

NNs

---

Henry Jones

---

---

---



## Prerequisites

- Basic Linear algebra
- Multivariable calculus
- Differential
- Jacobian's / gradients
- Basic ML knowledge

## Agenda

- Multivariable calculus review
- Neurons as a function
- Jacobians & NNs
- Gradient descent
- Back propagation
- BP as Matrix calculus

## Basic Notation

Define  $m = \#_{\text{data}}$  of training set

Define  $n = \# \text{ of input variables } x_1, x_2, \dots$

Define  $L = \# \text{ of layers in NN}$

Define  $l = \text{specific layer}$

$w^L = \text{weight}(l)$   $b^L = \text{bias}$  going into layer  $l$



Define  $x_i = \text{single input variable}$

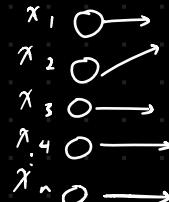
## Big picture

NNs are functions (Linear functions of weights & biases)

Big calculus problem (Linear approximations)

$$\rightarrow \frac{\partial \text{cost}}{\partial w} \quad (\text{want this})$$

Find how much each weight & bias contributes to the cost



with respect to every weight and bias

## Matrix Calculation Review

Gradients, Jacobians, Jacobian Chain Rule

Vector function

$$\text{eg. } f(x, y) = x^2 + \cos y$$

= Gradient" is a horizontal vector of partial derivatives

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

Notation: upside down triangle  
Triangle for change  
Gradient of  $f(x, y)$

$$f(x, y) = \begin{bmatrix} f_1 \\ f_2 \end{bmatrix} = \begin{bmatrix} 2x + y^3 \\ e^y - 13x \end{bmatrix} \rightarrow \text{Breaks up into 2 scalar functions}$$

$$\mathbb{R}^2 \rightarrow \mathbb{R}^2$$

$$\begin{aligned} \textcircled{1} f_1 &= 2x + y^3 \rightarrow \frac{\partial f_1}{\partial x} = 2 \quad \frac{\partial f_1}{\partial y} = 3y^2 \\ \textcircled{2} f_2 &= e^y - 13x \rightarrow \frac{\partial f_2}{\partial x} = -13 \quad \frac{\partial f_2}{\partial y} = e^y \end{aligned}$$

small function  
composing one large function

## Create Jacobian

denoted also as:

$$J = \begin{bmatrix} \frac{\partial f_1}{\partial x} & \frac{\partial f_1}{\partial y} \\ \frac{\partial f_2}{\partial x} & \frac{\partial f_2}{\partial y} \end{bmatrix} \begin{bmatrix} \nabla f_1 \\ \nabla f_2 \end{bmatrix} = \begin{pmatrix} 2 & 3y^2 \\ -13 & e^y \end{pmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

= Transposed  
defined rigorously from original function

Gradients:  
 $\mathbb{R}^n \rightarrow \mathbb{R}^1$  vector to scalar

$$x=2, y=0 \\ \begin{bmatrix} x \\ y \end{bmatrix} \xrightarrow{f(x, y)} s \text{ for scalar}$$

$$\text{eg. } \begin{bmatrix} 2 \\ 0 \end{bmatrix} = (5)$$

$$\mathbb{R}^2 \rightarrow \mathbb{R}^1$$

Jacobian:

$$\mathbb{R}^n \rightarrow \mathbb{R}^m$$

partial derivative,

$$\frac{\partial f}{\partial x} = 2x$$

$$\frac{\partial f}{\partial y} = -\sin(y)$$

vector to another vector

Since we can vary the input size and size of the Jacobian.

From last page

$$f(x, y) = \begin{pmatrix} 2x - y^3 \\ e^x - 13y \end{pmatrix}$$

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

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

$$\mathbb{R}^2 \Rightarrow \mathbb{R}^2$$

$$f_1 = 2x - y^3$$

$$f_2 = e^x - 13y$$

$$\left. \begin{array}{l} \frac{\partial f_1}{\partial x} = 2 \\ \frac{\partial f_1}{\partial y} = -3y^2 \\ \frac{\partial f_2}{\partial x} = e^x \\ \frac{\partial f_2}{\partial y} = -13 \end{array} \right\} = J = \begin{bmatrix} 2 & -3y^2 \\ e^x & -13 \end{bmatrix}^T = \begin{bmatrix} \nabla f_1^T \\ \nabla f_2^T \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x} & \frac{\partial f_1}{\partial y} \\ \frac{\partial f_2}{\partial x} & \frac{\partial f_2}{\partial y} \end{bmatrix}$$

Transpose the matrix to "factor" out  $x$  &  $y$ .

### Scalar Chain Rule

Let  $f(x) = \sin(x^2)$

$$f'(x) = 2x \cos(x^2)$$

Let  $g = x^2$

$f = \sin(g)$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial x} = 2x \cos g \quad \text{technical breakdown}$$

Now, extend to:  $f = g \cdot h(c(a \cdot x))$

as a general concept: use product of partial derivatives

① How does  $g$  change as  $\begin{pmatrix} x \\ y \end{pmatrix}$  change?

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

② How does  $f$  change as  $(g_1, g_2)$  change

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

③ Apply the Scalar chain rule

vector notation

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

Matrix multiplication

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

### Jacobian Chain Rule

Define the following:

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

$$g = \begin{pmatrix} x^2 + y \\ y^3 \end{pmatrix}_{g_1, g_2} \quad f = \begin{pmatrix} \sin(g_1) \\ \ln(g_2) \end{pmatrix}$$

