$

---

## 1.5.4 Jacobian chain rule

$$F(x,y) = \begin{bmatrix} \sin(x^2+y) \\ \ln(y^3) \end{bmatrix}$$

$$g = \begin{bmatrix} x^2+y \\ y^3 \end{bmatrix} \begin{matrix} \leftarrow g_1 \\ \leftarrow g_2 \end{matrix}$$

$$F(x,y) = \begin{bmatrix} \sin(g_1) \\ \ln(g_2) \end{bmatrix}$$

$$\frac{\partial \vec{g}}{\partial \vec{x}} = \begin{bmatrix} 2x & 1 \\ 0 & 3y^2 \end{bmatrix}$$

$$\frac{d\vec{F}}{\partial \vec{g}} = \begin{bmatrix} \cos(g_1) & 0 \\ 0 & \frac{1}{g_2} \end{bmatrix}$$

$$\frac{\partial \vec{F}}{\partial \vec{x}} = \frac{\partial \vec{F}}{\partial \vec{g}}\frac{\partial \vec{g}}{\partial \vec{x}}$$

$$\frac{\partial \vec{F}}{\partial \vec{x}} = \begin{bmatrix} \cos(g_1) & 0 \\ 0 & \frac{1}{g_2} \end{bmatrix}\begin{bmatrix} 2x & 1 \\ 0 & 3y^2 \end{bmatrix}$$

$$\frac{\partial \vec{F}}{\partial \vec{x}} = \begin{bmatrix} 2x\cos(g_1) & \cos(g_1) \\ 0 & \frac{3y^2}{g_2} \end{bmatrix}$$
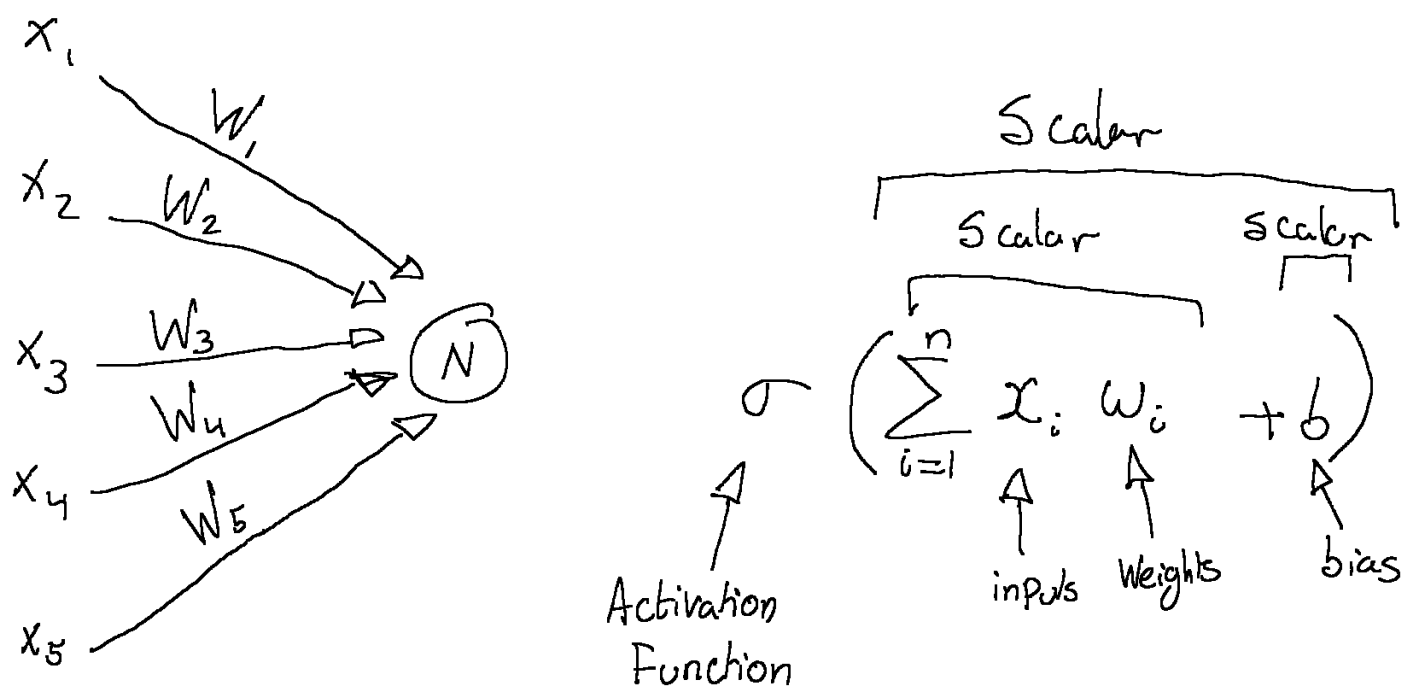
$$\frac{\partial \vec{F}}{\partial \vec{x}} = \begin{bmatrix} 2x\cos(x^2+y) & \cos(x^2+y) \\ 0 & \frac{3y^2}{y^3} \end{bmatrix}$$

# PART II:

## Forward Propogation

# 2.1 The Neuron Function



This can also be done by a dot product:

$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \qquad W = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}$$

$$\sum_{i=1}^{n} x_i w_i = X^T W = x_1 w_1 + x_2 w_2 + \ldots + x_n w_n$$

$$\therefore$$
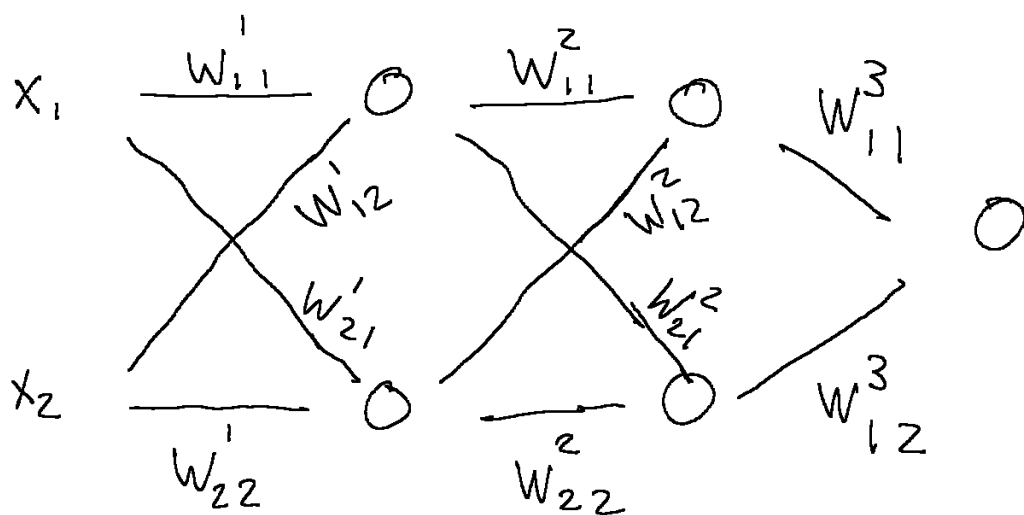
$$\sigma\left(\underbrace{X^T W + b}_{z}\right) \longrightarrow \quad z = X^T W + b$$
$$a = \sigma(z)$$

