

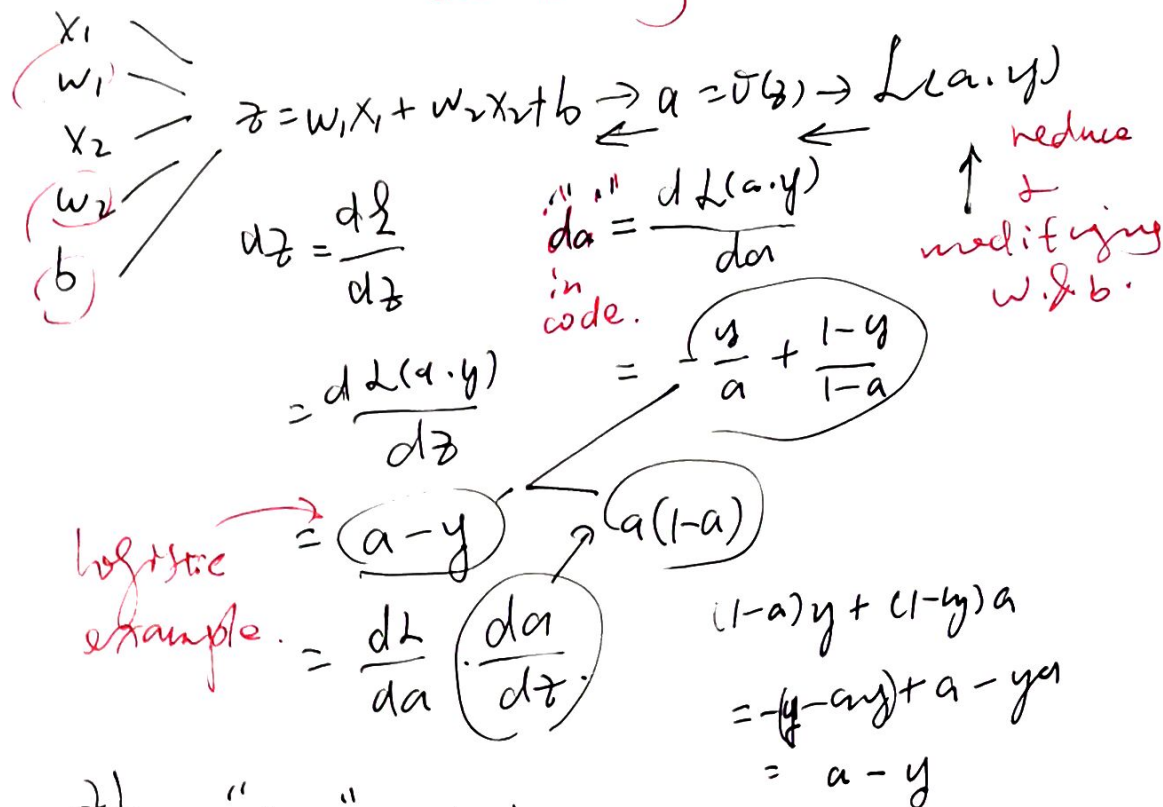
Computing Graph: For Logistic Regression.

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

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

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

One training example.



$$\begin{cases} \frac{\partial L}{\partial w_1} = "dw_1" = x_1 \cdot dz \\ \frac{\partial L}{\partial w_2} = "dw_2" = x_2 \cdot dz \\ db = \frac{\partial L}{\partial b} = dz \end{cases}$$

Backward propagation.

Going backwards to find out how much you need to change "dw1", "dw2" & "db" to ~~find out~~ change loss func a little bit.

Logistic Regression on m examples.

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

$$(x^{(i)}, y^{(i)})$$

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

$$dw_1^{(i)}, dw_2^{(i)}, db^{(i)}$$

$$\frac{\partial}{\partial w_1} J(w, b) = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial w_1} \mathcal{L}(a^{(i)}, y^{(i)}) \quad *$$

↓
overall
cost function derivative.

" $dw_1^{(i)}$ " → shown on single training example.