Inside functions

Outside functions

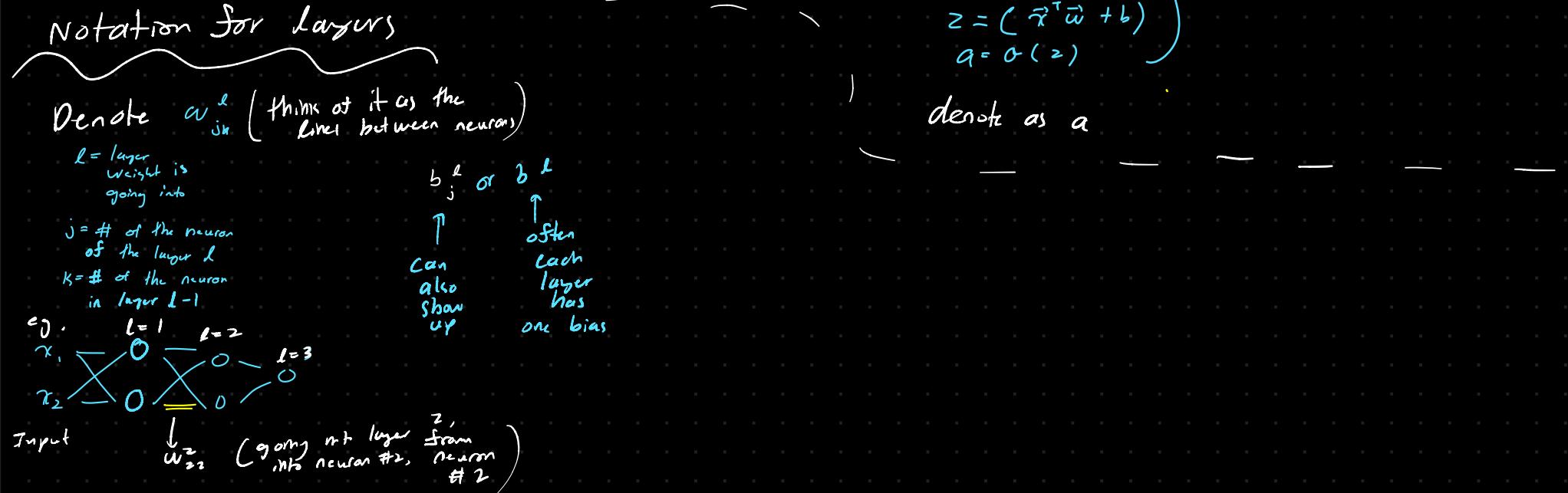
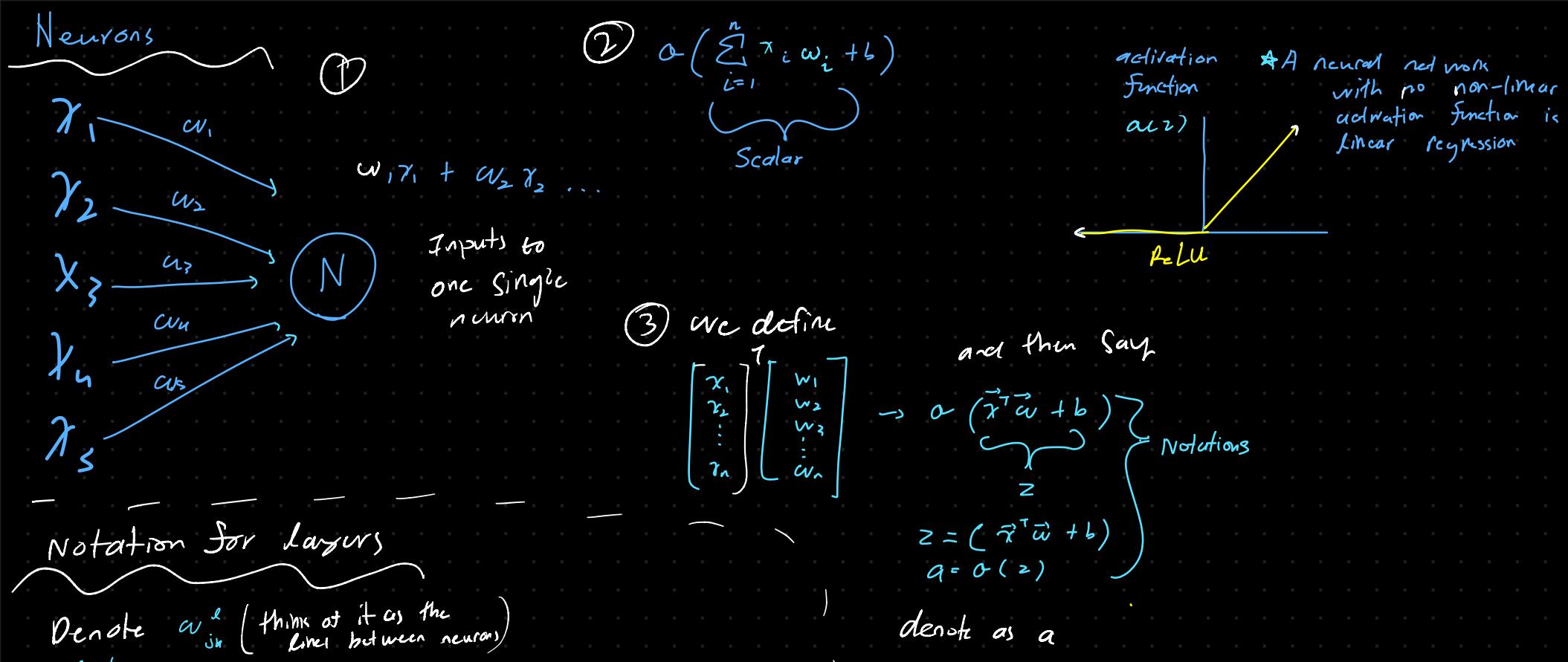
what if we cannot use chain rule?

$$f(x, y) = \begin{pmatrix} \sin(x^2 + y^2) \\ \cos(g_2^3) \\ xy^3 \end{pmatrix}$$

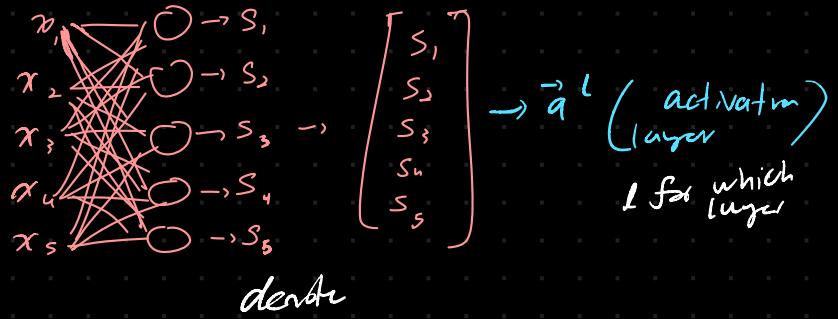
$$g = \begin{pmatrix} x^2 + y^2 \\ y^3 \\ 1 \end{pmatrix}$$

Leave as is if no intermediate

put all when you are not using chain rule aka / no intermediate function



For one layer: Define scalars



(a for activation )      input goes into

of weights

From single neuron to full layer:

Single node

$$a_l = \vec{x}^T \vec{w} + b$$

Entire layer

$a^l$

$a_{j,k}^l$  is equivalent to:

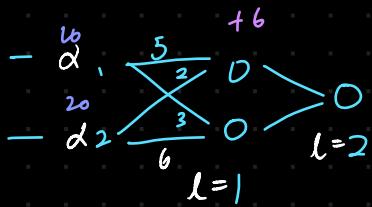
Define # of neurons in layer =  $k = j$  (usually square)  
in training examples

$$j \left[ \begin{array}{c} a^1 \\ \vdots \\ a^k \end{array} \right] \in \left[ \begin{array}{c} a_1 \\ a_2 \\ a_3 \\ a_4 \\ a_5 \end{array} \right]$$

comes from the  $k$ th neuron

and goes into the  $j$ th neuron

e.g.



