

NNs Lecture

Kenny Sue



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

Notation

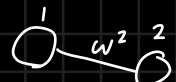
m = size of training Set

n = # of input Variables x_1, x_2, \dots

L = amount of layers in NN

L = specific layer

w^L_k = weight (s) k = layer weights going into



b^L - bias for layer L



x_i = single input variable

Big picture

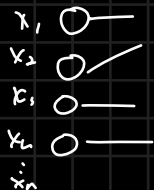
NN as functions (little functions of weights & biases)

Big calculus problem

$$\rightarrow \frac{\partial \text{Cost}}{\partial w}$$

Find how much each weight & bias contributes to the cost

with respect to every weight and bias



$$f_1 = 2x + 4y^3 \rightarrow \frac{\partial f_1}{\partial x} = 2 \quad \frac{\partial f_1}{\partial y} = 12y^2$$

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

$$f(x, y) = \begin{bmatrix} 2x + 4y^3 \\ e^y - 13x \end{bmatrix}$$

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

$$J = \begin{bmatrix} \frac{\partial f_1}{\partial x} & \frac{\partial f_1}{\partial y} \\ \frac{\partial f_2}{\partial x} & \frac{\partial f_2}{\partial y} \end{bmatrix} \begin{Bmatrix} \nabla f_1^T \\ \nabla f_2^T \end{Bmatrix} \text{ Transposed} = \begin{bmatrix} 2 & 12y^2 \\ -13 & e^y \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

Linear combination

Matrix calculation Review

Gradients, Jacobians, Jacobian chain rule

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

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

Vector function

$$f(x, y) = x^2 + \cos y$$

$$\begin{bmatrix} x \\ y \end{bmatrix} \xrightarrow{f(x, y)} (s)$$

$x=2, y=0$

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

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

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

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

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

Notation: upside down triangle

Jacobians

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

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

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

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

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

Recall: Jacobian Matrices

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

How does g change as (x, y) changes

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

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

Apply the Scalar chain rule

$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial u} \frac{\partial u}{\partial x} \quad \text{General notation}$$

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

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

what if we cannot use chain rule?

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

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

put a 1 when you are not using chain rule
aka / no intermediate function

Scalar chain rule

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

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

$$\text{let } g = x^2$$

$$f = \sin(g)$$

technical breakdown

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial x} = 2x \cos g \quad \text{Substitute back}$$

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

$$g(h(c(a(x))))$$

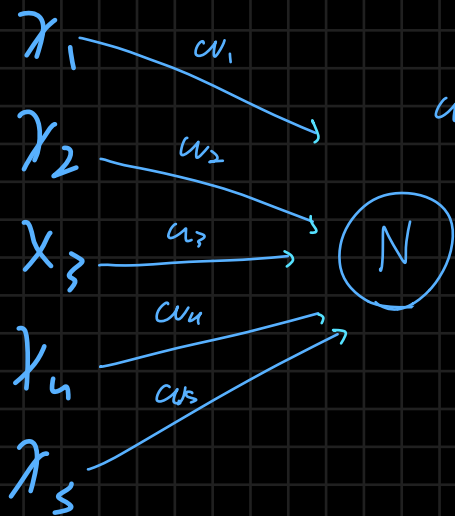
use product of partial derivatives

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

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

compute Jacobians of both

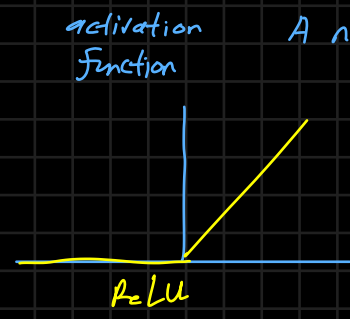
Neurons



$$w_1 x_1 + w_2 x_2 \dots$$

$$a \left(\sum_{i=1}^n x_i w_i + b \right)$$

Scalar



A neural network with no non-linear activation function is linear regression

take the dot product

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \cdot \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ \vdots \\ w_n \end{bmatrix}$$