* Procedure :

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

For $i=1$ to m → first for loop

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

refer to only one single training example

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

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

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

$$dw_j += x_j^{(i)} dz^{(i)}$$

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

accumulators

second for loop.

↓
mitigator by
vectorization
to reduce for loop.

$$J /= m$$

$$dw_1 /= m$$

$$dw_2 /= m$$

⋮

$$db /= m$$

get final answers.

For loop in calculating w_1, w_2, \dots, w_n .

Vectorization (numpy lib)

$$w = \begin{bmatrix} \\ i \\ \end{bmatrix} \quad x = \begin{bmatrix} \\ i \\ \end{bmatrix}$$

$$X, w \in \mathbb{R}^{n \times 1}$$

Non-Vec :

$$z = 0$$

$$\text{for } i \text{ in range}(n_x):$$

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

$$z += b.$$

slow.

$$u = Av$$

$$u_i = \sum_j A_{ij} v_j$$

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

for i ...

for j ...

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

$$v = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix} \rightarrow u = \begin{bmatrix} e^{v_1} \\ e^{v_2} \\ \vdots \\ e^{v_n} \end{bmatrix}$$

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

for i in range(n):

Vectorized:

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

$$\downarrow$$

$$w^T x.$$

GPU.

CPU

good at.

SIMD - single
instruction
multiple data.

Guideline:

avoid for loop
whenever necessary.

every time see
if it is possible
to call a numpy
command instead
of using for loop.

import numpy as np

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

$$\text{np.log}(v)$$

$$\text{np.abs}(v)$$

max(v, 0)

$$\left. \begin{array}{l} dw_1 = 0 \\ dw_2 = 0 \\ \vdots \end{array} \right\} \rightarrow dw = np.zeros((n_x, 1))$$

for $j = 1 \dots n_x$
 $dw_j += \dots$ (gets updated)
 gets optimised by
 $dw_j += x_j^{(i)} dz^{(i)}$

② $dw += x^{(i)} dz^{(i)}$ using
 up command.

③ $dw_1 = dw_1 / m \dots$
 $\rightarrow dw = dw / m$

For m
 trials examples:

$$\left\{ \begin{array}{l} z^{(1)} = w^T x^{(1)} + b \\ a^{(1)} = \sigma(z^{(1)}) \\ z^{(2)} = \dots \\ z^{(3)} = \dots \end{array} \right.$$

carry out forward
 propagation
 simultaneously with
 m examples without
 for loop.
 \rightarrow matrices

$$X = \left[\begin{array}{c|c|c|c} x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ \hline | & | & & | \\ \hline \end{array} \right]_{n_x}$$

m

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

$$[a_1, a_2, \dots, a_m] = \sigma(z)$$

how to 'sigmoid' all the
 vectorised z . shown in
 assignment.

\hookrightarrow using numpy `sigmoid`
`np.exp()`.

$$\begin{bmatrix} z^{(1)} & z^{(2)} & \dots & z^{(m)} \end{bmatrix} = w^T X + \begin{bmatrix} b & b & \dots & b \end{bmatrix}_{1 \times m} = \begin{bmatrix} w^T X + b & w^T X + b & \dots & w^T X + b \end{bmatrix}$$

$$z = np.dot(w.T, x) + b$$

$\mathbb{R}_{(1,1)} \rightarrow$ broadcasting

$$dz^{(1)} = a^{(1)} - y^{(1)} \dots$$

$$dz = [dz^{(1)} dz^{(2)} \dots dz^{(m)}]$$

$$A = [a^{(1)} \dots a^{(n)}]$$

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

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

$$db = \frac{1}{m} \sum_{i=1}^m dz^{(i)}$$

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

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

$$= \frac{1}{m} \begin{bmatrix} x^{(1)} & \dots & x^{(m)} \\ | & & | \\ 1 & & 1 \end{bmatrix} \begin{bmatrix} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{bmatrix}$$

