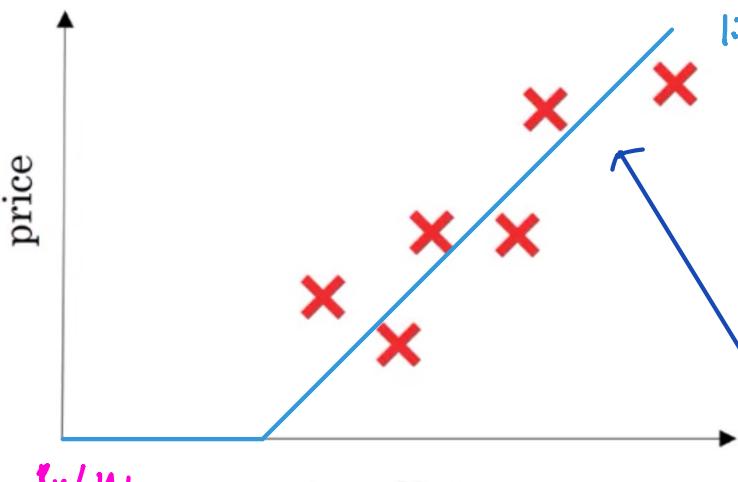


Course 1: Neural Network and Deep Learning

week 1 (Introduction)

what is Neural Network?

e.g. housing price prediction



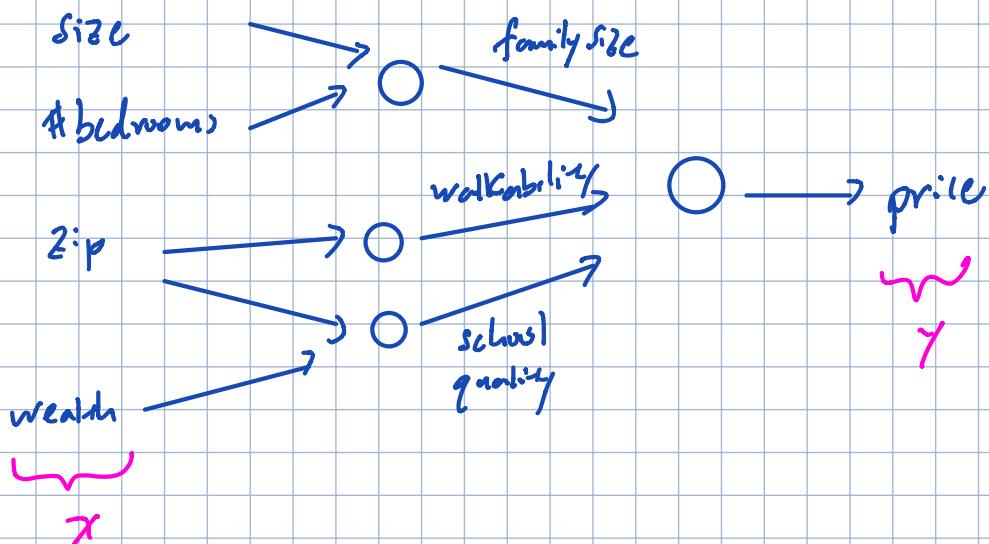
RuLu:
Rectified
Linear
Unit

linear regression by housing price
cannot be negative

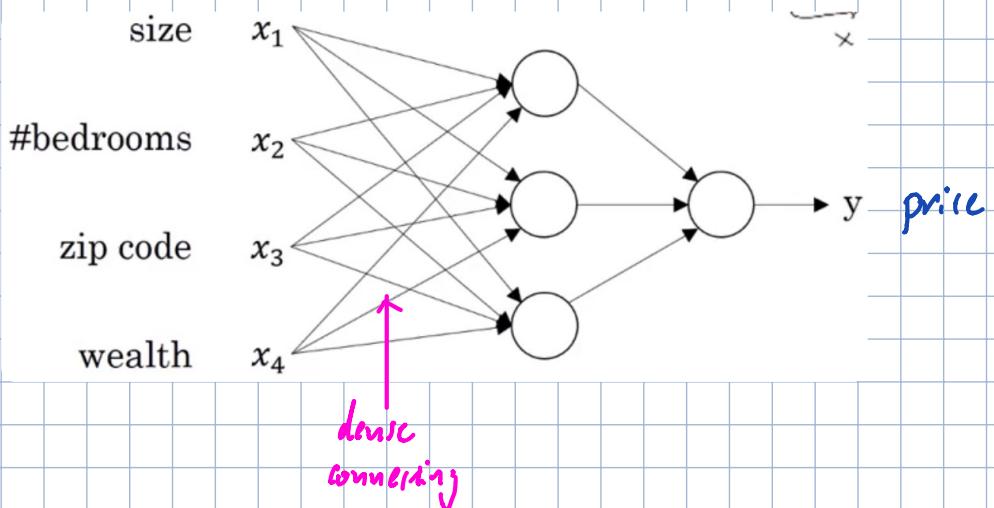
⇒ simplest neural network
input
size → ○ → output
X 'neuron' Y
implemented this function

All this neuron does is it inputs the size (x_1), computes this linear function, takes a max of 0, and then outputs the estimated price.

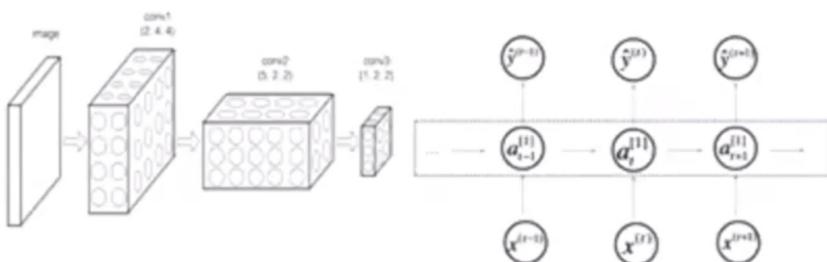
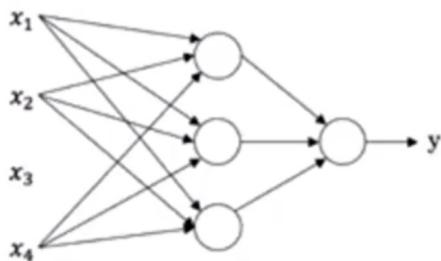
Say now we have multiple features: size, no. of bedrooms etc.



A neural network with 4 inputs



Neural Network Examples



Structured Data vs Unstructured Data

Structured Data

| Size | #bedrooms | ... | Price (1000\$) |
|------|-----------|-----|----------------|
| 2104 | 3 | | 400 |
| 1600 | 3 | | 330 |
| 2400 | 3 | | 369 |
| : | : | | : |
| 3000 | 4 | | 540 |

Unstructured Data

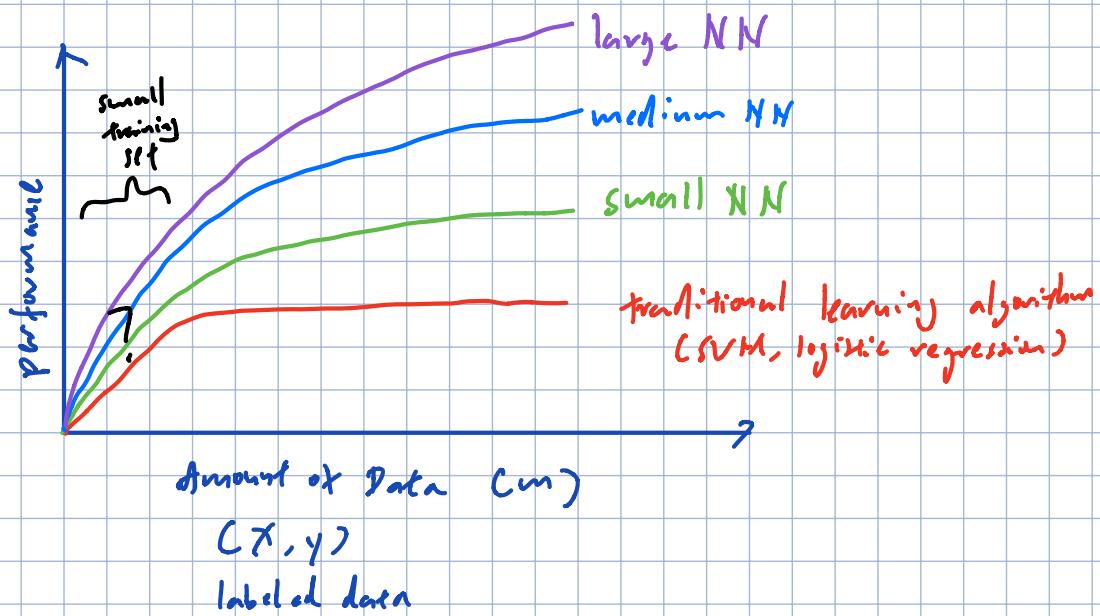


| User Age | Ad Id | ... | Click |
|----------|-------|-----|-------|
| 41 | 93242 | | 1 |
| 80 | 93287 | | 0 |
| 18 | 87312 | | 1 |
| : | : | | : |
| 27 | 71244 | | 1 |

Four scores and seven years ago...

Text

Scale drives deep learning progress



Week 2 (basics)

If we have a training set of m training examples, we might use a `for` loop, step through m training examples.

For neural network, usually want to process an entire training set, w/o using an explicit `for` loop to loop over the entire training set.

Binary Classification



y

→ 1 (cat) vs 0 (non cat)

| | | | | |
|-------|-----|-----|-----|-----|
| Blue | 255 | 134 | 93 | 22 |
| Green | 255 | 134 | 202 | 22 |
| Red | 255 | 231 | 42 | 22 |
| | 123 | 94 | 83 | 2 |
| | 34 | 44 | 187 | 92 |
| | 34 | 76 | 232 | 124 |
| | 67 | 83 | 194 | 202 |

unroll

feature vector

$X =$

$$\begin{bmatrix} 255 \\ 231 \\ \vdots \\ 255 \\ 134 \end{bmatrix}$$

if image is 64×64

$$n_x = 64 \times 64 \times 3$$

$$= 12288$$

Notation:

(x, y) : a single training example

$x \in \mathbb{R}^{n_x}$, $y \in \{0, 1\}$

m : no. of training examples: $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$

$\hookrightarrow m_{\text{train}}$: no. of training set examples

$\hookrightarrow m_{\text{test}}$: no. of test set examples

entire training set

compact representation of training examples:

$$X = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & | \end{bmatrix} \quad \begin{array}{c} \uparrow \\ n_x \\ \downarrow \end{array} \quad \text{i.e. } X \in \mathbb{R}^{n_x \times m}$$

$$Y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}] \quad \text{i.e. } Y \in \mathbb{R}^{1 \times m}$$

Logistic Regression

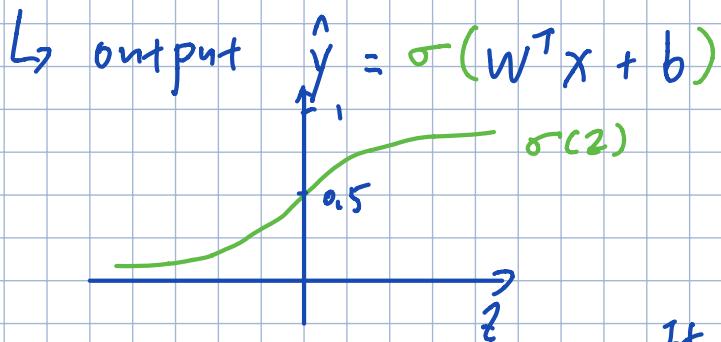
Given an input feature vector X , want output prediction \hat{y} ,

which is an estimate of y .

Formally: $\hat{y} = P(Y=1|X)$

$$X \in \mathbb{R}^{n_x}$$

Parameters of logistic regression: $w \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$



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

If z is large, $\sigma(z) \approx 1$

z is small, $\sigma(z) \approx 0$
large -ve

sometimes, the intercept, b , is incorporated in feature

i.e. $x_0 = 1$, $x \in \mathbb{R}^{n+1}$

↳ $\hat{y} = \sigma(\theta^T x)$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_{n_x} \end{bmatrix} \quad \left. \begin{array}{l} \text{b} \\ \text{w} \end{array} \right\}$$

For neural network, it is easier to separate these parameters.