A node is a Vector-to-Scalar function

# 2.2 Weight and Bias Indexing



Weights Notation: $\boxed{w_{jk}^{\ell}}$

where
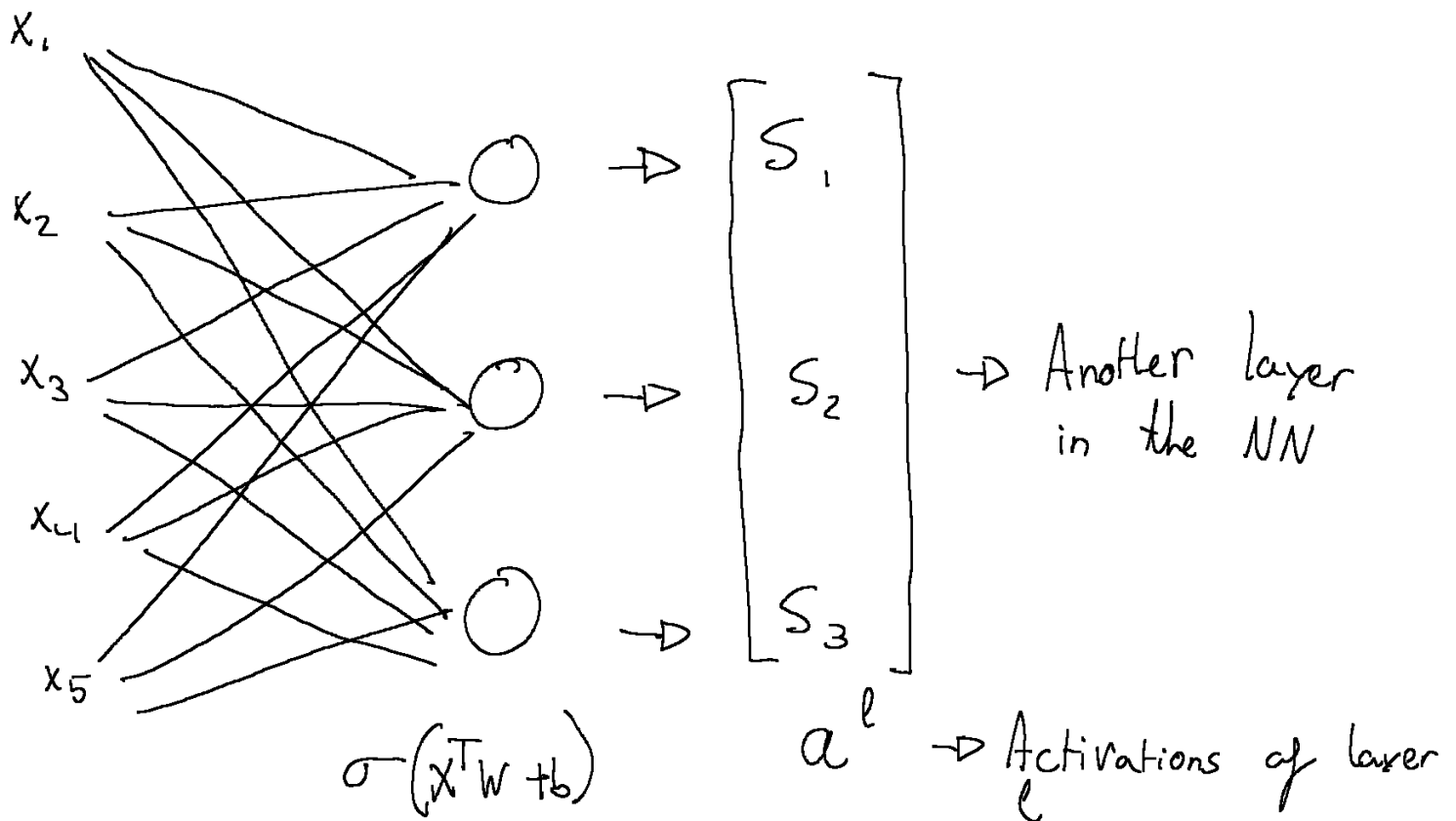
$\ell = $ Layer $w$ is going into

$j = $ Node in layer $\ell$

$k = $ Node in layer $\ell-1$

Bias Notation: $b^{\ell}$

note: Sometimes biases are attached to individual nodes (rather than an entire layer) where the notation would then be $b_j^{\ell}$

# 2.3 A layer of neurons



$$\sigma(X^T W + b)$$

$a^\ell \to$ Activations of layer $\ell$

$\to$ Another layer in the NN

~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~

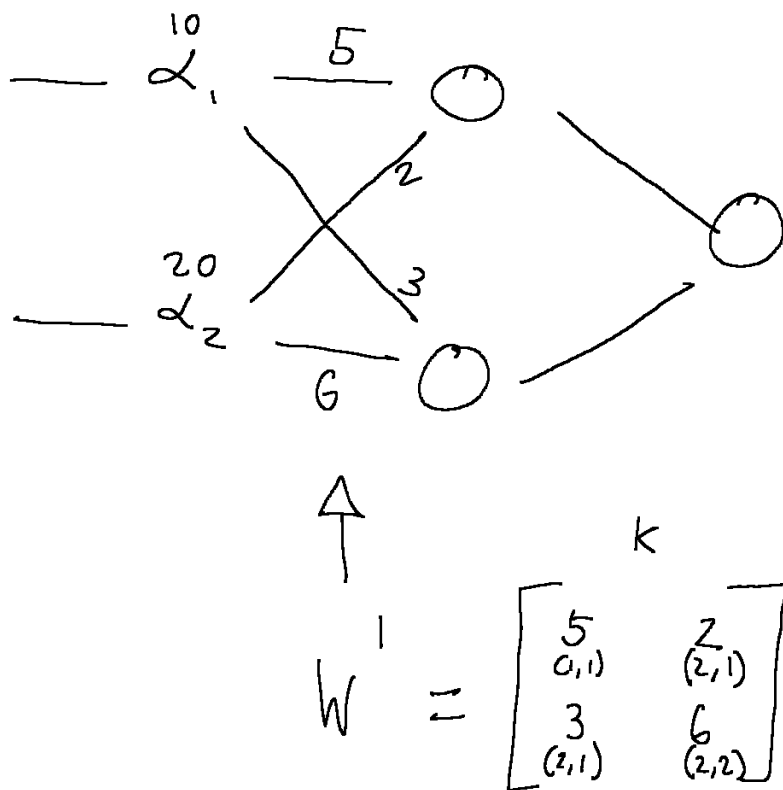$W^\ell \to$ Matrix of weights (Remember $w_{jk}^\ell$, $j^{th}$ row, $k^{th}$ column)

$$W = \begin{bmatrix} W_{11} & W_{12} & \dots & W_{1K} \\ W_{21} & W_{22} & \dots & W_{2K} \\ \vdots & \vdots & & \vdots \\ W_{j1} & W_{j2} & \dots & W_{jK} \end{bmatrix}$$

$\underbrace{\phantom{W_{11} \quad W_{12} \quad \dots \quad W_{1K}}}_{K}$

$j \nwarrow$ Going into neuron position 'j' in layer $\ell$

$\uparrow$ Coming from neuron position 'k' in layer $\ell-1$

Example of a weights matrix



$$W^1 = \begin{bmatrix} 5 \\ {\scriptstyle (1,1)} & \begin{matrix} 2 \\ {\scriptstyle (2,1)} \end{matrix} \\ \begin{matrix} 3 \\ {\scriptstyle (2,1)} \end{matrix} & \begin{matrix} 6 \\ {\scriptstyle (2,2)} \end{matrix} \end{bmatrix} \quad \cup$$

Output of a layer:

$$\sigma\left(W^\ell a^{\ell-1} + b^\ell\right)$$

↑ Activations of Prev layer

$$W^1 a^0 = \begin{bmatrix} 5 & 2 \\ 3 & 6 \end{bmatrix} \begin{bmatrix} 10 \\ 20 \end{bmatrix} = \begin{bmatrix} 90 \\ 150 \end{bmatrix}$$

Then add bias. Lets say $b^1 = 6$

$$z^1 = \begin{bmatrix} 90 \\ 150 \end{bmatrix} + 6 = \begin{bmatrix} 96 \\ 156 \end{bmatrix}$$

Then we have the activation function. Lets say $\sigma = ReLU$

$$a^1 = \sigma(z^1) = \sigma\left(\begin{bmatrix} 96 \\ 156 \end{bmatrix}\right) = \begin{bmatrix} 96 \\ 156 \end{bmatrix}$$

row ↓  ◁ col

In general, our weights matricies are $(n,m)$ where $n =$ nodes in layer $\ell$ and $m =$ nodes in layer $\ell-1$

Another way to represent a NN

$$a^0 \Rightarrow \sigma\left(\omega^1 a^0 + b^1\right) = a^1 \Rightarrow \sigma\left(\omega^2 a^1 + b^2\right) = a^2 \Rightarrow \dots$$

# PART

# III

Derivatives of Neural Networks and Gradient Descent

# 3.1 Motivation and Cost function

Cost function: A mathematical function that measures how well a neural network is performing on the given data.

