# Deep Learning Specialization
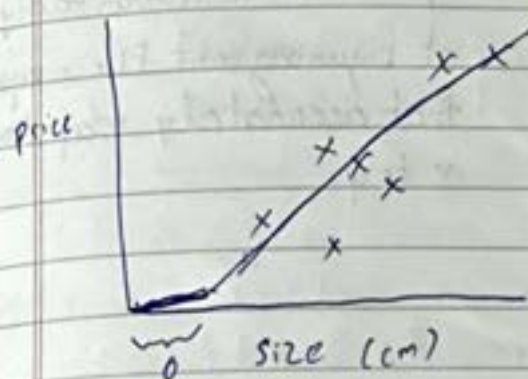
week-1 **Course-1: Neural Networks and Deep Learning**

Deep Learning : → Training very large NN's

## What is a neural Net?
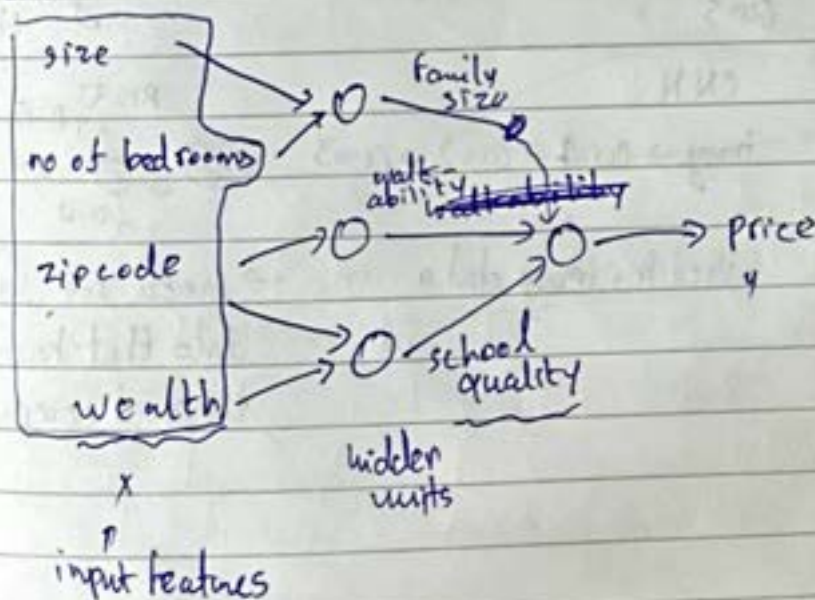


price

size (cm)

This ReLU

neuron

Size ——→ O ——→ Price
   x              y

S input - Neuron - this implements the function we drew on the left.
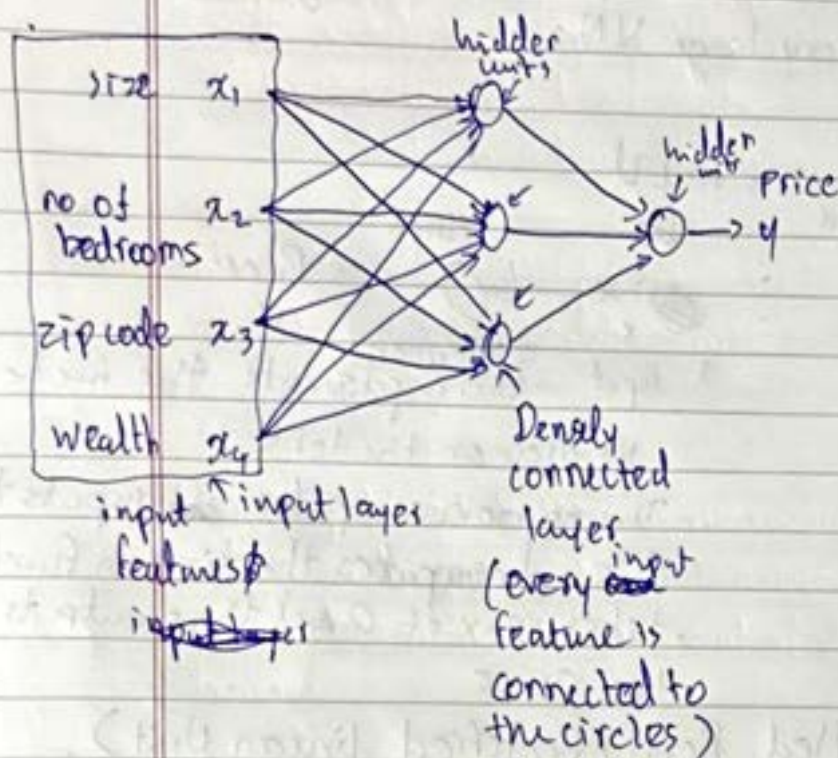
- The neuron computes its Inputs the Size and computes the linear function, takes max of 0 and then outputs the price

→ this function is called ReLU (Rectified linear Unit), the fn goes to 0 then takes off as a straight line. Rectify means taking max of 0 which is why we get a shape like this.

→ A single neuron NN is a very small NN, a large NN is formed by stacking taking many neurons and stacking them toghether.

ex.



size

no of bedrooms

zipcode

wealth

   x

input features

family size

walk ability

school quality

→ Price
    y

hidden units

## Denotation



hidden units

size $x_1$

no of bedrooms $x_2$

zip code $x_3$

wealth $x_4$

input input layer
features
input layer

hidden layer

Price → $y$

Densely
connected
layer
(every input
feature is
connected to
the circles)

Hidden units take in
input $(x_1, x_2, x_3, x_4)$

→ Given enough data about x and
y. NN's are remarkably good
at figuring out the function
that accurately map from
x to y.

## Supervised learning with NN

| input (x) | output (y) | Application | |
|---|---|---|---|
| Home features | Price | Real Estate | } Standard NN |
| Ad user info | Click on ad? (0/1) | Online Advertising | |
| Image | Object | Photo Tagging | } CNN RNN |
| English | Chinese | Machine Translation | |
| Image, Radar Info | Position of other Cars | Autonomous driving | } Custom / Hybrid |

Standard NN



$x_1$
$x_2$

CNN

image → conv1 → conv 2 → conv3

- Used for image data

RNN
$y^{[t-1]}$
$x^{[t-1]}$

→ Good for 1d seq
data that has a
temporal component

## Supervised Learning

We have applications of ML to both:
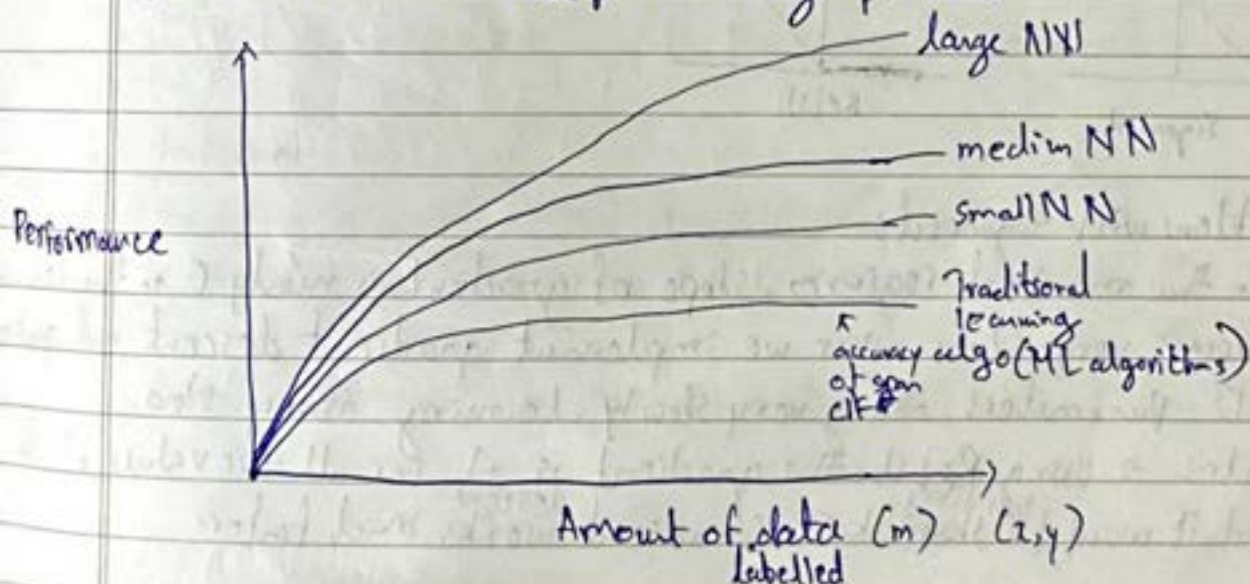- Structured data
- Unstructured data

## Structured data:

→ Databases of data

→ ex in housing price prediction, we have database or columns that tell us size of house, no of bedrooms

→ ex, whether or not you click on an ad, you might have info about the age, info about the ad and labels your trying to predict

## Unstructured data:

→ Refers to audio, text, images, where you might want to recognise
- ~~you might~~ the text or image

→ features can be pixel values of image or individual words in a piece of text.

→ With rise of NN's and DL, computers are much better in interpreting unstructured data.

Why is Deep learning taking off?
Scale drives the deep learning process



Performance (y-axis) vs Amount of data (m) (x,y) labelled (x-axis)

Curves from top to bottom:
- large NN
- medium NN
- small NN
- Traditional learning algo (ML algorithms)

accuracy of spam of cat

→ If we want a high level of performance, we need 2 things:
  - You need to train a big enough NN to take advantage of the huge amount of data
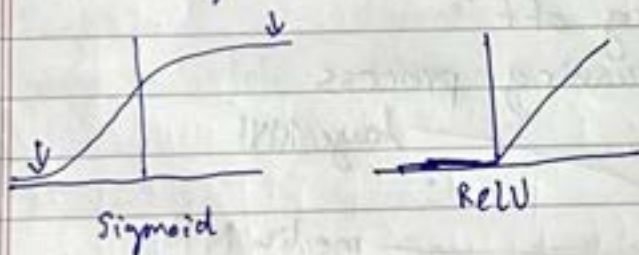  - We need a lot of data, ~~node~~

→ Scale has driven DL process. Scale refers to size of NN, just a neural network, a lot of hidden units, a lot of parameters, a lot of connections.

→ Best way to get better performance in a net:
  - Train a bigger net or throw more data it only works up to a point because we run out of data or the net is too big to train.

Scale drives DL progress:
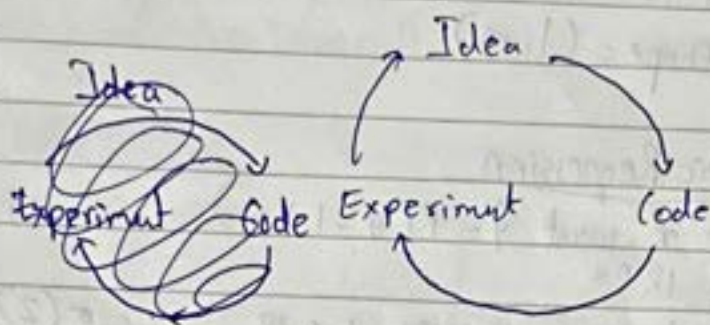→ Data
→ Computation
→ Algorithms

→ Breakthrough of NN: ex, switching from sigmoid to ReLU.