Logistic Regression Cost Function

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

we have $\hat{y}^{(i)} = \sigma(w^T x^{(i)} + b)$, where $\sigma(z) = \frac{1}{1+e^{-z}}$

Given $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$, we want $\hat{y}^{(i)} \approx y^{(i)}$

Loss (error) function:

option 1: squared error: $L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$

Generally Don't use this → optimisation problem becomes non-convex

i.e. gradient descent will find local optima (i.e. w)

option 2: (actual loss function for logistic regression)

$$L(\hat{y}, y) = -(y \log \hat{y} + (1-y) \log (1-\hat{y}))$$

for optimisation problems, we want L as small as possible

If $y=1$, $L(\hat{y}, y) = -\log \hat{y}$ ← want $\log \hat{y}$ large

↳ \hat{y} large (can never exceed 1)

If $y=0$, $L(\hat{y}, y) = -\log(1-\hat{y})$ ← want $\log(1-\hat{y})$ large

↳ \hat{y} small

(can never be below 0)

* Loss function was defined w.r.t. a single training example

* Cost function is w.r.t. the entire training set

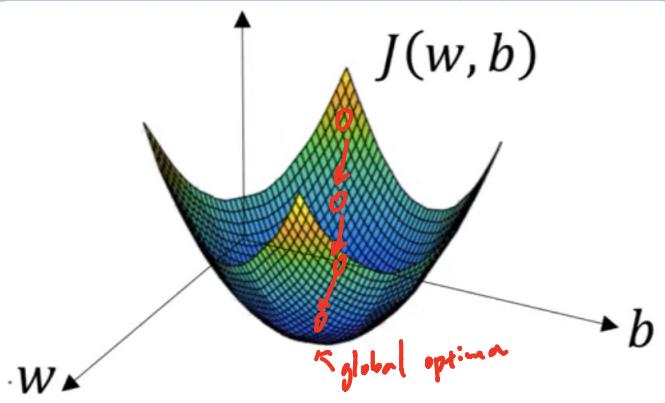
$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)}) \right]$$

gradient using a set
of w, b

Gradient Descent

From the above cost function, we want to find w, b that

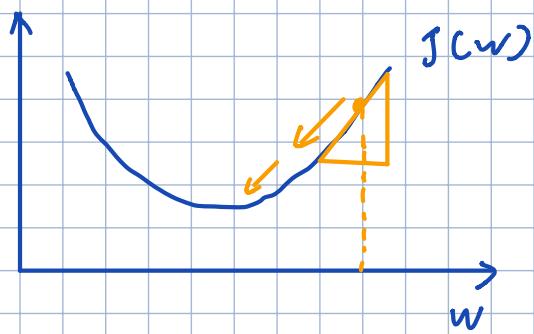
minimise $J(w, b)$



→ a convex function $J(w, b)$

↳ reason to use this J for logistic regression

Say we have some function $J(w)$ (ignore b for now)



Gradient Descent:

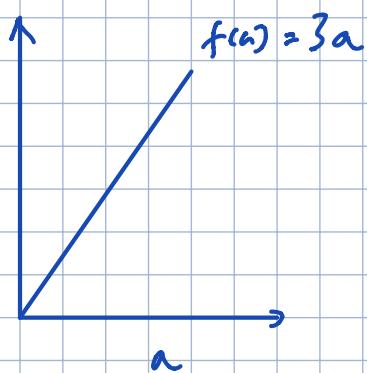
Repeat {
 $w := w - \alpha \frac{dJ(w)}{dw}$
} learning rate

till algorithm converges

For logistic regression, $J(w, b)$

Repeat {
 $w := w - \alpha \frac{\partial J(w, b)}{\partial w}$
 $b := b - \alpha \frac{\partial J(w, b)}{\partial b}$
}

Derivatives



$$a=2, f(a)=6$$

$$a=2.001, f(a)=6.003$$

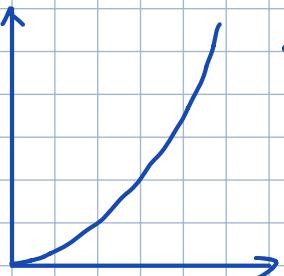
slope (derivatives) of $f(a)$
at $a=2$ is 3

$$a=5, f(a)=15$$

$$a=5.001, f(a)=15.003$$

slope at $a=5$, is also 3

$$\left. \frac{df(a)}{da} \right|_{a=2} = 3$$



$$f(a) = a^2$$

$$a=2 \quad f(a) = 4$$

$$a=2.001 \quad f(a) \approx 4.004$$

Slope at $a=2$ is 4 ($\because \frac{0.004}{0.001} = 4$)

$$\frac{d}{da} f(a) = 4 \text{ at } a=2$$

$$a=5 \quad f(a) = 25$$

$$a=5.001 \quad f(a) \approx 25.010$$

$$\frac{d}{da} f(a) = 10 \text{ at } a=5$$

$$\frac{d}{da} a^2 = 2a \quad \leftarrow$$

Computation Graph

say we have a function

$$J(a, b, c) = 3(a + bc)$$

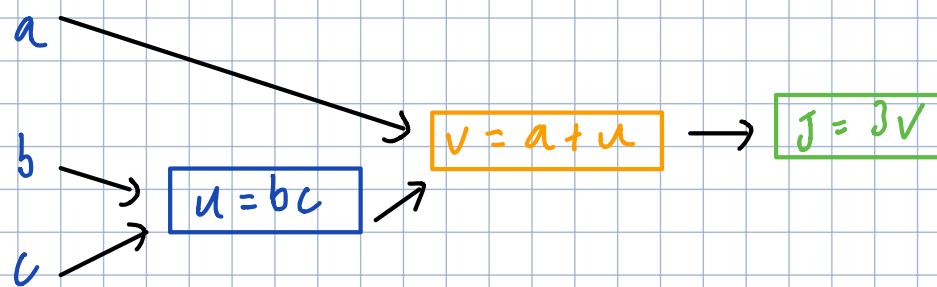
 u
 v
J

computation takes 3 steps:

$$u = bc$$

$$v = a + u$$

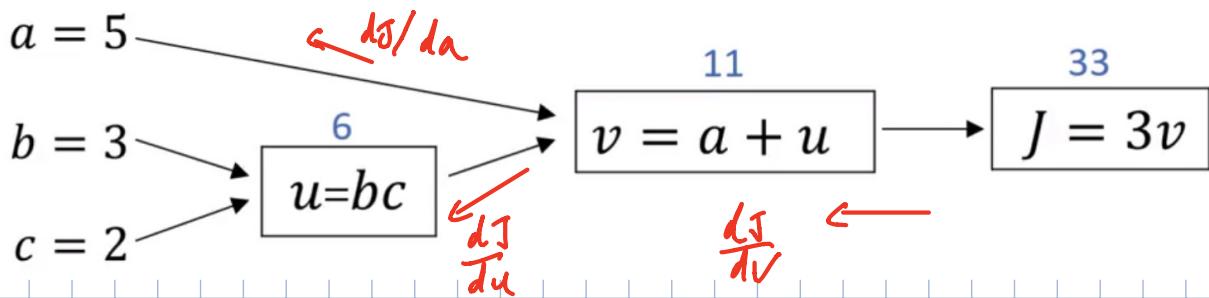
$$J = 3v$$



computation: left \longrightarrow right

derivatives: right \longrightarrow left

Computing Derivatives



Say we want to compute: $\frac{dJ}{dv}$

i.e. if we change value of V by a little bit, how will value of J change?

$$\text{Now: } J = 3v$$

$$v = 11 \rightarrow \text{if } v = 11.001 \\ \hookrightarrow J = 33.003 \quad \left. \begin{array}{l} \frac{dJ}{dv} = 3 \\ \because \frac{0.003}{0.001} = 3 \end{array} \right)$$

How about $\frac{dJ}{da}$?

$$\text{Now: } a = 5 \rightarrow \text{if } a = 5.001$$

$$\hookrightarrow v = a + u = 11.001 \quad \left. \begin{array}{l} \frac{dJ}{da} = 3 \\ J = 33.003 \end{array} \right.$$

$$\Rightarrow \frac{dJ}{da} = \frac{dJ}{dv} \frac{dv}{da} \quad \begin{array}{l} a \text{ cause } V \text{ to change} \\ v \text{ cause } J \text{ to change} \end{array}$$

$$\text{here } \frac{dv}{da} = 1 \quad (\because a \uparrow 0.001, v \uparrow 0.001)$$

How about $\frac{dJ}{du}$?

$$u = 6 \rightarrow \text{if } u = 6.001$$

$$\hookrightarrow v = 11.001$$

$$\hookrightarrow J = 33.003$$

$$\hookrightarrow \frac{dJ}{du} = 3 = \frac{dJ}{dv} \frac{dv}{du}$$

How about $\frac{dJ}{db}$?

$$\frac{dJ}{db} = \frac{dJ}{du} \frac{du}{db} = 6$$

if $b=3 \rightarrow 3.00$
 $\hookrightarrow u = b \times c = 6 \rightarrow 6.00$
i.e. $\frac{du}{db} = 2$

Logistic Regression Gradient Descent (for 1 training example)

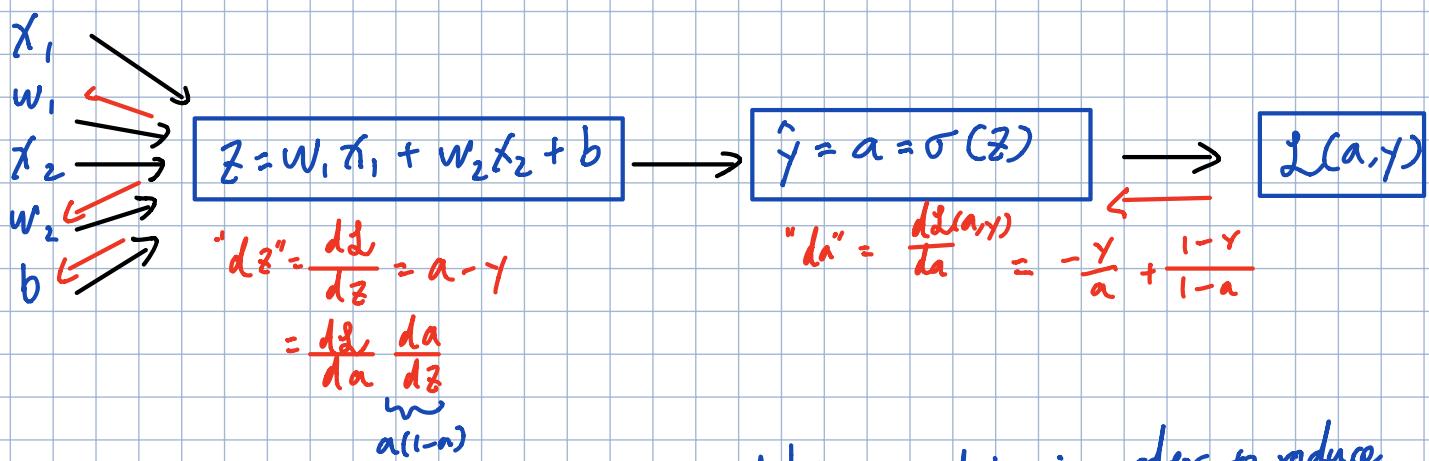
Recap:

$$Z = w^T x + b$$

$$\hat{y} = a = \sigma(Z)$$

$$L(a, y) = - (y \log(a) + (1-y) \log(1-a))$$

says we have 2 features: x_1, x_2



In logistic regression, we want to modify w and b , in order to reduce this loss L .

$$\frac{\partial L}{\partial w_1} = "dw_1" = x_1 \cdot dZ$$

$\leftarrow \frac{dL}{dZ}$
 $\leftarrow \frac{dZ}{dw_1}$

$$dw_2 = x_2 \cdot dZ$$

$\leftarrow \frac{dL}{dZ}$
 $\leftarrow \frac{dZ}{dw_2}$

$$db = dZ$$

i.e. compute dZ , then can get dw_1, dw_2, db

$$\hookrightarrow w_1 := w_1 - \alpha d w_1$$

Note: these are for 1 training example

$$w_2 := w_2 - \alpha d w_2$$

i.e. they are actually

$$b := b - \alpha d b$$

$d w_1^{(i)}, d w_2^{(i)}, d b^{(i)}$ for $(x^{(i)}, y^{(i)})$

Gradient Descent on m examples

Recall the cost function for logistic regression

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \ell(a^{(i)}, y^{(i)})$$

$$\text{where } a^{(i)} = \hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(w^T x^{(i)} + b)$$

For m training examples, the derivatives w.r.t the parameters (w_1, w_2, \dots, b) is also going to be the average of derivatives w.r.t to the parameters.

e.g.

$$\frac{\partial}{\partial w_1} J(w, b) = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{\partial}{\partial w_1} \ell(a^{(i)}, y^{(i)})}_{d w_1^{(i)}} \quad \text{for } (x^{(i)}, y^{(i)})$$

The algorithm (1 step of gradient descent)

$$J = 0; dw_1 = 0; dw_2 = 0; db = 0$$

For $i = 1 \text{ to } m$:

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

$$a^{(i)} = \sigma(z^{(i)})$$

$$J += -[y^{(i)} \log a^{(i)} + (1 - y^{(i)}) \log (1 - a^{(i)})]$$

$$dz^{(i)} = a^{(i)} - y^{(i)}$$

$$dw_1 += x_1^{(i)} dz^{(i)} \quad \left. \begin{array}{l} \\ n=2 \end{array} \right.$$

$$dw_2 += x_2^{(i)} dz^{(i)}$$

$$db += dz^{(i)}$$

$$J /= m$$

$$dw_1 /= m; dw_2 /= m; db /= m; \text{ # compute average}$$

Here:

$$dW_i = \frac{\partial J}{\partial w_i}$$

(accumulates to sum over $dW_i^{(i)}$ for the entire training set)

$$w_1 := w_1 - \alpha dW_1$$

$$w_2 := w_2 - \alpha dW_2$$

$$b := b - \alpha db$$

Everything here is just for 1 step of gradient descent

* In practice, avoid using explicit for loop, use vectorisation.