$$= \frac{1}{m} [x^{(1)} dz^{(1)} + \dots + x^{(m)} dz^{(m)}]$$

Single
step

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

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

$$A = \sigma(z)$$

$$dz = A - Y$$

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

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

still
need
a
for
loop to
do
multiple
iterations

Broadcasting example:

$$\textcircled{1} \text{cal} = A \cdot \text{sum}(\text{axis}=0)$$

want to sum up vertically
horizontally \rightarrow axis=1

$$\text{percentage} = 100 * A / \text{cal} \cdot \text{reshape}(1,4)$$

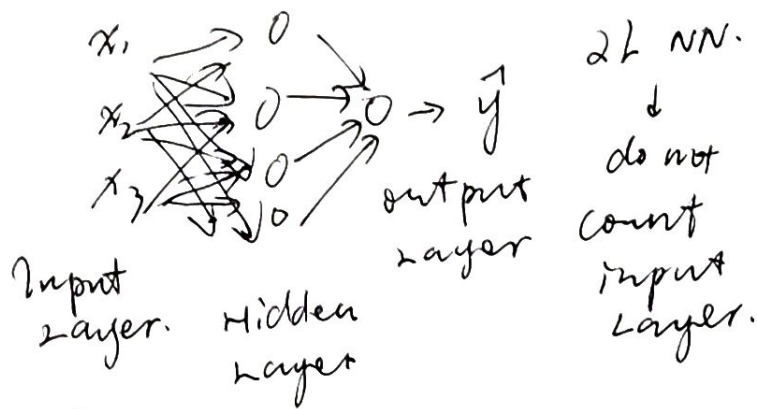
cheap to call
always confirm if
not sure.

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

$$\textcircled{3} \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \end{bmatrix} = \begin{bmatrix} 101 & 202 & 303 \\ 104 & 205 & 306 \end{bmatrix}$$

$$\textcircled{4} \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 \\ 200 \end{bmatrix} = \begin{bmatrix} 101 & 102 & 103 \\ 204 & 205 & 206 \end{bmatrix}$$

Neural Network Representation



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

2 L NN

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

(1.1) (1.4) (4.1) (1.1)

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

actual value of this layer is not observed.

Input activation $\rightarrow a^{[0]} = x$
Hidden: a_1, a_2, \dots

Given

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

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

$a_i^{[1]} \leftarrow$ node in layer.

$$\begin{cases} x \rightarrow a^{[2]} = \hat{y} \\ x^{(1)} \rightarrow a^{[2](1)} = \hat{y}^{(1)} \\ \vdots \\ x^{(m)} \rightarrow a^{[2](m)} = \hat{y}^{(m)} \end{cases}$$

layer 2 example i

For i from 1 to m:

$$\begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix} = \begin{bmatrix} - & w_1^{[1]T} & - \\ - & w_2^{[1]T} & - \\ & \vdots & \\ - & w_4^{[1]T} & - \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix}$$

4×3 $\downarrow W^{[1]}$ $\downarrow b^{[1]}$

vectorized to avoid for loop for training examples

$$z^{1} = w^{[1]} x^{(1)} + b^{[1]}, \quad z^{[1](2)} = w^{[1]} x^{(2)} + b^{[1]}, \quad z^{[1](3)} = w^{[1]} x^{(3)} + b^{[1]}$$

$$w^{[1]} = \begin{bmatrix} \equiv \\ \equiv \\ \equiv \end{bmatrix}, \quad w^{[1]} x^{(1)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}, \quad w^{[1]} x^{(2)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}, \quad w^{[1]} x^{(3)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

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

$$= \begin{bmatrix} z^{1} & z^{[1](2)} & z^{[1](3)} \\ | & | & | \end{bmatrix} = z^{[1]}$$

$$A = \sigma(z) \rightarrow A \begin{bmatrix} | & | & | \\ a^{1} & a^{[1](2)} & \dots \\ | & | & | \end{bmatrix} \left\{ \begin{array}{l} \text{hidden units} \\ \text{in 1st Layer} \end{array} \right\}$$