Sigmoid          ReLU

Problem with sigmoid:
→ In the marked regions, slope of gradient is nearly 0 so the learning becomes very slow when we implement gradient descent and gradient is 0, parameters change very slowly, learning is very slow
● Soln: → Using ReLU. The gradient is +1 for all +ve values of ile and it would less likely shrink 0. Gradient descent works much faster.

→ Just computation is very important, training process of net is very intuitive.
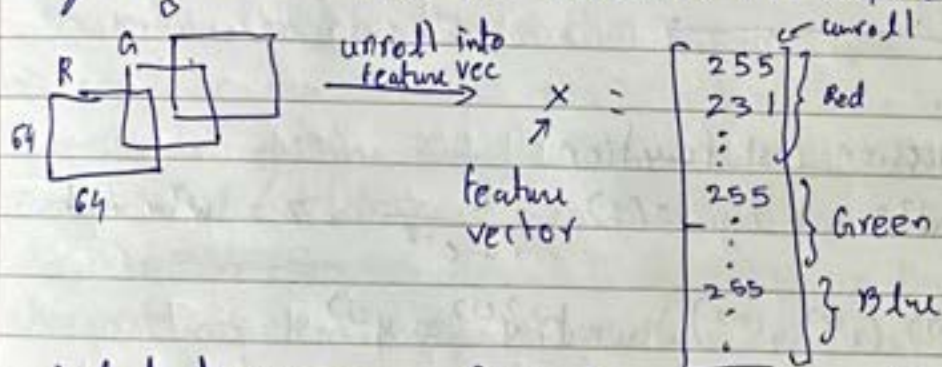


## Week-2

## Neural Network Basics

### Binary Classification:

- We might have an image of cat, we want to classify whether its cat (1) or not cat (0) → $y$.
- An img is represented as 3 seperate matrices for Red (R), Green (G), Blue (B) color channels - ex, if we have $64 \times 64$ (height × width) img we have 3 R, G, B matrices of $64 \times 64$.



total dimension of vector $x = 64 \times 64 \times 3 = 12288$

$n = n_x = 12288$      $n, n_x$: dimensions of $x$

### Notation

training example $(x, y)$  $x \in \mathbb{R}^{n_x}, y \in \{0, 1\}$

we will use $m$ to denote a train ex: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots (x^{(m)}, y^{(m)})\}$

$M = m_{train}, m_{test}$

Putting all train ex in compact notation

$$X = \begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & & | \end{bmatrix} \begin{array}{l} \uparrow \\ n_x \\ \downarrow \end{array}$$

$m$: columns
$n_x$: rows

$X \in \mathbb{R}^{n_x \times m}$
$x.shape: (n_x, m)$

$$y = [y^{(1)} \ y^{(2)} .. y^{(m)}]$$
$$y \in \mathbb{R}^{1 \times m}$$
$$y \cdot shape = (1, m)$$

## Logistic Regression

Given $x$, want $\hat{y} = P(y = 1 | x)$     $z = w^T x + b$

$x \in \mathbb{R}^{n_x}$
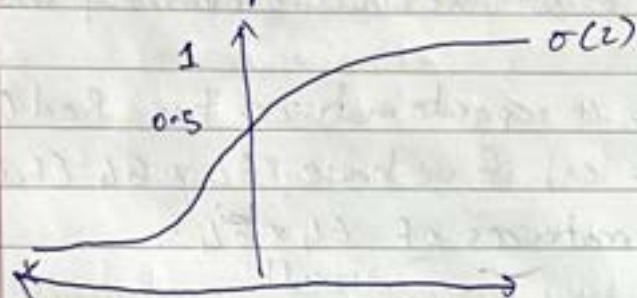
Parameters: $w \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$     $\sigma(z) = \dfrac{1}{1 + e^{-z}}$

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

Sigmoid

if $z$ is large:
$$\sigma(z) = \frac{1}{1 + 0} \approx 1$$



if $z$ is a large negative no
$$\sigma(z) = \frac{1}{1 + e^{-z}} \approx 0$$

## Logistic Regression Cost function

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

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

Loss (error) function: $L(\hat{y}, y) = -(y \log(\hat{y}) + (1-y) \log(1-\hat{y}))$

if $y = 1$: $L(\hat{y}, y) = -\log(\hat{y}) \leftarrow$ want $\log \hat{y}$ large, want $\hat{y}$ large

if $y = 0$: $L(\hat{y}, y) = -\log(1-\hat{y}) \leftarrow$ want $\log(1-\hat{y})$ large... want $\hat{y}$ small!

- Loss function was defined wrt single training ex

Cost function: measures how are you doing on the entire train set
$$J(w, b) = \frac{1}{m} \sum_{i=1}^{m} L(\hat{y}^{(i)}, y^{(i)})$$
$$= -\frac{1}{m} \sum_{i=1}^{m} \left[ y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log(1-y^{(i)}) \right]$$
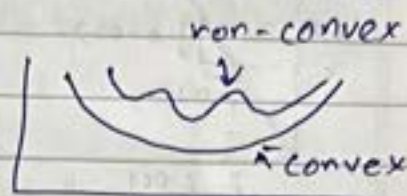
→ Cost function is the cost of training the parameters $(w, b)$, so in
→ ~~Using the~~ training the logistic regression model we will try finding
parameters $(w, b)$ that minimise overall cost function

## Gradient Descent (GD)

★ Recap on logistic regression,

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

$\hat{y} = \sigma(z) = \dfrac{1}{1 + e^{-z}}$


non-convex
convex

$$J(w,b) = \frac{1}{m} \sum_{i=1}^{m} L(\hat{y}^{(i)}, y^{(i)}) = \frac{-1}{m} \left[ \sum_{i=1}^{m} y^{(i)} \log(\hat{y}) + \sum_{i=1}^{m} (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}) \right]$$
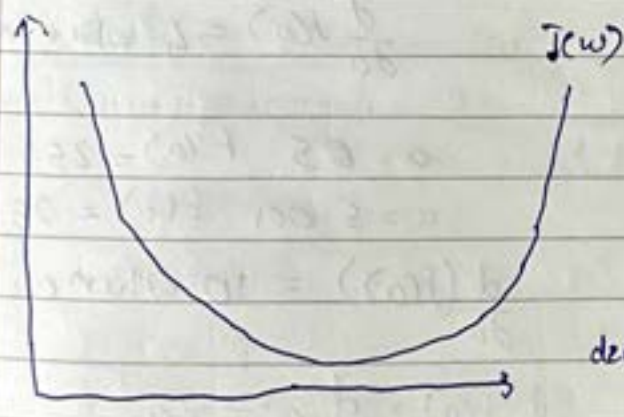
→ Find $w, b$ that minimise $J(w, b)$
- $J(w, b)$ is a convex function

→ To find good parameter values we initialize the values ~~to~~
of $w, b$ to $0$. (Note: we ~~can use~~ random initialization its not
recommended) for logistic regression ~~Bad~~

→ Since the function is convex, no matter where we initialize,
we should get the same point or roughly the same point.

→ ~~that~~ Gradient descent will start at the initial point and takes a
step ~~in the steepest~~ downhill ~~direction~~ in the direction of
steepest descent as quickly down as possible. This is one iteration
of GD

→ After a few iterations we converge, to the global optimum.
                                    close


$J(w)$

We repeatedly carry out this
update:
repeat {

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

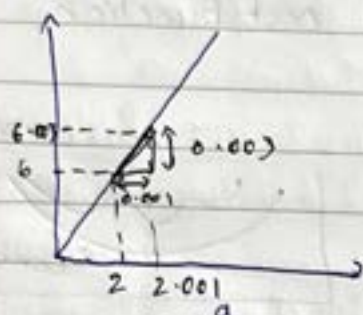$$b := b - \alpha \frac{\partial J(w,b)}{\partial b}$$

denotation in code
$\partial J(w,b)/\partial w \to dw$    $\partial J(w,b)/\partial b \to db$

$\alpha$: is the learning rate it controls how big a step we take in each step
of gradient descent

## Derivatives

Intuition about derivatives :

$f(a) = 3a$      $a = 2$      $f(\cdot 2) = 3(2) = 6$

$a = 2 \cdot 001$   $f(2 \cdot 001) = 6 \cdot 003$

- Slope (derivative) of $f(a)$

$slope = \dfrac{height}{width} = \dfrac{0.003}{0.001} = 3$

at $a = 2$, slope $= 3$

⇒Slope $= 3$ denotes the fact that when you, increase a to the right by 0.001.

➔ The amount at $f(a)$ goes up 3 times as high as you increased it
ir the horizontal direction

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

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

slope at $a = 5 = \dfrac{0.003}{0.001} = 3$

→ Derivatives are defined with an even smaller value of how much
you increase   $a$ to the right.

## More derivative examples

$f(a) = a^2$      $a = 2$,      $f(a) = 4$

$a = 2.001$      $f(a) = 4.004001$

slope at  $a = 2$ is 4

slope $= \dfrac{0.004}{0.001} = 4$

$\dfrac{d}{da} f(a) = 4$ when $a = 2$

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

$a = 5.001$  $f(0) = 25.010$

$\dfrac{d}{da}(f(a)) = 10$ when $a = 5$

$\dfrac{d}{da} f(a) = \dfrac{d}{da} a^2 = 2a$

More derivative examples

$f(a) = a^2$, $\frac{d}{da}(f(a)) = 2a$  $a = 2$, $f(a) = 4$

$a = 2.001, f(a) = 4.001$

$f(a) = a^3, \frac{d}{da} f(a) = 3a^2$, $a = 2$, $f(a) = 8$

$a = 2.001, f(a) = 8.012$

$f(a) = \log_e(a)$ & $\frac{d}{da} f(a) = \frac{1}{a}$

## computation Graph

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

Forward pass: Used to compute output of neural net.