$W_{ik}^l$  construct the weights matrix

$$W^l = \begin{bmatrix} (1,1) & (1,2) \\ (2,1) & (2,2) \end{bmatrix} = \begin{bmatrix} 5 & 2 \\ 3 & 6 \end{bmatrix}$$

$$\alpha(z) = \text{activation}$$

of the previous layer

$$= \begin{bmatrix} 5 & 2 \\ 3 & 6 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

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

$$\alpha(z) = \alpha\left(\begin{bmatrix} 96 \\ 156 \end{bmatrix}\right)$$

### Activation Function Calculation

General:  $a_n = \alpha(W^{l+1} a^n + b^l)$

$$a_0 \Rightarrow \alpha(w^1 a^0 + b^0)$$

from  $a_0$ ,  
 $a_1 = \alpha(w^2 a^1 + b^1)$

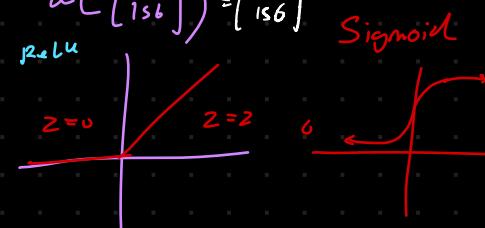
vector  $\rightarrow$  vector

output is  
a vector

$$= \begin{bmatrix} 96 \\ 156 \end{bmatrix}$$

$$\alpha\left(\begin{bmatrix} 96 \\ 156 \end{bmatrix}\right) = \begin{bmatrix} 96 \\ 156 \end{bmatrix}$$

relu



feed forward through  
an NN  
input through the network

$$O - O - O$$

$$O - O -$$

forward

$$O - O -$$

backward

$$\alpha(z) = \begin{bmatrix} 96 \\ 156 \end{bmatrix} \text{ stays the same}$$

## Jacobians of NNs

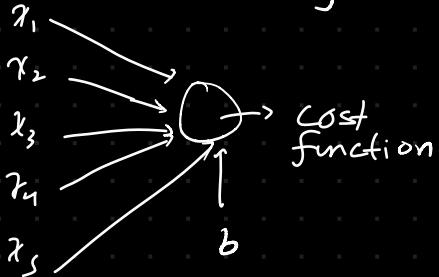
Minimize the cost of the output (error)

Mean Squared error

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

↓ answer we got  
real answer

$\pm$  does not matter (for the power rule)  
 $m$  is the # of training examples



## Jacobians of Operations

Binary element-wise operations

vector to vector

$F(\vec{v}, \vec{w}) \rightarrow \vec{z}$

$$\textcircled{1} \quad \begin{bmatrix} \vec{v} \\ \vec{w} \end{bmatrix} = \begin{bmatrix} \vec{z} \end{bmatrix}$$

element wise operators  
(can be any thing here)  
if general

Hadamard product ( $\otimes$ )  
is a binary element-wise  
operation.

$$\begin{bmatrix} v_1 \otimes w_1 \\ v_2 \otimes w_2 \\ v_3 \otimes w_3 \\ \vdots \\ v_n \otimes w_n \end{bmatrix} = \begin{bmatrix} b \end{bmatrix}$$

## ② $F(\vec{v}, \vec{w})$ (general)

$$= (f(\vec{v})) \circ g(\vec{w}) = \begin{bmatrix} f_1(\vec{v}) \circ g_1(\vec{w}) \\ f_2(\vec{v}) \circ g_2(\vec{w}) \\ \vdots \\ f_n(\vec{v}) \circ g_n(\vec{w}) \end{bmatrix}$$

operations on the  
functions before

Can get 2 Jacobians  
one for each vector

$$\frac{\partial F}{\partial \vec{v}} \quad \frac{\partial F}{\partial \vec{w}}$$

first element of  $v$ , second  
of  $w$

$$\begin{bmatrix} f_1(\vec{v}) \circ g_1(\vec{w}) \\ f_2(\vec{v}) \circ g_2(\vec{w}) \\ \vdots \\ f_n(\vec{v}) \circ g_n(\vec{w}) \end{bmatrix}$$

In the most general  
form (if  $F(\vec{v})$  has  
multiple element wise types  
of operations)

General representation  
increase  $f_i$  and  $g_i$ , don't change the vector

$$\textcircled{3} \quad J = \frac{\partial F}{\partial \vec{v}} = \begin{bmatrix} \frac{\partial}{\partial v_1} f_1(\vec{v}) \circ g_1(\vec{w}) & 0 & \dots & \frac{\partial}{\partial v_n} f_1(\vec{v}) \circ g_1(\vec{w}) \\ 0 & \frac{\partial}{\partial v_1} f_2(\vec{v}) \circ g_2(\vec{w}) & \dots & \frac{\partial}{\partial v_n} f_2(\vec{v}) \circ g_2(\vec{w}) \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \frac{\partial}{\partial v_n} f_n(\vec{v}) \circ g_n(\vec{w}) \end{bmatrix}$$

If  $f(\vec{v}) = \vec{v}$  and  $g(\vec{w}) = \vec{w}$

Scalar

$$\begin{bmatrix} \frac{\partial}{\partial v_1} (\vec{v}) \circ (\vec{w}) & 0 & \dots & 0 \\ 0 & \frac{\partial}{\partial v_2} (\vec{v}) \circ (\vec{w}) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \frac{\partial}{\partial v_n} (\vec{v}) \circ (\vec{w}) \end{bmatrix}$$

## Derivatives of Neuron Operators (Hadamard product)

→ Hadamard product

$$F(\vec{v}, \vec{w}) = \begin{bmatrix} f_1(\vec{v}) \otimes g_1(\vec{w}) \\ f_2(\vec{v}) \otimes g_2(\vec{w}) \\ \vdots \\ f_n(\vec{v}) \otimes g_n(\vec{w}) \end{bmatrix} = \begin{bmatrix} v_1 \otimes w_1 \\ v_2 \otimes w_2 \\ \vdots \\ v_n \otimes w_n \end{bmatrix} \Rightarrow \text{much better description}$$

$\frac{\partial F}{\partial \vec{v}} = \begin{bmatrix} \frac{\partial F_1}{\partial v_1} & \frac{\partial F_1}{\partial v_2} & \dots & \frac{\partial F_1}{\partial v_n} \\ \frac{\partial F_2}{\partial v_1} & \frac{\partial F_2}{\partial v_2} & \dots & \frac{\partial F_2}{\partial v_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial F_n}{\partial v_1} & \frac{\partial F_n}{\partial v_2} & \dots & \frac{\partial F_n}{\partial v_n} \end{bmatrix} = \begin{bmatrix} c v_1 & 0 & \dots & 0 \\ 0 & c v_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & c v_n \end{bmatrix}$

with respect to  $v$

\* For all element wise functions, the Jacobian will be a diagonal matrix

$$\begin{bmatrix} a & 0 & 0 \\ 0 & b & 0 \\ 0 & 0 & c \end{bmatrix} \rightarrow \text{diag}(a, b, c)$$