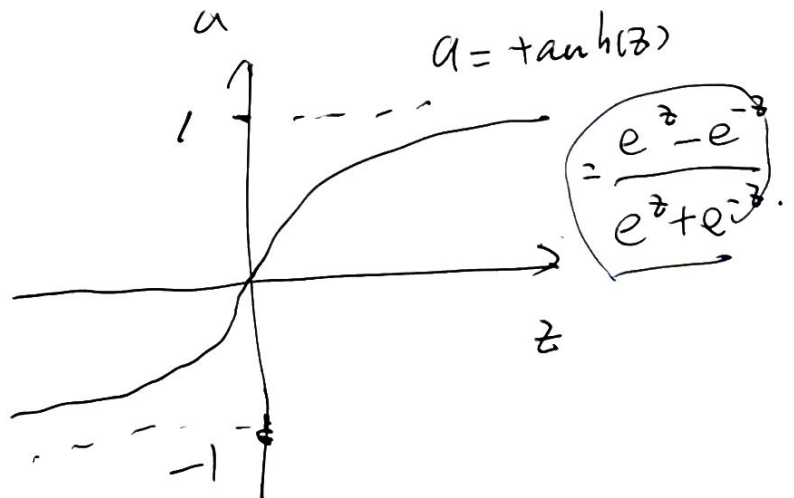
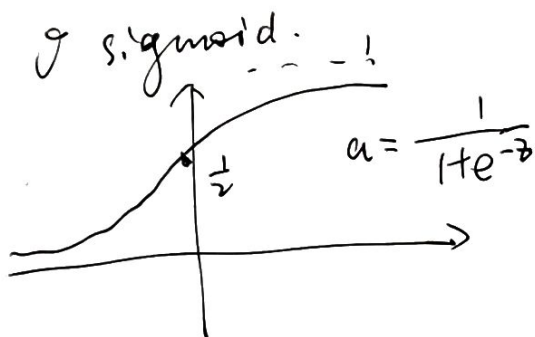
+ running example.

$$* z^{[i]} = w^{[i]} A^{[i-1]} + b^{[i]}$$

$$\left\{ \begin{array}{l} A^{[0]} = X \text{ (input layer)} \\ a^{[0](i)} = x^{(i)} \end{array} \right.$$

Forward propagation fully vectorized.

Activation Function:

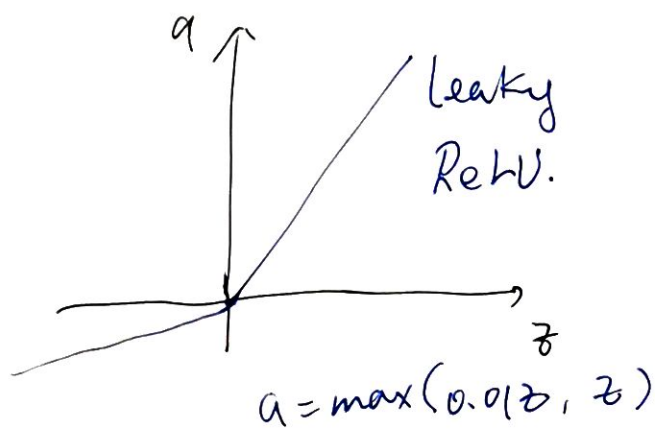


Tanh function almost always outperforms sigmoid function since it has 0 mean makes learning for 2nd layer easier.

Sigmoid function → used for binary classification ^{or} it requires the probability. So it is used for output layer