Back propagation: Used to compute gradients or derivatives

$u = bc$

$v = a + u$

$j = 3v$

computation graph

forward pass

$a = 5$  $a$

$b = 3$  $b$

$c = 2$  $c$

$u = b \cdot c$   $v = a + u \rightarrow J\ell = 3v$

$6$    $11$    $33$

back prop

inputs: $u, v, j$

forward pass: left to right we compute value of $j$
- In order to compute derivatives we go from right to left

## Derivatives with a computation Graph

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

$a = 5$

$v = a + u$     $J = 3v$     $\frac{dJ}{dv} = 3$   $\frac{dJ}{dv} = 3$

$\frac{du}{db} = c = 2$

$b = 3$    $u = bc$

$\frac{dJ}{dv} = 3$   $\frac{dJ}{da} = 3$   $\frac{dJ}{db} = 6$

$c = 2$

$= 6$

$\frac{dJ}{dc} = \frac{dJ}{da} \cdot \frac{du}{dc}$

$\frac{dJ}{da} = \frac{dJ}{dv} \cdot \frac{dv}{du} = 3 \cdot 1 = 3$  from $v = a + u$  $\frac{dv}{da} = 1, \frac{dv}{du} = 1$

$= 3 \cdot b$

$= 9$

$\frac{dJ}{da} = \frac{dJ}{dv} \cdot \frac{dv}{da} = 3 \cdot 1 = 3$

Derivatives ~~with a computal~~
Logistic Regression : Gradient Descent

$$z = w^T x + b$$
$$\hat{y} = a = \sigma(z)$$
$$L(a, y) = -(y \log(a) + (1-y) \log(1-a))$$

dz, da → are code variable names



$$z = w_1 x_1 + w_2 x_2 + b \rightarrow \boxed{\hat{y} = a = \sigma(z)} \rightarrow \boxed{L(a,y)}$$

$$dz = \frac{dL}{da} \cdot \frac{da}{dz} \nearrow^{a(1-a)} \qquad da = \frac{dL(a,y)}{da} \rightarrow da \text{ denotes the code representation}$$

$$= a - y \qquad \frac{dL}{da} = d\dot{a} = -\left(\frac{-y}{a} + \frac{1-y}{1-a}\right)$$

$$\frac{dL}{dw_1} = dw_1 = x_1 \cdot dz$$

$$w_1 := w_1 - \alpha \, dw_1$$
$$w_2 := w_2 - \alpha \, dw_2$$
$$\frac{dL}{dw_2} = dw_2 = x_2 \, dz \qquad b := b - \alpha \, db$$

$$db = dz$$

Gradient Descent on m train examples

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

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

$$\frac{\partial}{\partial w_1} J(w, b) = \frac{1}{m} \sum_{i=1}^{m} \underbrace{\frac{\partial}{\partial w_1} L(a^{(i)}, y^{(i)})}_{dw_1}$$

$$dw_1^{(i)} - (x^{(i)}, y^{(i)})$$

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

- for $i = 1$ to $m$

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

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

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

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

$dw_1 \;{+}{=}\; x_1^{(i)} dz^{(i)}$

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

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

$J /{=} m$

$dw_1 /{=} m, \; dw_2 /{=} m, \; db /{=} m$

-) Disadvantage with the above implementation:

- We are using 2 for loops to implement logistic regression.
  - → 1ˢᵗ for loop to iterate over all training ex.
  - → 2ʳᵈ for loop, the no of times we iterate is equal to the ro of features $(dw_1, dw_2)$
  - → This a problem because when we use bigger datasets, its important to implement your algorithms without explicitly "using for" loops is important and it helps scale to longer datasets.

-) Solution: Vectorization, this speeds up your code and we get rid of (or loops.

## Python and Vectorization

Vectorization in logistic regression

ex, $z = w^T x + b$    $\qquad w = \begin{bmatrix} : \end{bmatrix}, \; x = \begin{bmatrix} : \end{bmatrix} \quad \begin{matrix} w \in \mathbb{R}^{n_x} \\ x \in \mathbb{R}^{n_x} \end{matrix}$

| Non vectorized | Vectorized |
|---|---|
| $z = 0$ | $z = np.dot(w, x) + b$ |
| for $i$ in range $(n-x)$: | |
| $\quad z \;{+}{=}\; w[i] * x[i]$ | $\underbrace{\qquad\qquad}_{w^T \cdot x}$ |
| $z \;{+}{=}\; b$ | |

→ Scalable & DL implementations are done on GPU
→ GPU and CPU have parallelization instructions
→ They are called SIMD instructions (Single instruction multiple data)
→ SIMD: using built in functions. or other functions such np-dot.
we don't have to explicitly implement a for loop. It enables
Python numpy to take advantage of parallelism to do the computation
faster. These work great on GPU compared to CPU.

## More Vectorization examples   explicit

NN programming guideline: avoid for loops whenever possible

ex $u = Av$

$u_i = \sum_j A_{ij} v_j$

vectorized

$u = np.dot(A, v)$

non vectorized

u = np.zeros ((n,1))

for i...

for j...

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


Say you need to apply the exponential operation on every element of
a matrix/vector.

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

non vectorized

u = np.zeros((n,1))
for i in range (n):
    u[i] = math.exp (v[i])

vectorized

u = np.exp(v)

we can apply similar
operations:

np.log (v)
np.dot (v)
np.maximum(v, 0)

## logistic Regression derivatives

$J=0,$ $\boxed{dw1=0, dw2=0,}$ $\boxed{db=0}$

instead of initializing $dw_1, dw_2 = 0$ we can get rid of it and make $dw$ a vector

$dw = np.zeros((n_x, 1))$

→ for i=1 to m:
$$z^{(i)} = w^T x^{(i)} + b$$
$$a^{(i)} = \sigma(z^{(i)})$$

we want to eliminate the 2nd for loop

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

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

for i=1 to n_x
dw_j gets updated

$\boxed{\begin{array}{l} dw_1 += x_1^{(i)} dz^{(i)} \\ dw_2 += x_2^{(i)} dz^{(i)} \end{array}}$  $n_z = 2$ instead of doing this we can do

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

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

$J = J/m,$ $\boxed{dw_1 = dw_1/m, \; dw_2 = dw_2/m, \; db = db/m}$

we can have $dw /= m$

## Vectorizing Logistic Regression

Forward propagation in non-vectorized form
$$z^{(1)} = w^T x^{(1)} + b \qquad z^{(2)} = w^T x^{(2)} + b \qquad z^{(3)} = w^T x^{(3)} + b$$
$$a^{(1)} = \sigma(z^{(1)}) \qquad a^{(2)} = \sigma(z^{(2)}) \qquad a^{(3)} = \sigma(z^{(3)})$$

→ In order to carry the 4 propagation steps i.e. to compute preds on m train ex, there is a way to do it without using a single for loop

matrix of train i/p $X = \begin{bmatrix} x^{(1)} & x^{(2)} & \cdots & x^{(m)} \\ | & | & & | \end{bmatrix}$ $(n_x, m)$ $\mathbb{R}^{n_x, m}$

First we will compute $z^{(1)}, z^{(2)}, \ldots$ in one step.
- Construct $(1, m)$ matrix (row vector) while computing, $z^{(1)}, z^{(2)}, \ldots$
$$Z = [z^{(1)} \; z^{(2)} \ldots z^{(m)}] = w^T X + \underbrace{[b \; b \ldots b]}_{1 \times m}$$

$$w^T \begin{bmatrix} \frac{1}{x}^{(1)} & x^{(2)} & \cdots & \frac{1}{x}^{(m)} \\ 1 & 1 & & 1 \end{bmatrix}$$

$$Z = [\underbrace{w^T x^{(1)} + b}_{z^{(1)}} \underbrace{w^T x^{(2)} + b}_{z^{(2)}} \ldots w^T z^{(m)} + b]$$

$$\underset{1 \times m}{}$$

$$Z = np.dot(W.T, X) + b$$

$$\underset{(1,1)}{} \rightarrow \text{note that } b \text{ is a real}$$

Next we want to compute

$$A = [a^{(1)} \; a^{(2)} \; \ldots \; a^{(m)}]$$

no, not a vector. When
we add real no to vector
~~Python~~ automatically takes
real no $b$ and expands it
to $(1 \times m)$ row vector. This
operation is called Broadcasting
in Python

## Vectorizing Logistic Regressions Gradient Descent

$$dz^{(1)} = a^{(1)} - y^{(1)}, \quad dz^{(2)} = a^{(2)} - y^{(2)} \ldots$$

$$\boxed{dZ = [\underset{1 \times m}{dz^{(1)} \; dz^{(2)} \ldots dz^{(m)}}]} \quad \leftarrow \text{we can compute this at the same}$$
$$\text{time with 1 line of code}$$

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

$$dZ = A - Y = [a^{(1)} - y^{(1)} \; a^{(2)} - y^{(2)} \ldots]$$

from the prev implementation we got rid of 1 for loop but we still
~~from~~ have,

$$
\begin{array}{ll}
\text{for } j+1 \text{ to } n_x & \\
dw = 0 & db = 0 \\
dw \mathrel{+}= x^{(1)} dz^{(1)} & db \mathrel{+}= dz^{(1)} \\
dw \mathrel{+}= x^{(2)} dz^{(2)} & db \mathrel{+}= dz^{(m)} \\
\vdots & \\
dw / = m & db / = m
\end{array}
\left.\begin{array}{l} \\ \\ \\ \\ \\ \end{array}\right\} \text{nonvectorized}
$$

vectorized:

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

$$\boxed{\begin{aligned} dw &= \frac{1}{m} X \, dz^T \\ &= \frac{1}{m} \begin{bmatrix} x^{(1)} \ldots x^{(m)} \\ \vdots \end{bmatrix} \begin{bmatrix} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{bmatrix} \\ &= \frac{1}{m} \underset{(n \times 1)}{[x^{(1)} dz^{(1)} + \ldots + x^{(m)} dz^{(m)}]} \end{aligned}}$$

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

for $i=1$ to $m$:

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

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

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

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

$\quad dw_1 += x_1^{(i)} dz^{(i)} \left.\right\} dw^{(i)} += x^{(i)} * dz^{(i)}$

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

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

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

Right column:

for $i$ in range (1000):

$Z = w^T x + b$

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

$A = \sigma(Z)$

$dZ = A - y$

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

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

$w := w - \alpha \, dw$

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

---

## Broadcasting in Python

ex, Calories of carbs, proteins, fats in 100g of different food:

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

$\rightarrow A_{(3,4)}$  (3,4)

Calculate the percentage of calories from Carbs, proteins and fats
• if take apples. Can you do this without a for loop?

total no of calories = $56 + 1.2 + 1.8 = 59$

% of cal from carb = $(56/59) \times 100 = 94.91\%$

$\quad$ proteins = $(1.2/59) \times 100 = 2.03\%$

$\quad$ fats = $(1.8/59) \times 100 = 3.05\%$

Soln:- We will compute the column sum (we get total no of calorie from each food

- Next will divide each of 4 cols by their corresponding sum

Ans to the question

col-sum = A. sum(axis = 0)   # axis = 0, python sums vertically down
$\qquad\qquad$ 1, sums horizontally

percentage = (A / col-sum.reshape(1,4)) * 100

$\qquad\quad$ (3,4) $\uparrow$ $\qquad$ (1,4)

$\qquad\qquad$ Broadcasting

$\qquad\quad$ = (3,4) / (1,4)

more examples,

ex 1,

$$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + 100$$

Python will take 100 and auto expand it to a (4,1) vector

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

This type of broadcasting works for both col and row vectors.