Derivation of $\frac{dL}{dz} = \alpha - y$ ($\alpha \in \hat{J}$)

By chain rule:

$$\frac{dL}{dz} = \frac{dL}{da} \times \frac{da}{dz}$$

Step 1: $\frac{dL}{da}$

$$L = -(y \log(a)) + (1-y) \log(1-a)$$

$$\begin{aligned}\frac{dL}{da} &= -y \times \frac{1}{a} - (1-y) \times \frac{1}{1-a} \times (-1) \\ &= \frac{-y}{a} + \frac{1-y}{1-a} \\ &= \frac{a-y}{a(1-a)}\end{aligned}$$

Step 2: $\frac{da}{dz}$

$$\frac{da}{dz} = \frac{d}{dz} \sigma(z)$$

The derivative of a sigmoid has the form

$$\frac{d}{dz} \sigma(z) = \sigma(z) \times (1 - \sigma(z))$$

↳ since $a = \sigma(z)$

$$\frac{d}{dz} \sigma(z) = \frac{da}{dz} = a(1-a)$$

Step 3: $\frac{dL}{dz}$

$$\frac{dL}{dz} = \frac{dL}{da} \times \frac{da}{dz}$$

$$= \frac{a-y}{a(1-a)} \times a(1-a)$$

$$= a - y$$

Vectorisation

Recall in logistic regression, we need to compute

$$z = w^T x + b \quad \text{where } w = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad x = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

$$w \in \mathbb{R}^{n_x} \quad x \in \mathbb{R}^{n_x}$$

non-vectorised:

$$z = 0$$

for i in range ($N - n$):

$$z += w[i] * x[i]$$

$$z += b$$

vectorised

$$z = \underbrace{\text{np. dot}(w, x)}_{w^T x} + b$$

In neural network, whenever possible, avoid explicit for-loop.

e.g.

$$u = A v$$

\nearrow \uparrow \searrow
vector matrix vector

by definition: $u_i = \sum_j A_{ij} v_j$

non-vectorisation:

$$u = \text{np.zeros}(m, 1)$$

for i ...

for j ...

$$u[i] += A[i][j] * v[j]$$

vectorisation:

$$u = \text{np.dot}(A, v)$$

e.g. apply exponential on every element of a matrix/vector

$$v = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} \implies u = \begin{bmatrix} e^{v_1} \\ \vdots \\ e^{v_n} \end{bmatrix}$$

non-vectorisation

$$u = \text{np.zeros}(n, 1)$$

for i in range(n):

$$u[i] = \text{math.exp}(v[i])$$

Vectorisation

$$u = \text{np. exp}(v)$$

Logistic Regression Derivatives (not really covered, see the one below)

$$J = 0, dw_1 = 0, dw_2 = 0, db = 0$$

for i = 1 to m: ↪ 1 for loop

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

$$a^{(i)} = \sigma(z^{(i)})$$

$$J += -[y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})]$$

$$dz^{(i)} = a^{(i)}(1 - a^{(i)})$$

$$\left. \begin{array}{l} dw_1 += x_1^{(i)} dz^{(i)} \\ dw_2 += x_2^{(i)} dz^{(i)} \end{array} \right\}$$

$$\left. \begin{array}{l} dw_1 += x_1^{(i)} dz^{(i)} \\ dw_2 += x_2^{(i)} dz^{(i)} \end{array} \right\}$$

$$db += dz^{(i)}$$

$$J = J/m, dw_1 = dw_1/m, dw_2 = dw_2/m, db = db/m$$

2 features
also need for loop
for large n

for j = 1 ... n_x
dw_j += ...

Different from previous slide

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

$$\begin{aligned} &= \frac{a - y}{a(1-a)} \cdot a(1-a) \\ &= a - y \end{aligned}$$

make dw a vector: dw = np. zeros((n_x, 1))

$$\hookrightarrow dw += x^{(i)} dz^{(i)}$$

$$\hookrightarrow dw /= m$$

Vectorising Logistic Regression

If we have m training examples, then to make prediction on the first example:
predict on second example: think ...

$$z^{(1)} = w^T x^{(1)} + b$$

$$a^{(1)} = \sigma(z^{(1)})$$

$$z^{(2)} = w^T x^{(2)} + b$$

$$a^{(2)} = \sigma(z^{(2)})$$

$$z^{(3)} = w^T x^{(3)} + b$$

$$a^{(3)} = \sigma(z^{(3)})$$

So this is the forward propagation for m training examples.

We first try to compute all the z_s and a_s

Vectorisation

training inputs: $X = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & | \end{bmatrix} \quad \left\{ n_x \right\}$ i.e. $\mathbb{R}^{n_x \times m}$

construct: $Z = [z^{(1)} \ z^{(2)} \ \dots \ z^{(m)}]$

$$= w^T X + [b \ b \ \dots \ b]$$

w^T is a row vector
 $[\quad]$
 n_x

$1 \times m$

$[w^T x^{(1)} + b \ w^T x^{(2)} + b \ \dots \ w^T x^{(m)} + b]$

To implement this: $Z = \text{np. dot}(w.T, X) + b$

b
just a number
 \hookrightarrow broadcast to decent

$$A = [a^{(1)} \ a^{(2)} \ \dots \ a^{(m)}] = \sigma(Z)$$

Vectorizing Logistic Regression's Gradient Computation

Perform gradient computation on m training examples at the same time.

for the first training example:

$$\frac{dZ^{(1)}}{dZ^{(1)}} = a^{(1)} - y^{(1)}$$

2nd training example: --

$$dZ^{(2)} = a^{(2)} - y^{(2)} \quad \dots$$

Define a new variable

$$dZ = [dZ^{(1)} \ dZ^{(2)} \ \dots \ dZ^{(m)}] \quad 1 \times m$$

Recall we have computed

$$A = [a^{(1)} \dots a^{(m)}], Y = [y^{(1)} \dots y^{(m)}]$$

$$\hookrightarrow dZ = A - Y = [a^{(1)} - y^{(1)} \dots a^{(m)} - y^{(m)}]$$

In the for loop, we still has to loop through our training examples for dW, db :

$$dW = 0$$

$$dW += X^{(i)} dZ^{(i)}$$

$$dW += X^{(m)} dZ^{(m)}$$

⋮

$$dW / m$$

$$db = 0$$

$$db += dZ^{(1)}$$

$$db += dZ^{(2)}$$

⋮

$$db / m$$

here dW is already

a vector so no need
loops through dW_1, dW_2, \dots
(for each feature)

$$db = \frac{1}{m} \sum_{i=1}^m dZ^{(i)} = \frac{1}{m} \text{np. sum}(dZ)$$

$$dW = \frac{1}{m} X dZ^T = \frac{1}{m} \left[\begin{array}{c|c|c|c} X^{(1)} & \cdots & X^{(m)} \\ \hline & & & dZ^{(1)} \\ & & & \vdots \\ & & & dZ^{(m)} \end{array} \right]$$

$$= \frac{1}{m} [X^{(1)} dZ^{(1)} + \dots + X^{(m)} dZ^{(m)}]$$

$n \times 1$

Previously, inefficient non-vectorized version

$$J = 0, dW_1 = 0, dW_2 = 0, db = 0$$

for $i = 1$ to m :

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

$$a^{(i)} = \sigma(z^{(i)})$$

$$J += -[y^{(i)} \log a^{(i)} + (1 - y^{(i)}) \log(1 - a^{(i)})]$$

$$dz^{(i)} = a^{(i)} - y^{(i)}$$

$$\begin{bmatrix} dW_1 += x_1^{(i)} dz^{(i)} \\ dW_2 += x_2^{(i)} dz^{(i)} \\ db += dz^{(i)} \end{bmatrix}$$

$$J = J/m, dW_1 = dW_1/m, dW_2 = dW_2/m, db = db/m$$

get rid of this
for loop

on all our training example

$$Z = w^T X + b$$

$$= \text{np. dot}(w.T, X) + b$$

$$A = \sigma(Z)$$

$$dZ = A - Y$$

$$dW = \frac{1}{m} X dZ^T$$

$$db = \frac{1}{m} \text{np. sum}(dZ)$$

$$\hookrightarrow w := w - \alpha dW \quad \text{Iteration!}$$

$$b := b - \alpha db$$

Broadcasting in Python

Calories from Carbs, Proteins, Fats in 100g of different foods:

| | Apples | Beef | Eggs | Potatoes | |
|---------|--------|-------|------|----------|-----|
| Carb | 56.0 | 0.0 | 4.4 | 68.0 | = A |
| Protein | 1.2 | 104.0 | 52.0 | 8.0 | |
| Fat | 1.8 | 135.0 | 99.0 | 0.9 | |

\downarrow
59 cal : carbs: $56/59 \approx 95\%$

Say we want to calculate % of calories from Carbs, protein and fats for each of the 4 foods.

i.e. Sum up each col to get total number of calories in each food.

then divide throughout the matrix

$$cal = A.sum(axis=0) \quad \begin{matrix} 0 \\ \downarrow \\ 1 \end{matrix} \quad \text{actually here it is redundant}$$

$$\text{percentage} = 100 \times \frac{A / (cal.reshape(1, 4))}{(3 \times 4) / (1 \times 4)}$$

c-7.

$$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + 100 \Rightarrow \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + \begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \end{bmatrix} = \begin{bmatrix} 101 \\ 102 \\ 103 \\ 104 \end{bmatrix}$$

c-8.

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}_{(m,n)} + \begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix}_{(1,n) \rightsquigarrow (m,n)} \stackrel{\text{top}}{\Rightarrow} \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix} = \begin{bmatrix} 101 & 202 & 303 \\ 104 & 205 & 306 \end{bmatrix}$$

c-9.

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}_{(m,n)} + \begin{bmatrix} 100 \\ 200 \end{bmatrix}_{(m,1) \rightsquigarrow (m,n)} \Rightarrow \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 100 & 100 \\ 200 & 200 & 200 \end{bmatrix} = \begin{bmatrix} 101 & 102 & 103 \\ 204 & 205 & 206 \end{bmatrix}$$

General Principle

If we have a (m, n) matrix:

$$\begin{matrix} (m, n) \\ + \\ - \\ \times \\ / \end{matrix}$$

$$\begin{matrix} (m, 1) \\ (1, n) \end{matrix} \rightsquigarrow (m, n)$$

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

then apply the operation element wise

Numpy Vectors

$a = np.random.randn(5)$

$\hookrightarrow a.shape = (5,)$ \Rightarrow rank 1 array

} Do not use

$a = np.random.randn(5, 1)$

$b = np.random.randn(1, 5)$

$\text{assert}(a.shape == (5, 1))$

} consisting and easy to understand

$\hookrightarrow a = a.reshape((5, 1))$ convert to proper vector/matrix

Explanation of logistic regression cost function

To recap:

$$\hat{y} = \sigma(w^T x + b)$$

$$\text{where } \sigma(z) = \frac{1}{1+e^{-z}}$$

Interpret $\hat{y} = p(y=1|x)$

Another way of saying: if $y=1$, then $p(y|x) = \hat{y}$
if $y=0$, then $p(y|x) = 1 - \hat{y}$

$$\hookrightarrow p(y|x) = \hat{y}^y (1-\hat{y})^{1-y}$$

$$\text{If } y=1, \quad p(y|x) = \hat{y}^1 (1-\hat{y})^0 = \hat{y}$$

$$\text{If } y=0, \quad p(y|x) = 1-\hat{y}$$

Log function is a strictly monotonically increasing function, so maximise $\log p(y|x)$ is similar to maximise $p(y|x)$:

$$\begin{aligned}\log p(y|x) &= \log \hat{y}^y (1-\hat{y})^{1-y} = y \log \hat{y} + (1-y) \log (1-\hat{y}) \\ &= -\mathcal{L}(\hat{y}, y)\end{aligned}$$

minimise \mathcal{L} \Rightarrow maximise $\log p(y|x)$

Cost on an example

The probability of all the labels in the training set:

$$p(\text{labels in training set}) = \prod_{i=1}^m p(y^{(i)}|x^{(i)})$$

\hookrightarrow same as maximizing

$$\begin{aligned}\log p(\text{labels in training set}) &= \log \prod_{i=1}^m p(y^{(i)}|x^{(i)}) \\ &= \sum_{i=1}^m \log p(y^{(i)}|x^{(i)}) \\ &\quad - \mathcal{L}(\hat{y}^{(i)}, y^{(i)})\end{aligned}$$

Principle of maximum likelihood estimation

$$= - \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

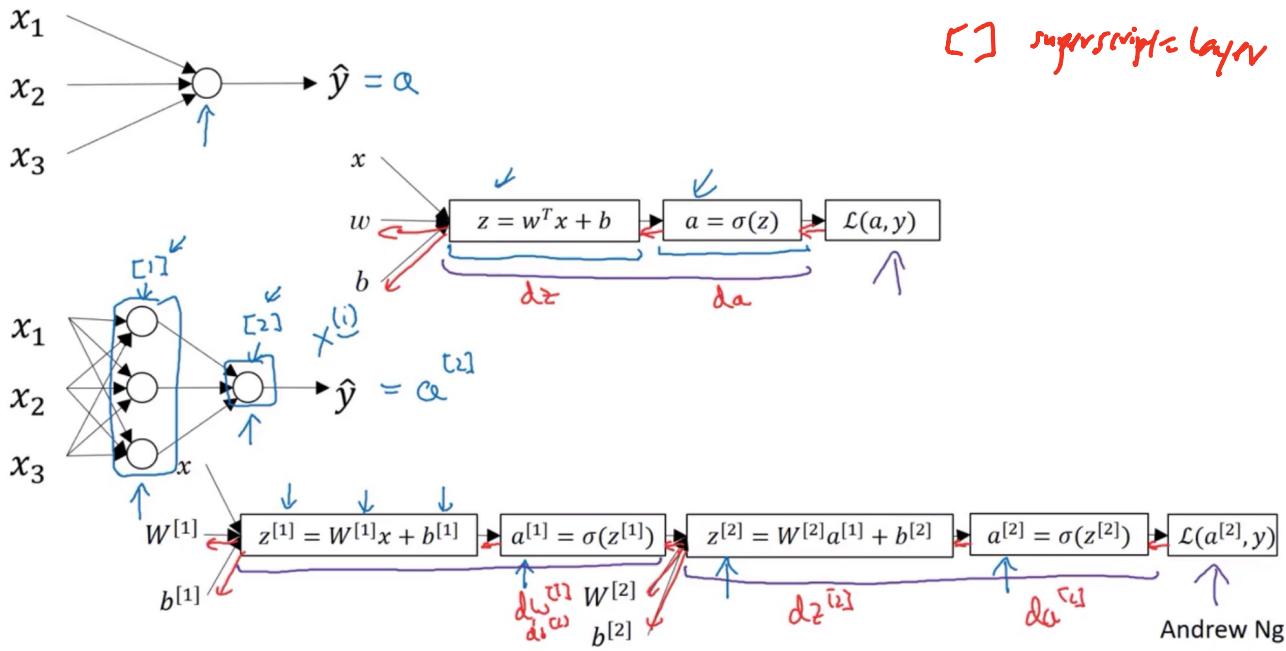
carry out maximum likelihood
on logistic regression model

