

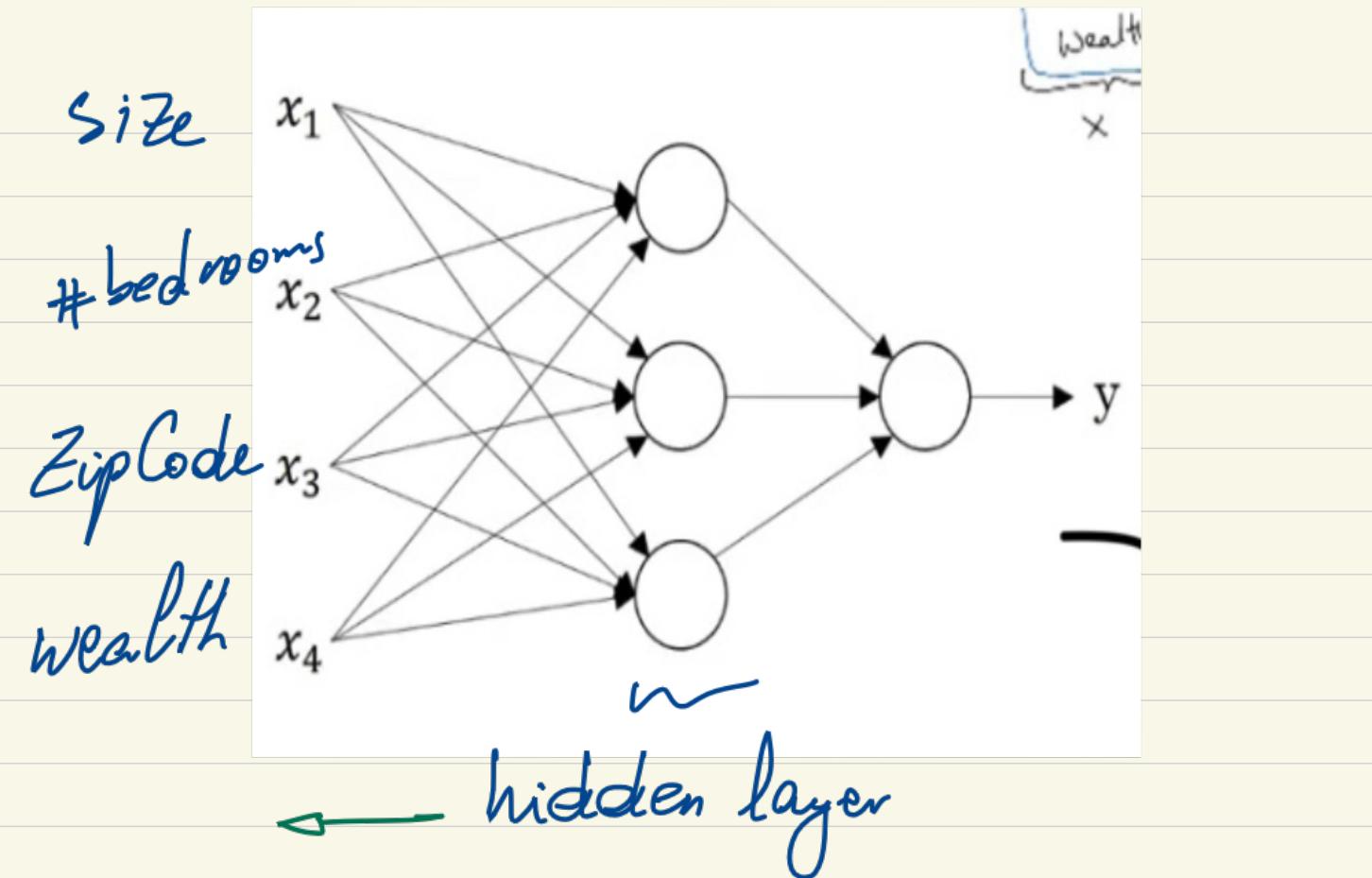
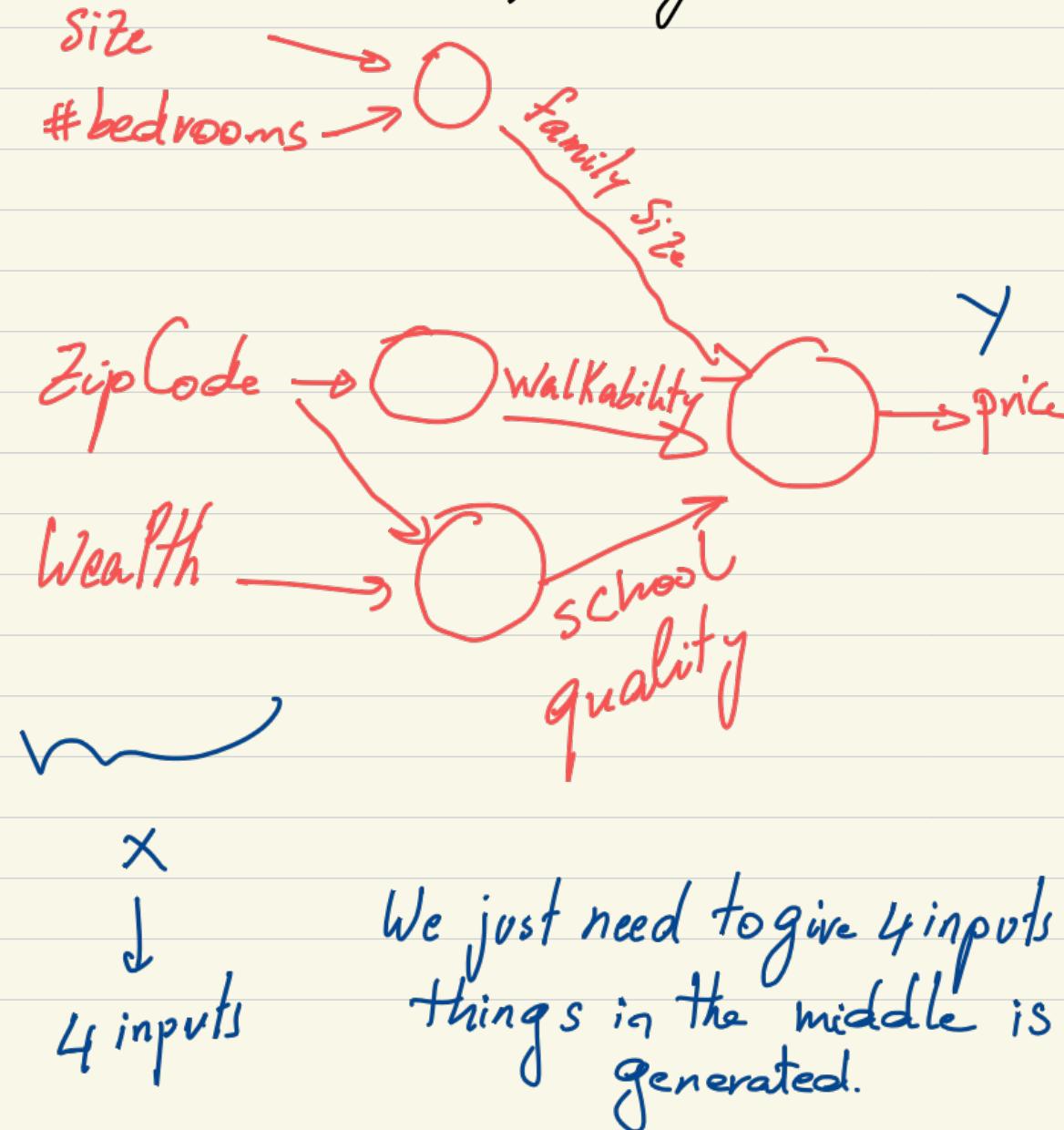
Deep learning specialisation by
Andrew Ng on Coursera