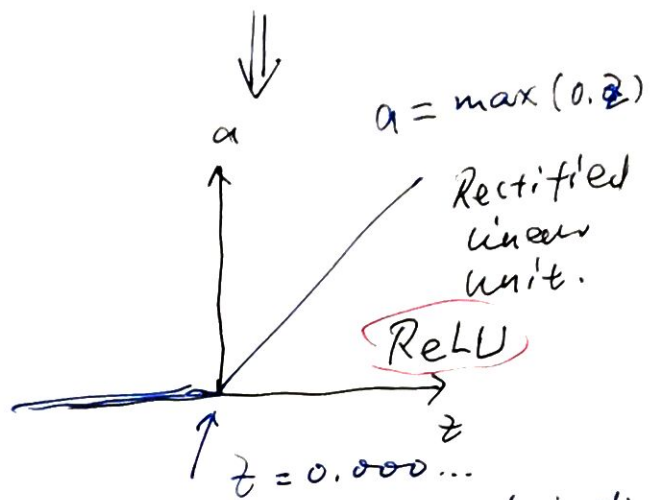
Different function for different layer.

g [i] ← Layer.



common drawback:

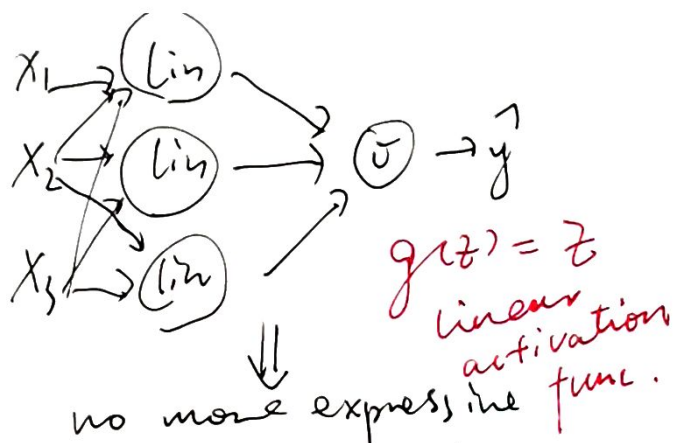
slope is small, z is either very large or very small



pretend the derivative when $z = 0$ is either 0 or 1. works just fine. Increasingly default activation function. (Hidden Layer)

If not sure what to use, ReLU is a great default choice.

ReLU: less effect the slope really small when going to zero → usually much faster than tanh & sigmoid



no more expressive than a standard Logistic Regression.

linear activation func always calculate the same \rightarrow Layers does not have their use all \rightarrow linear function.

one use: doing linear regression on housing price, etc.

use linear function for output layer.

For sigmoid func.

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

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

(clearly) ReLU:

$$g(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z > 0 \\ \text{undef} & z = 0. \end{cases}$$

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

$$\begin{cases} z \rightarrow \infty, g(z) \rightarrow 1 \\ z \rightarrow -\infty, g(z) \rightarrow 0 \\ z = 0, g(z) = \frac{1}{2} \times \frac{1}{2} = \frac{1}{4} \end{cases}$$

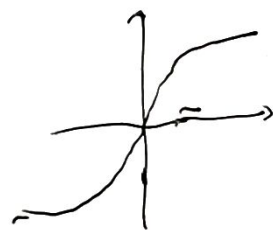
$$a' = a(1-a)$$

Tanh activation func.

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

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

$$a' = 1 - a^2$$



$$\begin{cases} g(z) \rightarrow 1, \text{ when } z \rightarrow \infty \\ g(z) \rightarrow 0, \text{ when } z \rightarrow -\infty \\ z = 0 \rightarrow g(z) = 0 \end{cases}$$

Gradient Descent for
One hidden layer NN.
2 L NN.

$$\frac{\partial L}{\partial z^{[2]}} = \frac{\partial L}{\partial z^{[1]}} \cdot \frac{\partial z^{[2]}}{\partial a^{[1]}} \cdot \frac{\partial a^{[1]}}{\partial z^{[1]}}$$

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

cost func: $J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{[i]}, y^{[i]})$

Gradient descent:

Repeat: {

→ Compute predictions ($\hat{y}^{[i]}, i=1 \dots m$)

$$dw^{[1]} = \frac{\partial J}{\partial w^{[1]}}, db^{[1]} = \frac{\partial J}{\partial b^{[1]}}, \dots$$

$$w^{[1]} = w^{[1]} - \alpha dw^{[1]}$$

⋮

}

troublesome.

element wise product.

until converging.