$\log p(y|x)$

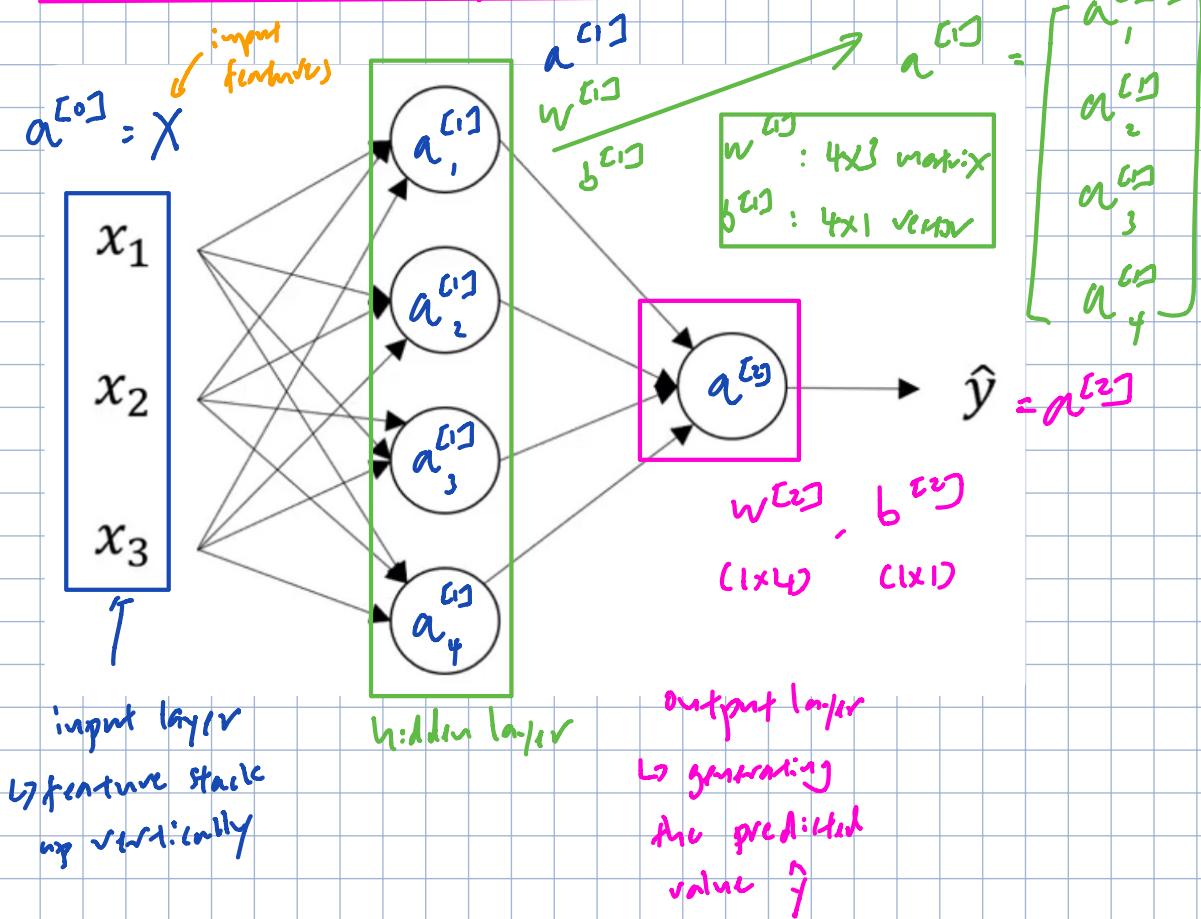
$$\hookrightarrow \text{cost: } J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

Week 3 (Shallow Neural Network)

Neural Network Overview



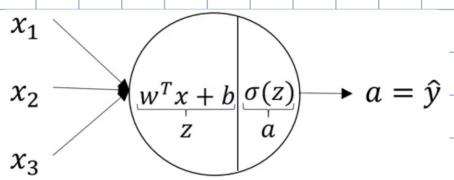
Neural Network Representation



- The term "hidden layer" refers to the fact that in the training set, the true values for these nodes are not observed.
- The above is called a 2-layer NN, because we usually don't count input layer. Input layer is sometimes called layer 0.
- The hidden layer and output layer will have parameters associated with them.

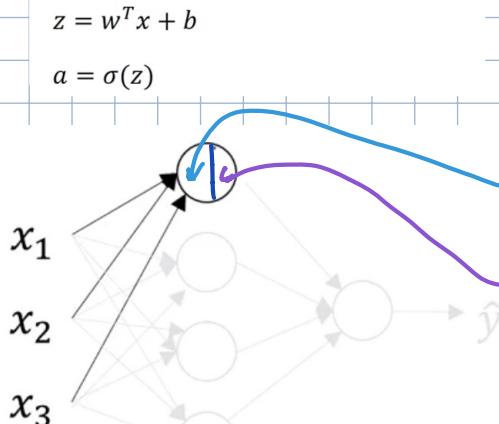
↳ w, b

Computing a NN's output



↙ logistic regression

In NN, we repeat this calculation many times.



↗ Neural network

↖ Let's focus on the first node in the hidden layer

$$z_1^{[1]} = w_1^{[1]T} x + b_1^{[1]}$$

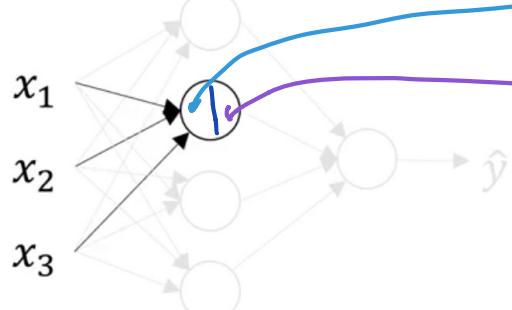
$$a_1^{[1]} = \sigma(z_1^{[1]})$$

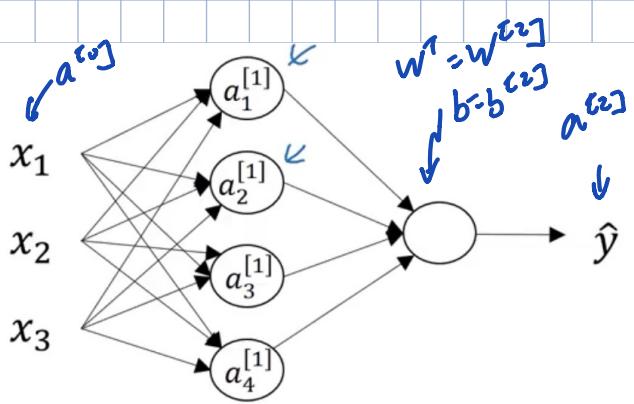
notation: $a_i^{[L]} \leftarrow$ layer
 $a_i^{[1]} \leftarrow$ node in that layer

Second node in layer 1

$$z_2^{[1]} = w_2^{[1]T} x + b_2^{[1]}$$

$$a_2^{[1]} = \sigma(z_2^{[1]})$$





$$z_1^{[1]} = w_1^{[1]T} x + b_1^{[1]}, \quad a_1^{[1]} = \sigma(z_1^{[1]})$$

$$z_2^{[1]} = w_2^{[1]T} x + b_2^{[1]}, \quad a_2^{[1]} = \sigma(z_2^{[1]})$$

$$z_3^{[1]} = w_3^{[1]T} x + b_3^{[1]}, \quad a_3^{[1]} = \sigma(z_3^{[1]})$$

$$z_4^{[1]} = w_4^{[1]T} x + b_4^{[1]}, \quad a_4^{[1]} = \sigma(z_4^{[1]})$$

Compute $z^{[1]}$ in vectorized manner:

$$\left[\begin{array}{c} w_1^{[1]T} \\ w_2^{[1]T} \\ w_3^{[1]T} \\ w_4^{[1]T} \end{array} \right] \left[\begin{array}{c} x_1 \\ x_2 \\ x_3 \\ x_4 \end{array} \right] + \left[\begin{array}{c} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{array} \right] = \left[\begin{array}{c} w_1^{[1]T} x + b_1^{[1]} \\ w_2^{[1]T} x + b_2^{[1]} \\ w_3^{[1]T} x + b_3^{[1]} \\ w_4^{[1]T} x + b_4^{[1]} \end{array} \right] = \left[\begin{array}{c} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{array} \right] = z^{[1]}$$

$w^{[1]T}$

$$a^{[1]} = \left[\begin{array}{c} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]} \end{array} \right] = \sigma(z^{[1]})$$

Given input x:

$$z^{[1]} = W^{[1]}x + b^{[1]}$$

(4,3) (4,3) (3,1) (4,1)

$$a^{[1]} = \sigma(z^{[1]})$$

(4,1) (4,1)

$$\hookrightarrow z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

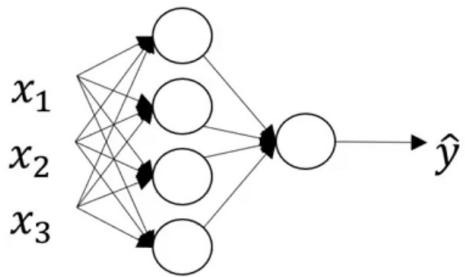
(1,4) (1,4) (4,1) (1,1)

$$\hookrightarrow a^{[2]} = \sigma(z^{[2]})$$

(1,1) (1,1)

This is for $\frac{1}{N}$ training example.

Vectorising across m training examples



$$z^{[1]} = W^{[1]}x + b^{[1]}$$

$$a^{[1]} = \sigma(z^{[1]})$$

$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

$$a^{[2]} = \sigma(z^{[2]})$$

$x \rightarrow a^{[2]} = \hat{y}$ using σ for single training sample.

If we have m training samples, we need to have

$$x^{(1)} \rightarrow a^{(1)(1)} = \hat{y}^{(1)}$$

$$x^{(2)} \rightarrow a^{(2)(2)} = \hat{y}^{(2)}$$

:

$$x^{(m)} \rightarrow a^{(2)(m)} = \hat{y}^{(m)}$$

notation: $a^{(2)(i)}$
 ↑
 training example i
 ↓
 layer 2

unvectorised version:

for $i = 1$ to m :

$$z^{(1)(i)} = W^{[1]}x^{(i)} + b^{[1]}$$

$$a^{(1)(i)} = \sigma(z^{(1)(i)})$$

$$z^{(2)(i)} = W^{[2]}a^{(1)(i)} + b^{[2]}$$

$$a^{(2)(i)} = \sigma(z^{(2)(i)})$$

for $i = 1$ to m :

$$z^{[1](i)} = W^{[1]}x^{(i)} + b^{[1]}$$

$$a^{[1](i)} = \sigma(z^{[1](i)})$$

$$z^{[2](i)} = W^{[2]}a^{[1](i)} + b^{[2]}$$

$$a^{[2](i)} = \sigma(z^{[2](i)})$$

recall we define the matrix X :

$$X = \begin{bmatrix} & | & | & | \\ & x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix}_{(n_r, m)}$$

To vectorise this for loop for m training samples:

$$z^{[1]} = W^{[1]}X + b^{[1]}$$

$$A^{[1]} = \sigma(z^{[1]})$$

$$z^{[2]} = W^{[2]}A^{[1]} + b^{[2]}$$

$$A^{[2]} = \sigma(z^{[2]})$$

$$z^{[1]} = \begin{bmatrix} | & | & | \\ z^{1} & z^{[1](2)} & \dots & z^{[1](m)} \end{bmatrix}$$

$$A^{[1]} = \begin{bmatrix} | & | & | \\ a^{1} & a^{[1](2)} & \dots & a^{[1](m)} \end{bmatrix}_{\text{nodes} \times \text{training sample}}$$

Explanation for vectorised implementation

Let's go through part of the propagation calculation for the few examples.

For the first train sample

$$z^{1} = w^{[1]} x^{(1)} + b^{[1]}$$

For second train sample:

$$z^{[1](2)} = w^{[1]} x^{(2)} + b^{[1]}$$

For the third training sample

$$z^{[1](3)} = w^{[1]} x^{(3)} + b^{[1]}$$

for simplification purpose

$$w^{[1]} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

same matrix

$$w^{[1]} x^{(1)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

column vector

$$w^{[1]} x^{(2)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

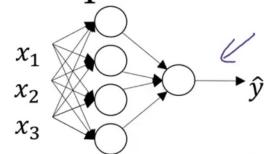
$$w^{[1]} x^{(3)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

Now consider training set with 3 samples:

$$\begin{aligned} w^{[1]} \begin{bmatrix} x^{(1)} & x^{(2)} & x^{(3)} \end{bmatrix} &= \begin{bmatrix} \cdot & \cdot & \cdot \end{bmatrix} \\ &= \begin{bmatrix} z^{1} & z^{[1](2)} & z^{[1](3)} \end{bmatrix} = z^{[1]} \end{aligned}$$

If we stack up x in cols, then z will also stack up, for my training examples.