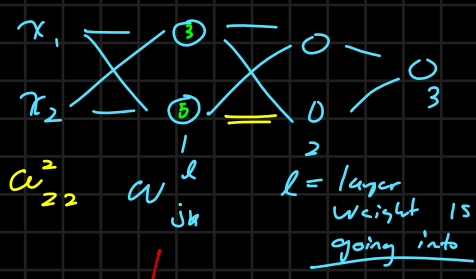
$$\rightarrow a(\underbrace{\vec{x}^T \vec{w}}_z + b)$$

Notations

$$z = (\vec{x}^T \vec{w} + b)$$

$$a = a(z)$$

LAYER of NN - NOTATION

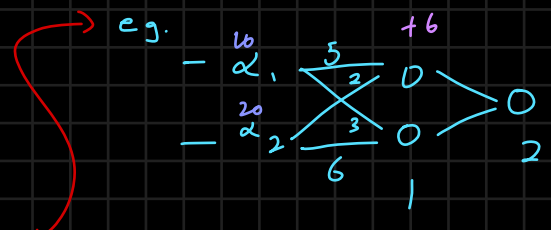


$$b_j^l$$

$$b_1^l = 3$$

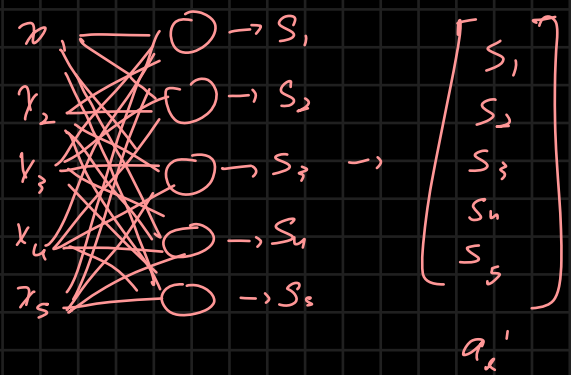
$$b_2^l = 3$$

l = layer weight is going into
 j = # of the neuron of the layer l
 k = # of the neuron in layer $l-1$

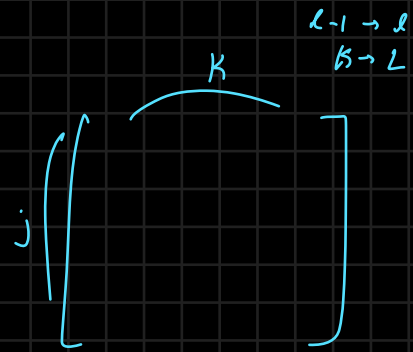


$$w^l = \begin{bmatrix} (1,1) & (1,2) \\ (2,1) & (2,2) \end{bmatrix} = \begin{bmatrix} 5 & 2 \\ 3 & 6 \end{bmatrix}$$

$$a(w^l \cdot a^{l-1} + b^l)$$



single node
 $a(\vec{x}^T \vec{w} + b)$
Entire layer
 w^l
 w_{jk}^l



activation of the previous layer

$$= \begin{bmatrix} 3 & 2 \\ 3 & 6 \end{bmatrix} \begin{bmatrix} 10 \\ 20 \end{bmatrix}$$

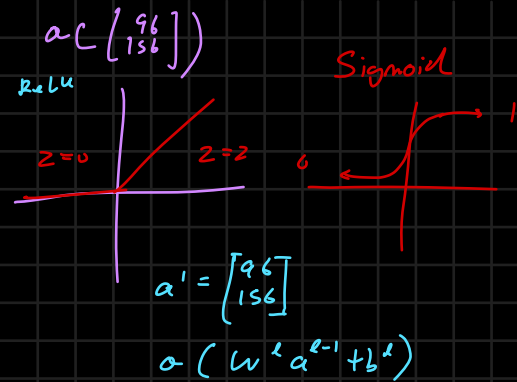
calculates the weighted sum

$$= \begin{bmatrix} 90 \\ 150 \end{bmatrix}$$

another dot product visualization

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

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



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

output is a vector (a_1)