$\frac{d}{dt}(xw) = x$  (easier way to understand)

## Derivative of a Scalar expansion

$$2 \begin{bmatrix} \vec{v} \\ 1 \end{bmatrix} = \begin{bmatrix} 2\vec{v}_1 \\ 2\vec{v}_2 \\ 2\vec{v}_3 \\ \vdots \\ 2\vec{v}_n \end{bmatrix}$$

$$F(\vec{v}, x) = f(\vec{v}) \circ g(x)$$

Defining for  $g(\vec{u}) = \vec{I} \cdot \vec{u}$  (expands into a vector)

The act of multiplying  $\vec{x}$  by the ones vector is an act (expands to a diag matrix or a Jacobian) of broadcasting itself

$$\begin{bmatrix} 2 \\ \vec{v} \\ 2 \end{bmatrix} \circ \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix}$$

When we get the Jacobian, we get a gradient if  $x$  is a scalar

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

$$\text{Gradient with respect to } x \quad \nabla F_x = \begin{bmatrix} \frac{\partial f_1}{\partial x} \\ \frac{\partial f_2}{\partial x} \\ \vdots \\ \frac{\partial f_n}{\partial x} \end{bmatrix} = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix}$$

don't forget, this is the scalar vector.

$$\begin{bmatrix} x & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & x \end{bmatrix}$$

$g$  is element wise

$$\text{Derivative of a Sum}$$

e.g.  $\vec{v} = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix} = \sum_{i=1}^n (g_i(\vec{v}))_i = \begin{bmatrix} \text{---} s \text{ ---} \end{bmatrix}$  for  $s$  to denote the sum

$$= \begin{bmatrix} \frac{\partial s}{\partial v_1} & \frac{\partial s}{\partial v_2} & \cdots & \frac{\partial s}{\partial v_n} \end{bmatrix}$$

$$J = \left[ \frac{\partial}{\partial v_i} \sum_{i=1}^n g_i(v) \right] = \frac{\partial}{\partial v_1} \sum_{i=1}^n g_i(v) \cdots \frac{\partial}{\partial v_n} \sum_{i=1}^n g_i(v)$$

Easy to generalize

\* The derivative of a sum is the sum of the derivatives

$$J = \left[ \sum_{i=1}^n \left( \frac{\partial}{\partial v_i} g_i(v) \right) \right] = \left[ \sum_{i=1}^n \frac{\partial}{\partial v_1} g_i(v) \cdots \sum_{i=1}^n \frac{\partial}{\partial v_n} g_i(v) \right]$$

If  $g(\vec{v}) = 2\vec{v}$

Consider  $g(\vec{v}) = \vec{v}$

$$J = \left[ \sum_{i=1}^n \frac{\partial}{\partial v_i} (v_i) \sum_{i=1}^n \frac{\partial}{\partial v_2} (v_i) \cdots \sum_{i=1}^n \frac{\partial}{\partial v_n} (v_i) \right] = \begin{bmatrix} 1 & 1 & \dots & 1 \end{bmatrix}$$

Basically, a horizontal line of what  $g$  does if it does anything

# Derivative of a neuron

activation  
 $a = \sigma(\vec{w}^T \vec{x} + b)$

$\vec{z}$  dot product  
 = Hadamard product

wrt  
 $\frac{\partial a}{\partial w} = \frac{\partial a}{\partial z} \frac{\partial z}{\partial w}$

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

Define an activation slice:

$a = \sigma(\text{sum}(w \otimes x) + b)$

$\vec{z}$  Hadamard product

$\vec{w}$  SCH

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

either  
 $0$  or  
 $1$

previous  
 slide

$$[-] \begin{bmatrix} 1 \\ | \\ | \\ m \end{bmatrix}, [-] \begin{bmatrix} 1 \\ | \\ | \\ n \end{bmatrix}$$

$$1 \begin{bmatrix} 1 \\ | \\ | \\ m \end{bmatrix}$$

$$\sigma(\sum(w \otimes x) + b)$$

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

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

By Jacobian convention,  $H = w \otimes x$

①  $\frac{\partial H}{\partial w} = \begin{bmatrix} x_1 & \dots & 0 \\ \vdots & x_2 & \vdots \\ 0 & \dots & x_n \end{bmatrix}$

diag ( $x_1, \dots, x_n$ )

②  $Z = S(H) + b$

$\frac{\partial Z}{\partial H} = \vec{1}$

③  $Z = S^{-1}$

$\frac{\partial Z}{\partial S} = \vec{1}$

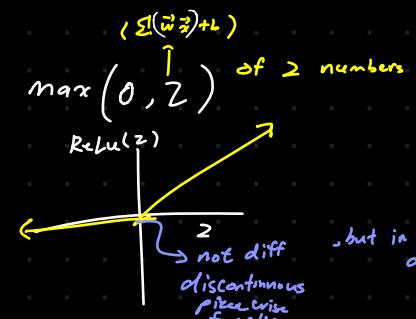
(Scalar)

$\frac{\partial S}{\partial H} = (1 \ 1 \ 1 \ 1 \ 1)$

④ Derivative of ReLU wrt the activation

$$\frac{\partial a}{\partial z} = \frac{\partial a}{\partial z} [x^T]$$

Simplifies  
 after all  
 3 substitutions



Need to find this

To solve, we split up into piecewise functions

taking the derivative

$$\frac{\partial a}{\partial w} = \max \begin{cases} z \geq 0, 0 \rightarrow \vec{0} \times \frac{\partial z}{\partial w} = \vec{0} = [\vec{0}] \\ z > 0, z \rightarrow 1 \times \frac{\partial z}{\partial w} = \frac{\partial z}{\partial w} = [x^T] \end{cases}$$

For some  $w$ ,  $\frac{\partial a}{\partial w}$  gives the best approximation for the instant rate of change of the activation, wrt that specific  $w$ .  
 also putting all derivatives together

important so that the vector operations make sense dimensionally

$$\frac{\partial a}{\partial w} = [0 \ 0 \ 0 \ \dots \ 0]$$

$$\text{or} \\ \frac{\partial a}{\partial w} = [x_1 \ x_2 \ x_3 \ \dots \ x_n]$$

depending on  $\vec{w}^T \vec{x} + b$

\* The derivative of the cost wrt a weight is directly proportional to the input  $\vec{x}$

# The Gradient of the Loss Function

Cost function

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

$\downarrow$  answer activation  $a^L$



We want to find:

$$\frac{\partial C}{\partial w} = \frac{\partial C}{\partial a} \frac{\partial a}{\partial w} \quad (\text{Chain rule})$$