Vectorized:

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

$$dw^{[2]} = \frac{1}{m} dz^{[2]} A^{[1]T}$$

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

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

$$dw^{[1]} = dz^{[1]} X^T$$

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

Forward & Backward Graph: $w^{[1]}, b^{[1]}$

Similar to
logistics.

$$X \rightarrow z^{[1]} = w^{[1]} X + b^{[1]} \rightarrow a^{[1]} = \sigma(z^{[1]}) \rightarrow z^{[2]} = w^{[2]} a^{[1]} + b^{[2]} \rightarrow a^{[2]} = \sigma(z^{[2]}) \rightarrow L(y, \hat{y})$$

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

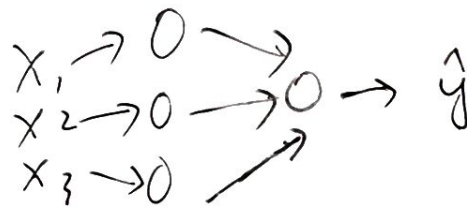
$$dz^{[2]} = a^{[2]} - y$$

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

$$dw^{[1]} = dz^{[1]} X^T$$

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

$$L_1, L_2, db^{[2]} = dz^{[2]}$$



If weights are initialised as zero.

$$W^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}. \quad dw = \begin{bmatrix} u & v \\ u & v \end{bmatrix}$$

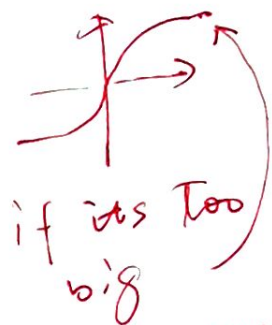
no matter how many neurons are in the layer. they are essentially computing the same thing. \rightarrow as one neuron.

$$W^{[l]} = \begin{bmatrix} - & - \\ - & - \end{bmatrix} \text{ equal.}$$

Random Initialisation:

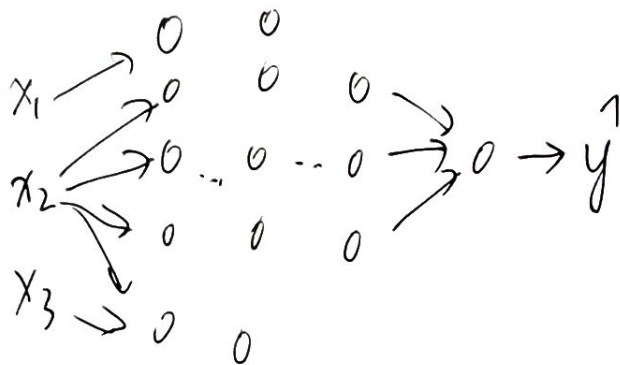
$$\begin{aligned} W^{[1]} &= \text{np.random.randn}(2,2) * 0.01 \\ b^{[1]} &= \text{np.zeros}(2,1) \rightarrow \text{no symmetry} \\ W^{[2]} &= \dots \text{rand} \dots (1,2) \rightarrow \text{problem} \\ b^{[2]} &= \text{zero} \dots \end{aligned}$$

works for the shallow neural network.
deep one different mech.
Initial with small value.



if its too big
learning will be very slow.

Deep Neural Networks.



$L=4$ (# layer)

$n^{[L]}$ # units in layer L .

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

$a^{[L]} \rightarrow$ activation in layer L .

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

$W^{[L]}$ weights for $z^{[L]}$.

$$n^{[2]} = n_x = 3.$$

* Forward propagation in deep NN

$$\begin{cases} z^{[l]} = w^{[l]} A^{[l-1]} + b^{[l]} \\ A^{[l]} = g^{[l]}(z^{[l]}) \end{cases}$$

$$\begin{cases} b^{[l]} : (n^{[l]}, 1) \\ dw^{[l]} : (n^{[l]}, n^{[l-1]}) \\ db^{[l]} : (n^{[l]}, 1) \end{cases}$$

z & a — same dimension

↓
vectorized implementation
 $(n^{[l]}, 1)$

$$z^{[l]} = w^{[l]} X + b^{[l]}$$

$(n^{[l]}, m) \quad (n^{[l]}, n^{[l-1]}) \quad (n^{[l]}, m)$

bd. $(n^{[l]}, m)$
broadcast

For loop for different layer is perfectly suitable.

