# Deep Neural Network:

→ We've seen frwd, bckwd propogation for a single layer, log reg & vectorization. We saw why its important to initiate vectors randomly. We'll use these together for a deeper NN.

## Deep NN:
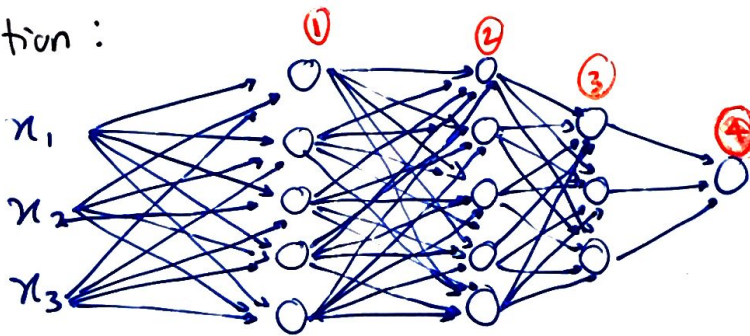
log reg: kinda shallow

1 hidden layer: better but basic

Deeper NNs? functions that can be learned here, that shallower Networks are unable to recognize.

## New Problems to think?

→ depth needed for a NN to work optimally : New hyper parameter

## Notation:



→ 4 layed NN    $L=4$
→ 3 hidden layers
→ 5, 5, 3

$n^{[l]}$ = #units in layer $l$

$n^{[1]} = 5$       $n^{[3]} = 3$
$n^{[2]} = 5$       $n^{[4]} = 1$
$n^{[0]} = n_x = 3$

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

$a^{[l]} = g^{[l]}(z^{[l]})$ , $W^{[l]}$, $b^{[l]}$ = weights for $z^{[l]}$

## ✳ Forward propagation:

for layer 1,

$x$:  $z^{[1]} = W^{[1]} x + b^{[1]}$
$a^{[1]} = g^{[1]}(z^{[1]})$  $\longrightarrow$

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

$z^{[4]} = W^{[4]} a^{[3]} + b^{[4]}$
$a^{[4]} = g^{[4]}(z^{[4]}) = \hat{y}$

## General form?

Going by the pattern $l: 1 \to 4$, we see,

$$z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}$$
$$a^{[l]} = g^{[l]}(z^{[l]})$$

# Vectorized form:

→ We'll basically need to stack together the computed $z, a$ values.

$$Z = \begin{bmatrix} | & | & & | \\ z^{[1](1)} & z_2^{[2](2)} & \cdots & z^{[n](n)} \\ | & | & & | \end{bmatrix} \qquad \hat{y} = g(z^{[4]}) = A^{[4]}$$

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

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

→ need for loop to compute activations for layer $1, 2, 3 \ldots n$

for $l = 1 \ldots 4$:

$$\begin{matrix} z, \\ a \end{matrix}$$  ↑↓ Calculation needed.

# Note on dimentions:

→ helpful to write the dims on paper before implementation.