$\frac{\partial C}{\partial b}$  know this already

Recall that

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

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

$$\frac{1}{m} \left\{ \begin{array}{l} \vec{0}^T \\ \sum_{i=1}^m (\vec{w}^T \vec{x} + b - y_i) \vec{x}^T \end{array} \right\}$$

$$\text{one training example} \quad \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix}$$

$$\frac{1}{2m} \sum_{i=1}^m (y_i - a^L)^2$$

$$\text{Several training examples } C_m \quad \begin{bmatrix} x_{11} & x_{21} & \dots & x_{n1} \\ x_{12} & x_{22} & \ddots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ x_{1n} & x_{2n} & \dots & x_{nn} \end{bmatrix}$$

$$\text{Define and say: } v = \underset{\text{correct answer}}{y} - a^L$$

$$\frac{1}{2m} \sum_{i=1}^m (y_i - a^L)^2 = \frac{1}{2m} \sum_{i=1}^m (v)^2$$

$$\frac{\partial C}{\partial w} = \frac{\partial C}{\partial v} \frac{\partial v}{\partial a} \frac{\partial a}{\partial w}$$

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

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

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

$$\text{Find } \frac{\partial C}{\partial w} \quad \frac{\partial v}{\partial w} = -(\mathbf{1})(\frac{\partial a}{\partial w})$$

don't forget that  
 $L$  is not  
actually a  
constant in  
the equation

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

$$\text{Take out the constant} \quad = \frac{1}{2m} \left( \sum_{i=1}^m \frac{\partial}{\partial w} (v^2) \right) \left( \frac{\partial v^2}{\partial v} = \frac{\partial v^2}{\partial v} \frac{\partial v}{\partial w} \right)$$

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

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

$$\text{what does this mean? taking } \rightarrow \text{ at a piecewise} \\ \downarrow \quad \left\{ \begin{array}{l} [0]^T \text{ if...} \\ -v \frac{\partial a}{\partial w} \text{ if...} \end{array} \right. \quad \text{make the options negative} \\ \text{if... } \cancel{\text{DOES not actually change}} \quad \text{the piecewise}$$

$$\text{Subbing in the variables,} \quad = \frac{1}{m} \sum_{i=1}^m \left\{ \begin{array}{l} [0]^T \text{ if...} \\ -(y_i - a^L) \left[ \frac{1}{x_i} \right]^T \end{array} \right.$$

From ReLU,

$$\theta(w^T x + b) \rightarrow \max(0, w^T x + b)$$

$$\Rightarrow \frac{1}{m} \sum_{i=1}^m \left\{ \begin{array}{l} [0]^T \text{ if...} \\ -\{y_i - (\max(0, w^T x + b))\} \left[ \frac{1}{x_i} \right]^T \end{array} \right. \\ = -\{y_i - (\max(0, w^T x + b))\} \left[ \frac{1}{x_i} \right]^T$$

if  $w^T x + b > 0$   
always the case  
if it went down this path.

max function  
is redundant

$$\frac{\partial C}{\partial w} = \frac{1}{m} \sum_{i=1}^m (\vec{w}^T \vec{x} + b - y_i) \vec{x}^T$$

= scalar

$$| \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} | \quad ]$$

What did we find out?  
 $\frac{\partial C}{\partial w} = \frac{1}{m} \sum_{i=1}^m (w^T x + b - y_i) x^i$  if  $w^T x + b \leq 0$   
 Want a Jacobian for back propagation

$$\left[ \begin{array}{c} \frac{\partial C}{\partial w_1} \\ \vdots \\ \frac{\partial C}{\partial w_n} \end{array} \right] \text{ called total weights}$$

$$\Rightarrow \frac{1}{m} \sum_{i=1}^m e_i(x)$$

where  $e_i = w^T x - y_i$  (scalar)

Consider  $m=1$  and 5 inputs, transposing our  $e_i | x^T$  so it becomes a column,

$$\frac{\partial C}{\partial w} = \begin{bmatrix} e_1 x_1 \\ e_2 x_2 \\ \vdots \\ e_5 x_5 \end{bmatrix} \frac{\partial C}{\partial w_1} \text{ has } C \text{ changes wrt } w_1, \dots \text{ how } C \text{ changes wrt } w_2, \dots$$

When you change  $w_i$  by some amount, the corresponding change is ...

$$\begin{aligned} \text{consider } m \sum_{i=1}^m e_i x^T \\ = \frac{1}{m} e_1 x^T + e_2 x^T + e_3 x^T \dots e_5 x^T \\ = \frac{1}{m} \left( \begin{bmatrix} e_1 x_1 \\ e_2 x_2 \\ \vdots \\ e_5 x_5 \end{bmatrix} + \begin{bmatrix} e_1 x_1 \\ e_2 x_2 \\ \vdots \\ e_5 x_5 \end{bmatrix} + \dots + \begin{bmatrix} e_1 x_1 \\ e_2 x_2 \\ \vdots \\ e_5 x_5 \end{bmatrix} \right) \end{aligned}$$



Consider a general case

$$\frac{1}{m} \begin{bmatrix} e_1 x_1 + e_2 x_1 \dots e_5 x_1 \\ e_1 x_2 + e_2 x_2 \dots e_5 x_2 \\ e_1 x_3 + e_2 x_3 \dots e_5 x_3 \\ e_1 x_4 + e_2 x_4 \dots e_5 x_4 \\ e_1 x_5 + e_2 x_5 \dots e_5 x_5 \end{bmatrix}$$

random averages

$$\Rightarrow \begin{bmatrix} g_1 \\ g_2 \\ g_3 \\ g_4 \\ g_5 \end{bmatrix} = \frac{\partial J}{\partial w} \text{ all } w's, \text{ averaged out over } \underline{\text{all training examples}}$$

## Derivative of the bias

$$\frac{\partial C}{\partial b} = \frac{1}{2m} \sum_{i=1}^m (\gamma - a^L)^2$$

$$v = \gamma - a^L \quad (-1) \quad \checkmark$$

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

$$= \frac{1}{m} \sum_{i=1}^m v \begin{cases} 0 & \text{if } \dots \\ -1 & \text{if } \dots \end{cases} \quad (\text{plug into our piecewise})$$

$$= \frac{1}{m} \sum_{i=1}^m \begin{cases} 0 & \text{if } \dots \\ -(v - a^L) & \text{if } \dots \end{cases}$$

$\downarrow$   
max redundancy again, can remove

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

\* scalar \*

$$= \frac{1}{m} \sum_{i=1}^m \begin{cases} (w^T x + b - \gamma) & \dots \end{cases}$$

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

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

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

$$= \frac{1}{m} \sum_{i=1}^m v \frac{\partial v}{\partial b} \rightarrow \text{piecewise}$$

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

$$= (-1) \begin{cases} 0 & \text{if } \dots \\ 1 & \text{if } \dots \end{cases}$$

Now generally,

$$\frac{\partial C}{\partial b} = \frac{1}{m} \sum_{i=1}^m e_i$$

$$\boxed{\frac{\partial C}{\partial b} = \frac{1}{m} [e_1 + e_2 + \dots + e_n]}$$