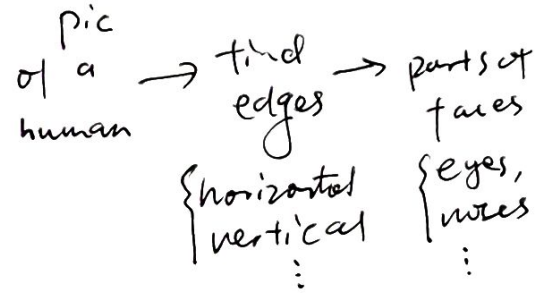
("not avoidable" — Andrew Ng)

Dimensions Check

$$z^{[l]} = w^{[l]} X + b^{[l]}$$

$(3,1) \quad (3,2) \quad (2,1) \quad (3,1)$

$\Rightarrow (n^{[l]}, n^{[l-1]})$



→ different type of faces.
compose hierarchically.
another eg: sound recog.

Some functions require exponentially more hidden ~~neurons~~ units if layers are limited.

Circuit theory

$$x_1 \text{ XOR } x_2 \text{ XOR } x_3 \dots \text{ XOR } x_n$$

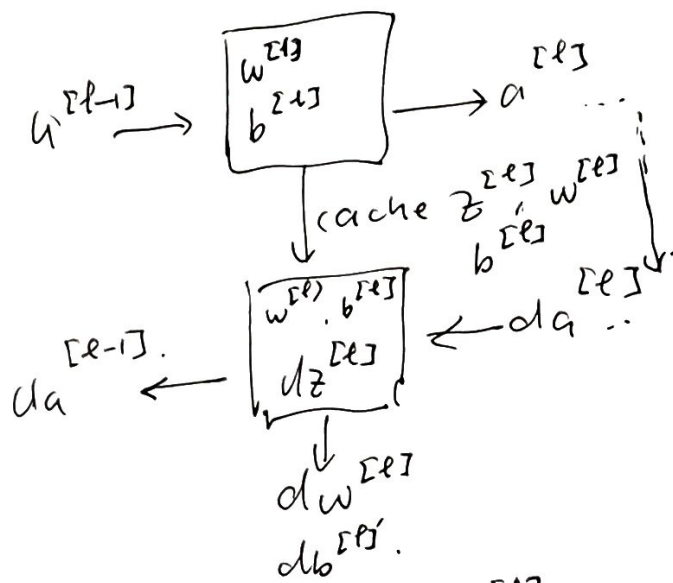
$(\log n)$
elongate exponentially
 2^n

Backward functions.

Input: $da^{[l]}$
 $cache(z^{[l]})$

Output: $da^{[l-1]}$
 $dw^{[l]}, db^{[l]}$

layer l .



$$\begin{cases} w^{[l]} = w^{[l]} - \alpha dw^{[l]} \\ b^{[l]} = b^{[l]} - \alpha db^{[l]} \end{cases}$$

Backward Propagation for layer l .

$$\begin{cases} dz^{[l]} = da^{[l]} * g'(z^{[l]}) \\ dw^{[l]} = dz^{[l]} \cdot a^{[l-1]} \\ db^{[l]} = dz^{[l]} \\ da^{[l-1]} = W^{[l]T} \cdot dz^{[l]} \end{cases}$$

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

Parameters & Hyperparameters

$$W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}, \dots$$

Hyperparameters:

Learning rate α
 # iterations
 # hidden layers L
 # hidden unites.
 $\eta^{[1]}, \eta^{[2]}, \dots$
 choice of activation func.
 many more ...

empirical process
 try to tune the parameters

Hyperparameter might be different in different industries.