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

$(3,1) \quad \therefore (3,2) \quad (2,1) \quad \mathbf{(3,1)}$

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

$[n^{[1]}, 1] \qquad [n^{[1]}, n^{[0]}] \left( n^{[0]}, 1 \right)$

② $W^{[1]} : (n^{[1]}, n^{[0]})$

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

③ $Z^{[2]} = W^{[2]} \cdot a^{[1]} + b^{[2]}$

$(5,1) \quad \sim (5,3) \quad (3,1) \quad \mathbf{(5,1)}$

④ Similarly,

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

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

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

# Generic form:

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

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

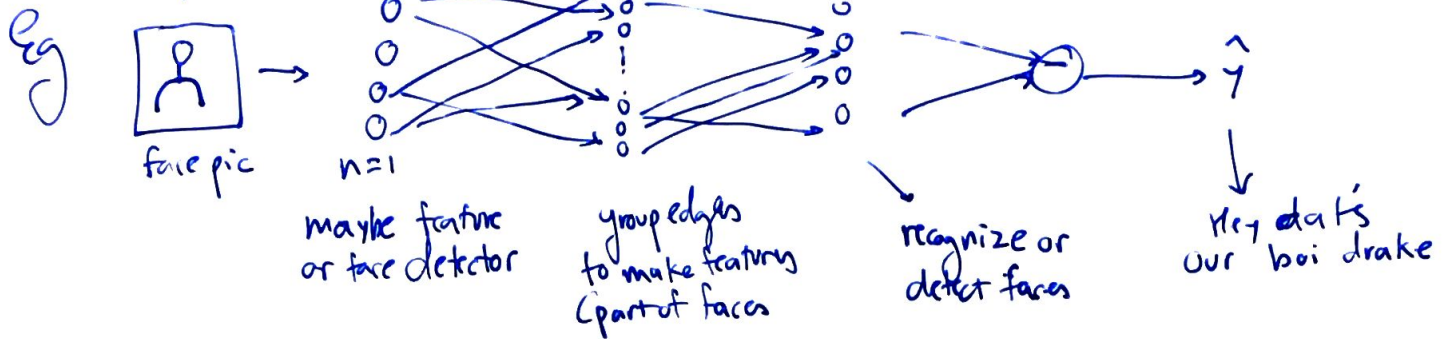$$z^{[l]}, a^{[l]} : (n^{[l]}, 1)$$

$$Z^{[l]}, A^{[l]} : (n^{[l]}, m)$$

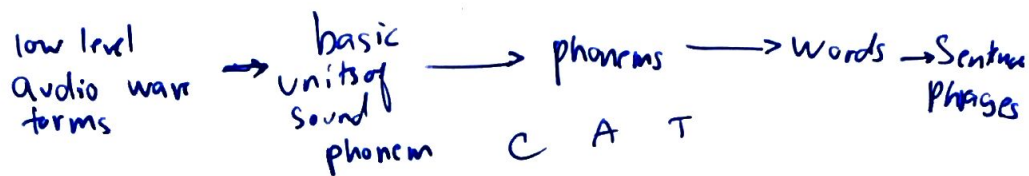— dw has same dims

db has same dims

$dZ^{[l]}, dA^{[l]}$ have same dims

# Why deep networks work:

deep network is computing

Eg  [face pic]  →  n=1 ⚬⚬⚬⚬ → ⚬⚬⚬⚬⚬ → ⚬⚬⚬⚬ → ◯ → $\hat{y}$

maybe feature
or face detector

group edges
to make features
(part of faces)

recognize or
detect faces

Hey dats
our boi drake

→ audio recognition

low level
audio wave
forms

→ basic
units of
sound
phonem

→ phonems
C  A  T

→ words → Sentence
Phrases

## Circuit thory:

There are fns you can compute with a 'small' L-layer network
that shallower networks require exponentially more hidden units to compute

$$x_1 \; XOR \; x_2 \; XOR \; x_3 \; XOR \; x_4 \; \dots \dots \; x_n$$

deep N

$x_1$ ⟩ XOR
$x_2$ ⟩

$x_3$ ⟩ XOR
$x_4$

⋮ XOR
⋮
⋮ XOR → $y$

$x_n$

XOR

shallow N

$x_1$
$x_2$
$x_3$
⋮
⋮
$x_n$

[hidden layer column] → $y$

$O(2^{...})$

$\approx 2^{(n-1)}$

hidden layer
~ Expo. large

→ basically, to compute these XORs, a shallower network's hidden layer would
to compute expo. larger computations.

# Building blocks of a deep NN

→ visualizing ideas behind fwd, bckwd propagation

for layer $l$, $\quad W^{[l]}, b^{[l]}$

forward prg$^n$ → $Z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}$, $\quad a^{[l]} = g^{[l]}(Z^{[l]})$
& cache $Z^{[l]}$ for bckwd p.

Backward propagation: Input $da^{[l]}$, $\left[ \text{output } dg^{[l-1]} \quad dw^{[l]} \quad db^{[l]} \right]$
Cache $Z^{[l]}$,