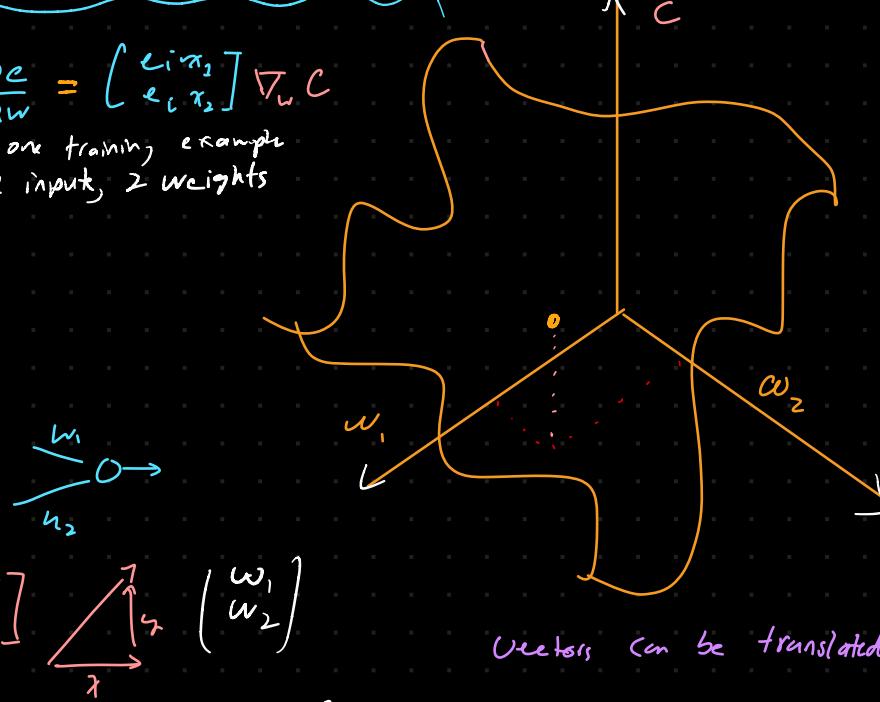
Just a scalar

derivative of the cost wrt to the bias, averaged over all training examples.

## Gradient Descent Intuition

$$\frac{\partial C}{\partial w} = \begin{bmatrix} e_i x_1 \\ e_i x_2 \end{bmatrix} \nabla_w C$$

one training example  
2 inputs, 2 weights



The gradient of a function points in the direction of steepest ascent

The value of the input decides the derivative  
want to change the weights of the corresponding input that causes the greatest effect on the cost function

$\begin{bmatrix} e_i x_1 \\ e_i x_2 \end{bmatrix}$  will point in the direction of the weights that have the most impact

$\nabla_w C$  goes in a higher direction of cost (evaluate  $e_i x_i$  first)  
error is magnified

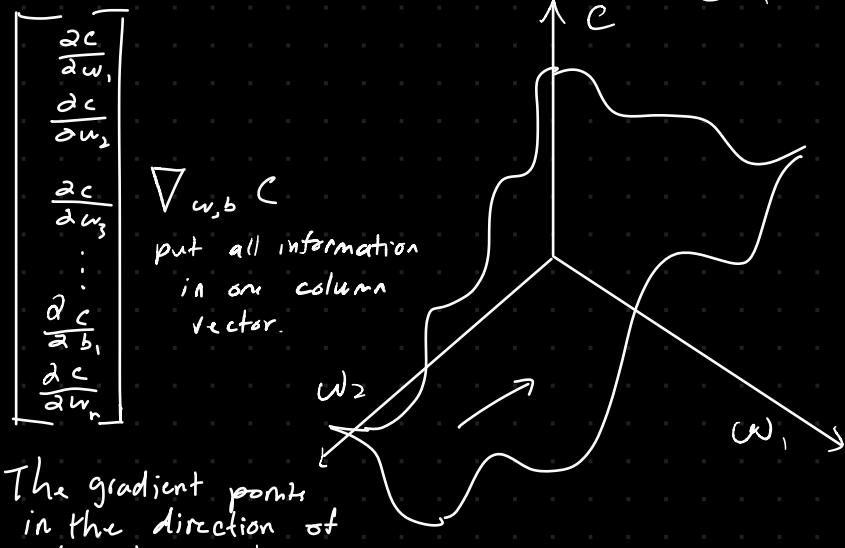
- $\nabla_w C$  tells us how to most quickly decrease the cost

The vector points in the direction that have the most impact (cost activation is  $x_i$ , error is magnified by  $x_i$ )

This we want to do  
gradient descent in the opposite direction

All machine learning models optimize to point in the direction of steepest descent

## Gradient Descent & Stochastic Gradient Descent



$\theta$  (defined to be all weights and biases)

Initially starts as randomly assigned

$$\theta = \theta - \alpha \nabla_{w,b} C$$

↓  
Learning rate (Step size)

$$\begin{bmatrix} w_1 - \alpha \frac{\partial C}{\partial w_1}, \\ w_2 - \alpha \frac{\partial C}{\partial w_2}, \\ \vdots \\ w_n - \alpha \frac{\partial C}{\partial w_n}, \\ b_1 - \alpha \frac{\partial C}{\partial b_1}, \\ b_2 - \alpha \frac{\partial C}{\partial b_2}, \\ \vdots \\ b_r - \alpha \frac{\partial C}{\partial b_r} \end{bmatrix}$$

$\alpha$  (Learning rate)  
depends on research

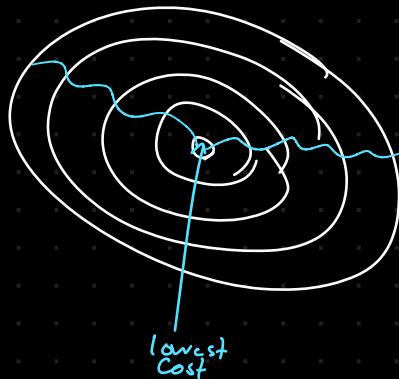
Training batch: m  
↳ Can be very computationally expensive.

$$\begin{bmatrix} 1 \\ \frac{1}{m} \\ \vdots \\ 1 \end{bmatrix}$$

over m examples

Mini-batch Gradient descent:

SGD → batching smaller datasets and then, running gradient descent on each batch  
Shuffle the data well. Batch should represent data as a whole.