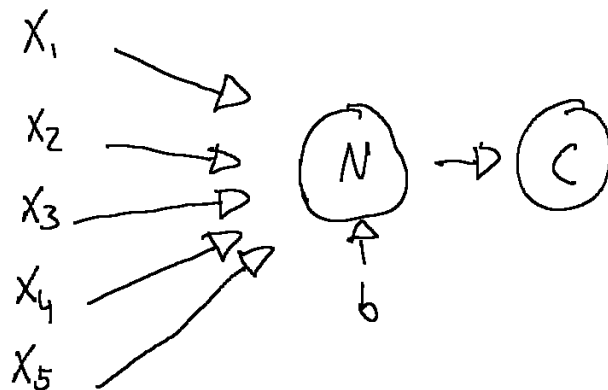## Mean Square Error (MSE)

$$\frac{1}{2m} \sum_{i=1}^{m} \left( y - \hat{y} \right)^2$$

Training Examples → $m$

The 2 makes calculations easier w/ the derivative

Actual training example

Output of network

## Simple Model

$X_1$
$X_2$
$X_3$
$X_4$
$X_5$

$N \rightarrow C$

$b$

# 3.2 Differentiating a neuron's operations

## 3.2.1 Derivative of a binary element-wise operation

Binary element-wise operation: $F(\vec{v}, \vec{w}) \rightarrow \vec{b}$

i.e. function that takes two vectors and returns a single vector where an operation is done 'element wise'.

e.g.

$$\begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix} \oplus \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} = \begin{bmatrix} v_1 + w_1 \\ v_2 + w_2 \\ \vdots \\ v_n + w_n \end{bmatrix}$$

Element-wise addition

$$F(\vec{v}, \vec{w}) = f(\vec{v}) \bigcirc g(\vec{w})$$

Does not need to be elementwise but can have a function to the vectors

Element-wise

✱ An element wise operation can be a variety of things such as addition, multiplication and can even be a comparison (e.g '<') where a vector of binary elements is outputted

A lot of the time in element wise functions,

$$f(\vec{v}) = \vec{v}$$

$$g(\vec{w}) = \vec{w}$$

We can now take the jacobian of an element-wise operation

We are now passing through two vectors, and can get two jacobians out of this. One w.r.t. the elements of $\vec{v}$ and one w.r.t. the elements of $\vec{w}$

$\therefore$ we can find: $\dfrac{\partial F}{\partial \vec{w}}$ and $\dfrac{\partial F}{\partial \vec{v}}$

We will find the jacobian w.r.t. $\vec{v}$ $\left(\dfrac{\partial F}{\partial \vec{v}}\right)$

$$F(\vec{v}, \vec{w}) = \begin{bmatrix} f_1(\vec{v}) \odot g_1(\vec{w}) \\ f_2(\vec{v}) \odot g_2(\vec{w}) \\ \vdots \\ f_n(\vec{v}) \odot g_n(\vec{w}) \end{bmatrix} \begin{matrix} F_1 \\ F_2 \\ \\ F_n \end{matrix}$$

$$\frac{\partial F}{\partial \vec{v}} = \begin{bmatrix} \frac{\partial}{\partial v_1} f_1(\vec{v}) \odot g_1(\vec{w}) & \frac{\partial}{\partial v_2} f_1(\vec{v}) \odot g_1(\vec{w}) & \cdots & \frac{\partial}{\partial v_r} f_1(\vec{v}) \odot g_1(\vec{w}) \\ \frac{\partial}{\partial v_1} f_2(\vec{v}) \odot g_2(\vec{w}) & \frac{\partial}{\partial v_2} f_2(\vec{v}) \odot g_2(\vec{w}) & \cdots & \frac{\partial}{\partial v_n} f_2(\vec{v}) \odot g_2(\vec{w}) \\ & & \vdots & \\ \frac{\partial}{\partial v_1} f_n(\vec{v}) \odot g_n(\vec{w}) & \frac{\partial}{\partial v_2} f_n(\vec{v}) \odot g_n(\vec{w}) & \cdots & \frac{\partial}{\partial v_n} f_n(\vec{v}) \odot g_n(\vec{w}) \end{bmatrix}$$

if $f(\vec{v}) = \vec{v}$ and $g(\vec{w}) = \vec{w}$ then:

$$\frac{\partial F}{\partial \vec{v}} = \begin{bmatrix} \frac{\partial}{\partial v_1} v_1 \odot w_1 & \frac{\partial}{\partial v_2} v_1 \odot w_1 & \cdots & \frac{\partial}{\partial v_n} v_1 \odot w_1 \\ \frac{\partial}{\partial v_1} v_2 \odot w_2 & \frac{\partial}{\partial v_2} v_2 \odot w_2 & \cdots & \frac{\partial}{\partial v_n} v_2 \odot w_2 \\ & & \vdots & \\ \frac{\partial}{\partial v_1} v_n \odot w_n & \frac{\partial}{\partial v_2} v_n \odot w_n & \cdots & \frac{\partial}{\partial v_n} v_n \odot w_n \end{bmatrix}$$

The Jacobian is very often a diagonal as $\frac{\partial}{\partial v_j}\left(F_i(\vec{v}) \odot g_i(\vec{w})\right) = 0$

where $j \neq i$

as $F_i$ and $g_i$ are __not__ functions of $v_j$

$\therefore \frac{\partial F}{\partial \vec{v}}$ Simplifies to

$$\frac{\partial F}{\partial \vec{v}} = \begin{bmatrix} \frac{\partial}{\partial v_1} V_1 \odot w_1, & 0 & , \cdots, & 0 \\ 0 & , \frac{\partial}{V_2} V_2 \odot w_1 & , \cdots, & 0 \\ 0 & , & 0 & , \cdots, \frac{\partial}{V_n} V_n \odot w_n \end{bmatrix}$$

So for all elementwise functions, the jacobian will be diagonal.