Recap of vectorizing across multiple examples



$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix}$$

↑ stackup

$$A^{[1]} = \begin{bmatrix} a^{1} & a^{[1](2)} & \dots & a^{[1](m)} \end{bmatrix}$$

for $i = 1$ to m

$$z^{[1](i)} = W^{[1]} x^{(i)} + b^{[1]}$$

$$a^{[1](i)} = \sigma(z^{[1](i)})$$

$$z^{[2](i)} = W^{[2]} a^{[1](i)} + b^{[2]}$$

$$a^{[2](i)} = \sigma(z^{[2](i)})$$

non-vectorized

$$Z^{[1]} = W^{[1]} X + b^{[1]}$$

$$A^{[1]} = \sigma(Z^{[1]})$$

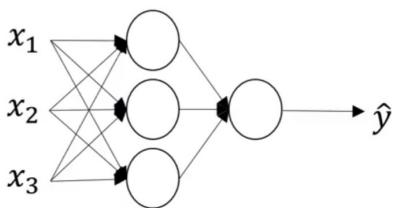
$$Z^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$$

$$A^{[2]} = \sigma(Z^{[2]})$$

vectorize w/ids
m samples

vectorized
w/o for loop

Activation Functions



Given x :

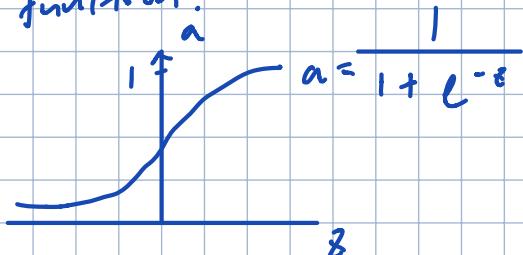
$$z^{[1]} = W^{[1]}x + b^{[1]}$$

$$a^{[1]} = \sigma(z^{[1]})$$

$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

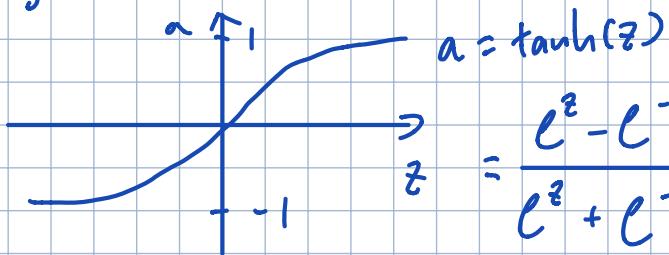
$$a^{[2]} = \sigma(z^{[2]})$$

In forward propagation, there are steps we used sigmoid function as the activation function:



more generally, we can have a function $g(z^{[1]})$, which is a non-linear function.

e.g. \tanh



$$z = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

(Shifted sigmoid)

For hidden layer, tanh is almost always better than the sigmoid function, because the values between $+1$ and -1 , the mean of the activations are closer to having a 0 mean. When we pre-process data, we might center the data.

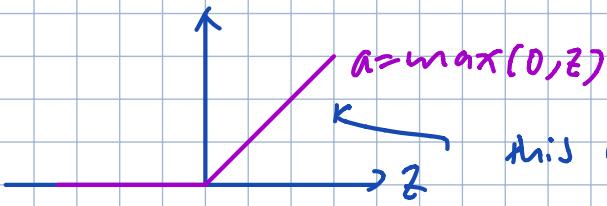
↳ for output layer, if y is either 0 or 1, then sigmoid make sense as $0 \leq \hat{y} \leq 1$.

✗ activation function can differ from layer to layer: $g^{[l]}(z^{[l]})$

If z is either very large or very small, then the gradient of the function is small.

layer 1

ReLU:



this slope makes ReLU learning faster than tanh

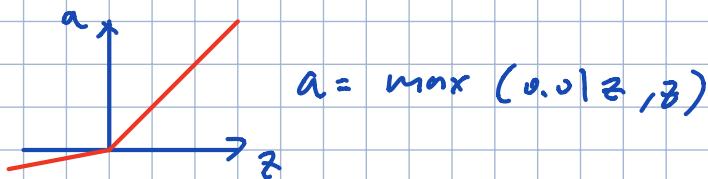
Some guidelines:

- * If output is 0/1, binary classification, sigmoid function for output layer
- * All other units, use ReLU

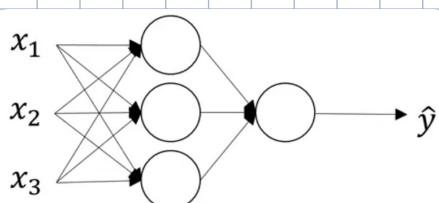
Disadvantage of ReLU:

→ Gradient = 0 if z is negative

↳ leaky ReLU:



why activation function?



In forward propagation, we use activation function.

If we use linear activation function:

$$g(z) = z \quad (\text{practically do nothing})$$

Given x :

$$z^{[1]} = W^{[1]}x + b^{[1]}$$

$$a^{[1]} = g^{[1]}(z^{[1]})$$

$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

$$a^{[2]} = g^{[2]}(z^{[2]})$$

If we do this, then this model just computing \hat{y} as a linear function of input feature.

$$\text{i.e. } a^{[1]} = z^{[1]} = W^{[1]}x + b^{[1]}$$

$$a^{[2]} = z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

$$= W^{[2]}(W^{[1]}x + b^{[1]}) + b^{[2]}$$

$$= (w^{[2]} w^{[1]}) X + (w^{[2]} b^{[1]} + b^{[2]})$$

$$= w' X + b'$$

i.e. neural network just outputting a linear function of the input.

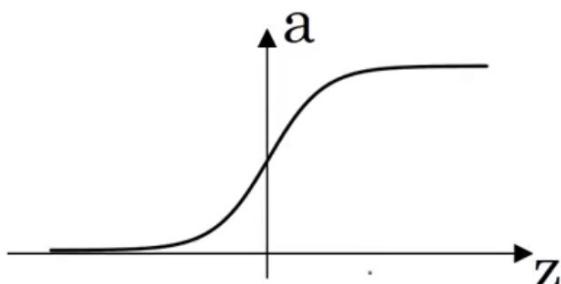
\hookrightarrow no matter how many layers \Rightarrow output is just linear function of the input.

For Regression problem, the output layer may be can use linear activation function, to output a real value.

Derivatives of activation functions

when implementing back propagation, need slope of the activation function.

Sigmoid function



$$g(z) = \frac{1}{1 + e^{-z}}$$

$$\text{slope of } g(z) \text{ at } z : \frac{d}{dz} g(z) = \frac{1}{1 + e^{-z}} \left(1 - \frac{1}{1 + e^{-z}} \right)$$

$$= g(z) (1 - g(z))$$

Sanity check: if z is large e.g. $z = 10$, $g(z) \approx 1$

$$\hookrightarrow \frac{d}{dz} g(z) \approx 1 (1 - 1) \approx 0$$

if z is small, $z = -10$: $g(z) \approx 0$

$$\hookrightarrow \frac{d}{dz} g(z) \approx 0 (1-0) \approx 0$$

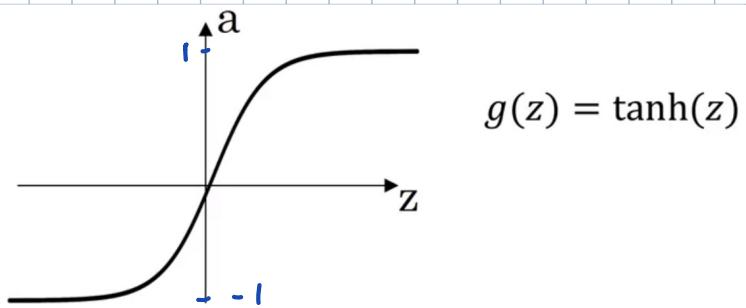
if $z=0$, $g(z) = \frac{1}{2}$ $\Rightarrow \frac{d}{dz} g(z) = \frac{1}{2}(1-\frac{1}{2}) = \frac{1}{4}$

Notation: $g'(z) = \frac{d}{dz} g(z)$

In neural networks : $a = g(z) = \frac{1}{1+e^{-z}}$

$$\hookrightarrow \frac{da}{dz} = a(1-a) = g'(z)$$

Tanh function



$$g(z) = \tanh(z) \\ = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$\frac{d}{dz} g(z) = 1 - (\tanh(z))^2 = g'(z)$$

sanity check: $z=10$ $\tanh(10) \approx 1$

$$\hookrightarrow g'(10) \approx 0$$

$z=-10$ $\tanh(-10) \approx -1$

$$\hookrightarrow g'(-10) \approx 0$$

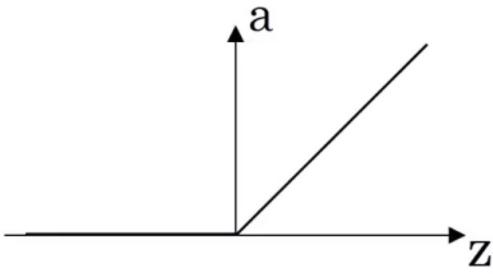
$z=0$ $\tanh(0) = 0$

$$\hookrightarrow g'(0) = 1$$

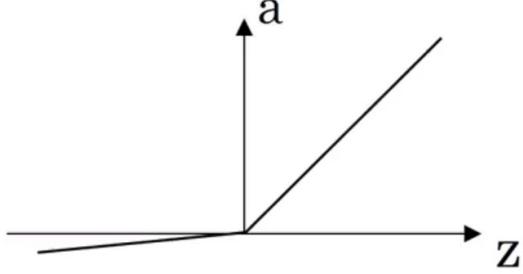
If $a = g(z)$

$$g'(z) = 1 - a^2$$

ReLU and Leaky ReLU



ReLU



Leaky ReLU

$$\text{ReLU: } g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z > 0 \\ \text{undefined if } z=0 \end{cases} \quad (\text{technically})$$

$$g(z) = \max(0.01z, z)$$

$$g'(z) = \begin{cases} 0.01 & \text{if } z < 0 \\ 1 & \text{if } z > 0 \\ \text{undefined if } z=0 \end{cases}$$

Gradient Descent for Neural Networks

A neural network with a single hidden layer, will have parameters

$$w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$$

Recall that we have :

$$n_x = n^{[0]}, n^{[1]}, n^{[2]} = 1$$

input features hidden unit output unit

then $w^{[1]}$ is $(n^{[0]}, n^{[1]})$

$b^{[1]}$ is $(n^{[1]}, 1)$

$w^{[2]}$ is $(n^{[1]}, n^{[2]})$

$b^{[2]}$ is $(n^{[2]}, 1)$

Cost function: (assume binary classification)

$$J(w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}) = \frac{1}{m} \sum_{i=1}^m \underbrace{L(\hat{y}, y)}_{a^{(2)}}$$

To train these parameters, need to perform gradient descent:

* need to initialize parameters randomly (rather than to all 0s)

Repeat [

compute predict ($\hat{y}^{(i)}$, for $i = 1 \dots m$)

$$dw^{(1)} = \frac{\partial J}{\partial w^{(1)}}, db^{(1)} = \frac{\partial J}{\partial b^{(1)}}, dw^{(2)}, db^{(2)}$$

$$w^{(1)} := w^{(1)} - \alpha dw^{(1)}$$

$$b^{(1)} := b^{(1)} - \alpha db^{(1)}$$

$$w^{(2)} := w^{(2)} - \alpha dw^{(2)}$$

$$b^{(2)} := b^{(2)} - \alpha db^{(2)} \quad \}$$

\downarrow true value
 $Y = [y^{(1)} \ y^{(2)} \dots \ y^{(m)}]$

To get all the derivatives:

forward propagation:

$$z^{(1)} = w^{(1)} X + b^{(1)}$$

$$A^{(1)} = \sigma(z^{(1)})$$

$$z^{(2)} = w^{(2)} A^{(1)} + b^{(2)}$$

$$A^{(2)} = j^{(2)}(z^{(2)}) = \sigma(z^{(2)})$$

for binary classification

$$(n^{(1)}, m)$$

Backpropagation:

$$dz^{(2)} = A^{(2)} - Y$$

$$dw^{(2)} = \frac{1}{m} dz^{(2)} A^{(1)T}$$

$$db^{(2)} = \frac{1}{m} \text{np.sum}(dz^{(2)}, \text{axis}=1, \text{keepdims=True})$$

[so far similar to logistic regression

$$dz^{(1)} = w^{(2)T} dz^{(2)} \times j^{(1)'}(z^{(1)})$$

$$(n^{(1)}, m)$$

$$dw^{(1)} = \frac{1}{m} dz^{(1)} X^T$$

$$db^{(1)} = \frac{1}{m} \text{np.sum}(dz^{(1)}, \text{axis}=1, \text{keepdims=True})$$

\downarrow assume signal

\downarrow horizontally

keep dims = True)