ex 2,

$$\underset{(m,n)}{\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}} + \underset{(1,n)}{[100 \quad 200 \quad 300]}$$

Python will copy matrix m times and turns it into (m, n) matrix

$$\underset{(m,n) \; (2,3)}{\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}} + \underset{(1,n)\leadsto (m,n) \; (2,3)}{\begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix}} = \begin{bmatrix} 101 & 202 & 303 \\ 104 & 205 & 306 \end{bmatrix}$$

$$\underset{(m,n)}{\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}} + \begin{bmatrix} 100 \\ 200 \end{bmatrix}$$

(m, 1) ~ Python copies n times horizontally

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

General Principle to Broadcasting:

$$(m,n) \quad \begin{array}{c} + \\ - \\ * \\ / \end{array} \quad \underset{\text{vec}}{\overset{\text{Vector}}{(1,n)}} \leadsto (m,n)$$

matrix

$$/ \overset{\text{col}}{\underset{\text{vec}}{(m,1)}} \leadsto (m,n)$$

$$(m,1) + \text{Real no } \mathbb{R}$$

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

$$[1 \quad 2 \quad 3] + 100 = [101 \quad 102 \quad 103]$$

# A note on numpy vectors

Note :

generates 5 random no with $\mu = 0, \sigma = 1$

When you use np.random.randn(5), → Don't use
- we get a rank 1 vector and its shape is (5,)
   ↳ This is not a row or column vector
   ↳ a.transpose will be the same as a
   ↳ dot product between a and a.transpose gives us a no . instead of matrix

Use

$a = np.random.randn(5,1)$ → column vector
$a = np.random.randn(1,5)$ → row vector

when we are not sure about dimensions of one of the vectors use assert statement.

assert(a.shape == (5,1))

Important functions in numpy
~~Explanation of logistic Regression cost function~~

☞ reshape : np.reshape((m.n)) reshapes an array
norm : ~~np.linalg~~ =norm Normalizes an array either row-wise or column wise . Normalizing data leads to better performance because gradient
~~∂x ⟶~~ = descent converges faster after normalization.
- While normalizing $\frac{x}{||x||}$ divide each row of vector by norm

$$x = \begin{bmatrix} 0 & 3 & 4 \\ 2 & 6 & 4 \end{bmatrix}$$

then $||x|| = np.linalg.norm(x, axis = 1, keepdims = True) = $

$$||x|| = \begin{bmatrix} 5 \\ \sqrt{56} \end{bmatrix}$$

$$x\_norm = \frac{x}{||x||} = \begin{bmatrix} 0 & 3/5 & 4/5 \\ 2/\sqrt{56} & 6/\sqrt{56} & 4/\sqrt{56} \end{bmatrix}$$

Note we can divide matrices by diff sizes
- Broadcasting

keepdims = True will broadcast correctly against original .
axis = 1, we will get the norm row-wise
axis = 0, we will get col wise norm
ord : type of norm (degree for root)

np.absolute : computes abs value for each element

np.sum: computes element wise sum

during the exercise

Note: when they ask to flatten array of shape $(209, \overset{h}{64}, \overset{w}{64}, \overset{c}{3})$

to flatten this we are using:

$X = X.reshape(X.shape[0], -1).T$

transpose will change shape to $(12288, 209)$

tells numpy to automatically calculate 2nd dimension which is :

$64 \times 64 \times 3 = 12288$

np.zeros(shape = (m,n)) : you get (m,n) arr of zeros

Note while vectorizing gradient descent for logistic regression.

we can write $\frac{dJ}{dw} = \frac{1}{m} X (A-y)^T$ as:
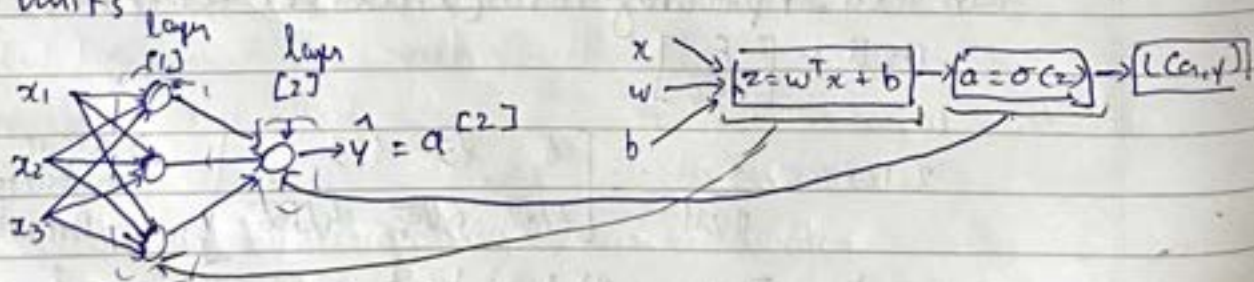
$dz = A - y$

$dw = (1/m) * np.dot(x, dz.T))$

## Week-3

### Shallow Neural Networks

### Neural Nets Overview :

What is a Neural Net ?

- You can form an NN by stacking a lot of little sigmoid units



layer
[1]     layer
        [2]

$x_1$

$x_2$

$x_3$

$\hat{y} = a^{[2]}$

$x$
$w$ $\rightarrow$ $z = w^T x + b$ $\rightarrow$ $a = \sigma(z)$ $\rightarrow$ $L(a, y)$
$b$

i/p feature   $x$

parameters  $w^{[1]}$
             $b^{[1]}$
$z^{[1]} = W^{[1]} x + b^{[1]}$ $\rightarrow$ $a^{[1]} = \sigma(z^{[1]})$ $\rightarrow$ $z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$

$L(a^{[2]}, y)$

## Neural Network Representation



$a^{[0]} = 2$    $a^{[1]} \rightarrow$    $w^{[1]}, b^{[1]}$
$(4,3)$  $(4,1)$

$x_1$

$x_2$

$x_3$

$w^{[2]}_{(1,4)}, b^{[2]}_{(1,1)}$

$\hat{y} = a^{[2]}$

input layer     hidden layer     output layer

in logistic regression we used $\hat{y} = a$ because we only had one output. In NN we will use ~~that~~ superscript square brackets to distinguish the layer it came from

Note: This is a 2 layer net. we don't count i/p layer, so we only have hidden and output ~~out~~ layer.

representing the hidden layer

$$a^{[1]} = \begin{bmatrix} a^{[1]}_1 \\ a^{[1]}_2 \\ \vdots \\ a^{[1]}_4 \end{bmatrix} \leftarrow \text{ in this we got 4 hidden units in the first layer}$$

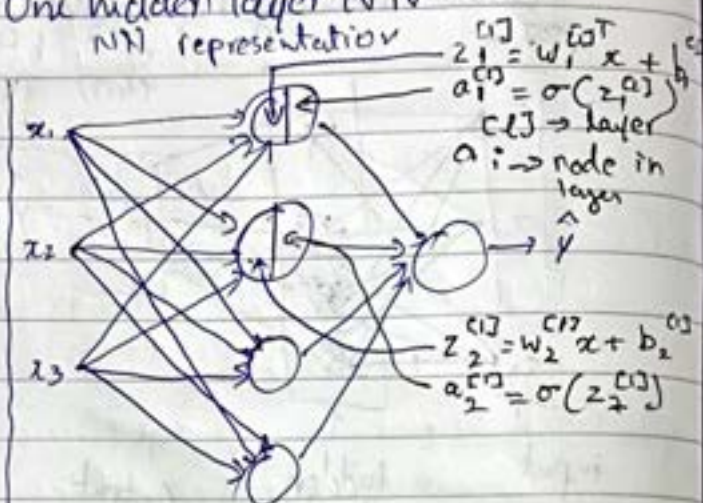→ Hidden and output layers have parameters associated with them. Hidden layer is associated with params $w, b$. @
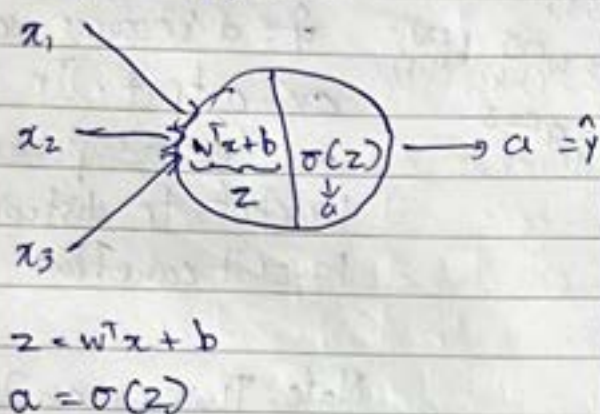$a^{[1]}$ has $w^{[1]}_{(4,3)}, b^{[1]}_{(4,1)}$, w & is a $(4,3)$ matrix and b is $(4,1)$ vector

$w^{[1]} \rightarrow (4,3)$ → We have 3 i/p features in input layer
↳ we have 4 nodes in hidden layer

→ Output layer $a^{[2]}$ is associated with parameters $w^{[2]}_{(1,4)}, b^{[2]}_{(1,1)}$
$w^{[2]} \rightarrow (1,4)$ dims $(1,4)$ because we got 1 node and 4 i/p features from prev layer (hidden)

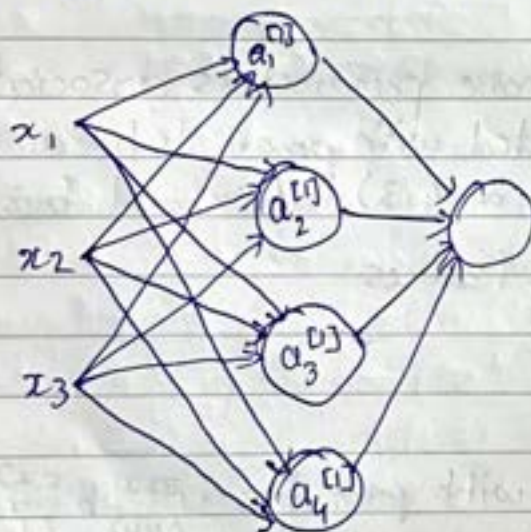Computing an NN's output : One hidden layer NN

Logistic Regression

NN representation

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

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

$[1] \to$ layer

$i \to$ node in layer



$z = w^T x + b$

$a = \sigma(z)$

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

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

—In each node there are 2 steps of computation:

hidden layer

$z = w^T x + b$

$a = \sigma(z)$

Lets draw the NN with the notation



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

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

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

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

We will vectorize these eqns because doing it with a for loop is inefficient

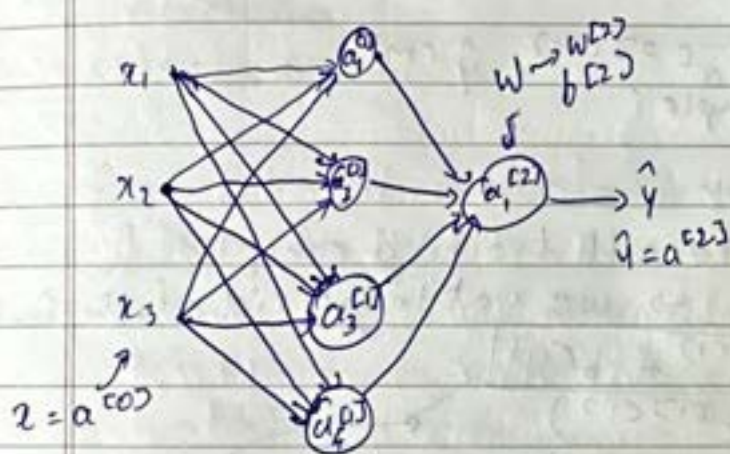→ We will start by computing z as a vector