This can simplify to

$$\frac{\partial F}{\partial \vec{v}} = \text{diag}\left(\frac{\partial}{\partial v_1} V_1 \odot w_1, \frac{\partial}{\partial v_2} V_2 \odot w_2, \cdots \frac{\partial}{\partial v_n} V_n \odot w_n\right)$$

and

$$\frac{\partial F}{\partial \vec{w}} = \text{diag}\left(\frac{\partial}{\partial w_1} V_1 \odot w_1, \frac{\partial}{\partial w_2} V_2 \odot w_2, \cdots \frac{\partial}{\partial w_n} V_n \odot w_n\right)$$

# 3.2.2. Derivative of a Hadamard Product

Hadamard Product: Element wise multiplies two vectors

$$F(\vec{V}, \vec{w}) = \begin{bmatrix} f_1(\vec{v}) \otimes g_1(\vec{w}) \\ f_2(\vec{v}) \otimes g_2(\vec{w}) \\ \vdots \\ f_n(\vec{v}) \otimes g_n(\vec{w}) \end{bmatrix} = \begin{bmatrix} V_1 \otimes w_1 \\ V_2 \otimes w_2 \\ \vdots \\ V_n \otimes w_n \end{bmatrix} = \vec{V} \otimes \vec{w}$$

Since we are not manipulating $\vec{v}$ and $\vec{w}$ before element wise multiplying them, we can remap $f_n()$ and $g_n()$

$$\therefore \quad \frac{\partial F}{\partial \vec{V}} = \begin{bmatrix} \frac{\partial F_1}{\partial V_1} & \frac{\partial F_1}{\partial V_2} & \cdots & \frac{\partial F_1}{\partial V_n} \\ \frac{\partial F_2}{\partial V_1} & \frac{\partial F_2}{\partial V_2} & \cdots & \frac{\partial F_2}{\partial V_n} \\ \vdots & \vdots & & \vdots \\ \frac{\partial F_n}{\partial V_1} & \frac{\partial F_n}{\partial V_2} & \cdots & \frac{\partial F_n}{\partial V_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial V_1} & 0 & \cdots & 0 \\ 0 & \frac{\partial f_2}{\partial F_2} & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & \frac{\partial f_n}{\partial V_n} \end{bmatrix}$$

The derivative, $\frac{\partial F_n}{\partial V_n}(V_n \otimes W_n) = W_n$, so:

$$\frac{\partial F}{\partial \vec{V}} = \begin{bmatrix} W_1 & 0 & \cdots & 0 \\ 0 & W_2 & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & W_n \end{bmatrix}$$

Equivalently:

$$\frac{\partial F}{\partial \vec{w}} = \begin{bmatrix} V_1 & 0 & \cdots & 0 \\ 0 & V_2 & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & V_n \end{bmatrix}$$

# 3.2.3 Derivative of a Scalar Expansion

Scalar Expansion: Multiplying a scalar w/ a vector

$$2 \cdot \begin{bmatrix} V_1 \\ V_2 \\ \vdots \\ V_n \end{bmatrix} = \begin{bmatrix} 2 \\ 2 \\ \vdots \\ 2 \end{bmatrix} \begin{bmatrix} V_1 \\ V_2 \\ \vdots \\ V_n \end{bmatrix} = \begin{bmatrix} 2V_1 \\ 2V_2 \\ \vdots \\ 2V_n \end{bmatrix}$$

We broadcast/expand the scalar to be a vector of the same size

$$F(\vec{V}, x) = F(\vec{V}) \; O \; g(x) \qquad \overset{\text{scalar}}{\nearrow}$$

where $g(x) = \vec{1} \cdot x$   (The act of multiplying x by the ones vector is an act of broadcasting itself)

$\uparrow$

Expands x to be a vector by multiplying it by the ones vector

$$F(\vec{V}, x) = \begin{bmatrix} f_1(\vec{V}) \; O \; g_1(x) \\ f_2(\vec{V}) \; O \; g_2(x) \\ \vdots \\ f_n(\vec{V}) \; O \; g_n(x) \end{bmatrix}$$

$$\frac{\partial F}{\partial \vec{V}} = \begin{bmatrix} \frac{\partial f_1}{\partial V_1} & \frac{\partial f_1}{\partial V_2} & \cdots & \frac{\partial f_1}{\partial V_n} \\ \frac{\partial f_2}{\partial V_1} & \frac{\partial f_2}{\partial V_2} & \cdots & \frac{\partial f_2}{\partial V_n} \\ \vdots & \vdots & & \vdots \\ \frac{\partial f_n}{\partial V_1} & \frac{\partial f_n}{\partial V_2} & \cdots & \frac{\partial f_n}{\partial V_n} \end{bmatrix} = \begin{bmatrix} x & 0 & \cdots & 0 \\ 0 & x & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & x \end{bmatrix}$$

Since $x$ is a scalar, taking the derivative w.r.b. $x$ will give us a gradient not a jacobian.

If elementwise is multiplication then:

$$\nabla f_x = \begin{bmatrix} \frac{\partial f_1}{\partial x} \\ \frac{\partial f_2}{\partial x} \\ \vdots \\ \frac{\partial f_n}{\partial x} \end{bmatrix} = \begin{bmatrix} V_1 \\ V_2 \\ \vdots \\ V_n \end{bmatrix}$$

If elementwise is addition then

$$\nabla f_x = \begin{bmatrix} \frac{\partial f_1}{\partial x} \\ \frac{\partial f_2}{\partial x} \\ \vdots \\ \frac{\partial f_n}{\partial x} \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$$

If elementwise is subtraction then

$$\nabla f_x = \begin{bmatrix} \frac{\partial f_1}{\partial x} \\ \frac{\partial f_2}{\partial x} \\ \vdots \\ \frac{\partial f_n}{\partial x} \end{bmatrix} = \begin{bmatrix} -1 \\ -1 \\ \vdots \\ -1 \end{bmatrix}$$

# 3.2.4. Derivative of a Sum

$$\vec{V} = \begin{bmatrix} V_1 \\ V_2 \\ \vdots \\ V_n \end{bmatrix} = \sum_{i=1}^{n} g_i(v)$$

A summation is a vector-to-scalar function

The derivative of the sum:

$$\frac{\partial S}{\partial \vec{V}} = \begin{bmatrix} \frac{\partial S}{V_1} & \frac{\partial S}{V_2} & \cdots & \frac{\partial S}{V_n} \end{bmatrix}$$

$$= \begin{bmatrix} \frac{\partial}{\partial v_1} \sum_{i=1}^{n} g_i(v) & \frac{\partial}{\partial v_2} \sum_{i=1}^{n} g_i(v) & \cdots & \frac{\partial}{\partial v_n} \sum_{i=1}^{n} g_i(v) \end{bmatrix}$$

The derivative of a Sum is equivalent to the sum of a derivative

$$\therefore = \begin{bmatrix} \sum_{i=1}^{n} \frac{\partial}{\partial v_1} g_i(v) & \sum_{i=1}^{n} \frac{\partial}{\partial v_2} g_i(v) & \cdots & \sum_{i=1}^{n} \frac{\partial}{\partial v_n} g_i(v) \end{bmatrix}$$

If $g(v) = v$

$$= \begin{bmatrix} \sum_{i=1}^{n} \frac{\partial}{\partial v_1} V_i & \sum_{i=1}^{n} \frac{\partial}{\partial v_2} V_i & \cdots & \sum_{i=1}^{n} \frac{\partial}{\partial v_n} V_i \end{bmatrix}$$
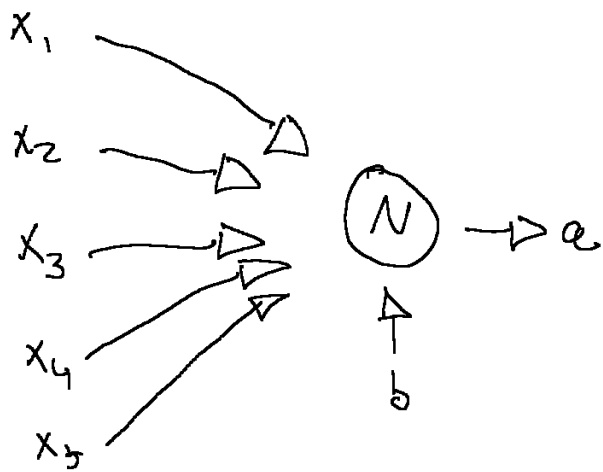
The derivative of $v$, for every summation, will be 0 everywhere except on the element, $v_i$, being derived where it would be 1. So

$$= \begin{bmatrix} 1 & 1 & ... & 1 \end{bmatrix}$$

If instead $g(v) = z \cdot v$, then

$$= \begin{bmatrix} z & z & ... & z \end{bmatrix}$$

# 3.3 Derivative of a neuron's activation



$$a = \sigma(W^T x + b) = \sigma(z)$$

where $z = W^T x + b$

We can further break this down

$$a = \sigma(W^T x + b) = \sigma\left(Sum(W \otimes X) + b\right)$$

$\underbrace{W \otimes X}_{H}$

$\underbrace{Sum(W \otimes X)}_{S(H)}$

We will now investigate the derivative of the weights and bias using the chain rule

weights: $\dfrac{\partial a}{\partial w} = \dfrac{\partial a}{\partial z} \dfrac{\partial z}{\partial w} = \dfrac{\partial a}{\partial z} \dfrac{\partial z}{\partial S} \dfrac{\partial S}{\partial H} \dfrac{\partial H}{\partial w}$