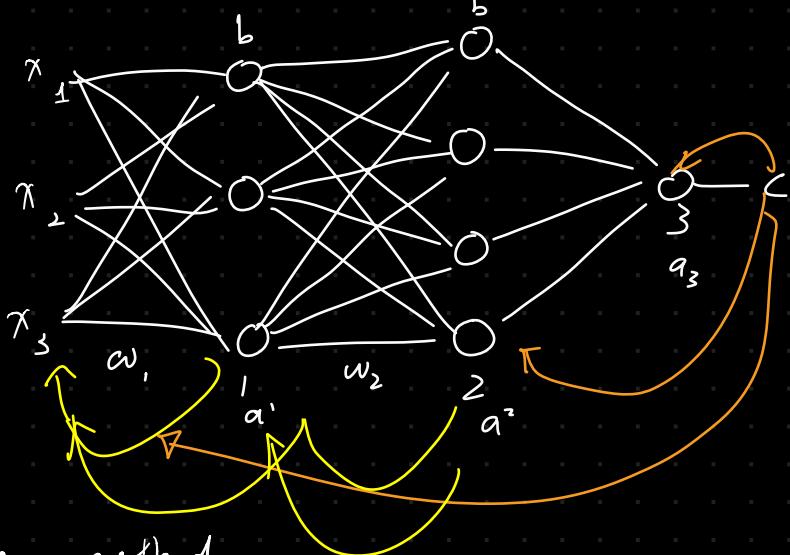


$$\left\{ x_1 \right\} \quad \left\{ x_2 \right\}$$

↓      ↓  
SGD   SGD

then, you can take as many steps, as you have batches instead of one step over all examples.

Finding derivatives of entire layers (why it is also different from one neuron)



Back propagation - saves efficiency by not having to go back and forth

Accumulate our "error"

One method

$$\frac{\partial C}{\partial w_1} = \left( \frac{\partial a'}{\partial w_1} \right) \left( \frac{\partial a^2}{\partial a_1} \right) \left( \frac{\partial a^3}{\partial a^2} \right) \left( \frac{\partial C}{\partial a^3} \right)$$

$$\left( \frac{\text{first layer}}{\text{weights}} \right) \left( \frac{\text{second layer}}{\text{first layer}} \right) \left( \frac{\text{third layer}}{\text{second layer}} \right) \left( \frac{\text{cost}}{\text{third layer}} \right)$$

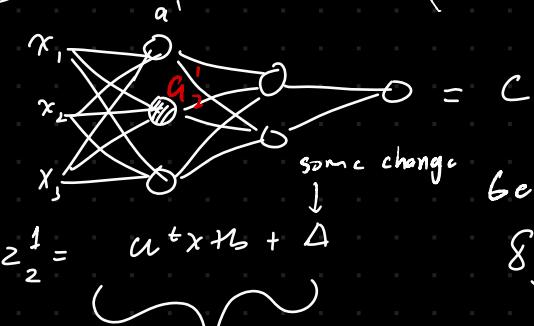
$$\frac{\partial C}{\partial w_2} = \frac{\partial a^2}{\partial w_2} \frac{\partial a^3}{\partial a^2} \frac{\partial C}{\partial a^3}$$

Doing these calculations to compute the derivatives going forward is resource-intensive  
Not often done, but possible

error of node  $j$ , of layer  $l$

$$\delta_j^l$$

The error of a node



If we add  $\Delta$  to a node,  
how does the cost  
function change?

General form error  
of a node (impact of a node on  
the cost)

$$z_2^1 = a^t x + b + \Delta$$

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

Activation of node  $j$  of  
layer 1, before ReLU

$(z_2^1 + \Delta)$  is the "total error"

1st equation - error of Last node  $a^L$

$$a^L = \sigma(z^L)$$

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

$$\delta_j^1 = \left( \frac{\partial C}{\partial a_j^1} \right) \left( \frac{\partial a_j^1}{\partial z_j^1} \right) \text{ chain rule again}$$

General form error  
of a node (impact of a node on  
the cost)

if close to 0,  
this means that changing the value of  
 $z$  has little effect on the cost  
Vice versa.

How much "error" is being caused  
by a node  
"Bang for your buck"

Now, do the entire last layer:

$$\nabla_{a^L} C \quad (\text{gradient of cost wrt } a^L)$$

$$\delta^L = \nabla_{a^L} C \odot \sigma'(z^L)$$

element-wise derivative  
operation

$$\begin{bmatrix} \frac{\partial C}{\partial a_1^L} \\ \frac{\partial C}{\partial a_2^L} \\ \vdots \\ \frac{\partial C}{\partial a_n^L} \end{bmatrix} \otimes \begin{bmatrix} \sigma'(z_1^L) \\ \sigma'(z_2^L) \\ \vdots \\ \sigma'(z_n^L) \end{bmatrix}$$

$$\delta_j^1 = \frac{\partial C}{\partial a_j^1} \cdot \sigma'(z_j^1) \quad \text{general form}$$

2nd equation: error of any node:  $\delta^L = ((w^{L+1})^T \otimes \delta^{L+1}) \otimes \sigma'(z^L)$

= weights going into layer  $L+1$

↓ weight matrix connecting  $L, L+1$

↓ total error of  $L+1$  layer

derivatives of the activations

$$\delta^L = \begin{bmatrix} \vdots \\ \delta^L \end{bmatrix}$$

why does this work:

If we know errors at layer  $L+1$

$(w^{L+1})^T$  is like going back and getting that error

e.g. Find  $\frac{\partial C}{\partial z_1}$ , given  $\frac{\partial C}{\partial z_2}$

$$\delta(w_1 a_0 + b_1) + \delta(w_2 a_1 + b_2)$$

given

using the chain rule:

$$\frac{\partial C}{\partial z_1} = \frac{\partial z_2}{\partial z_1} \frac{\partial C}{\partial z_2}$$

$$\frac{\partial C}{\partial z_1} = \frac{\partial a_1}{\partial z_1} \frac{\partial z_2}{\partial a_1} \frac{\partial C}{\partial z_2}$$

$\delta(z_2) \otimes w_2^T \otimes \delta^2$

Generally, derivative of matrix-vector product wrt the vector is,

$$\frac{\partial Ab}{\partial b} = A^T$$

works for any layer

$$\delta^{L+1} \quad C = \begin{bmatrix} \vdots \\ \delta^{L+1} \end{bmatrix}_j^K \quad u^T = \begin{bmatrix} \vdots \\ u^T \end{bmatrix}_K^J \quad \left. \right\} \text{Reversing through the networks}$$

$$z_1 = w_2 a_1$$

$$\frac{\partial z_1}{\partial a_1} = w_2^T \leftarrow \text{Transposed}$$

Equation 3 - derivative of cost wrt any bias

$$\frac{\partial C}{\partial b_1} = \delta^L$$

$$\begin{bmatrix} 4 \\ 4 \\ 4 \\ 4 \\ 4 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} \end{bmatrix}$$

Given  $\sigma(w, a_0 + b_1) \Rightarrow \sigma(w_1 a_1 + b_1)$

$$\delta^L = \frac{\partial C}{\partial z_1}$$

given  $\delta'$  (error of  $z_1$ )