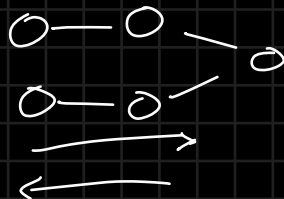
vector \rightarrow vector

$\rightarrow \sigma(w^2 a^1 + b^2)$

$a_2 \rightarrow \dots$

feed forward through an NN

input through the network



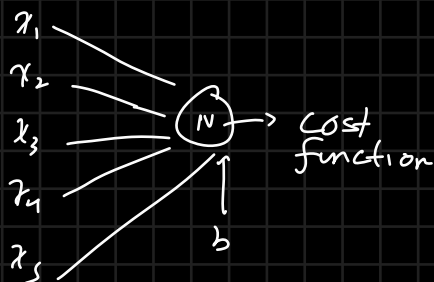
Jacobians & NNs

Minimize the cost of the output (error)

Mean Squared error

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

↓ does not matter (for the power rule)



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

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

Derivatives of Neuron Operators (Hadamard product)

→ Hadamard product

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

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

with respect to \vec{v}

$$\frac{d}{dw}(xw) = x$$

Jacobians of operations

Binary element-wise operations
Hadamard product (\odot)

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

$$\begin{bmatrix} \vec{v} \\ \vec{w} \end{bmatrix} \xrightarrow{\odot} \begin{bmatrix} \vec{z} \end{bmatrix}$$

element wise operators

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

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

operations on the functions before F

Can get 2 Jacobians

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

General representation

$$J = \frac{\partial F}{\partial \vec{v}} = \begin{bmatrix} \frac{\partial}{\partial v_1} f_1(\vec{v}) \odot g_1(\vec{w}) & \frac{\partial}{\partial v_2} f_1(\vec{v}) \odot g_1(\vec{w}) & \dots \\ \frac{\partial}{\partial v_1} f_2(\vec{v}) \odot g_2(\vec{w}) & \frac{\partial}{\partial v_2} f_2(\vec{v}) \odot g_2(\vec{w}) & \dots \\ \vdots & \vdots & \ddots \end{bmatrix}$$

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

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

Derivative of a Scalar expansion

$$2 \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 2v_1 \\ 2v_2 \\ 2v_3 \\ \vdots \\ 2v_n \end{bmatrix}$$

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

$$g(\vec{v}) = \vec{1} \cdot x \text{ (expands into a vector)}$$

The act of multiplying x by the ones vector is an act of broadcasting itself

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

x is a scalar

$$= \begin{bmatrix} f_1(\vec{v}) \circ g_1(x) & F_1 \\ f_2(\vec{v}) \circ g_2(x) & F_2 \\ \vdots & \vdots \\ f_n(\vec{v}) \circ g_n(x) & F_n \end{bmatrix} \frac{\partial F}{\partial \vec{v}} = \begin{bmatrix} \frac{\partial f_1}{\partial v_1} & \frac{\partial f_1}{\partial v_2} & \dots & \frac{\partial f_1}{\partial v_n} \\ \frac{\partial f_2}{\partial v_1} & \frac{\partial f_2}{\partial v_2} & \dots & \frac{\partial f_2}{\partial v_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_n}{\partial v_1} & \frac{\partial f_n}{\partial v_2} & \dots & \frac{\partial f_n}{\partial v_n} \end{bmatrix} = \begin{bmatrix} x & \dots & 0 \\ \vdots & x & \vdots \\ 0 & \dots & x \end{bmatrix}$$

$$\nabla F_x = \begin{bmatrix} \frac{\partial f_1}{\partial x} \\ \frac{\partial f_2}{\partial x} \\ \vdots \\ \frac{\partial f_n}{\partial x} \end{bmatrix} = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix}$$

with respect to x

don't forget, this is the scalar vector.