Bias: $\dfrac{\partial a}{\partial b} = \dfrac{\partial a}{\partial z} \dfrac{\partial z}{\partial b}$

We will solve for this on the next 2 pages

We can now find $\dfrac{\partial a}{\partial w}$

$$\frac{\partial a}{\partial w} = \frac{\partial a}{\partial z} \frac{\partial z}{\partial w} = \frac{\partial a}{\partial z} \frac{\partial z}{\partial s} \frac{\partial s}{\partial H} \frac{\partial H}{\partial w}$$

$$\frac{\partial H}{\partial w} = \frac{\partial}{\partial w}(W \otimes x) = \begin{bmatrix} x_1 & 0 \dots 0 \\ 0 & x_2 \dots 0 \\ \vdots & \vdots \ddots \vdots \\ 0 & 0 \dots x_n \end{bmatrix} = diag(x_1 \dots x_n)$$

So:

$$\frac{\partial a}{\partial w} = \frac{\partial a}{\partial z} \frac{\partial z}{\partial s} \frac{\partial s}{\partial H}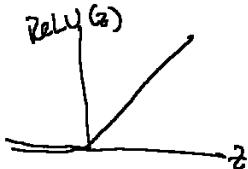 \frac{\partial H}{\partial w} = \frac{\partial a}{\partial z} \frac{\partial z}{\partial s} \frac{\partial s}{\partial H} \, diag(x_1 \dots x_n)$$

$$\frac{\partial s}{\partial H} = \frac{\partial}{\partial H} Sum(x \otimes W) = \vec{1}^T$$

note:
$$\vec{1}^T \cdot diag(x_1 \dots x_n) = x^T$$

So:

$$\frac{\partial a}{\partial w} = \frac{\partial a}{\partial z} \frac{\partial z}{\partial s} \frac{\partial s}{\partial H} \, diag(x_1 \dots x_n) = \frac{\partial a}{\partial z} \frac{\partial z}{\partial s} \vec{1}^T \, diag(x_1 \dots x_n) = \frac{\partial a}{\partial z} \frac{\partial z}{\partial s} x^T$$

$$\frac{\partial z}{\partial s} = \frac{\partial}{\partial s}(s(H) + b) = 1$$

So

$$\frac{\partial a}{\partial w} = \frac{\partial a}{\partial z} \frac{\partial z}{\partial s} x^T = \frac{\partial a}{\partial z} \cdot 1 \cdot x^T = \frac{\partial a}{\partial z} x^T = \frac{\partial a}{\partial z} [x_1, x_2, \dots, x_n]$$

Using $ReLU = max(0, z)$

ReLU(z)


The derivative of ReLU is:

$$\frac{\partial a}{\partial z} = \begin{cases} 0 & \text{if } z \leq 0 \\ 1 & \text{if } z > 0 \end{cases}$$

$$\frac{\partial a}{\partial w} = \frac{\partial a}{\partial z} x^T = \begin{cases} \vec{0}^T & \text{if } w^T x + b \leq 0 \\ \vec{x}^T & \text{if } w^T x + b > 0 \end{cases}$$

We can now find $\dfrac{\partial a}{\partial b}$

$$\frac{\partial a}{\partial b} = \frac{\partial a}{\partial z} \frac{\partial z}{\partial b}$$

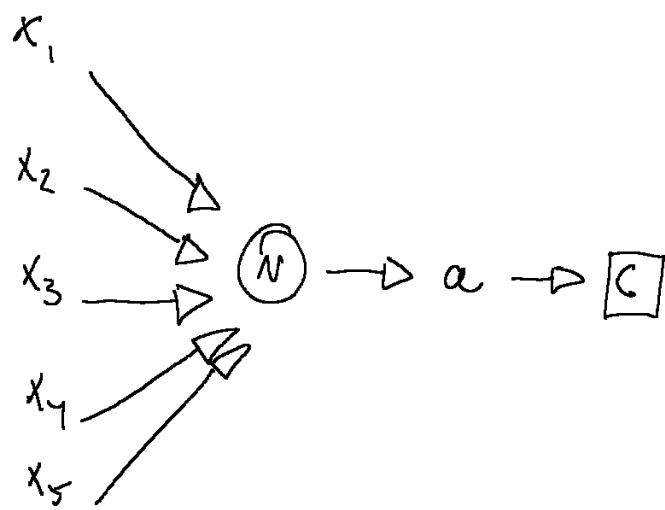$$\frac{\partial z}{\partial b} = 1$$

so

$$\frac{\partial a}{\partial b} = \frac{\partial a}{\partial z} \frac{\partial z}{\partial b} = \frac{\partial a}{\partial z} \cdot 1 = \frac{\partial a}{\partial z}$$

Thus

$$\frac{\partial a}{\partial b} = \frac{\partial a}{\partial z} = \begin{cases} 0 & \text{if } w^T x + b \leq 0 \\ 1 & \text{if } w^T x + b > 0 \end{cases}$$

# 3.4 Derivative of the cost for a simple neural network



Cost (c) is Mean Square Error

$$MSE = \frac{1}{2m} \sum_{i=1}^{m} (y - \hat{y})^2$$

where

$\hat{y}$ = activations = last layer of activations

We want to find $\frac{\partial c}{\partial w}$ and $\frac{\partial c}{\partial b}$

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

where $\frac{\partial a}{\partial w} = \begin{cases} \vec{0}^T & \text{if } w^T x + b \leq 0 \\ \frac{\partial z}{\partial w} & \text{if } w^T x + b > 0 \end{cases}$

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

where $\frac{\partial a}{\partial b} = \begin{cases} 0 & \text{if } w^T x + b \leq 0 \\ 1 & \text{if } w^T x + b > 0 \end{cases}$

Training example is a matrix where each training example is a column

$$X = \begin{bmatrix} x_1^1 & x_1^2 & \dots & x_1^m \\ x_2^1 & x_2^2 & \dots & x_2^m \\ \vdots & \vdots & & \vdots \\ x_n^1 & x_n^2 & \dots & x_n^m \end{bmatrix}$$

$$\frac{1}{2m} \sum_{i=1}^{m} \left(Y - a^L\right) = \frac{1}{2m} \sum_{i=1}^{m} V^2$$

where

$$\left(Y - a^L\right) = V$$

The chain rule then becomes:

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

$$\frac{\partial V}{\partial w} = \frac{\partial}{\partial w} \left(Y - a^L\right) = \frac{\partial}{\partial w} \left(0 - a^L\right) = - \frac{\partial a^L}{\partial w} = \frac{\partial a}{\partial w}$$