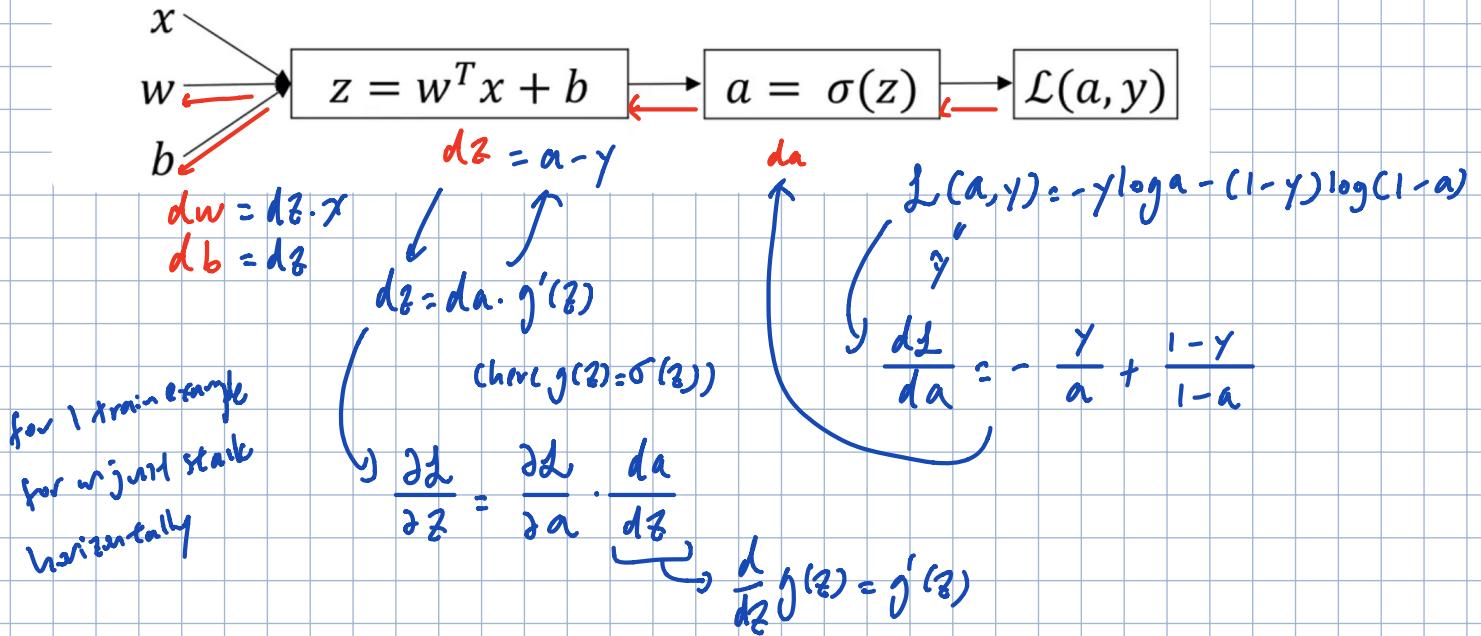
here this is 1x1 number

element-wise derivative of what the activation function on layer 1

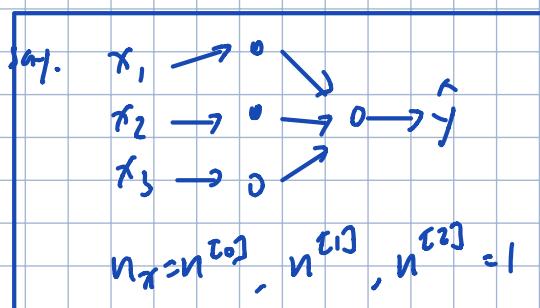
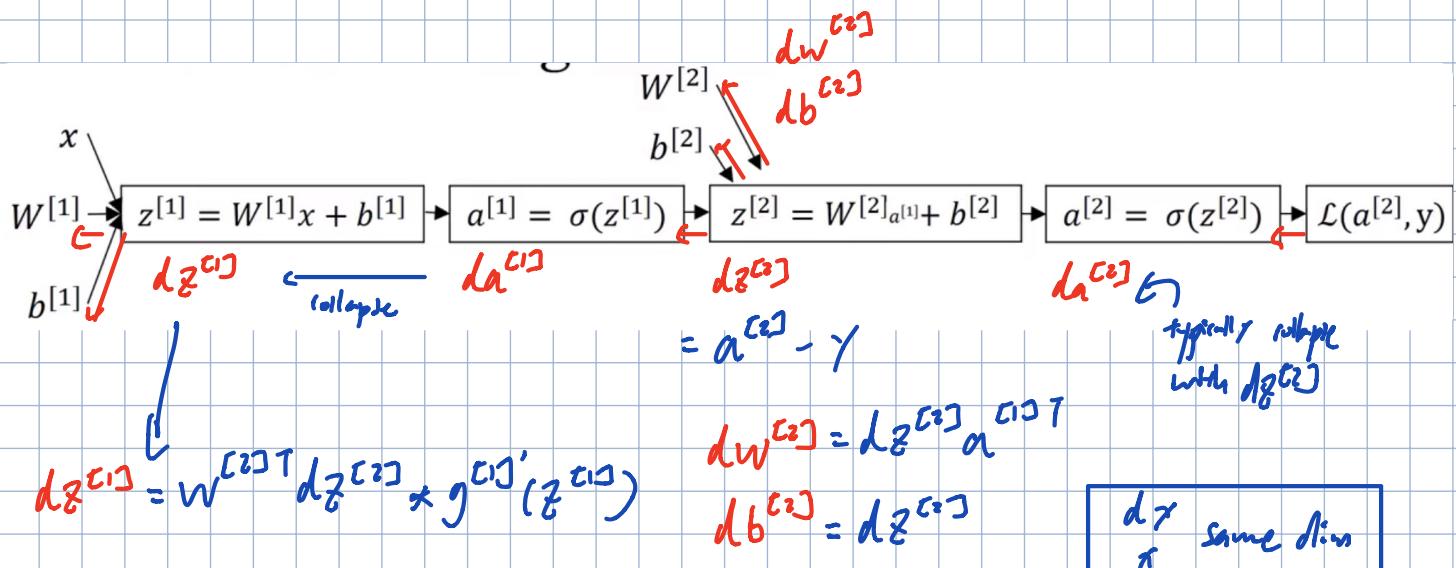
activation function on layer 1

Back propagation

Logistic regression



Neural Network Gradient



for 1 training example

$$\begin{aligned}
 & \text{w}^{[2]} : (n^{[2]}, n^{[1]}) \\
 & z^{[2]} : (n^{[2]}, 1) \\
 & z^{[1]}, dz^{[2]} : (n^{[1]}, 1) \\
 & dz^{[1]} = W^{[2]T} dz^{[2]} \cdot \sigma'(z^{[1]}) \\
 & (n^{[0]}, 1) \quad (n^{[1]}, n^{[2]}) \quad (n^{[2]}, 1)
 \end{aligned}$$

$$dW^{[1]} = dz^{[1]} \cdot X^T$$

$$db^{[1]} = dz^{[1]}$$

Summary of gradient descent

$$dz^{[2]} = a^{[2]} - y \quad | \text{ training}$$

$$dW^{[2]} = dz^{[2]} a^{[1]T}$$

$$db^{[2]} = dz^{[2]}$$

$$dz^{[1]} = W^{[2]T} dz^{[2]} * g^{[1]'}(z^{[1]})$$

$(n^{[1]}, 1)$

$$dW^{[1]} = dz^{[1]} x^T$$

$$db^{[1]} = dz^{[1]}$$

$$dZ^{[2]} = A^{[2]} - Y$$

$$dW^{[2]} = \frac{1}{m} dZ^{[2]} A^{[1]T}$$

$$db^{[2]} = \frac{1}{m} np.sum(dZ^{[2]}, axis=1, keepdims=True)$$

$$dZ^{[1]} = \underbrace{W^{[2]T} dZ^{[2]} * g^{[1]'}(z^{[1]})}_{\substack{\text{elementwise product} \\ (n^{[1]}, m) \quad (n^{[1]}, m) \quad (n^{[1]}, m)}}$$

$$dW^{[1]} = \frac{1}{m} dZ^{[1]} X^T$$

$$db^{[1]} = \frac{1}{m} np.sum(dZ^{[1]}, axis=1, keepdims=True)$$

in training, stack vertically

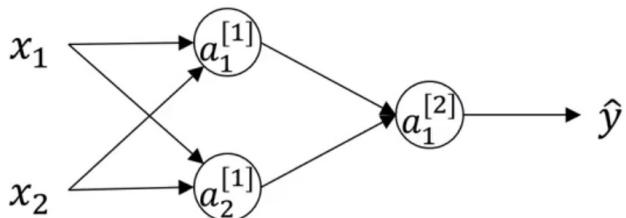
$$\mathcal{J}(\dots) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}_i, y_i)$$

\approx

Random Initialisation

when training NN, it is important to initialize the weights randomly.

What happens if you initialize weights to zero?



$$n^{[1]} = 2 \quad n^{[2]} = 2$$

update same for both nodes.

$$W^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad \begin{matrix} \uparrow \\ \text{problem} \end{matrix}$$

$$b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \begin{matrix} \uparrow \\ \text{To K} \end{matrix}$$

$$\Rightarrow a_1^{[1]} = a_2^{[1]} \quad \Rightarrow dz_1^{[1]} = dz_2^{[1]}$$

$$W^{[1]} := W^{[1]} - \alpha dW$$

1)

$$dW = \begin{bmatrix} u & v \\ u & v \end{bmatrix}$$

each row is same

Need to initialize the weight randomly

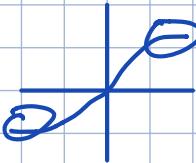
initialize to small random values

$$W^{[1]} = \text{np. random. random}((2, 2) \times 0.01)$$

$$b^{[1]} = \text{np. zeros}(2, 1)$$

$$W^{[2]} = \text{np. random. random}(1, 2) \times 0.01$$

$$b^{[2]} = 0$$



$$z^{[1]} = W^{[1]} x + b^{[1]}$$

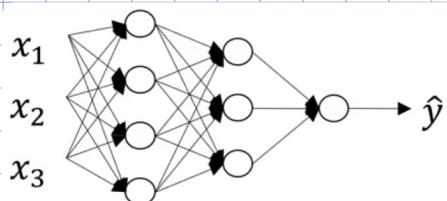
$$a^{[1]} = g^{[1]}(z^{[1]})$$

If not, gradient descent will be slow

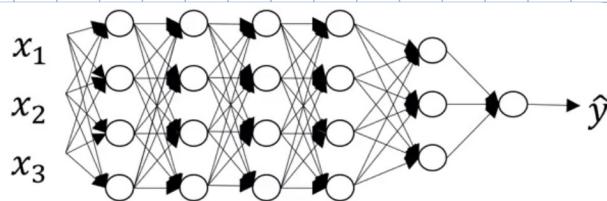
For deep network, we want choose different parameter 0.01.

Week 4 (Deep Neural Networks)

Deep L-layer neural network



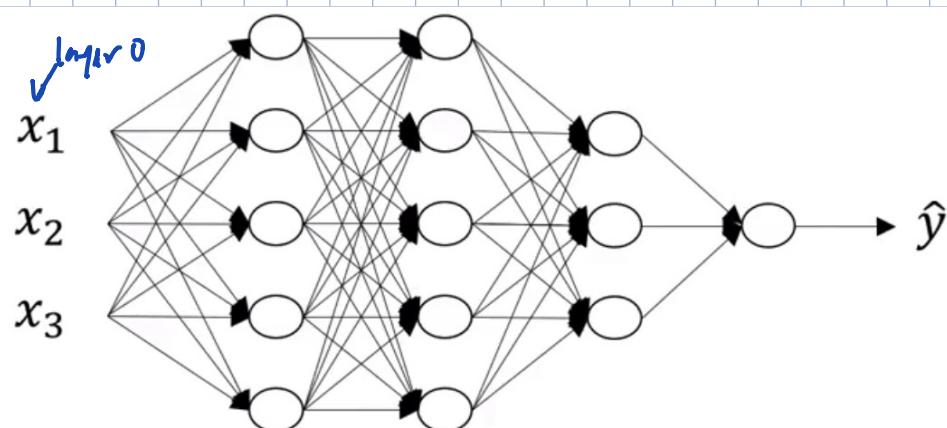
2 hidden layers



5 hidden layers

Deep neural network notation

4 layer NN (3 hidden)



$$L = \text{no. of layers} = 4$$

$a^{[l]}$ = activation in layer l

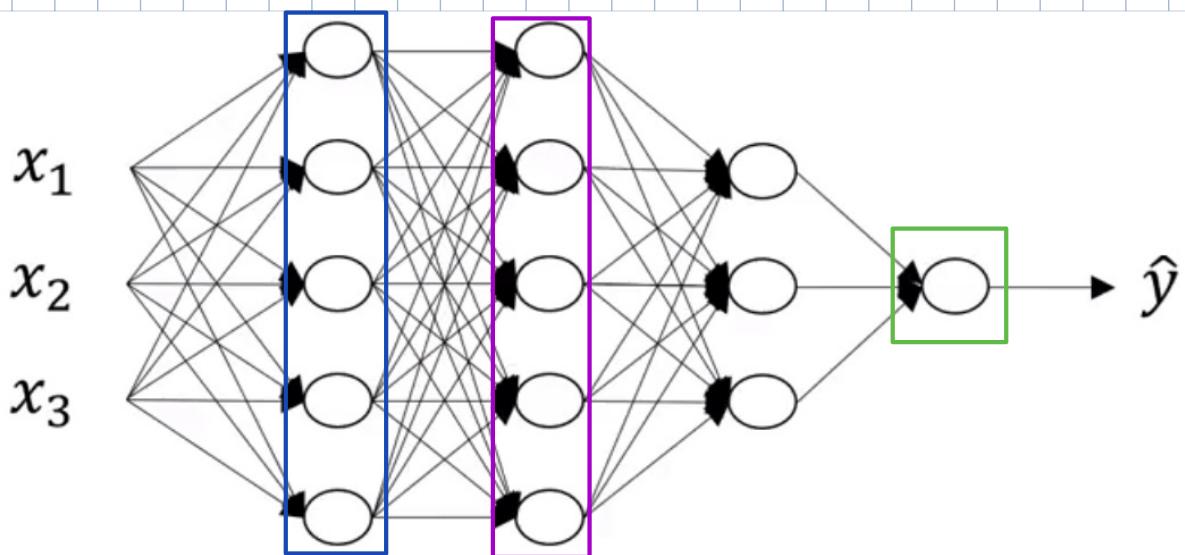
$$n^{[l]} = \text{no. of nodes in layer } l$$