Notes by Mina Rafla

Course 1

Week 1

let's say instead of predicting the price of a house only from the size, we can other features as # bedrooms. These 2 features tell us if the house fit the family size. The Zip Code can tell us the walkability. And maybe size & wealth tell us the school quality



It densely Connected since each i/p feature is connected to each neuron in the hidden layer.

the NN decides the weights linking between each i/p feature & each hidden unit.
Giving enough data, NN can figure out function that accurately map X & Y.

In supervised learning, we have Input X & tries to estimate the price Y .

Many applications

↳ Online advertising to estimate ad clicks

Photo tagging

Speech Recognition

Machine Translation

Autonomous driving

Different types of data

↳ Structured Data (like tabular data)

↳ Unstructured Data →

Audio

Images

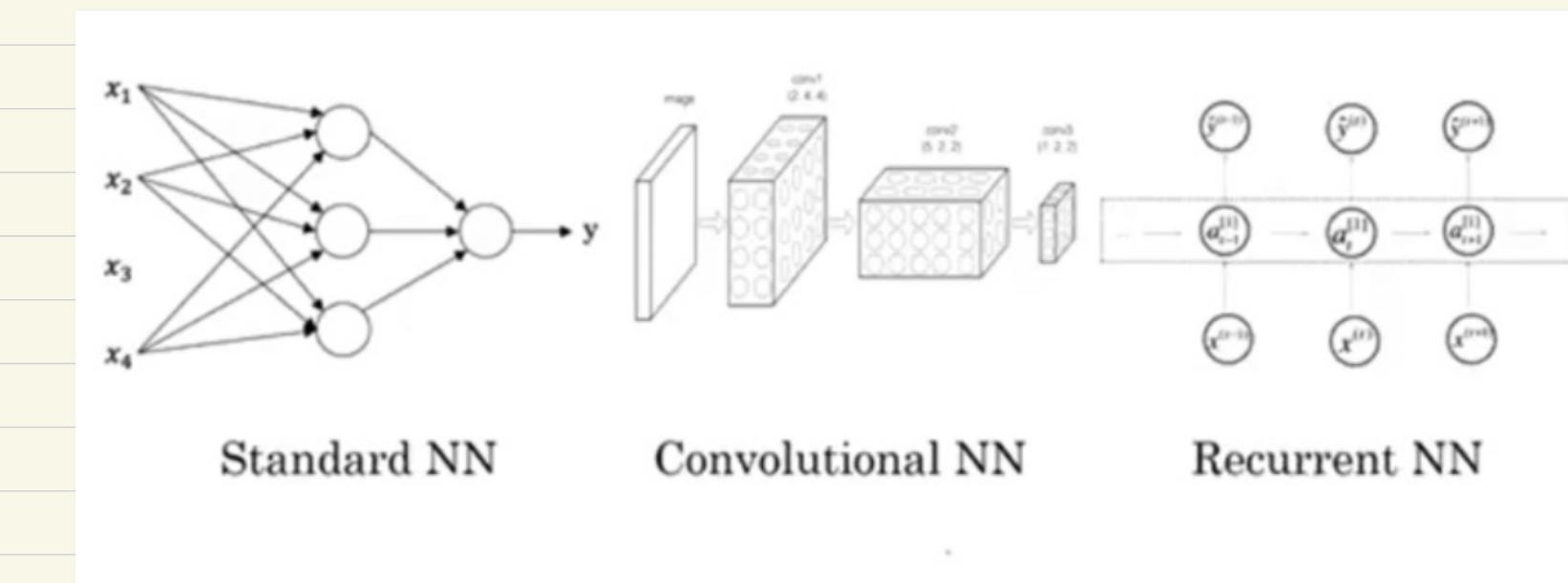
Text

Depending on the application, a different type of NN can be used.

Photo tagging → CNN (Convolutional NN)

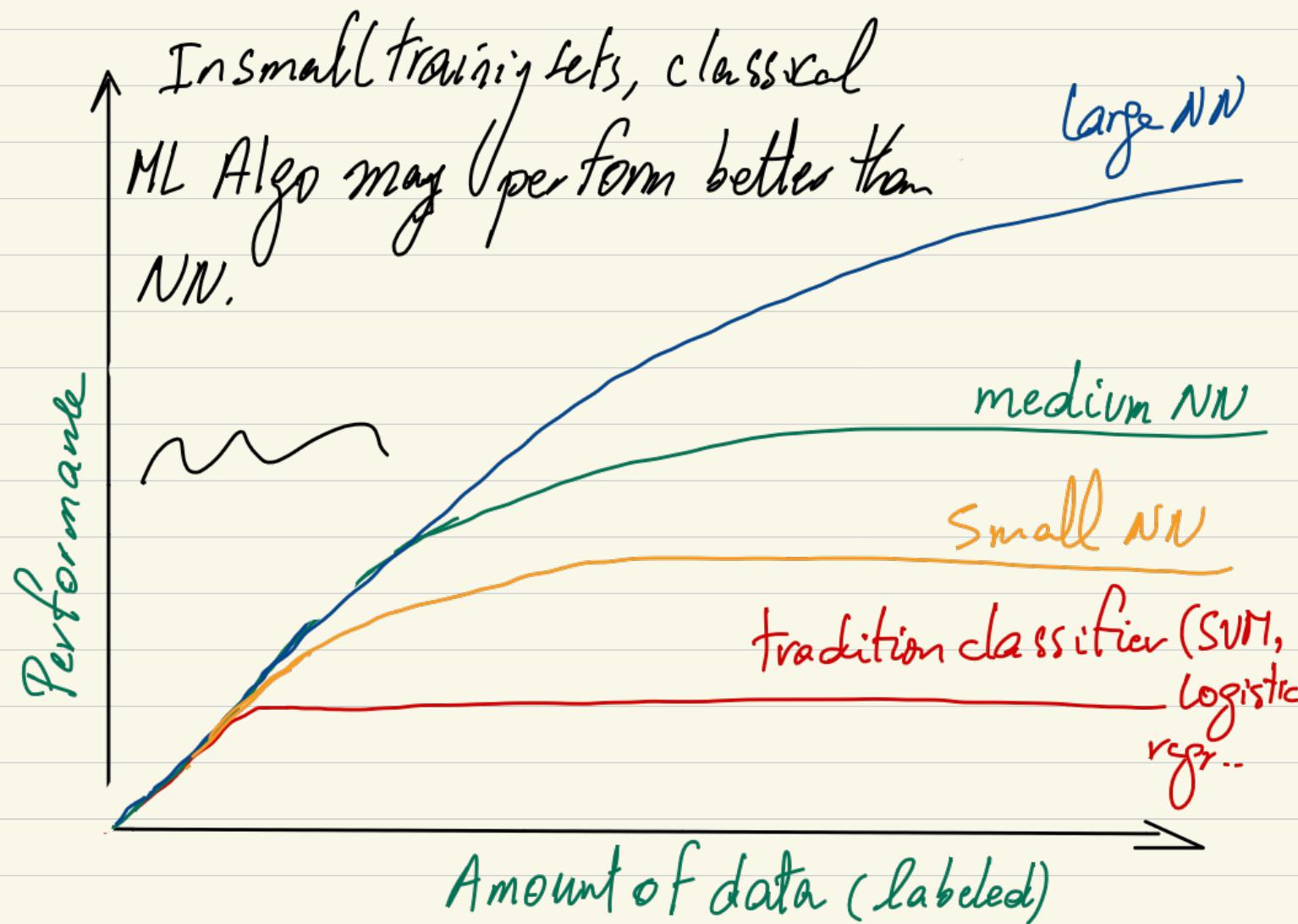
Speech Recognition → is a sequence data so RNN
Machine Translation are used.