Derivative of a Sum

$$\vec{v} = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix} = \sum_{i=1}^n (g(v))_i = \text{---S---}$$

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

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

The derivative of a sum is the sum of the derivatives

$$J = \begin{bmatrix} \sum_{i=1}^n \frac{\partial}{\partial v_1} g_i(v) & \sum_{i=1}^n \frac{\partial}{\partial v_2} g_i(v) & \dots & \sum_{i=1}^n \frac{\partial}{\partial v_n} g_i(v) \end{bmatrix} \xrightarrow{\text{If } g(\vec{v}) = 2\vec{v}} [2 \ 2 \ 2 \ \dots \ 2]$$

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

$$J = \begin{bmatrix} \sum_{i=1}^n \frac{\partial}{\partial v_1} v_i & \sum_{i=1}^n \frac{\partial}{\partial v_2} v_i & \dots & \sum_{i=1}^n \frac{\partial}{\partial v_n} v_i \end{bmatrix} = [1 \ 1 \ \dots \ 1]$$

1+0+... 0+1+0... 0+0+...+1

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

dot product = Hadamard product

pulling all concepts together

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

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

Does not change

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

$$a = \sigma(\text{Sum}(w \odot x) + b)$$

$$H = \text{Hadamard product } (\sum w \odot x + b)$$

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

$$\frac{\partial a}{\partial w} = \frac{\partial a}{\partial z} \begin{bmatrix} 1 \\ x^T \end{bmatrix}$$

Simplifies after all 3 substitutions

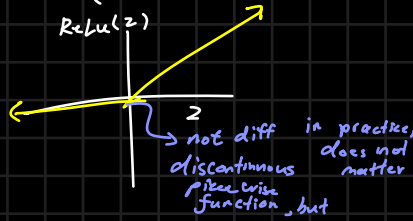
$$\frac{\partial H}{\partial w} = \begin{bmatrix} x_1 & \dots & 0 \\ \vdots & x_2 & \vdots \\ 0 & \vdots & x_n \end{bmatrix}$$

diag(x_1, \dots, x_n)

$$\frac{\partial S}{\partial H} = \vec{1} \text{ (no g function)}$$

$$\max(0, z)$$

of 2 numbers



taking the derivative

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

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

$$a = \sigma(\sum (w \odot x) + b)$$

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

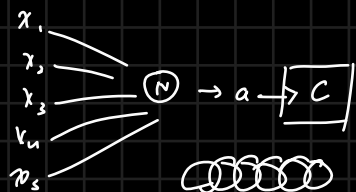
$$\frac{\partial a}{\partial b} = \begin{cases} 0 \times 1 = \vec{0} & \text{if } z \leq 0 \\ \vec{1} \times 1 = \vec{1} & \text{if } z > 0 \end{cases} = [\vec{0}]$$

(important so that the vector operations make sense dimensionally)

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

The Gradient of the Loss Function

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



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

$$\frac{\partial C}{\partial b}$$

ReLU

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

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

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

$$\begin{bmatrix} x_{11} & x_{21} & \dots & x_{n1} \\ x_{12} & x_{22} & \dots & x_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ x_{1m} & x_{2m} & \dots & x_{nm} \end{bmatrix}$$

let then, $(y - a^i) = v$

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

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

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

M training examples not just one

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

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

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

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

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

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

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

what does this mean? taking C as a piecewise

$$= \frac{1}{m} \sum_{i=1}^m \begin{cases} [-0]^T & \text{if } \dots \\ -v \frac{\partial v}{\partial w} & \text{if } \dots \end{cases}$$

make the options negative
DOES not actually change the piecewise

$$= \frac{1}{m} \sum_{i=1}^m \begin{cases} -[0]^T & \text{if } \dots \\ -v x^T & \text{if } \dots \end{cases}$$

$$= \frac{1}{m} \sum_{i=1}^m \begin{cases} [0]^T & \text{if } \dots \\ -(y - a^i) x^T & \text{if } \dots \end{cases}$$

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

$$\begin{cases} 0 & \text{if } w^T x + b \leq 0 \\ -(y - (\max(0, w^T x + b))) x^T & \text{if } w^T x + b > 0 \end{cases}$$

always the case

$$= -(y - (w^T x + b)) x^T$$

$$= (w^T x + b - y) (x^T)$$

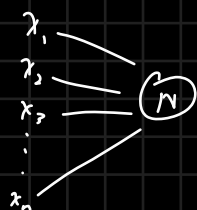
$$\frac{1}{m} \begin{cases} 0^T \\ \sum_{i=1}^m (w^T x + b - y) (x^T) \end{cases}$$