$$\Rightarrow a^{[l]} = g^{[l]}(z^{[l]})$$

$$\hookrightarrow n^{[1]} = 5, \quad n^{[2]} = n_x = 3 \\ n^{[3]} = 3, \quad n^{[4]} = n^{[L]} = 1$$

$$w^{[l]} = \text{weights for } z^{[l]}$$

Forward Propagation in a Deep Network



for a single training example x
for the first layer

$$x: z^{[1]} = W^{[1]} x + b^{[1]}$$

$$a^{[1]} = g^{[1]}(z^{[1]})$$

$$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$$

$$a^{[2]} = g^{[2]}(z^{[2]})$$

:

$$z^{[4]} = W^{[4]} a^{[3]} + b^{[4]}$$

$$a^{[4]} = g^{[4]}(z^{[4]}) = \hat{y}$$

in general:

$$z^{[l]} = W^{[l]} A^{[l-1]} + b^{[l]}$$

$$A^{[l]} = g^{[l]}(z^{[l]})$$

vectorized (in training example):

$$z^{[l]} = W^{[l]} X + b^{[l]} \quad (X = 1^{100})$$

$$A^{[l]} = g^{[l]}(z^{[l]})$$

$$z^{[l]} = W^{[l]} A^{[l-1]} + b^{[l]}$$

$$A^{[l]} = g^{[l]}(z^{[l]})$$

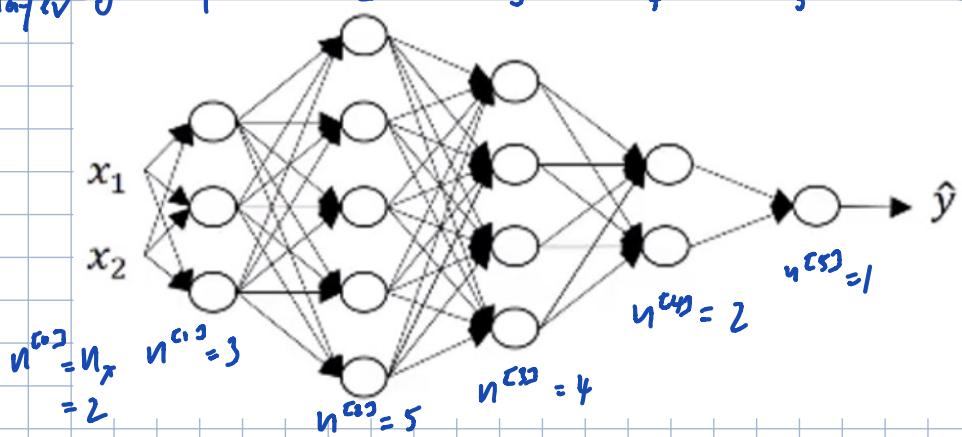
for
loop

$$\hat{Y} = g(z^{[4]}) = A^{[4]}$$

$$\text{e.g. } z^{[2]} = \begin{bmatrix} 1 \\ z^{[2]_{1,1}} & \dots & z^{[2]_{1,100}} \\ 1 & \dots & 1 \end{bmatrix}$$

Getting the matrix dimension right

layer 0 1 2 3 4 5



$$L = 5$$

\hookrightarrow 4 hidden layers + 1 output layer

independent of m

Forward propagation!

$$z^{[1]} = W^{[1]} \cdot X + b^{[1]}$$

$$\begin{matrix} (3,1) & (3,2) & (2,1) & (3,1) \\ \downarrow & \downarrow & \downarrow & \downarrow \\ (n^{[0,1]}, 1) & (n^{[0,2]}, n^{[0,0]}) & (n^{[0,0]}, 1) & (n^{[0,1]}, 1) \end{matrix}$$

$$W^{[1]} : (n^{[0,0]}, n^{[0,1]})$$

$$W^{[2]} : (5, 3) \rightarrow (n^{[1,0]}, n^{[1,1]})$$

$$z^{[2]} = W^{[2]} \cdot a^{[1]} + b^{[2]}$$

$$\begin{matrix} (5,1) & (5,3) & (3,1) & (5,1) \rightarrow (n^{[1,0]}, 1) \\ \downarrow & \downarrow & \downarrow & \downarrow \\ (5,1) & (5,3) & (3,1) & (5,1) \rightarrow (n^{[1,0]}, 1) \end{matrix}$$

$$\hookrightarrow W^{[3]} : (4, 5)$$

$$\hookrightarrow W^{[4]} : (2, 4)$$

$$\hookrightarrow W^{[5]} : (1, 2)$$

$$b^{[l]} : (n^{[l,0]}, 1)$$



$$d_b^{[l]} : (n^{[l,0]}, 1)$$

$$W^{[6]} : (n^{[0,0]}, n^{[0,1]})$$



$$d_W^{[l]} : (n^{[l,0]}, n^{[l,1]})$$

Note that $a^{[l]} = g^{[l]}(z^{[l]}) \Rightarrow z^{[l]}$ and $a^{[l]}$ same dimension
in this type of network.

$$a^{[l]} : (n^{[l,0]}, 1)$$

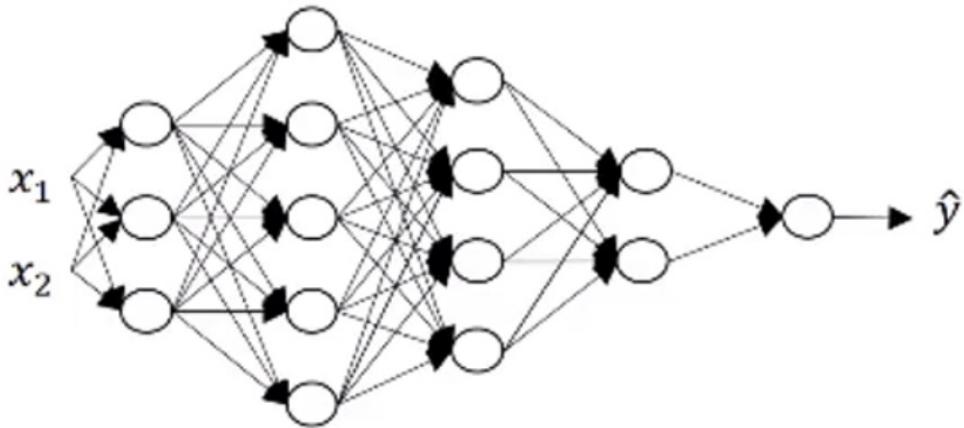
$$z^{[l]} : (n^{[l,0]}, 1)$$

1 training example!

vectorised version (in training example)

dimension of w , dw and db remains the same for vectorised version,

but dimension of z , a , x will change



single example:

$$z^{(l)} = W^{(l)} \cdot x + b^{(l)}$$

$$(n^{(l)}, 1) \quad (n^{(l)}, n^{(l)}) \quad (n^{(l)}, 1) \quad (n^{(l)}, 1)$$

m examples:

$$z^{(l)} = W^{(l)} \cdot X + b^{(l)}$$

Diagram illustrating the dimensions of $z^{(l)}$ for m examples:

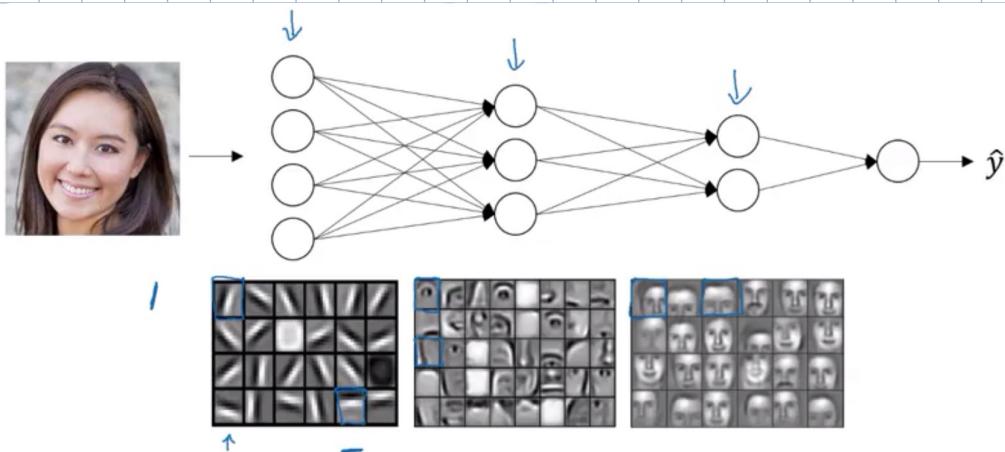
- $z^{(l)}$ is represented as a column vector: $[z^{(l)(1)} \dots z^{(l)(m)}]$ with dimensions $(n^{(l)}, m)$.
- X is represented as a matrix: $(n^{(l)}, n^{(l)})$.
- $b^{(l)}$ is represented as a column vector: $(n^{(l)}, 1)$.
- A note indicates: "Python broadcast" under the $b^{(l)}$ term.
- $W^{(l)}$ is represented as a matrix: $(n^{(l)}, m)$.

$$z^{(l)}, A^{(l)} : (n^{(l)}, m) \quad \text{when } l=0, A^{(0)} \cdot X = (n^{(l)}, m)$$

$$dz^{(l)}, dA^{(l)} : (n^{(l)}, m)$$

why deep representation?

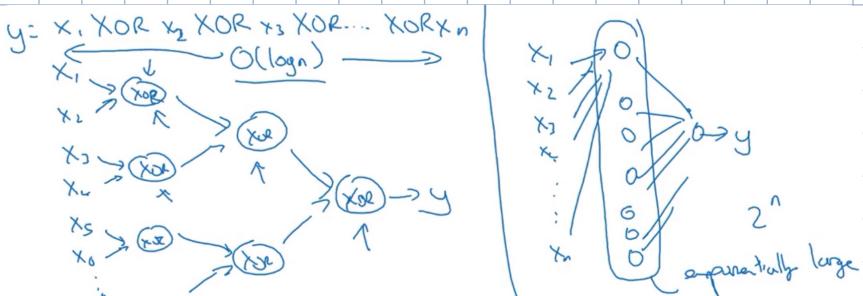
useful networks have many hidden layers, why?



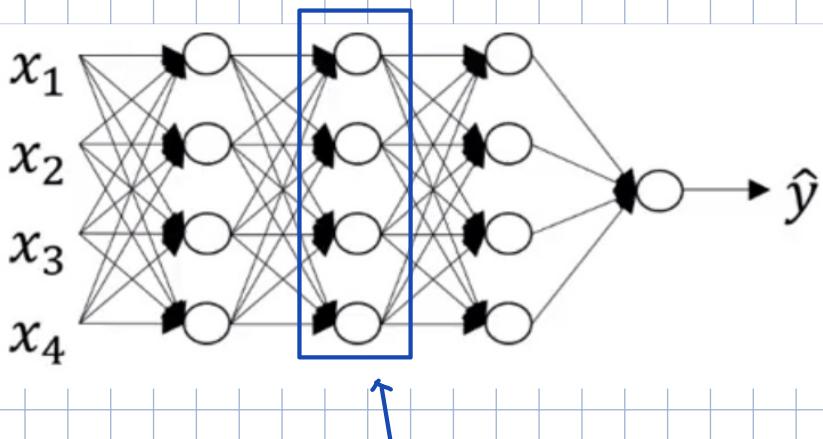
Circuit theory and deep learning

Informally: There are functions you can compute with a "small" L-layer deep neural network that shallower networks require exponentially more hidden units to compute.

less hidden units



Building blocks of deep neural networks



Say we focus on this particular layer

→ for layer ℓ : we have some parameters $W^{[\ell]}, b^{[\ell]}$

Forward: Input $a^{[\ell-1]}$ output $a^{[\ell]}$

$$\text{Using: } z^{[\ell]} = W^{[\ell]} a^{[\ell-1]} + b^{[\ell]} \quad \text{cache } z^{[\ell]} \text{ for later use}$$