→ Stack the w's into matrix

$$z^{[1]} = \begin{bmatrix} -W_1^{[1]T} \rightarrow \\ -W_2^{[1]T} - \\ -W_3^{[1]T} - \\ -W_4^{[1]T} - \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} = \begin{bmatrix} W_1^{[1]T} x + b_1^{[1]} \\ W_2^{[1]T} x + b_2^{[1]} \\ W_3^{[1]T} x + b_3^{[1]} \\ W_4^{[1]T} x + b_4^{[1]} \end{bmatrix} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}$$

$W^{[1]}$ $(4,3)$ $(3,1)$ $b^{[1]}$ $(4,1)$ $(4,1)$ $(4,1)$

→ We will stack a's toghter

$$a = \begin{bmatrix} a_1^{[1]} \\ \vdots \\ a_1^{[4]} \end{bmatrix} \rightarrow \sigma(z^{[1]})$$

Clear representation of each layer and the dimensions of each layer



$W \rightarrow \begin{matrix} W^{[2]} \\ b^{[2]} \end{matrix}$

$\hat{y} = a^{[2]}$

$x = a^{[0]}$

$z = a^{[0]}$

Given input $x$

$z_1^{[1]} = W^{[1]} a^{[0]} + b^{[1]}$
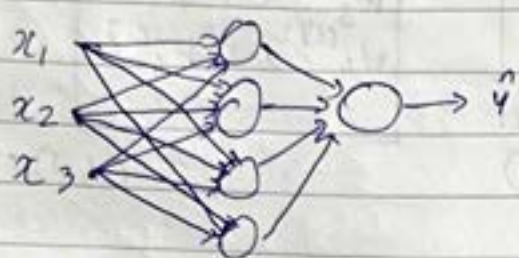$(4,1)$ $(4,3)(3,1)$ $(4,1)$

$a^{[1]} = \sigma(z^{[1]})$
$(4,1)$ $(4,1)$

$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$
$(1,1)$ $(1,4)(4,1)$ $(1,1)$

→ we can think of the last unit as being analogous to logistic regression

# Vectorizing across multiple examples

Aim: compute the outputs for all examples in NN at the same time



$$\begin{cases} z^{[1]} = W^{[1]} x + b^{[1]} \\ a^{[1]} = \sigma(z^{[1]}) \\ z^{[2]} = W^{[2]} a^{[1]} + b^{[2]} \\ a^{[2]} = \sigma(z^{[2]}) \end{cases}$$

* These eqns tell how given ~~these eqn~~ an input feature $x$. You can use them to generate $a2 = \hat{y}$ hat for a single train ex.

$$x \longrightarrow a^{[2]} = \hat{y}$$

if we have $m$ train ex we need to repeat the process

$$x^{(1)} \longrightarrow a^{[2](1)} = \hat{y}^{(1)}$$
$$x^{(2)} \longrightarrow a^{[2](2)} = \hat{y}^{(2)}$$

.
.
.

$$x^{(m)} \longrightarrow a^{[2](m)} = \hat{y}^{(m)}$$

$a^{[2](i)} \longrightarrow$ example $i$
$\qquad \searrow$ layer 2

Computing predictions for all training examples

for $i = 1$ to $m$: $\rightarrow$ we want to get rid of the for loop

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

X is matrix of train ex stacked in columns

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

$(n_x, m)$

We can do it similarly for $Z, a$ as we did for $X$.

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

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

horizontal index corresponds to diff train ex, when we sweep left to right through train ex.

vertical index corresponds to diff roles in NN, as we scan down we index into hidden unit no.

$A^{[0]} \; x^{(i)} = a^{[0]}$
To simplify

$$Z^{[1]} = W^{[1]} \boxed{X} + b^{[1]} \quad z^{(i)}, a^{[0](i)}$$

$\longrightarrow W^{[1]} A^{[0]} + b^{[1]} \longleftarrow$ train ex

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

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

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

This notation is used for matrix $X, Z$ as well

Explanation for vectorized implementation
Justification for vectorized implementation:
Forward propagation calculation for a few ex

$\cancel{Z^{[1](1)} = W^{[1](1)} + b^{[1]}}$

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

$W^{[1]} = \begin{bmatrix} — \\ — \\ — \end{bmatrix} \qquad W^{[1]} x^{(1)} = \begin{bmatrix} . \\ . \\ . \end{bmatrix} \qquad W^{[1]} x^{(2)} = \begin{bmatrix} : \\ : \end{bmatrix} \qquad W^{[1]} x^{(3)} = \begin{bmatrix} : \\ : \end{bmatrix}$

col vector
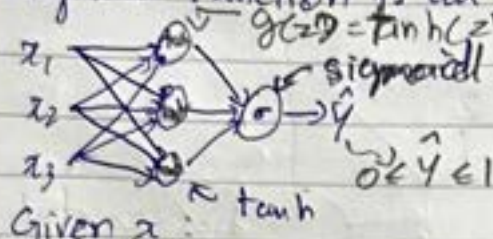
$$W^{[1]} \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & x^{(3)} \\ | & | & | \end{bmatrix} = \begin{bmatrix} : & : & : \\ . & . & . \\ : & : & : \end{bmatrix} = \begin{bmatrix} | & | & | \\ Z^{[1](1)} & Z^{[1](2)} & Z^{[1](3)} \\ | & | & | \end{bmatrix} = Z^{[1]}$$

$X$

$W^{[1]} x^{(i)} = Z^{[1](i)}$

$Z^{[1]} = W^{[1]} X + b^{[1]}$ — this line allows you vectorize all m examples at the

## Activation Functions

-> Sigmoid function is an activation function

$g(z) = \tanh(z^{[1]})$
sigmoid



$x_1$
$x_2$
$x_3$

tanh

$0 \leq \hat{y} \leq 1$

sigmoid
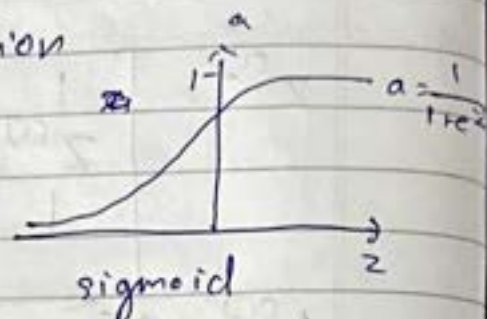
Given x:

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

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

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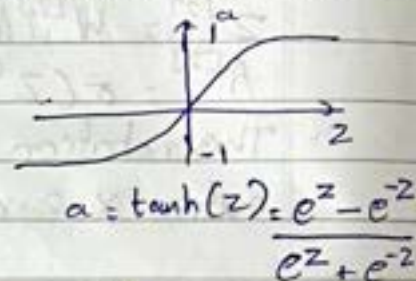
$a^{[2]} = \sigma(z^{[2]}) \quad g(z^{[2]})$

-> We can have a different function $g(z)$ where $g$ could be non-linear function that may not be a sigmoid function.
A tangent or hyperbolic tangent (tanh) almost always works better than sigmoid.

-> tanh goes between $-1$ and $+1$.

$a = \tanh(z) = \dfrac{e^z - e^{-2}}{e^z + e^{-2}}$



$a = \tanh(z) = \dfrac{e^z - e^{-2}}{e^z + e^{-2}}$

-> If hidden units were tanh it ~~almost~~ its better than sigmoid because with values between $-1$ and $1$, the mean of the activation that come out of your hidden ~~units~~ layers, are closer to $0$ mean

-> When training a learning algorithm, you might center the data and have your data have $0$ mean using tanh instead of sigmoid

-> The effect of centering your data so that the mean of your data is close to $0$ rather than maybe $0.5$; this makes learning for the next layer a little easy.

-> tanh is almost strictly superior.

-> One exception is for the output layer where it will be sigmoid because the output should bet $0 \leq \hat{y} \leq 1$, while binary classification

-> We can have tanh activation fn for hidden units and a sigmoid fn for output layer.

-> Sometimes activation functions can be different for different layers

$$g^{[1]}(z^{[1]}) = \tanh(z^{[1]})$$
$$g^{[2]}(z^{[2]}) = \sigma(z^{[2]})$$

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

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

→ Downside of both sigmoid and tanh is $z^{[1]}$, $z$ is very small or very large if, the gradient of derivative of the slope is very small. If $z$ is very small it ends up, being close to 0 and it slows down gradient descent. Another

→ Another popular choice is ReLU (Rectified linear unit)
- The derivative is 1 as long as slope $z$ is +ve       $a = \max(0,z)$
- Derivative is 0. When $z$ is -ve                      $z$ +ve slope
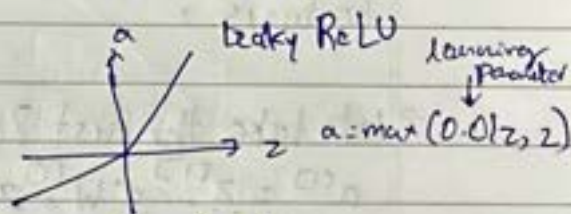                                                        $z$ -ve slope → $z$ ReLU

Rule of thumb for choosing activations func.
- if output is bet 0 and 1, we are using Binary classification then sigmoid is the choice for o/p layer.
- For all other units ReLU is the default choice. Even for hidden unit

- Disadvantage of ReLU: Derivative is equal to zero when $z$ is -ve

Leaky ReLU:
- When $z$ is -ve, instead of being 0 it takes a small gradient for -ve $z$ values       leaky ReLU    learning parameter
                                                        $a = \max(0.01z, z)$
Note: Using either is fine R. Usually we ca Use ReLU

Advantage of ReLU and Leaky ReLU: → A lot of space for $z$, the slope derivative of activation fn. is very different from 0.
→ Using ReLU the NN learns faster than tanh or sigmoid. The main reason is that there is less effect of slope of fn going down to 0 that slows down learning

Note: if your not sure about which activation fn to use try them all and evaluate them on a holdout set or development set.      validation

## Why do we need non-linear activation func?
## Activation functions



Given x

We can get rid of g and set $a^{[1]} = z^{[1]}$

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

$$a^{[1]} = g^{[1]}(z^{[1]}) = z^{[1]}$$ (We can say $g(z) = z$. This is called the
$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$ linear activation fn / identity activation
$$a^{[2]} = g^{[2]}(z^{[2]}) = z^{[2]}$$ fn - since it outputs whatever was the
input)

- What if $a^{[2]} = z^{[2]}$ ?