Gradient Descent Intuition

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

Want:

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



$$\frac{1}{n} \sum_{i=1}^m e_i x^T$$

Consider a general case

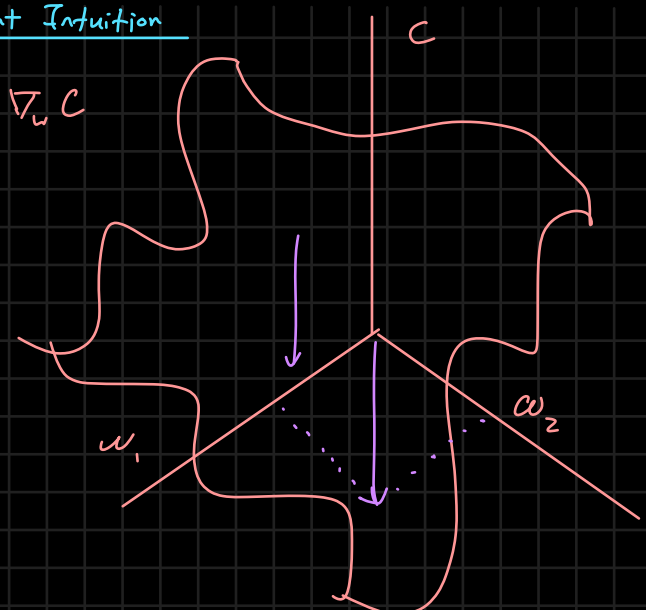
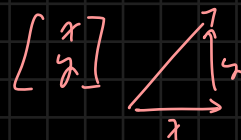
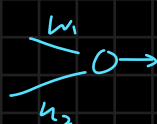
$$\frac{1}{n} \begin{bmatrix} e_1 x_1 + e_2 x_1 \dots e_5 x_1 \\ e_1 x_2 + e_2 x_2 \dots e_5 x_2 \\ e_1 x_3 + e_2 x_3 \dots e_5 x_3 \\ e_1 x_4 + e_2 x_4 \dots e_5 x_4 \\ e_1 x_5 + e_2 x_5 \dots e_5 x_5 \end{bmatrix} \Rightarrow \begin{bmatrix} g_1 \\ g_2 \\ g_3 \\ g_4 \\ g_5 \end{bmatrix}$$

Consider $m=1$ and 5 inputs

$$\frac{\partial C}{\partial \vec{w}} = \begin{bmatrix} e_1 x_1 \\ e_1 x_2 \\ \vdots \\ e_1 x_5 \end{bmatrix} \begin{matrix} \frac{\partial C}{\partial w_1} \\ \frac{\partial C}{\partial w_2} \\ \vdots \\ \frac{\partial C}{\partial w_5} \end{matrix}$$

derivative of the cost with respect to all w 's, averaged out over all training examples

$$\frac{\partial C}{\partial w} \begin{bmatrix} e_1 x_1 \\ e_1 x_2 \end{bmatrix} \nabla_w C$$



Vectors can be translated

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

Derivative of the bias

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

$$y - a^i = v$$

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

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

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

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

max redundancy
 $=(w^T x + b - y)$

$$\begin{cases} 0 & \text{if } \dots \\ \frac{1}{n} \sum_{i=1}^m (w^T x + b - y) & \text{if } \dots \end{cases} e_i$$

$$\frac{\partial C}{\partial b} = \frac{1}{n} [e_1 + e_2 + e_3]$$

$\begin{bmatrix} e_1 x_1 \\ e_1 x_2 \end{bmatrix}$ 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