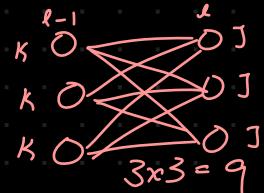
$$\frac{\partial C}{\partial b_1} = \left( \frac{\partial C}{\partial z_1} \right) \left( \frac{\partial z_1}{\partial b_1} \right)$$

$\delta^L$       1

$$\boxed{\frac{\partial C}{\partial b_1} = \delta^L}$$

Vectorized form:

Example:



$$w_{lk} = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \end{bmatrix}_{j,k} \quad (3,3)$$

$$\frac{\partial C}{\partial w_{lk}} \begin{bmatrix} \frac{\partial C}{\partial w_{11}} & \frac{\partial C}{\partial w_{12}} & \frac{\partial C}{\partial w_{13}} \\ \vdots & \vdots & \vdots \end{bmatrix}$$

$$\frac{\partial C}{\partial w_{lk}} = \begin{bmatrix} a'_1 \delta_1^2 & a'_2 \delta_1^2 & a'_3 \delta_1^2 \\ a'_1 \delta_2^2 & a'_2 \delta_2^2 & a'_3 \delta_2^2 \\ a'_1 \delta_3^2 & a'_2 \delta_3^2 & a'_3 \delta_3^2 \end{bmatrix}$$

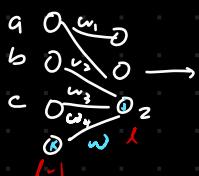
Calculate  
like:

$$\delta^2 = \begin{bmatrix} \delta_1^2 \\ \delta_2^2 \\ \delta_3^2 \end{bmatrix} \quad a'_1 = \begin{bmatrix} 1 \\ a'_1 \\ a'_2 \\ a'_3 \end{bmatrix}$$

$$\boxed{(\vec{\delta}^2)(\vec{a}'_1)^\top = \frac{\partial C}{\partial w_{lk}}}$$

Equation 4 - find the derivative of the cost wrt any weight

$$\frac{\partial C}{\partial w_{jk}} = a_{k-1}^{L-1} \delta_j^L \quad (\text{Scalar})$$



have  $\frac{\partial C}{\partial z_j^L}$   
know  $\frac{\partial C}{\partial w_{jk}} = \frac{\partial C}{\partial z_j^L} \frac{\partial z_j^L}{\partial w_{jk}}$

$$z_j^L = (a_{k-1} + b_{k-1} + c_{k-1} + \dots + w_k) \quad \text{dot product}$$

$$z_j^L = (a_{k-1} + b_{k-1} + c_{k-1} + \dots + w_k) \quad \text{dot product}$$

$$\frac{\partial z_j^L}{\partial w} = k$$

What is  
k?  
activation  
of node  
l-1 from  
node k

Tying everything together

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

$$\frac{\partial c}{\partial w} \begin{cases} \vec{e} & \\ \frac{1}{m} \sum_{i=1}^m (w^T x + b - y_i) x^T & \end{cases}$$

$$\frac{\partial c}{\partial b} = e^i = \frac{\partial c}{\partial b} = \delta^L$$

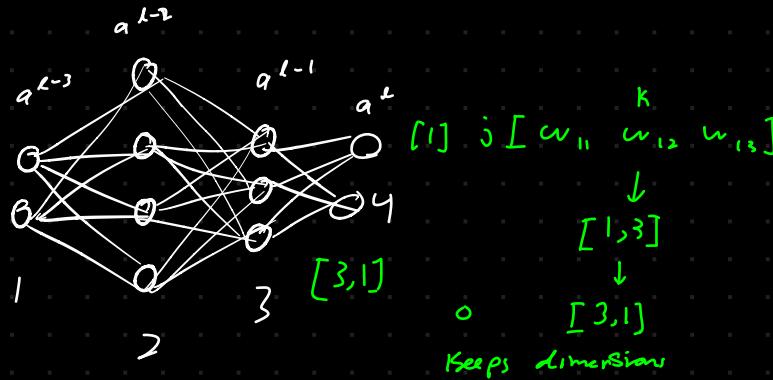
$$\frac{\partial c}{\partial w} = e^i x^T = \frac{\partial c}{\partial w^i} = \delta^L (a^{L-1})^T$$

\* These are equivalent

$e^i$   
 $e^i x^T$  are wrt only one node



① Feed forward and store  $z^l$  and  $a^l$  values at all  $L$  layers



② compute the cost

③ use eq 1 to calculate  $\delta^L$

④ use eq 2 to compute all  $\delta^L$  terms

⑤ use eq 3 to find bias derivatives  $\frac{\partial c}{\partial b^L}$

⑥ use eq 4 to find weight derivatives  $\frac{\partial c}{\partial w^i}$

$$\theta = \theta - \alpha \nabla_{w,b} C$$

forward pass  $\xrightarrow{\hspace{2cm}}$   
 back propagation

$\xleftarrow{\hspace{2cm}}$   
 back propagation

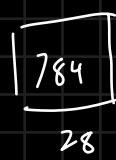
$$\nabla_{w,b} C = \begin{bmatrix} & \\ & \\ & \end{bmatrix}$$

unroll all matrices / flatten

$$\begin{bmatrix} a \\ b \\ c \\ d \\ e \\ f \end{bmatrix} = \begin{bmatrix} a \\ b \\ c \\ d \\ e \\ f \end{bmatrix}$$

C U

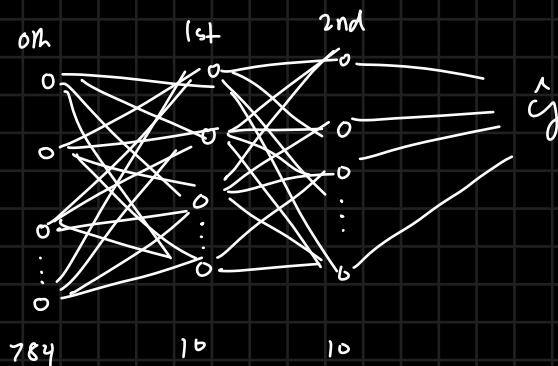
Start



28      n training images

$$dz' =$$

$$x = \begin{bmatrix} x^1 & x^2 & \dots & x^n \end{bmatrix} \Rightarrow \begin{matrix} 0, 1, 2, 3, \dots, 9 \\ 10 \text{ classes} \end{matrix}$$



### ① Forward Propagation

$$A^0 = x \quad (\text{first layer})$$

$$z^1 = w^1 A^0 + b^1$$

$10 \times m$      $10 \times 784$      $784 \times m$      $10 \times m$

Apply activation function (otherwise layers are just linear combination of previous layers)

use ReLU ( $Z'$ )

adds complexity rather than a linear model



$$Z^2 = w^2 A^1 + b^2$$

$$\frac{e^{z_0}}{\sum_{i=1}^k e^{z_i}}$$

$A^2 = \text{softmax}(Z^2)$   
our outputs  
are converted  
between 0 and 1  
to a probability