→ The model is computing y or $\hat{y}$ as a linear function of the input features

→ Lets take the first 2 eqns:
$$a^{[1]} = z^{[1]} = W^{[1]}x + b^{[1]} \quad -- ①$$

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

Subs ① in ②

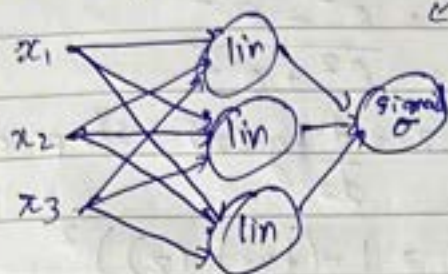$$a^{[2]} = W^{[2]}(\underbrace{W^{[1]}x + b^{[1]}}_{a^{[1]}}) + b^{[2]}$$

$$= \underbrace{(W^{[2]}W^{[1]})}_{W'}x + \underbrace{(W^{[2]}b^{[1]} + b^{[1]})}_{b'}$$

$$= W'x + b'$$

→ If we use a linear activation fn or alternatively if we don't have an activation fn, then no matter how many layers the NN has all its doing is computing a linear activation function. So we might as well not have any

hidden layers :
→ Lets say we had :



← ⟶ This model is no more expressive
⟶ $\hat{y}$ than a standard logistic regression without
any hidden layer.
→ Imp: Linear hidden layer is more-less
useless because the composition of 2 linear
functions is itself a linear function.

→ Unless there is non-linearity there is nothing interesting we are
computing as you go deeper in the network.
→ We can use linear activation only if we are doing regression
problem. We can use this in the output layer

Derivatives of Activation functions
Sigmoid activation function



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

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

$\frac{d}{dz}(g(z)) = $ slope of $g(z)$ at $z$

$= \frac{1}{1+e^{-z}}\left(1 - \frac{1}{1+e^{-z}}\right) \quad - ②$

subs ① in ②

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

$= a(1-a)$

if $z = 10, g(z) \approx 1$

② $\frac{d}{dz}g(z) = 1(1-1) = 0$

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

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

if $z = 0, g(z) = 1/2$

$\frac{d}{dz}(g(z)) = \frac{1}{2} \ast \frac{1}{2}\left(1 - \frac{1}{2}\right) = \frac{1}{4}$

## tanh activation function :



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

$$g'(z) = \frac{d}{dz} g(z) = \text{slope of } g(z) \text{ at } z$$

$$= 1 - (\tanh(z))^2$$

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

if $z = 10$, $\tanh(z) \approx 1$

$\quad g'(z) \approx 0$

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

$\quad g'(z) \approx 0$

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

$\quad g'(z) \approx 1$

## ReLU and Leaky ReLU :



ReLU

$$g(z) = \max(0, z)$$

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

Leaky ReLU

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

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

## Gradient Descent for Neural Networks

Implementing Gradient Descent for a layer

Parameters: $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}$

$(n^{[1]}, n^{[0]}) (n^{[1]}, 1) (n^{[2]}, n^{[1]}) (n^{[2]}, 1)$

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

input features   hidden units   output unit

Cost function: $J(W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^{n} L(\hat{y}, y)$

Gradient descent:

Repeat {

Compute predictions $(\hat{y}^{(i)}, i=1....m)$

$dw^{[1]} = \frac{dJ}{dw^{[1]}}, db^{[1]} = \frac{dJ}{db^{[1]}}, ...$

~~$W^{[1]} - \frac{dJ}{}$~~

$W^{[1]} = \alpha \frac{dJ}{dw^{[1]}}, b^{[1]} -= \alpha \frac{dJ}{db^{[1]}}$

$W^{[2]} -= \alpha \frac{dJ}{dw^{[2]}}, b^{[2]} -= \alpha \frac{dJ}{db^{[2]}}$

Formulas to compute derivatives:

forward propagation

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

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

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

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

Back Prop

~~Backword~~

$dz^{[2]} = A^{[2]} - y \qquad y = [y^{(1)} y^{(2)} .. y^{(n)}]$

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

sum horizontal

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

Prevents python from outputting one rank arrays where dimension are $(n,)$. Setting it to True will output $(n,1)$

$dz^{[1]} = W^{[2]T} dz^{[2]} * g^{[1]'}(Z^{[1]}) \quad (n^{[1]}, m)$

$\underbrace{\qquad}_{(n^{[1]}, m)}$   ↑ elementwise prodt.

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

$db^{[1]} = \frac{1}{m}$ np-sum ~~( ... keepdims=True)~~

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

~~(....)~~

Backpropagation intuition (Optional)
Forward pass on computation graph for logistic regression
Computing gradients
~~logistic~~ Logistic regression

$x$

$w \Longrightarrow$ $\boxed{z = w^T x + b} \rightleftarrows \boxed{a = \sigma(z)} \rightleftarrows \boxed{L(a,y)}$

$b$

$dz = a - y$

$\boxed{dz = dz \cdot g'(z)}$

$dw = dz \cdot x$

$g(z) = \sigma(z)$

$db = dz$

$\boxed{\dfrac{dL}{dz} = \dfrac{dL}{da} \cdot \dfrac{da}{dz}}$

$dz = da$

$da = -\dfrac{y}{a} + \dfrac{1-y}{1-a}$

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

$\dfrac{dL(a,y)}{da} = -\dfrac{y}{a} + \dfrac{1-y}{1-a}$

$\boxed{\dfrac{da}{dz} = \dfrac{d}{dz} g(z) = g'(z)}$

### Neural network gradients:

$dw^{[2]}$ $\dfrac{dw}{dw^{[2]}} = \dfrac{dL}{dz^{[2]}} \cdot \dfrac{dz}{dw}$

$W^{[2]}$ $dw^{[2]}$
$db^{[2]}$

$b^{[2]}$ $dw^{[2]} = dz^{[2]} a^{[1]\top}$ $= (a - y)x$

$db^{[2]} = dz^{[2]}$ $dw^{[2]} = dz^{[2]} a^{[1]\top}$

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

$x$

$w^{[1]} \Rightarrow \boxed{z^{[1]} = W^{[1]}x + b^{[1]}} \to \boxed{a^{[1]} = \sigma(z^{[1]})} \to \boxed{z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}} \to \boxed{a^{[2]} = \sigma(z^{[2]})} \rightleftarrows \boxed{L(a^{[2]},y)}$

$b^{[1]}$

$da^{[1]}$

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

$da^{[2]} = -\dfrac{y}{a} + \dfrac{1-y}{1-a}$

$dz^{[1]} = W^{[2]\top} dz^{[2]} * g^{[1]'}(z^{[1]})$

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

$z^{[1]}, dz^{[1]} \quad (n^{[2]}, 1) - (1,1)$

$z^{[1]}, dz^{[1]} \quad (n^{[1]}, 1)$

$dz^{[1]} = w^{[2]\top} dz^{[2]} * g^{[1]'}(z^{[1]})$

$(n^{[1]}, 1) \ (n^{[1]}, n^{[2]})(n^{[2]}, 1) * (n^{[1]}, 1)$

$dz^{[2]} \dfrac{dL}{dz^{[2]}} = \dfrac{dL}{da} \cdot \dfrac{da}{dz}$

$= \left(-\dfrac{y}{a} + \dfrac{1-y}{1-a}\right) a(1-a)$

$= a(1-a) \cdot -\dfrac{y}{a} + \dfrac{a(1-a)(1-y)}{(1-a)}$

$= -y(1-a) + a(1-y)$

$= -y + ay + a - ay$

$da^{[2]} \dfrac{dL}{da^{[2]}} = -(y \log a + (1-y) \log(1-a))$

$= -\dfrac{y}{a} + \dfrac{(1-y)}{1-a}$

$dw^{[1]} = dz^{[1]} x^\top$

$db^{[1]} = dz^{[1]}$ $a^{[0]\top}$

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

$dw^{[2]} = dz^{[2]} a^{[1]\top}$ dw is row vector thats we ou
using $a^{[1]}$.

Summary of gradient descent:

Normal

$$dz^{[2]} = a^{[2]} - y$$
$$dW^{[2]} = dz^{[2]} a^{[1]T}$$
$$db^{[2]} = dz^{[2]}$$
$$dz^{[1]} = W^{[2]T} dz^{[2]} \cdot g^{[1]'}(z^{[1]})$$
$$dW^{[1]} = dz^{[1]} x^T$$
$$db^{[1]} = dz^{[1]}$$

Vectorized implementation

$$dz^{[2]} = A^{[2]} - y$$
$$dW^{[2]} = \frac{1}{m} dz^{[2]} A^{[1]T}$$
$$db^{[2]} = \frac{1}{m} np.sum(dz^{[2]}, axis=1, keepdims=True)$$
$$dz^{[1]} = W^{[2]T} dz^{[2]} \cdot g^{[1]'}(z^{[1]})$$
$$dW^{[1]} = \frac{1}{m} dz^{[1]} x^T$$
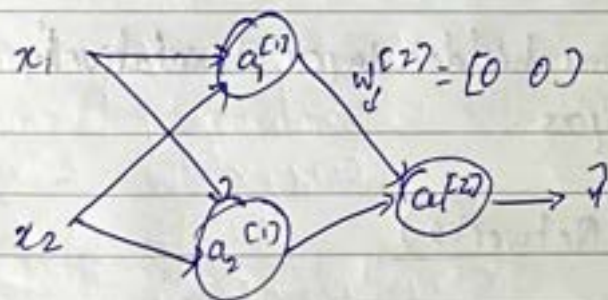$$db^{[1]} = \frac{1}{m} np.sum(dz^{[1]}, axis=1, keepdims=True)$$

## Random Initialization

→ In logistic regression we could ~~initialut~~ initialize the weights to ~~parameto~~ to 0. In NN's initializing weight parameters to 0 wo'nt work.