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

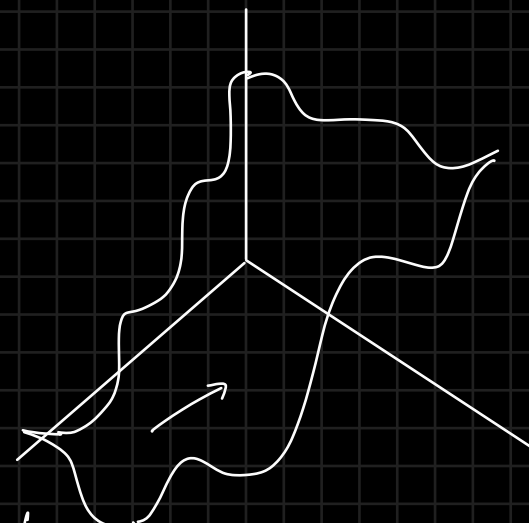
$$= \frac{1}{2n} \sum_{i=1}^m \frac{\partial}{\partial b} (v^2) \left(\frac{dv^2}{dv} \frac{dv}{db} \right)$$

$$= \frac{1}{2n} \sum_{i=1}^m \cancel{v} \frac{dv}{db}$$

$$= \frac{1}{2n} \sum_{i=1}^m v \frac{dv}{db}$$

Gradient Descent & Stochastic Gradient Descent

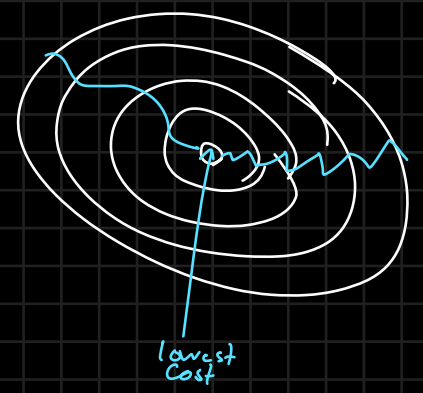
$$\begin{bmatrix} \frac{\partial C}{\partial w_1} \\ \frac{\partial C}{\partial w_2} \\ \frac{\partial C}{\partial w_3} \\ \vdots \\ \frac{\partial C}{\partial b_1} \\ \frac{\partial C}{\partial w_n} \end{bmatrix} \nabla_{w,b} C$$



The gradient points in the direction of steepest descent

Mini-batch GD

SGD \rightarrow batching smaller datasets and then, running gradient descent on each batch \rightarrow then, you can take as many steps, as you have batches
Shuffle the data well



Θ (all weights & biases)

Initially starts as randomly assigned

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

learning rate

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

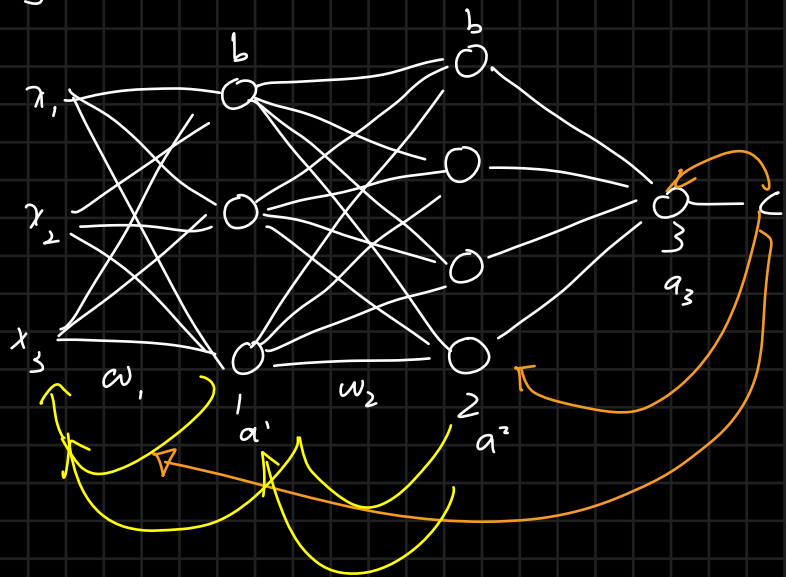
α (learning rate)
depends on research

Training batch: m

\hookrightarrow Can be very computationally expensive

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



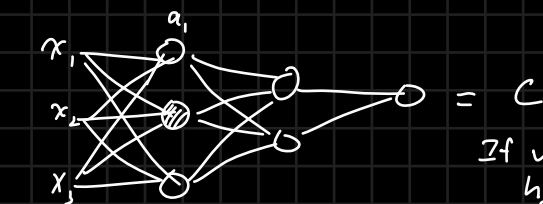
$$\frac{\partial C}{\partial w_1} = \frac{\partial a^1}{\partial w_1} \frac{\partial a^2}{\partial a^1} \frac{\partial a^3}{\partial a^2} \frac{\partial C}{\partial a^3}$$