We now get

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

Lets find $\frac{\partial c}{\partial V}$

$$\frac{\partial c}{\partial w} = \frac{\partial}{\partial (..)} \frac{1}{2m} \sum_{i=1}^{m} V^2 = \frac{1}{2m} \sum_{i=1}^{m} \frac{\partial}{\partial w} V^2$$

$$= \frac{1}{2m} \sum_{i=1}^{m} \frac{\partial V^2}{\partial V} \frac{\partial V}{\partial w} = \frac{1}{2m} \sum_{i=1}^{m} 2V \cdot - \frac{\partial a}{\partial w}$$

$$= \frac{1}{m} \sum_{i=1}^{m} V \cdot - \frac{\partial a}{\partial w} = \frac{1}{m} \sum_{i=1}^{m} V \begin{cases} \vec{0}^T & \text{if } w^T x + b \leq 0 \\ -\frac{\partial z}{\partial w} & \text{if } w^T x + b > 0 \end{cases}$$

$$= \frac{1}{m} \begin{cases} \vec{0}^T & \text{if } w^T x + b \leq 0 \\ -V \frac{\partial z}{\partial w} & \text{if } w^T x + b > 0 \end{cases} = \frac{1}{m} \begin{cases} \vec{0}^T & \text{if } w^T x + b \leq 0 \\ -V \cdot x^T & \text{if } w^T x + b > 0 \end{cases}$$

$$= \frac{1}{m} \sum_{i=1}^{m} \begin{cases} \vec{0}^T & \text{if } x^T w + b \leq 0 \\ -(y - a^L) x^T & \text{if } x^T w + b > 0 \end{cases}$$

$$= \frac{1}{m} \sum_{i=1}^{m} \begin{cases} \vec{0}^T & \text{if } x^T w + b \leq 0 \\ -(y - \underbrace{\max(0, w^T x + b)}_{\text{This represents}}) x^T & \text{if } x^T w + b > 0 \end{cases}$$

ReLU

The piece wise and the $\max(0, w^T x + b)$ do the same thing.

$$\therefore = \frac{1}{m} \sum_{i=1}^{m} \begin{cases} \vec{0}^T & \text{if } x^T w + b \leq 0 \\ -(y - (w^T x + b)) \cdot x^T & \text{if } x^T w + b > 0 \end{cases}$$

$$= \frac{1}{m} \sum_{i=1}^{m} \begin{cases} \vec{0}^T & \text{if } x^T w + b \leq 0 \\ (w^T x + b - y) \cdot x^T & \text{if } x^T w + b > 0 \end{cases}$$

$$\frac{\partial c}{\partial w} = \frac{1}{m} \begin{cases} \vec{0}^T & \text{if } x^T w + b \leq 0 \\ \sum_{i=1}^{m} (w^T x + b - y) x^T & \text{if } x^T w + b > 0 \end{cases}$$

$$\frac{\partial C}{\partial w} = \frac{1}{m} \begin{cases} \vec{0}^T & \text{if } W^Tx+b \leq 0 \\ \sum_{i=1}^{m} (W^Tx + b - y) x^T & \text{if } W^Tx+b \geq 0 \end{cases}$$

We want :
$$\begin{bmatrix} \frac{\partial C}{\partial w_1} \\ \frac{\partial C}{\partial w_2} \\ \vdots \\ \frac{\partial C}{\partial w_n} \end{bmatrix}$$

i.e. how does the cost change when we change one weight. We use this for gradient descent to know how much/what to change the specific weight by

$W^Tx + b - y = a^L - y_i$  is the error term $(e_i)$

$\therefore$

$$\frac{\partial C}{\partial w} = \frac{1}{m} \begin{cases} \vec{0}^T & \text{if } W^Tx+b \leq 0 \\ \sum_{i=1}^{m} e_i x^T & \text{if } W^Tx+b > 0 \end{cases}$$

The error term, $e_i$, is how incorrect the answer is.

$\therefore$ For every training example, $i$, we get:

$$\frac{\partial c}{\partial \vec{w}} = e_i X^T = \begin{bmatrix} e_i X_1 \\ e_i X_2 \\ \vdots \\ e_i X_n \end{bmatrix} \begin{array}{l} \to \frac{\partial c}{\partial w_1} \\ \to \frac{\partial c}{\partial w_2} \\ \\ \to \frac{\partial c}{\partial w_n} \end{array}$$

The $\frac{\partial c}{\partial \vec{w}}$ vector represents the ratios of change in $C$ when changing the weights by some amount

---

For multiple training examples:

For each training example, $X^T$ is the same but there is a different error, $e_i$

$$\frac{\partial c}{\partial \vec{w}} = \frac{1}{m} \sum_{i=1}^{m} e_i X^T = \frac{1}{m} \cdot \begin{bmatrix} e_1 X_1 + e_2 X_1 + \ldots e_m X_1 \\ e_1 X_2 + e_2 X_2 + \ldots e_m X_2 \\ \vdots \quad \vdots \quad \vdots \\ e_1 X_n + e_2 X_n + \ldots e_m X_n \end{bmatrix}$$

Each row is an approximate derivative of the cost w.r.t. all weights averaged $(\frac{1}{m})$ over all training example

## 3.6 Differentiating the Bias

$$C = \frac{1}{2m} \sum_{i=1}^{m} (y-a^L)^2 = \frac{1}{2m} \sum_{i=1}^{m} (v)^2 \qquad \text{where } v = y-a^L$$

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

---

We know :

$$\frac{\partial a^L}{\partial b} = \begin{cases} 0 & \text{if } w^T x + b \leq 0 \\ 1 & \text{if } w^T x + b > 0 \end{cases}$$

$$\frac{\partial v}{\partial a^L} = \frac{\partial}{\partial a^L} (y - a^L) = -1$$

$$\therefore \frac{\partial v}{\partial b} = \frac{\partial v}{\partial a^L} \cdot \frac{\partial a^L}{\partial b} = -1 \cdot \begin{cases} 0 & \text{if } w^T x + b \leq 0 \\ 1 & \text{if } w^T x + b > 0 \end{cases} = \begin{cases} 0 & \text{if } \dots \\ -1 & \text{if } \dots \end{cases}$$

---

We could find $\frac{\partial c}{\partial v}$ and multiply it by $\frac{\partial v}{\partial b}$, but it is easier to find $\frac{\partial c}{\partial b}$ and substitute derivatives along the way

$$\frac{\partial c}{\partial b} = \frac{\partial}{\partial b}\left(\frac{1}{2m}\sum_{i=1}^{m}(v^2)\right) = \frac{1}{2m}\sum_{i=1}^{m}\frac{\partial}{\partial b}v^2$$

$$= \frac{1}{2m}\sum_{i=1}^{m}\frac{\partial v^2}{\partial v}\frac{\partial v}{\partial b} = \frac{1}{2m}\sum_{i=1}^{m}2v\frac{\partial v}{\partial b}$$

$$= \frac{1}{m}\sum_{i=1}^{m}v\cdot\begin{cases}0 & \text{if } w^Tx+b \leq 0 \\ -1 & \text{if } w^Tx+b\end{cases} = \frac{1}{m}\sum_{i=1}^{m}\begin{cases}0 & \text{if} \dots \\ -v & \text{if} \dots\end{cases}$$

$$= \frac{1}{m}\sum_{i=1}^{m}\begin{cases}0 & \text{if } w^Tx+b \leq 0 \\ -(y-a^L) & \text{if } w^Tx+b > 0\end{cases}$$

$$= \frac{1}{m}\sum_{i=1}^{m}\begin{cases}0 & \text{if } w^Tx+b \leq 0 \\ -(y-\max(0, w^Tx+b) & \text{if } w^Tx+b > 0\end{cases}$$