⊗

What happens if we initialize the weights to zero?



$$w_i^{[2]} = [0 \ 0]$$

→ Initializing bias terms to 0 is fine. But not ok to do it for weights

$$n^{[0]} = 2 \qquad n^{[1]} = 2$$

$$W^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \qquad b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$a_1^{[1]} = a_2^{[1]}, \qquad dz_1^{[1]} = dz_2^{[1]}$$

Both activations in the same

Hidden units are completely identical when we compute
$$dz_1^{[1]} = dz_2^{[1]} \quad (\text{outgoing wts are symmetrical})$$
$$W^{[2]} = [0 \ 0]$$

After a few iterations of training the hidden units compute the same function.

$$dw = \begin{bmatrix} u & v \\ u & v \end{bmatrix} \quad \text{every row takes on the same value. so we perform a weight update.}$$

$$w = w^{[1]} - \alpha dw$$
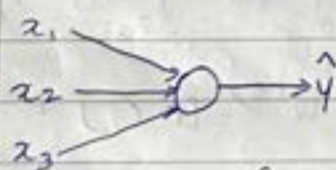
- After every iteration 1st row equals 2nd row
- By the proof of induction if we initialized all in the values w to 0 then both hidden units start off computing the same function and both hidden units have the same influence on output unit.
- 2 hidden units are still symmetric
- By induction no matter how ~~many times we compute~~ long we train the N N, the hidden units compute the same function.
- In this case we don't need more than one hidden unit.
- If we have a large N.N. with 3 features and many ~~weights~~ hidden units, if initialized ~~to~~ weights to 0, all hidden units are symmetric. No matter how long we run gradient descent they compute the same function

→ We want different hidden units to compute different ~~activation~~ functions. Solution: Initialize parameters randomly and multiply with small no.

$$W^{[1]} = np.random.randn((2,2)) * 0.01 \quad \text{we small no because we wants weights to have small random values}$$
$$b^{[1]} = np.zeros((2,1)) \quad (b \text{ doesn't have the symmetry breaking problem})$$

→ When we train with just one hidden layer its relatively shallow N N without too many hidden layers.

Week - 4   (Deep Neural Networks)
Deep L layer NN
What is a deep N N ?
ex lets take logistic regression, 1 hidden layer, 5 hidden layers



$x_1$
$x_2$ →$\hat{y}$
$x_3$

logistic ("shallow")
regression

1 layer NN

$x_1$
$x_2$ →0→$\hat{y}$
$x_3$

1 hidden
layer
(2 layer N N)
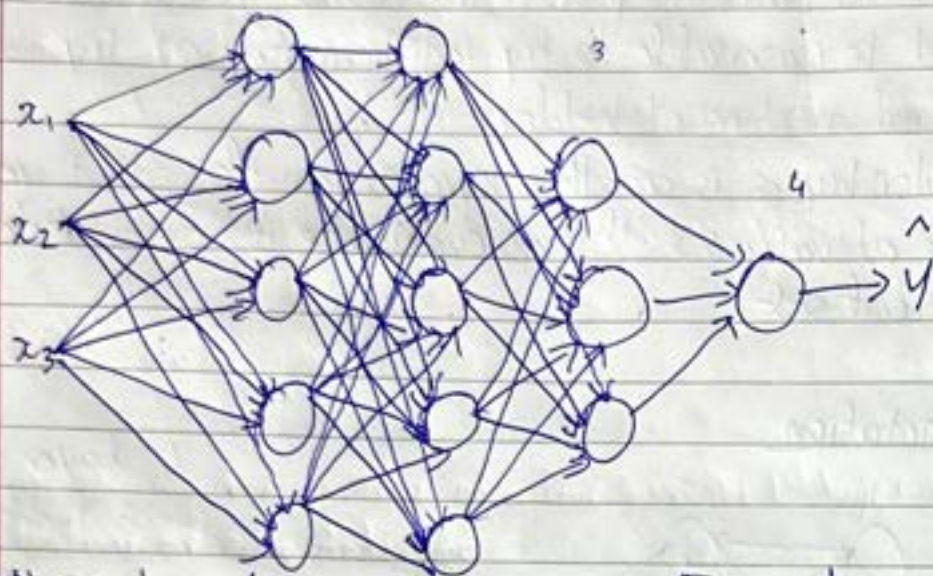we dont include
i/p layer

$x_1$
$x_2$
$x_3$

5 hidden layers
(Deep)

→ Very deep NN can learn, ~~that~~ functions that shallower models are unable to.

→ For any problem its hard to predict in advance how deep an NN should be. So it would be reasonable to try logistic regression, try 1 and then 2 hidden layer; ~~and view no of hidden~~

→ The no of hidden layers is another hyperparameter that you could try a variety of values with and evaluate it on validation data, or on development set.

## Deep NN notation



4 layer NN
3 hidden layers
1 output layer

layer→0

$\hat{y} = a^{[4]}$

$n^{[4]} = 1$

$n^{[3]} = 3$

$n^{[0]} = n_x = 3$
$x = a^{[0]}$

$n^{[1]} = 5 \quad n^{[2]} = 5$

$L = 4$ (no of layers)
$n^{[\ell]} = $ no of nodes in layer $\ell$
$a^{[\ell]} = $ activations in layer $\ell$
$a^{[\ell]} = g(z^{[\ell]})$
$w^{[\ell]} = $ weights for $z^{[\ell]}$
$b^{[\ell]} = $ bias

# Forward Propagation in a Deep Network



## Non vectorized

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

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

$$\vdots$$

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

## Forward prop eqns

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

## Vectorized implementation

$$Z^{[\ell]} = W^{[\ell]} \overset{A^{[0]}}{X} + b^{[1]}$$

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

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

$$\hat{y} = g(Z^{[4]}) = A^{[4]}$$
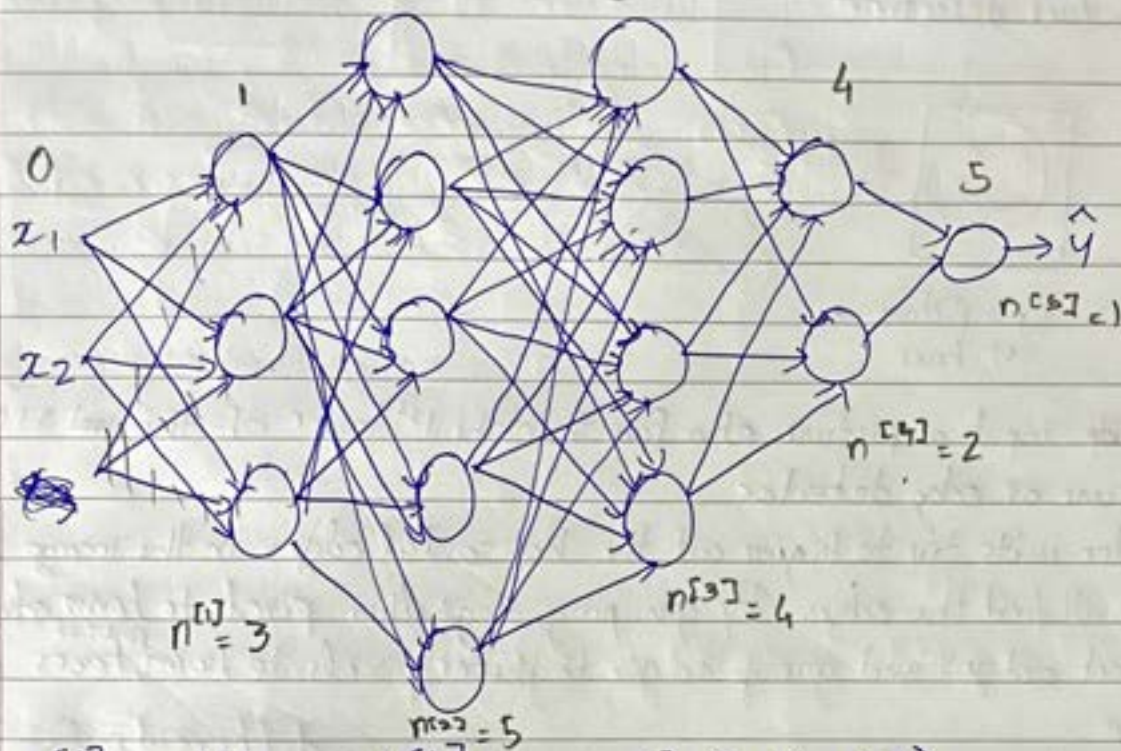
(We will have a for loop)
to compute activations
for $\ell = 1$ to $4$

$$X = A^{[0]}$$

$$\left[ \underset{1}{Z}^{[2](1)} \quad \underset{1}{Z}^{[2](2)} \cdots \underset{1}{Z}^{[2](2)} \right]$$

Getting your dimensions right (matrix)

Parameters are $W^{[l]}$ and $b^{[l]}$

$L = 5$ (ignore $x$)



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

$$z^{[1]} = W^{[1]} x + b^{[1]}$$
$(3,1)$   $(3,2)$  $(2,1)$   $(3,1)$
$(n^{[1]},1)$  $(n^{[1]}, n^{[0]})(n^{[0]},1)$ $(n^{[1]},1)$

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

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

Note: a and z should have dimensions $(n^{[l]}, 1)$

Vectorized implementation:

$$z^{[1]} = W^{[1]} x + b^{[1]}$$
$(n^{[1]},1)$ $(n^{[1]}, n^{[0]})$ $(n^{[0]},1)$  $(n^{[1]},1)$

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

$$Z^{[1]} = W^{[1]} X + b^{[1]}$$
$(n^{[1]},m)$ $(n^{[1]},n^{[0]})$ $(n^{[0]},m)$ $(n^{[1]},1)$ → $(n^{[1]},m)$

$$W^{[1]} = (n^{[1]}, n^{[0]})$$
$$W^{[2]} = (5,3)  (n^{[2]}, n^{[1]})$$
$$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]} \leftarrow (n^{[2]},1)$$
$(5,1)$  $(5,3)$ $(3,1)$   $(5,1)$
$$W^{[3]}: (4,5)$$
$$W^{[4]}: (2,4)$$
$$W^{[5]}: (1,2)$$

$$\boxed{\begin{array}{l} W^{[l]}: (n^{[l]}, n^{[l-1]}) \\ b^{[l]}: (n^{[l]}, 1) \\ dW^{[l]}: (n^{[l]}, n^{[l-1]}) \\ db^{[l]}: (n^{[l]}, 1) \end{array}}$$

$$z^{[l]}, a^{[l]}: (n^{[l]}, 1)$$
$$Z^{[l]}, A^{[l]}: (n^{[l]}, m)$$
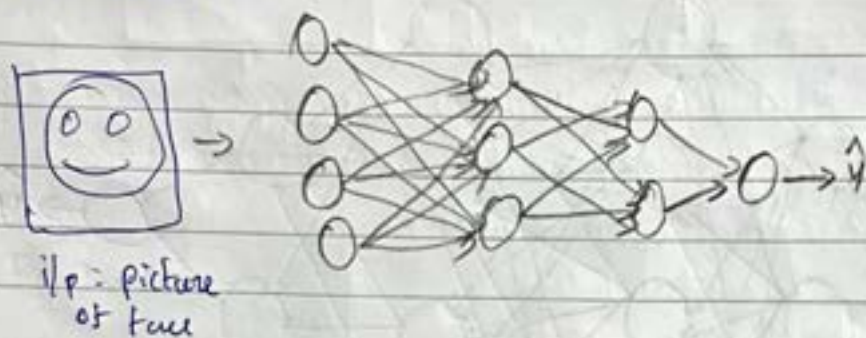$$l = 0, \quad A^{[0]} \cdot X - (n^{[0]}, m)$$
$$dZ^{[l]}, dA^{[l]}: (n^{[l]}, m)$$

Why deep representations ?
Intuition about deep representation
ex, Face detection



i/p : picture
of face

- We ~~it~~ input a picture of a face then the 1st layer of the ~~red~~ NN can be a feature or edge detector
- Hidden units try to figure out the horizontal edges in the image
- We will find the edges by grouping toghether pixels to form edges. It can detect edges and group edges togheter pixels to form parts of faces ex, nose, , eye
- By putting a lot of edges it can detect $\overset{different\ pats}{faces}$ . Then it can try to detec different types of faces.
- Earlier layers are of ~~a~~ NN detect simple functions like edges. and composing them toghethen in the later layers of neuralnet so that it can learn more and more complex function.

ex, Speech $\overset{recognition}{\cancel{Detection}}$ system
- i/p audio clip
- Lower level layers learn ~~the~~ detect the low level audio waveform features such as tone going up or down, pitch, etc
- By composing low level waveforms will learn to detect the basic units of sound.