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

Doing these calculations to compute the derivatives going forward is resource-intensive

$$\delta_j^k$$

The error of a node



If we add Δ to a node,
how does the cost
function change

$$z'_2 = w^L x + b + \Delta$$

z

$$z'_2 + \Delta$$

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

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

$$\delta_2^L = \frac{\partial C}{\partial z_2^L}$$

How much error is being caused
by a node

1st equation - error of last node a^L

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

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

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

Now, do the entire layer

$$\nabla_{a^L} C$$

derivative

$$\begin{bmatrix} \frac{\partial C}{\partial z_1^L} \\ \frac{\partial C}{\partial z_2^L} \\ \vdots \end{bmatrix}$$

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

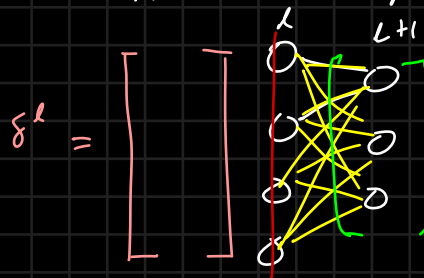
element-wise
operation

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

$$\frac{\partial C}{\partial z_1^L} = \frac{\partial C}{\partial a_1^L} \frac{\partial a_1^L}{\partial z_1^L}$$

δ_1^L

$$\text{error of any node: } \delta^L = ((w^{L+1})^T \delta^{L+1}) \cdot \sigma'(z^L)$$



weight
matrix
connecting
 $L, L+1$

derivatives of
the activations

total
error
of $L+1$
layer

$$w_{jk}^L$$

$$w = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}_j$$

K

$$w^T = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}_k$$

j

K

proof

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

$$\frac{\partial C}{\partial z_1^L} = \frac{\partial C}{\partial a_1^L} w_{1j}^L \delta_j^{L+1}$$

$L=1$

works for any layer

Equation 3 - derivative of cost wrt any bias

$$\frac{\partial c}{\partial b^l} = \delta^l$$

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

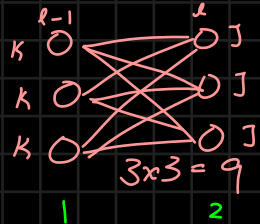
$$\sigma(\underbrace{w_1 a_0 + b_1}_{z_1}) \Rightarrow \sigma(\underbrace{w_2 a_1 + b_2}_{z_2})$$

$$\delta^l = \frac{\partial c}{\partial z_1}$$

$$\frac{\partial c}{\partial b_1} = \frac{\partial z_1}{\partial b_1} \frac{\partial c}{\partial z_1}$$

$$\frac{\partial c}{\partial b_1} = 1$$

Vectorized form:



$$w_2 = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \end{bmatrix} \quad (3,3)$$

$$\frac{\partial c}{\partial w_2} = \begin{bmatrix} \frac{\partial c}{\partial w_{11}} & \frac{\partial c}{\partial w_{12}} & \frac{\partial c}{\partial w_{13}} \\ \vdots & \vdots & \vdots \end{bmatrix}$$

$$\frac{\partial c}{\partial w_2} = \begin{bmatrix} a_1^1 \delta_1^2 & a_1^2 \delta_1^2 & a_1^3 \delta_1^2 \\ a_1^1 \delta_2^2 & a_1^2 \delta_2^2 & a_1^3 \delta_2^2 \\ a_1^1 \delta_3^2 & a_1^2 \delta_3^2 & a_1^3 \delta_3^2 \end{bmatrix}$$