Autonomous driving → Custom or Hybrid NN



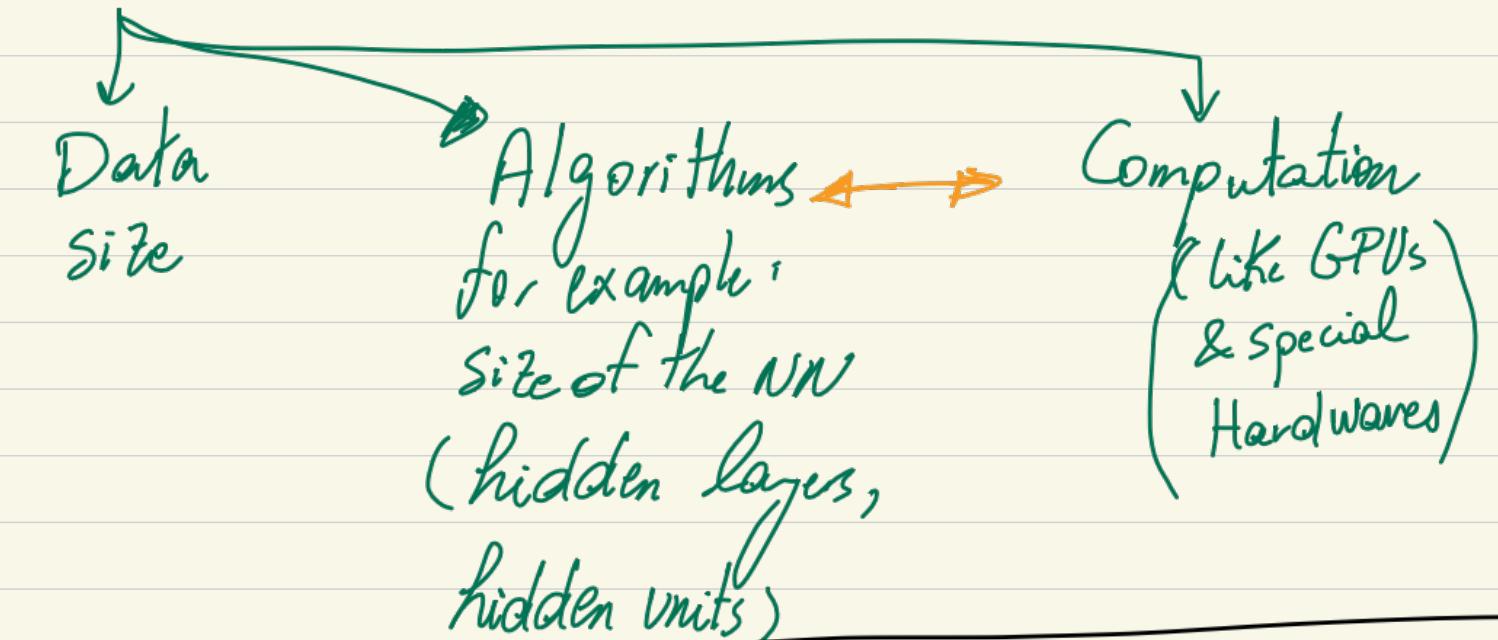
Why is deep learning taking off

Scale drives deep learning progress



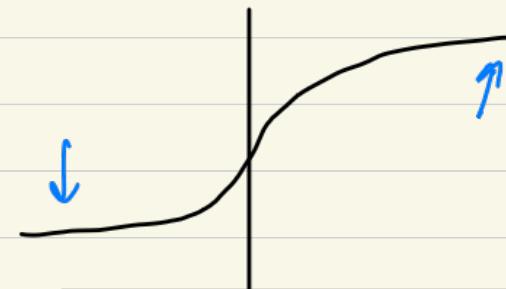
* We now have a fairly large amount of data!!
because of digitization