Piecewise makes $\max(0, w^Tx+b)$ redundant so we can remove

$$= \frac{1}{m}\sum_{i=1}^{m}\begin{cases}0 & \text{if } w^Tx+b \leq 0 \\ -(y-(w^Tx+b)) & \text{if } w^Tx+b > 0\end{cases}$$

$$= \frac{1}{m}\sum_{i=1}^{m}\begin{cases}0 & \text{if } w^Tx+b \leq 0 \\ w^Tx+b-y & \text{if } w^Tx+b > 0\end{cases}$$

$$\frac{\partial c}{\partial b} = \frac{1}{m}\begin{cases}0 & \text{if } w^Tx+b \leq 0 \\ \sum_{i=1}^{m}(w^Tx+b-y) & \text{if } w^Tx+b > 0\end{cases}$$
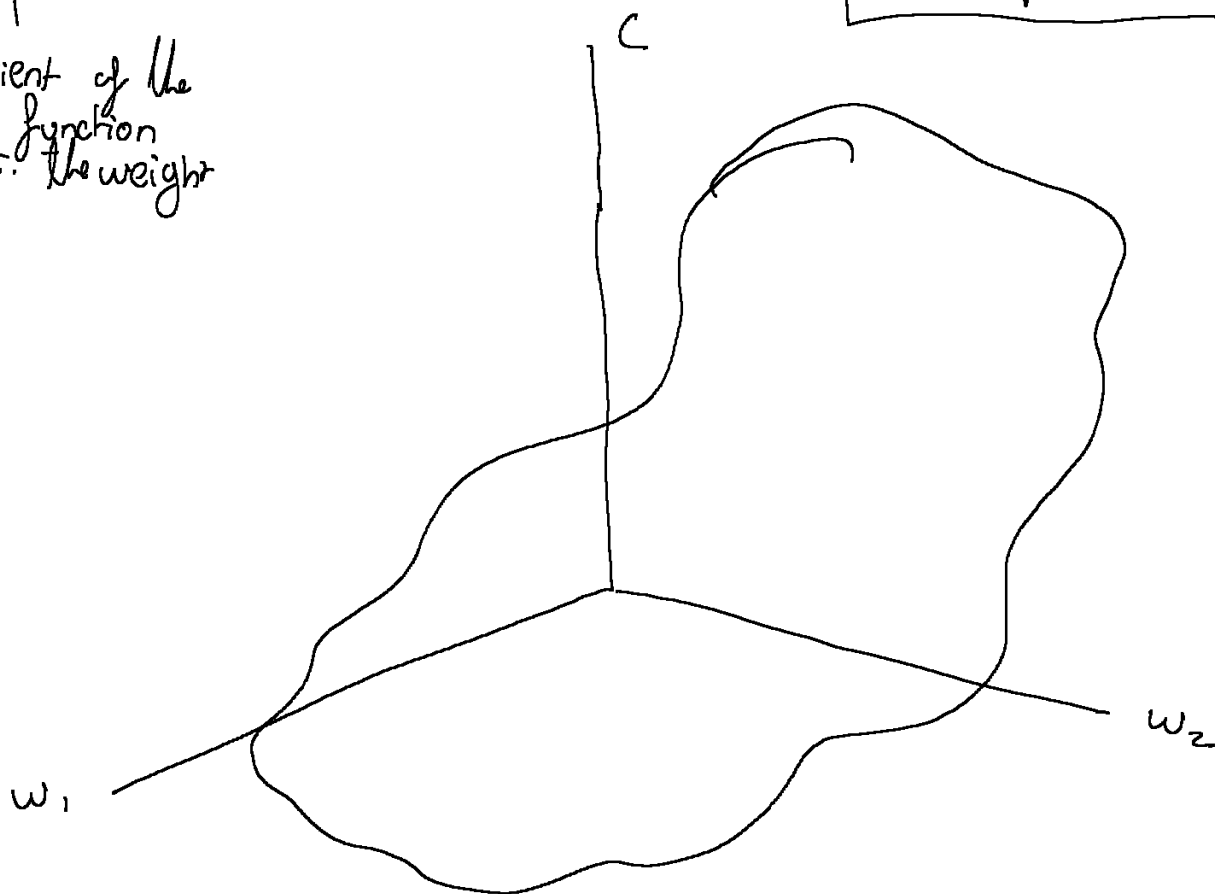
This is a scalar

$$\frac{\partial c}{\partial \vec{w}} = \begin{bmatrix} e_i\, x_1 \\ e_i\, x_2 \\ \vdots \\ e_i\, x_n \end{bmatrix}$$

Lets take a network w/ 2 weights and graph it in 3-D

$$\nabla_w C = \begin{bmatrix} e_i\, x_1 \\ e_i\, x_2 \end{bmatrix} \begin{matrix} \rightarrow w_1 \\ \rightarrow w_2 \end{matrix}$$

↑

Gradient of the
cost function
w.r.t. the weight

> ☆ Remember, the gradient
> of a function points
> in the direction of
> steepest ascent

We can then take the negative of the steepest gradient $(-\nabla_w C)$ as this will tell us how to most quickly decrease the cost

# 3.8 Gradient Descent Algorithm and SGD

$$\nabla_{w,b} C = \begin{bmatrix} \dfrac{\partial c}{\partial w_1} \\[6pt] \dfrac{\partial c}{\partial w_2} \\[6pt] \vdots \\[6pt] \dfrac{\partial c}{\partial w_n} \end{bmatrix}$$

We can "unravel" all the weights and biases into a vector that is the same length as $\nabla_{w,b} C$. We can call this vector $\Theta$

Gradient descent is an iterative process. It takes many iterations to lower the cost by adjusting the weights and bias

$$\therefore \; \Theta = \Theta - \alpha \, \nabla_{w,b} C = \begin{bmatrix} w_1 - \alpha \dfrac{\partial c}{\partial w_1} \\[6pt] w_2 - \alpha \dfrac{\partial c}{\partial w_2} \\[6pt] \vdots \\[6pt] w_n - \alpha \dfrac{\partial c}{\partial w_n} \\[6pt] b - \alpha \dfrac{\partial c}{\partial b} \end{bmatrix}$$

$\uparrow$

Learning rate

The learning rate plays a critical role in determining how big the "steps" of adjusting the weights and biases are

If the learning rate is too big, it could overshoot the minima and thus it may not converge