Summary: Deep NN $\overset{with\ multiple\ hidden\ layers}{will}$ ~~learn lower level~~ simple features ~~and~~ then later might ~~be~~ be able to have earlier layers ~~and team~~ these the loss

Deep NN was having multiple hidden layers, in the earlier layer they learn lower level simple features and in the deeper layers put together the simpler things learned like specific words or phrases or sentences.
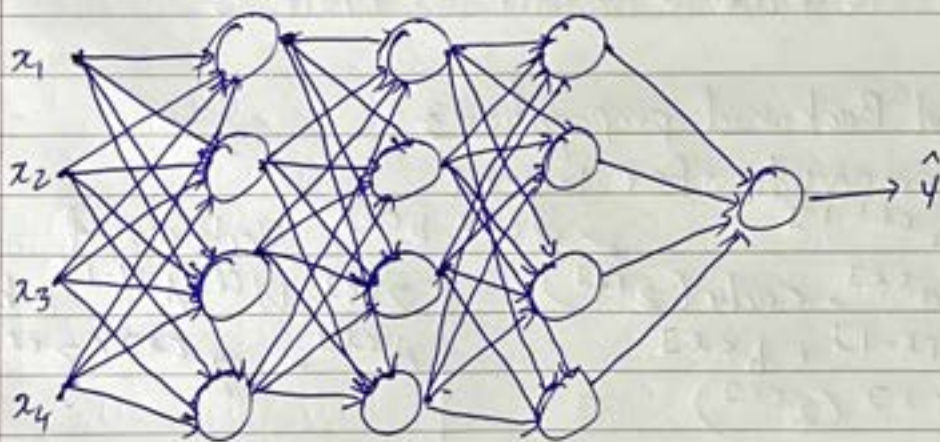
Why deep networks work well?
- This comes circuit theory
circuit theory and deep learning :
⇒ Informally: There are functions you can compute with a small L layer network that shallower networks require exponentially more hidden units to compute.

Building Blocks of Deep Neural Networks :
Forward and Backward functions :

$x_1$

$x_2$ → $\hat{y}$

$x_3$

$x_4$

Layer $l$: $W^{[l]}, b^{[l]}$
Forward: Input $a^{[l-1]}$, output $a^{[l]}$
$Z^{[l]} : W^{[l]} a^{[l-1]} + b^{[l]}$       cache $Z^{[l]}$
$a^{[l]} : g^{[l]} (Z^{[l]})$

Backward: Input $\rightarrow da^{[l]}$, output $da^{[l-1]}$
                  $\rightarrow$ cache $(Z^{[l]})$          $dw^{[l]}$
                                                        $db^{[l]}$

layer $l$

$a^{[l-1]} \rightarrow \boxed{W^{[l]}, b^{[l]}} \rightarrow a^{[l]}$

↓ cache $Z^{[l]}$

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

↓
$dw^{[l]}, db^{[l]}$

Diagram showing: $a^{[0]} \rightarrow [W^{[1]}, b^{[1]}] \xrightarrow{a^{[1]}} [W^{[2]}, b^{[2]}] \xrightarrow{a^{[2]}} [\ ] \rightarrow \cdots \rightarrow [\ ] \rightarrow$

with cache $z^{[1]}$, cache $z^{[1]}$, $z^{[2]}$, $z^{[3]}$, $z^{[l]}$

Backward boxes with $da^{[1]}$, $da^{[2]}$, $da^{[l-1]}$, $\leftarrow da$

$dw^{[1]}$, $db^{[1]}$, $dw^{[2]}$, $db^{[2]}$, $dw^{[3]}$, $db^{[3]}$, $dw^{[l]}$, $db^{[l]}$

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

## Forward and Backward propagation:

Forward propagation for layer $l$:

Input: $a^{[l-1]}$

output: $a^{[l]}$, cache $(z^{[l]})$

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

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

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

Backward propagation for layer $l$:

Input: $da^{[l]}$

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

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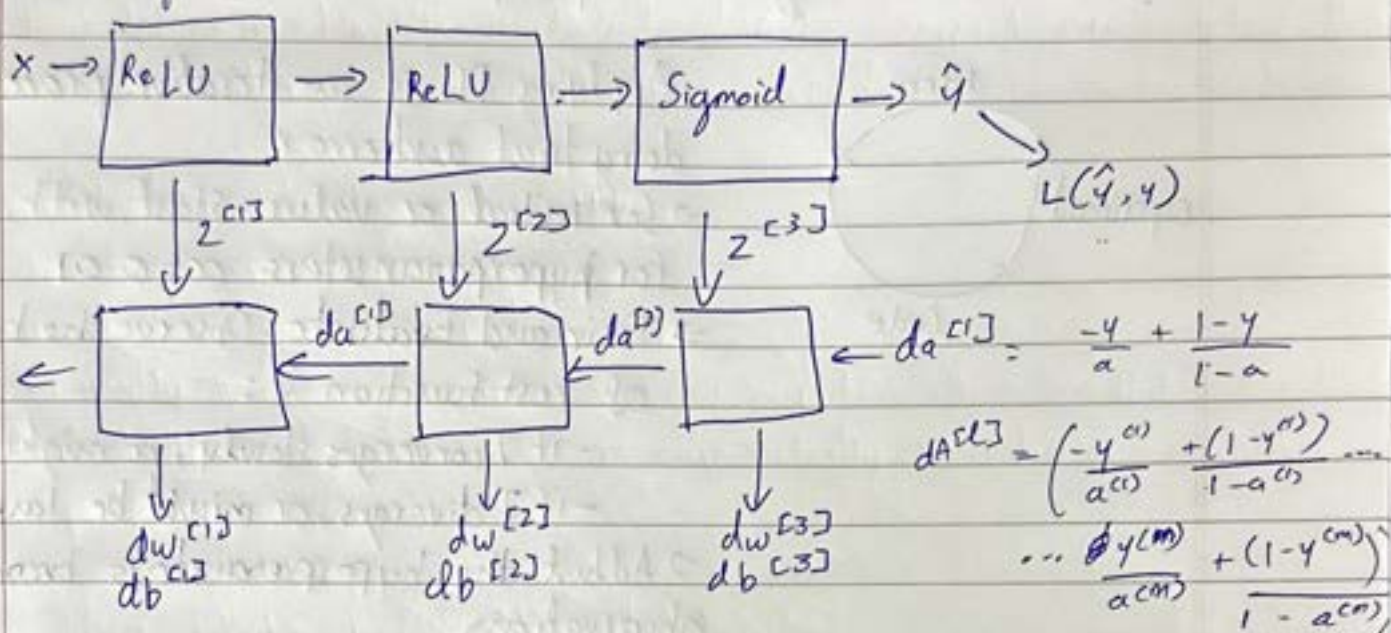$$dz^{[l]} = W^{[l+1]T} dz^{[l+1]} * g^{[l]'}(z^{[l]})$$

Vectorized

$$dZ^{[l]} = dA^{[l]} * g^{[l]'}(Z^{[l]})$$
$$dW^{[l]} = \frac{1}{m} dZ^{[l]} \cdot A^{[l-1]T}$$
$$db^{[l]} = \frac{1}{m} np\text{-}sum(dZ^{[l]}, axis=1, keepdims=True)$$
$$dA^{[l]} = W^{[l]T} \cdot dZ^{[l]}$$

## Summary

$$X \rightarrow \boxed{ReLU} \rightarrow \boxed{ReLU} \rightarrow \boxed{Sigmoid} \rightarrow \hat{y}$$
$$\searrow$$
$$L(\hat{y}, y)$$

$z^{[1]} \downarrow \qquad z^{[2]} \downarrow \qquad z^{[3]} \downarrow$

$$\boxed{\phantom{xxx}} \xleftarrow{da^{[1]}} \boxed{\phantom{xxx}} \xleftarrow{da^{[2]}} \boxed{\phantom{xxx}} \xleftarrow{} da^{[L]} = \frac{-y}{a} + \frac{1-y}{1-a}$$

$\downarrow \qquad\qquad \downarrow \qquad\qquad \downarrow$

$$dw^{[1]} \qquad dw^{[2]} \qquad dw^{[3]}$$
$$db^{[1]} \qquad db^{[2]} \qquad db^{[3]}$$

$$dA^{[L]} = \left( \frac{-y^{(1)}}{a^{(1)}} \quad \frac{+(1-y^{(1)})}{1-a^{(1)}} \cdots \right.$$
$$\left. \cdots \frac{-y^{(m)}}{a^{(m)}} \quad \frac{+(1-y^{(m)})}{1-a^{(m)}} \right)$$

## Parameters vs Hyperparameters
### What are 'hyperparameters'?

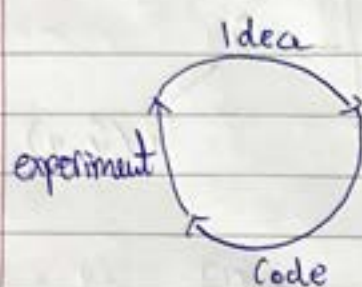We have parameters: $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}, \ldots$

hyperparameters : → learning rate $\alpha$
- → no of iterations of gradient descent
- → ~~no of~~ hidden layers $L$
- → no of hidden units $n^{[1]}, n^{[2]}, \ldots$
- → Choice of activation function

→ Hyperparameters are parameters that control $w, b$. They determine the final value of $w, b$.

→ other hyperparameters like momentum term, mini batch size and various regularization parameters.

## Applied deep learning is a very emperical process

Idea

experiment

Code

Applying DL is an iterative process by doing trial and error :

→ Set initial $\alpha$ value. Start with guesses for hyperparameters $\alpha = 0.01$

→ Train and Evaluate: Observe the behavior of cost function $J$ :

- if $J$ converges slowly, $\alpha$ might be too small
- if $J$ diverges, $\alpha$ might be large

→ Adjust the hyperparameters based on observations

## What does this have to do with the brain?