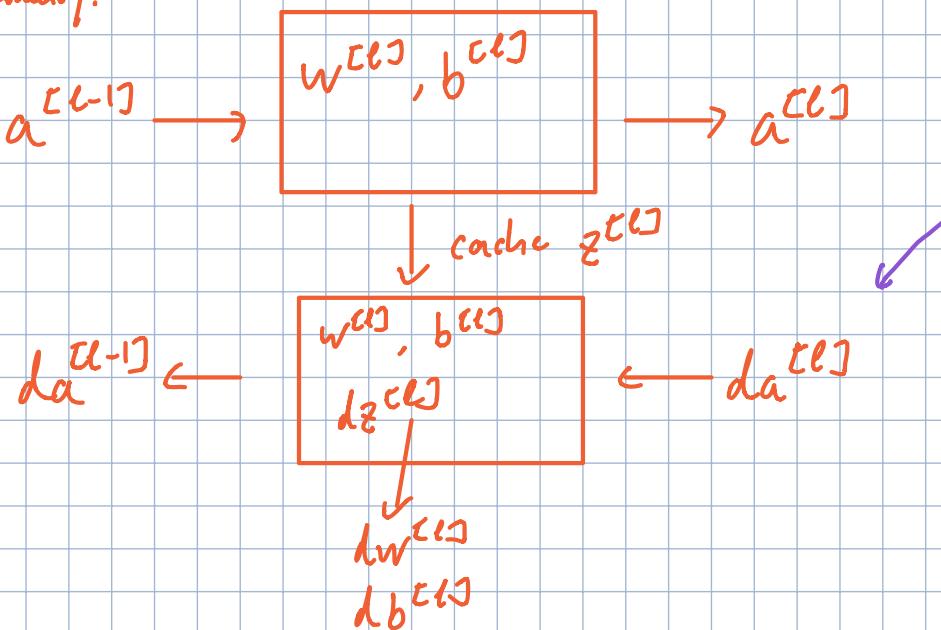
$$a^{[\ell]} = g^{[\ell]}(z^{[\ell]})$$

Backward: Input: $da^{[\ell]}$ output $da^{[\ell-1]}$

$$\text{cache}(z^{[\ell]}) \quad \frac{dw^{[\ell]}}{da^{[\ell]}} \quad \frac{db^{[\ell]}}{da^{[\ell]}}$$

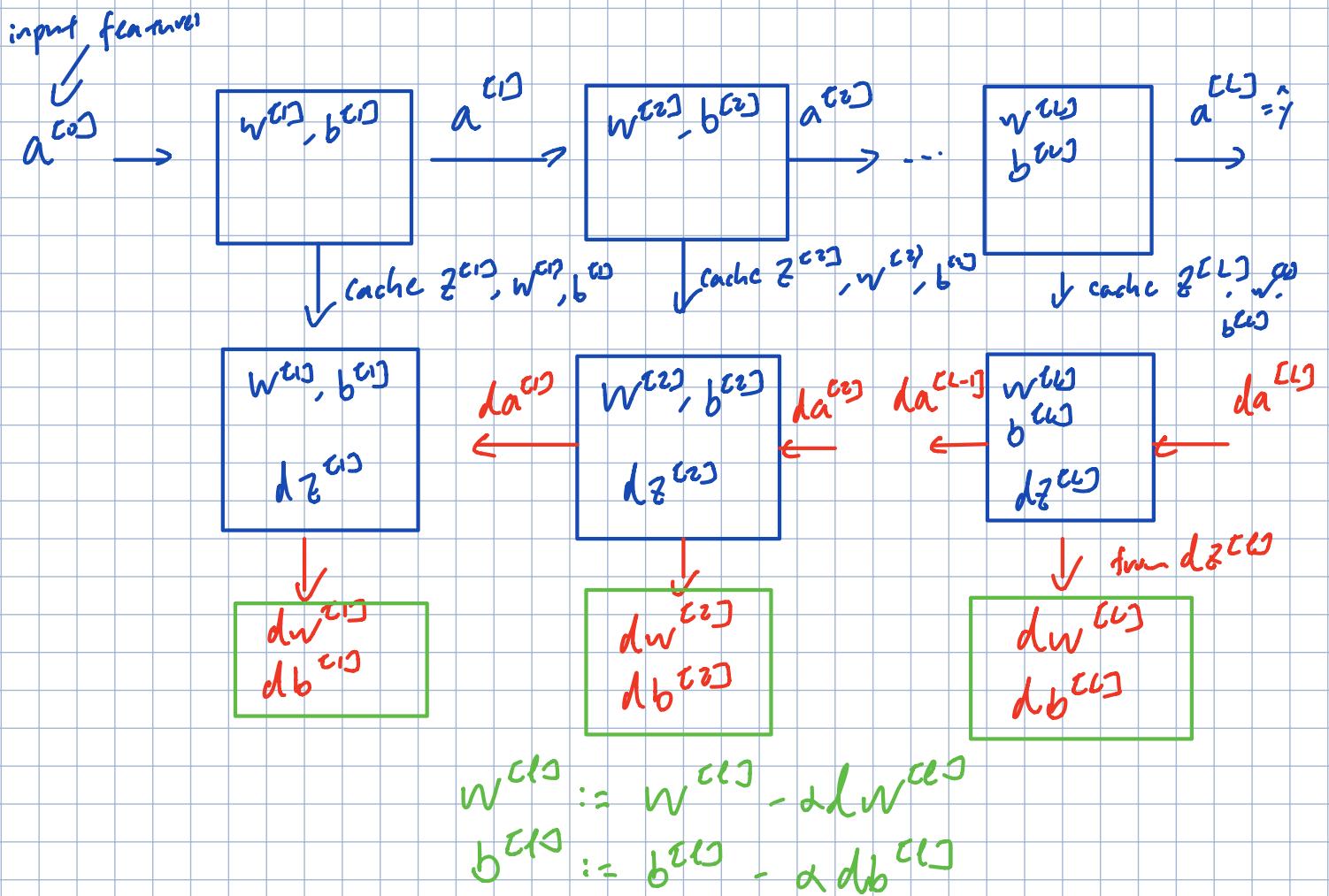
layer ℓ

Summary:



gives the derivatives w.r.t.
those activations (i.e. da)

Basic computation of NN



1 forward, 1 backward, update $w, b \leftarrow$ 1 iteration of gradient descent

Forward and Backward propagation

Forward for layer l :

input $a^{(l-1)}$

output $a^{(l)}, \text{cache}(z^{(l)}, w^{(l)}, b^{(l)})$

$$\begin{aligned} z^{(l)} &= w^{(l)} \cdot a^{(l-1)} + b^{(l)} \\ a^{(l)} &= g^{(l)}(z^{(l)}) \end{aligned}$$

implementation

with:

vectorize:

$$\begin{aligned} z^{(l)} &= w^{(l)} \cdot A^{(l-1)} + b^{(l)} \\ A^{(l)} &= g^{(l)}(z^{(l)}) \end{aligned}$$

backward for layer l :

Input $da^{[l]}$

output $da^{[l-1]}, dW^{[l]}, db^{[l]}$

$$dZ^{[l]} = da^{[l]} \cdot g^{[l]'}(Z^{[l]})$$

dimmt mit

$$dW^{[l]} = dZ^{[l]} \cdot a^{[l-1]T}$$

$$db^{[l]} = dZ^{[l]}$$

$$da^{[l-1]} = W^{[l]T} \cdot dZ^{[l]}$$

vectorized:

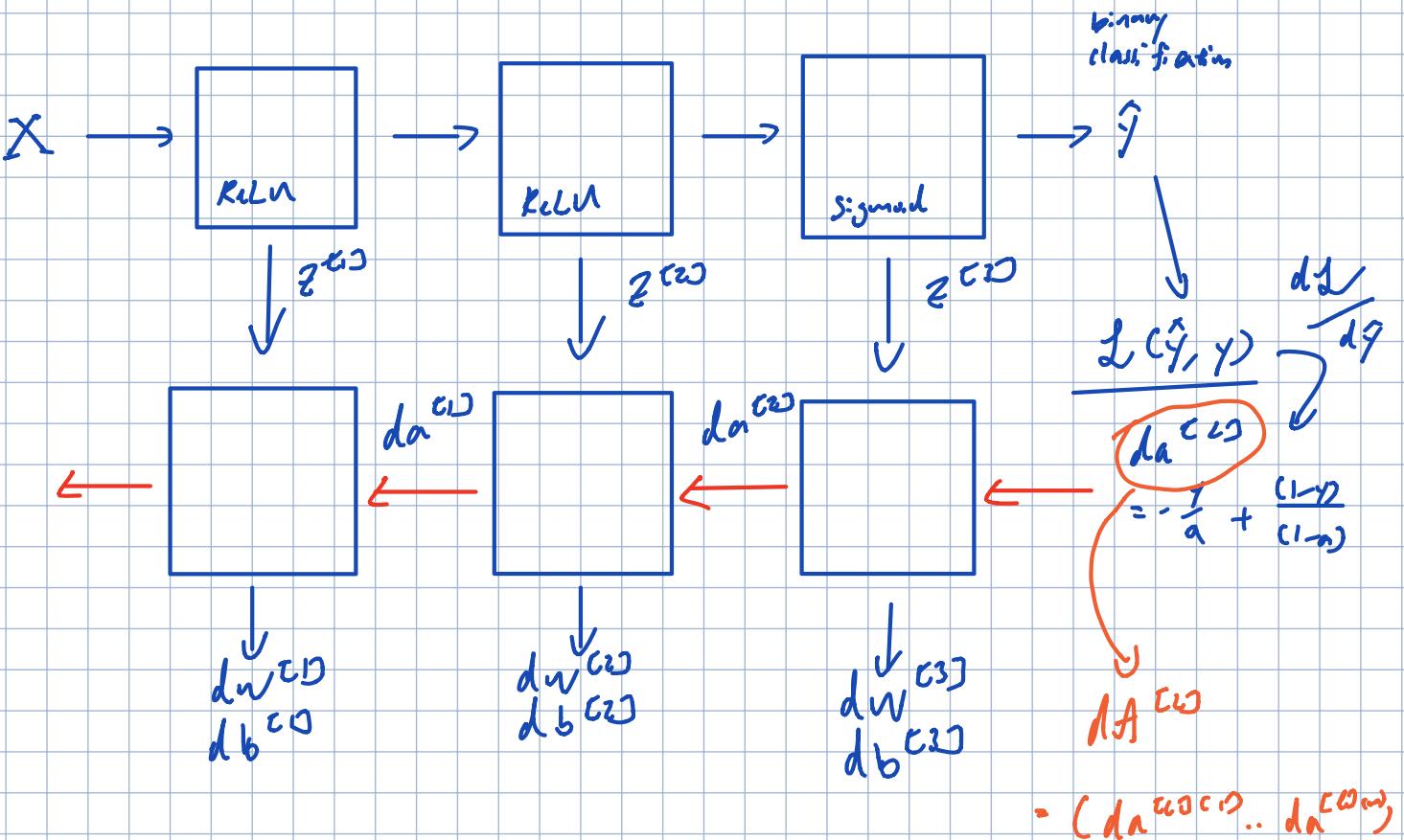
$$dZ^{[l]} = dA^{[l]} \cdot g^{[l]'}(Z^{[l]})$$

$$dW^{[l]} = \frac{1}{m} dZ^{[l]} \cdot A^{[l-1]T}$$

$$db^{[l]} = \frac{1}{m} \text{up. sum}(dZ^{[l]}, \text{axis}=1, \text{keep dim}) = \text{Tr}(dZ^{[l]})$$

$$dA^{[l-1]} = W^{[l]T} \cdot dZ^{[l]}$$

Summary



Parameters vs Hyperparameters

Parameters: $W^{(1)}, b^{(1)}, W^{(2)}, b^{(2)}, \dots$

Hyperparameters: learning rate α

iterations for gradient descent

hidden layer L

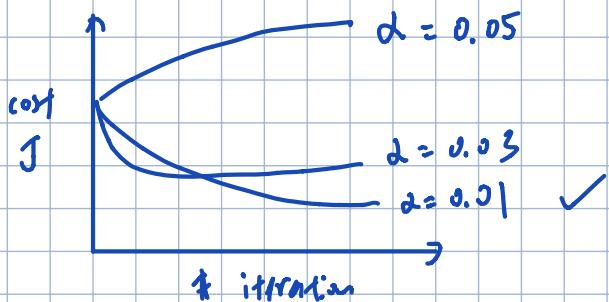
hidden units: $n^{(1)}, n^{(2)}, \dots$

choice of activation function

These hyperparameters can ultimately decide parameters, and they cannot be learned explicitly.

More advanced hyperparameters:

momentum, minibatch size, regularisation parameters...



Summary

Forward and backward

$$\begin{aligned} Z^{[1]} &= W^{[1]}X + b^{[1]} \\ A^{[1]} &= g^{[1]}(Z^{[1]}) \\ Z^{[2]} &= W^{[2]}A^{[1]} + b^{[2]} \\ A^{[2]} &= g^{[2]}(Z^{[2]}) \\ &\vdots \\ A^{[L]} &= g^{[L]}(Z^{[L]}) = \hat{Y} \end{aligned}$$

$$dZ^{[L]} = A^{[L]} - Y$$

$$dW^{[L]} = \frac{1}{m} dZ^{[L]} A^{[L-1]^T}$$

$$db^{[L]} = \frac{1}{m} np.sum(dZ^{[L]}, axis = 1, keepdims = True)$$

$$dZ^{[L-1]} = W^{[L]^T} dZ^{[L]} * g'^{[L-1]}(Z^{[L-1]})$$

Note that * denotes element-wise multiplication

⋮

$$dZ^{[1]} = W^{[2]} dZ^{[2]} * g'^{[1]}(Z^{[1]})$$

$$dW^{[1]} = \frac{1}{m} dZ^{[1]} A^{[0]^T}$$

Note that $A^{[0]^T}$ is another way to denote the input features, which is also written as X^T

$$db^{[1]} = \frac{1}{m} np.sum(dZ^{[1]}, axis = 1, keepdims = True)$$