layer $l$

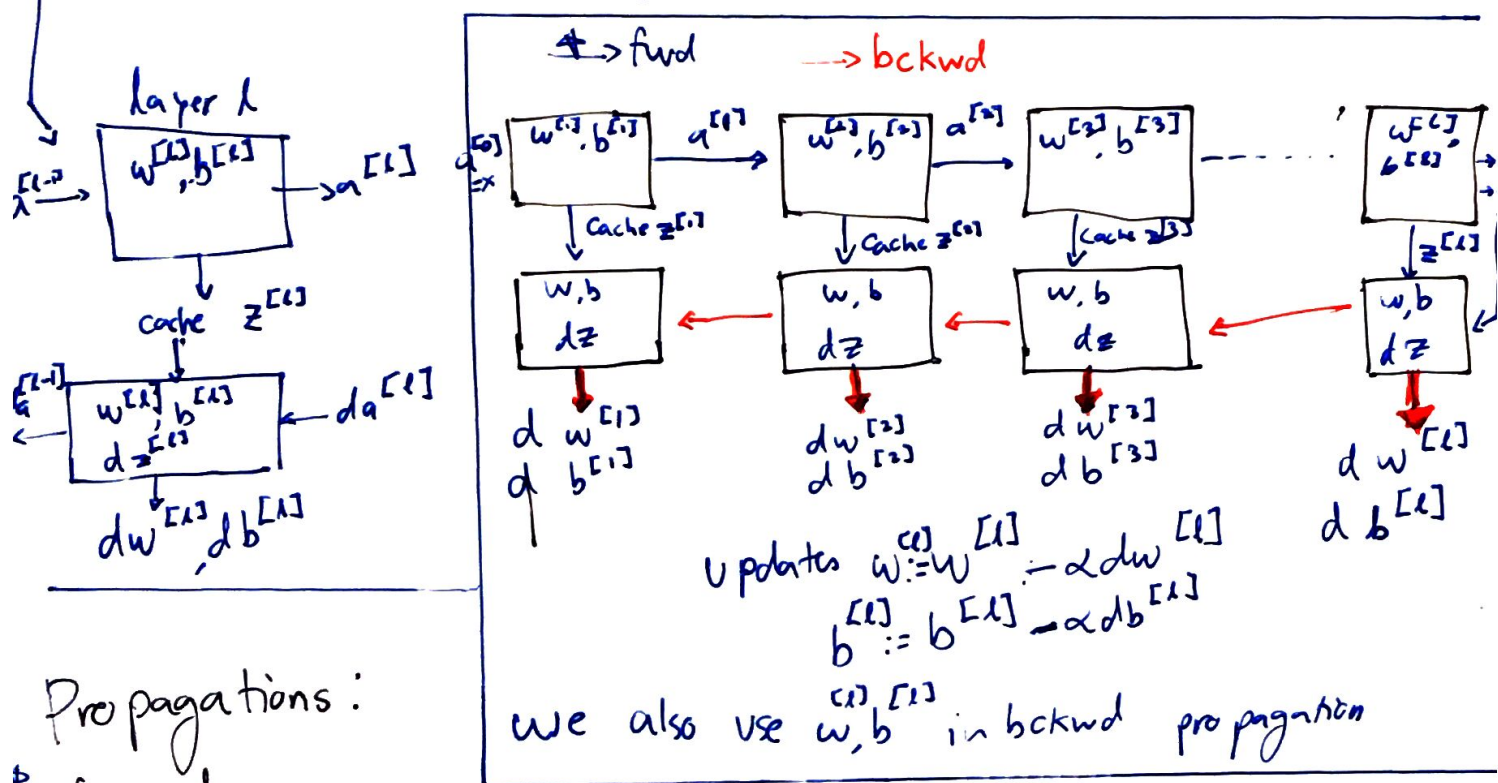$a^{[l-2]} \rightarrow \boxed{W^{[l]}, b^{[l]}} \rightarrow a^{[l]}$

cache $Z^{[l]}$

$a^{[l-1]} \leftarrow \boxed{\begin{array}{c} W^{[l]}, b^{[l]} \\ dz^{[l]} \end{array}} \leftarrow da^{[l]}$

$dw^{[l]}, db^{[l]}$

↑→ fwd    →bckwd

$a^{[0]}=x \rightarrow \boxed{W^{[1]}, b^{[1]}} \xrightarrow{a^{[1]}} \boxed{W^{[2]}, b^{[2]}} \xrightarrow{a^{[2]}} \boxed{W^{[3]}, b^{[3]}} - - - \boxed{\begin{array}{c} W^{[l]} \\ b^{[l]} \end{array}}$

Cache $Z^{[1]}$ ↓  Cache $Z^{[2]}$ ↓  cache $Z^{[3]}$ ↓   $Z^{[l]}$ ↓

$\boxed{\begin{array}{c} w,b \\ dz \end{array}} \leftarrow \boxed{\begin{array}{c} w,b \\ dz \end{array}} \leftarrow \boxed{\begin{array}{c} w,b \\ dz \end{array}} \leftarrow \boxed{\begin{array}{c} w,b \\ dz \end{array}}$

$d w^{[1]}$       $dw^{[2]}$       $dw^{[3]}$        $dw^{[l]}$
$d b^{[1]}$       $db^{[2]}$       $db^{[3]}$        $db^{[l]}$

updates $w^{[l]} := w^{[l]} - \alpha \, dw^{[l]}$
$\qquad b^{[l]} := b^{[l]} - \alpha \, db^{[l]}$

we also use $w^{[l]}, b^{[l]}$ in bckwd propagation

## Propagations:

$\rho$ – forward
inputs $a^{[l-1]}$, outputs $a^{[l]}$ & caches $Z^{[l]}$

$\left. \begin{array}{l} Z^{[l]} = W^{[l]} A^{[l-1]} + b^{[l]} \\ A^{[l]} = g^{[l]}(Z^{[l]}) \end{array} \right\}$ vectorized version

– backward
inputs $da^{[l]}$, outputs $da^{[l-1]}, dW^{[l]}, db^{[l]}$    $* = $ pairwise mult.