Scale drives Neural Network



for example:

When using Sigmoid function the learning may be very slow



Since the slope of the Sigmoid in the marked regions is

nearly zero so the gradient becomes very small & the learning (with the gradient descent) becomes very slow

Unlike the ReLu function that made the gradient descent much faster.



Course 1

Week 2

Logistic regression as a neural network

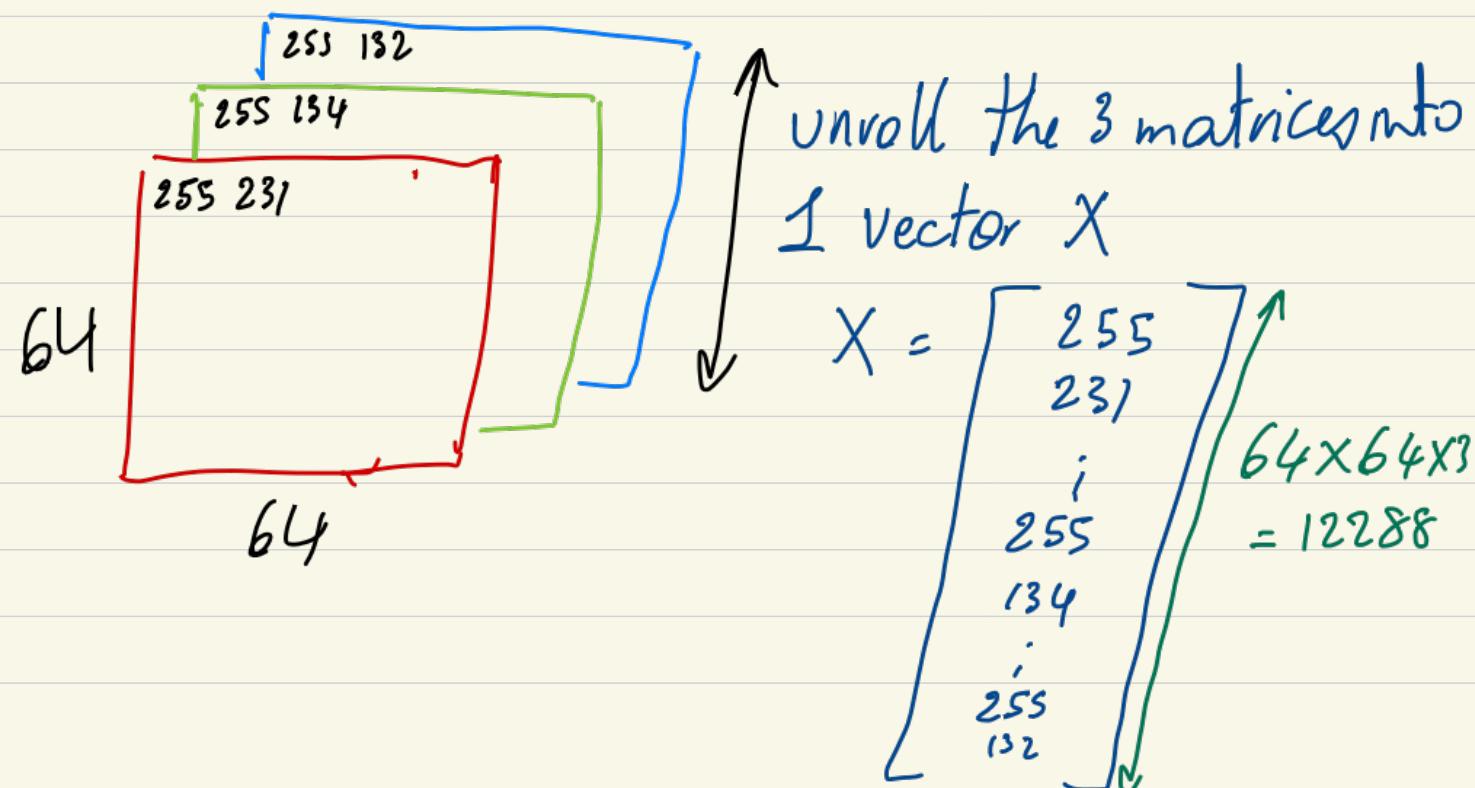
Basic Classification

Logistic Regression is an algo for binary classification



How to present an i/p image in a computer?