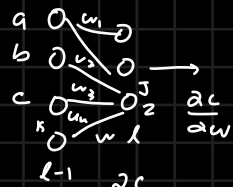
$$\delta^2 = \begin{bmatrix} \delta_1^2 \\ \delta_2^2 \\ \delta_3^2 \end{bmatrix} \quad a_1 = \begin{bmatrix} a_1^1 \\ a_1^2 \\ a_1^3 \end{bmatrix}$$

$$= \vec{\delta}^2 (\vec{a}^{l-1})^T$$

\downarrow
 $(n,1)$ $(1,m)$ m n

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

$$\frac{\partial c}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l$$



$$\frac{\partial c}{\partial w_{jk}^l} = \frac{\partial z_j^l}{\partial w_{jk}^l} \frac{\partial c}{\partial z_j^l}$$

$$z_j^l = (a w_1 + b w_2 + c w_3 + k w + b)$$

$$\frac{\partial z_j^l}{\partial w} = k = a_k^{l-1}$$

Tying everything together

$$\frac{\partial C}{\partial b} \begin{cases} 0 & w^L x \leq 0 \\ \frac{1}{n} \sum_{i=1}^n (w^L x + b - y_i) & \text{if } w^L x > 0 \end{cases}$$

$$\frac{\partial C}{\partial w} \begin{cases} 0 \\ \frac{1}{n} \sum_{i=1}^n (w^L x + b - y_i) x^T \end{cases}$$

$$\frac{\partial C}{\partial b} = e_i \quad \frac{\partial C}{\partial b} = \delta^L$$

$$\frac{\partial C}{\partial w} = e_i x^T \quad \frac{\partial C}{\partial w^L} = \delta^L (a^{L-1})^T$$

BP Equations

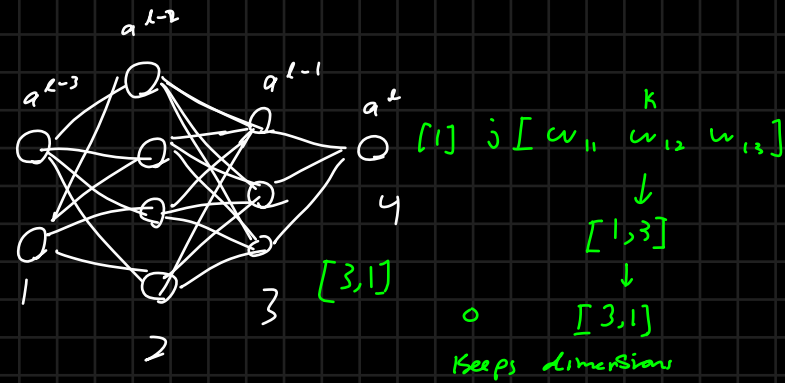
$$\delta^L = \nabla_a C \circ \sigma'(z^L)$$

$$\delta^L = ((w^{L+1})^T \delta^{L+1}) \circ \sigma'(z^L)$$

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

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

① Feed forward and store z's and a's values



② Compute the cost

③ Use eq 1 to calculate δ^L

④ Use eq 2 to compute all δ^L terms

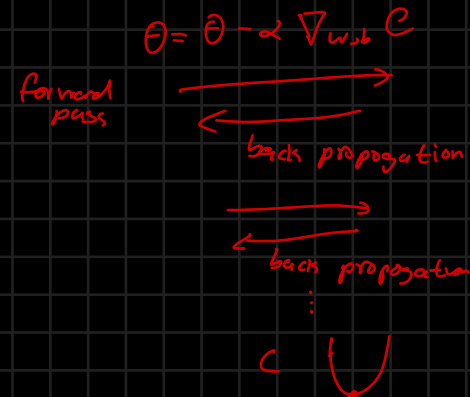
⑤ Use eq 3 to find bias derivatives $\frac{\partial C}{\partial b^L}$

⑥ Use eq 4 to find weight derivatives $\frac{\partial C}{\partial w^L}$

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

Unroll all matrices

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



Start

$$\begin{matrix} \boxed{784} & 28 \\ 28 \end{matrix}$$

n training images

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

10 classes