$dZ^{[l]} = da^{[l]} * g^{[l]}(Z^{[l]})$
$dW^{[l]} = dz^{[l]} g * a^{[l-1]T}$       $dz^{[l]} = w^{[l+1]T} dz^{[l+1]} * g^{[l]}(Z^{[l]})$
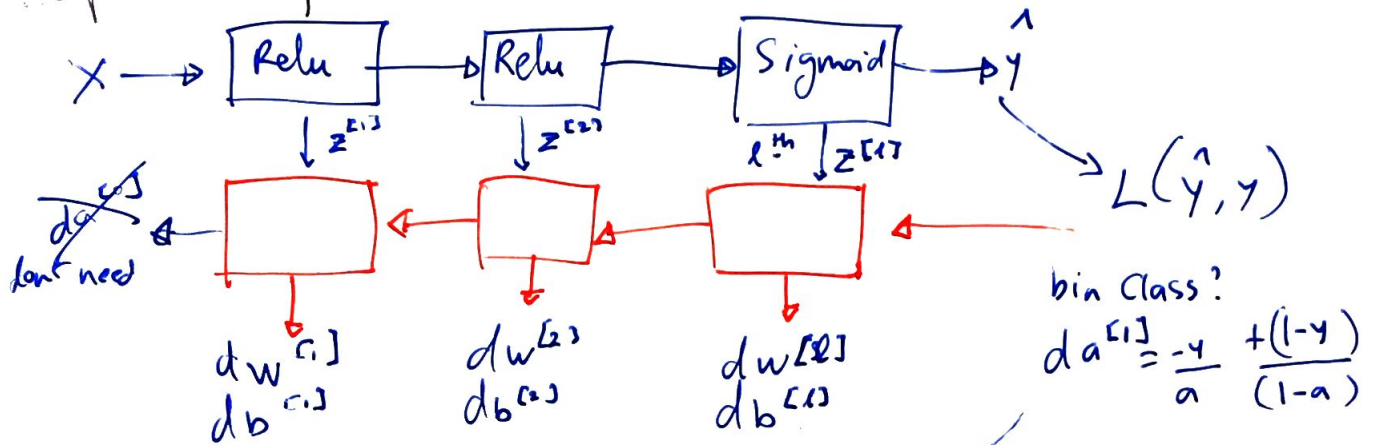$db^{[l]} = dz^{[l]}$; $da = w^{[l]T} \cdot dz^{[l]}$

# Vectorized version:

cached $dz^{[l]} = dA^{[l]} * g^{[l]}(z^{[l]})$

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

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

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

# Update Cycle:



$X \rightarrow$ [Relu] $\rightarrow$ [Relu] $\rightarrow$ [Sigmoid] $\rightarrow \hat{y} \rightarrow L(\hat{y}, y)$

$z^{[1]}$   $z^{[2]}$   $l^{th}$ $z^{[l]}$

$da^{[l]}$ don't need

$dw^{[1]}$   $dw^{[2]}$   $dw^{[l]}$
$db^{[1]}$   $db^{[2]}$   $db^{[l]}$

bin Class?
$da^{[l]} = \frac{-y}{a} + \frac{(1-y)}{(1-a)}$

vectorized $v^n =$
$dA^{[l]} = \left( \frac{-y^{[1]}}{a^{[1]}} + \frac{1-y^{[1]}}{1-a^{[1]}} \cdots \right.$
$\cdots \frac{y^{[m]}}{a^{[m]}} + \frac{1-y^{[m]}}{1-a^{[m]}}$

# Hyper parameters:

Hyper params values impact learning algo's params

Params: $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}, W^{[3]}, b^{[3]}$ ---

controls

Hyper params:   Learning rate $\alpha$
                # iterations
                # hidden layer L
                # hidden units $h^{[1]}, n^{[2]}$ ---
                choice of activation $f^n$

later will also see   Momentum, mini batch size, regularizations, ---

$\rightarrow$ Applied deep learning is a very emperical process
                                    gotta try a lot of values