↳ Three matrices → Red / Green / Blue



Since the image is 64×64 pixels so the i/p X

will have a size $N_X = 64 \times 64 \times 3 = 12288$

$n = n_X = 12288 \rightarrow$ the size of the feature vector

So the goal is to learn a classifier that
has as i/p our feature vector X & outputs Y
 $(1 \text{ or } 0) \Rightarrow$ Binary classification

Notation for the rest of the Course :



a single training example is presented by a pair

$$(X, Y) \quad X \in \mathbb{R}^{n_X}, Y \in \{0, 1\}$$

m training examples : $\{(X^{(1)}, Y^{(1)}), (X^{(2)}, Y^{(2)}), \dots, (X^{(m)}, Y^{(m)})\}$

$M = M_{\text{train}}$ # of train examples $M_{\text{test}} =$ # of test examples

$\mathbb{R}^{n_X \times m}$

$$\hat{X} = \begin{bmatrix} | & | & | \\ X^{(1)} & X^{(2)} & \dots & X^{(m)} \\ | & | & \dots & | \end{bmatrix}_{n_X \times m}$$

Examples are stacked in columns

$$Y = \begin{bmatrix} | & | & | \\ Y^{(1)} & Y^{(2)} & \dots & Y^{(m)} \\ | & | & \dots & | \end{bmatrix}_{1 \times m}$$

Logistic regression

Given an i/p vector \underline{x} (for example an image)

We want to estimate $\hat{y} = P(Y=1 | \underline{x})$

$$\underline{x} \in \mathbb{R}^{n_x}$$

\hookrightarrow What is the chance
that \underline{x} is a cat

Parameters of a logistic regression picture ??

$$\hookrightarrow w \in \mathbb{R}^{n_x}, b \in \mathbb{R}$$

Output $\hat{y} = w^T \underline{x} + b$ \leftarrow One solution
that doesn't

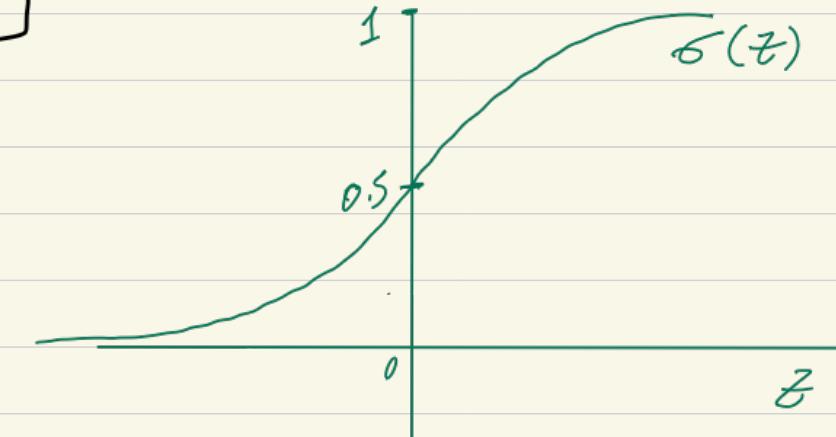
I want \hat{y} to be between 0 & 1

In logistic Regression

$$\hat{y} = \underbrace{\sigma}_{\text{Sigmoid function}}(w^T \underline{x} + b)$$

Logistic Regression is used for classification

Sigmoid function



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

\hookrightarrow If z very large: $\sigma(z) = \frac{1}{1+0} = 1$

\hookrightarrow If z is very small

$$\sigma(z) = \frac{1}{1 + e^{-z}} = \frac{1}{1 + \text{Big Number}} \approx 0$$

When implementing Logistic Regression we try to find
 $z = w^T \underline{x} + b$ // so we try to find params
 w & b

Logistic regression cost function

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

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

We have m training examples

Of course $\hat{y}^{(i)} = \sigma(w^T x^{(i)} + b)$ Notations here
 $\& z^{(i)} = w^T x^{(i)} + b$

Loss (Error) function: bad idea
↳ one thing I could do, squared Error

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

In logistic Regression
we don't do this since the optimization Problem becomes non-convex, so we can fall into local optimum

 ← Loss fn

In logistic Regression we use other loss function

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

Intuition:

We want the loss as small as possible

If $y=1 \Rightarrow L(\hat{y}, y) = -\log \hat{y} \Rightarrow$ So we want $\log \hat{y}$ to be large so \hat{y} large ($=1$)

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

so \hat{y} we want it small as possible ($=0$)

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

Cost function: to show how the algo is doing on

the entire dataset $J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$

Expanding the Cost Function

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log 1-\hat{y}^{(i)}$$

 the average of all loss functions for all the examples.