If the learning rate is too small, it could get "stuck" and not converge to the optimal point as it moves too slowly
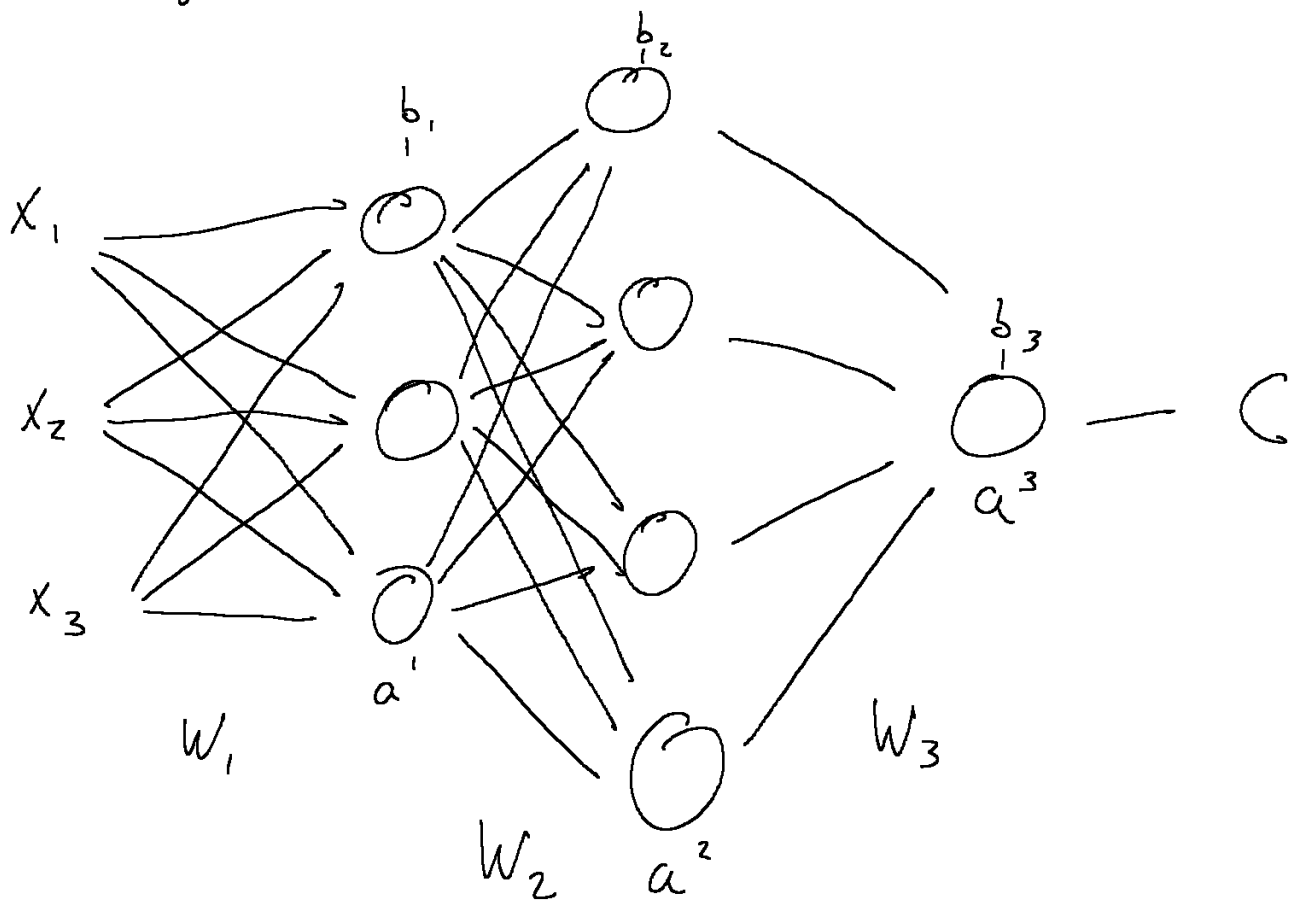
**Stochastic gradient descent (SGD):** A modification of gradient descent where you calculate the gradient using just a small part of the observations instead of all. Can reduce computation time

**Batch stochastic gradient descent:** The gradients are calculated and the decision variables are updated iteratively with a subset of observations, called minibatches

Lets say we have 1,000,000 training examples, $m$, and we take batches of size 120. Then we would get 8,333 batches (with the last one being partially full) all with 120 examples in the batch

We could then run gradient descent on each batch

# 3.9 Finding Derivatives of an Entire Layer (and why it doesn't work well)



If we want to see how changing $w_1$ affects the cost, we get:

$$\frac{\partial c}{\partial w_1} = \frac{\partial c}{\partial a^3} \frac{\partial a^3}{\partial a^2} \frac{\partial a^2}{\partial a^1} \frac{\partial a^1}{\partial w_1}$$

If we want to see how $w_2$ affects the cost:

$$\frac{\partial c}{\partial w_2} = \frac{\partial c}{\partial a^3} \frac{\partial a^3}{\partial a^2} \frac{\partial a^2}{\partial w_2}$$
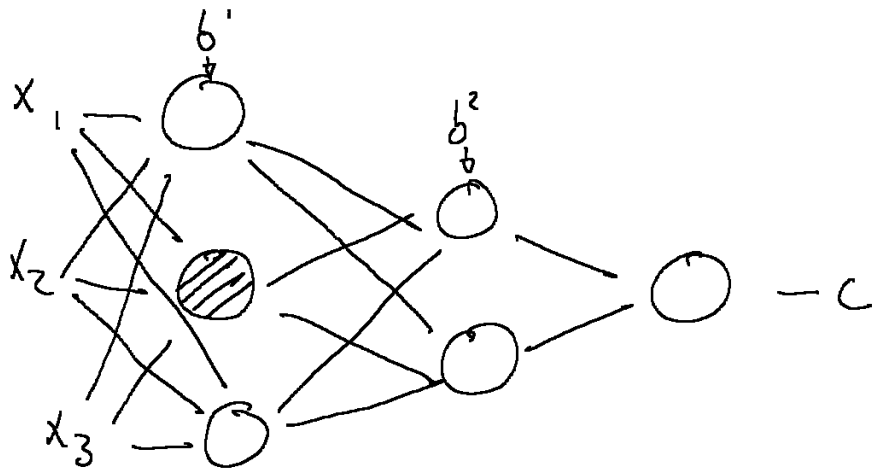
This is quite difficult to calculate everything "forward", so instead we use "back propogation" so that we don't have to recalculate the derivatives when we change a weight

# PART IV

# Backpropagation

We want to examine how the cost changes when we add a Value, $\Delta$, to a node

$$a'_2 = \sigma \underbrace{(W^T x + b + \Delta)}_{z}$$

We would actually be adding the value to the $z$ (before it goes through the activation function so it would be

$$z'_2 + \Delta$$

Therefore we can say that the error of a node, $\delta$, is given by

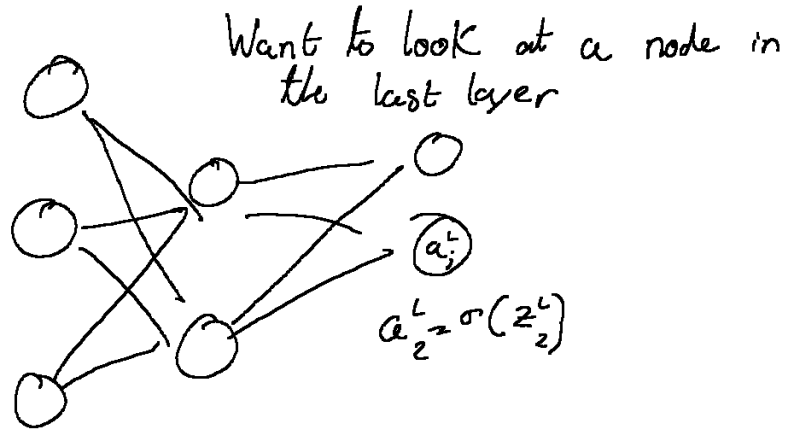$$\delta^\ell_j = \frac{\partial c}{\partial z^\ell_j}$$

It shows us how much the cost will change by changing a value of a node

# 4.2 The Four Equations of Backpropagation

## 4.2.1 Equation 1: Finding the error of all nodes in the last layer

$$\delta_j^l = \frac{\partial C}{\partial z_j^l}$$

Error of $j^{th}$ node in layer $l$

Want to look at a node in the last layer

$$a_2^L = \sigma(z_2^L)$$

Error for ONE node in the last layer:

$$\delta_j^L = \frac{\partial C}{\partial a_j^L} \cdot \frac{\partial a_j^L}{\partial z_j^L}$$

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

→ Derivative of activation function w.r.t. $z_j^L$

Error for EVERY node in the last layer:

$$\delta^L = \nabla_{a^L} C \cdot \sigma'(z_j^L)$$

This will be a vector of errors for each node in the last layer