Loss function is applied on a single training example

Cost function is the cost of the parameters



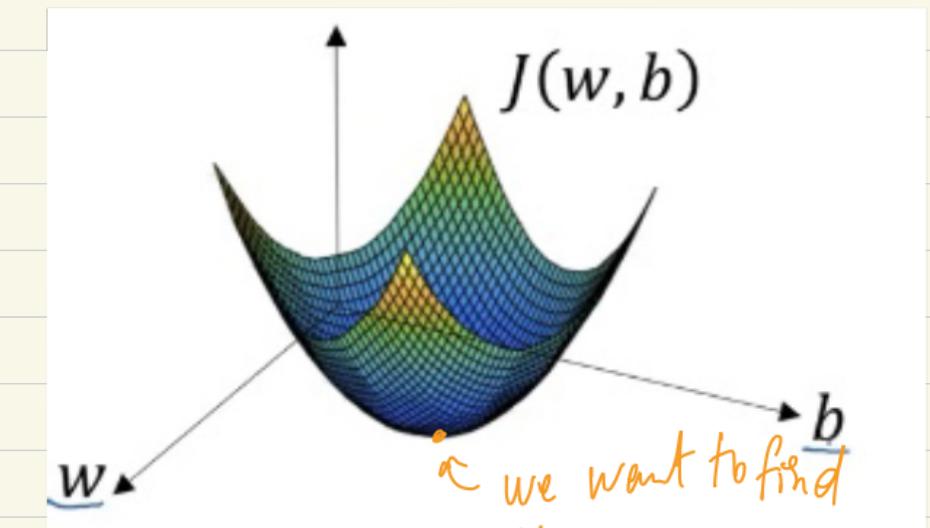
Logistic Regression is a very small Neural Network

Gradient descent

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

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log 1-\hat{y}^{(i)}$$

We want to find w, b that minimizes $J(w, b)$

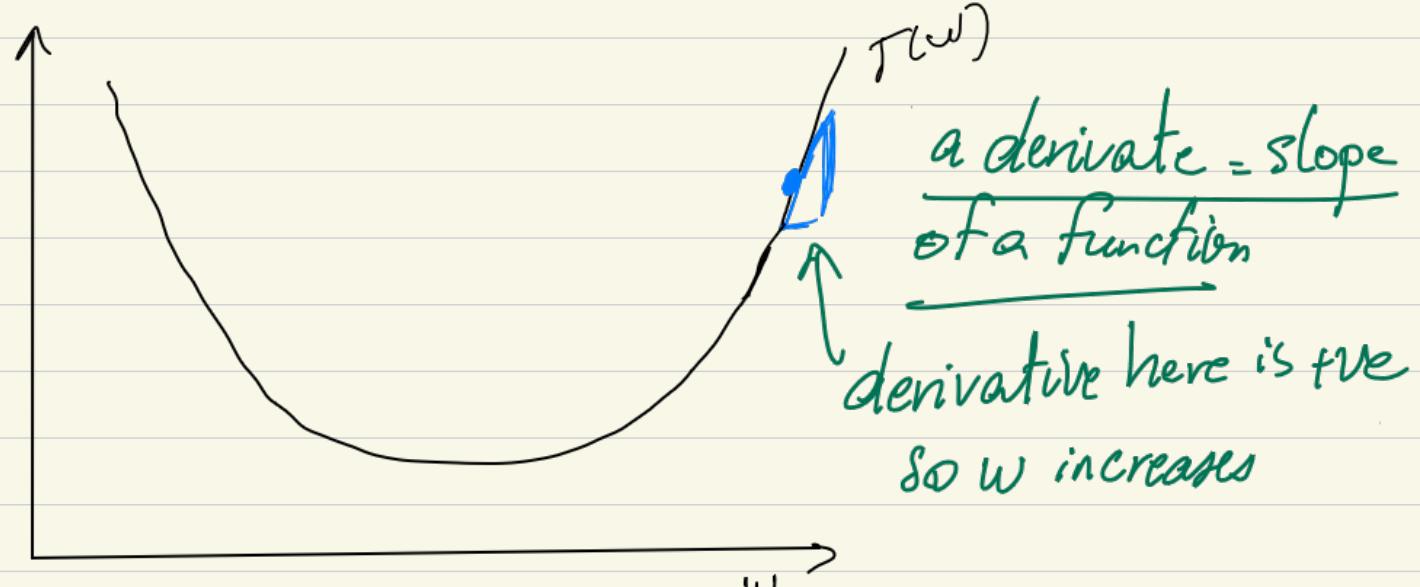


$J(w, b)$ is a surface
w can be multidimensional

$J(w, b)$ is a convex function \cup contrary to non-convex functions that contain several local optima

To find a good value for the params we initialize w, b to some initial value.

In Logistic Regression, no matter what the value initial is it should find the minimum cost function since it's a convex function



Repeat {

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

learning rate
Controls how big the steps I am taking

$$\left\{ \frac{dJ(w)}{dw} = ? \right.$$

The updates of the parameters are then

$$w := w - \alpha \frac{dJ(w, b)}{dw}$$

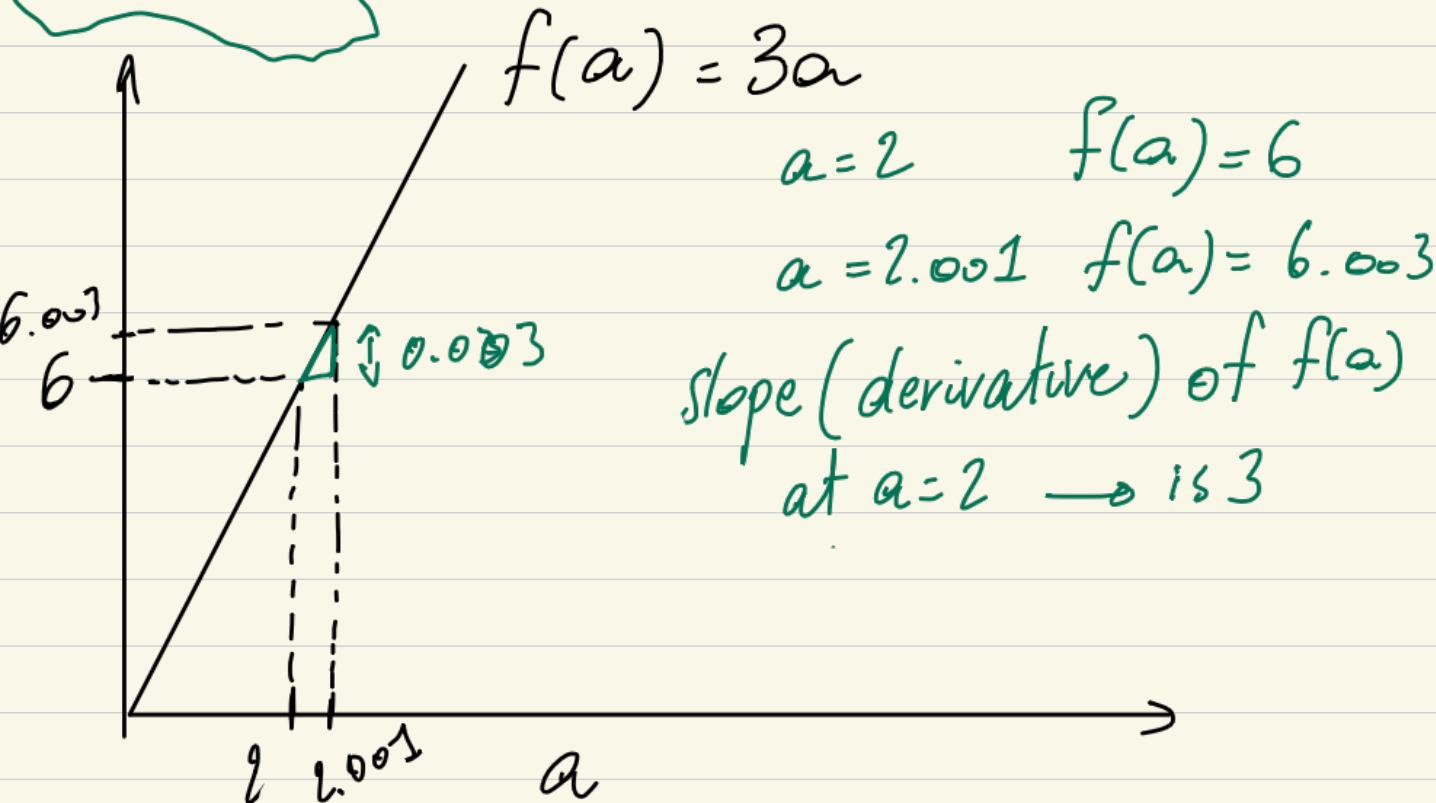
$$b := b - \alpha \frac{dJ(w, b)}{db}$$

To write the symbol of the derivative of a function of 2 or more variables \Rightarrow we write $\frac{\partial J(w, b)}{\partial w}$ partial derivative symbol

$$\frac{\partial J(w, b)}{\partial b}$$

Derivatives

Intuition of derivatives



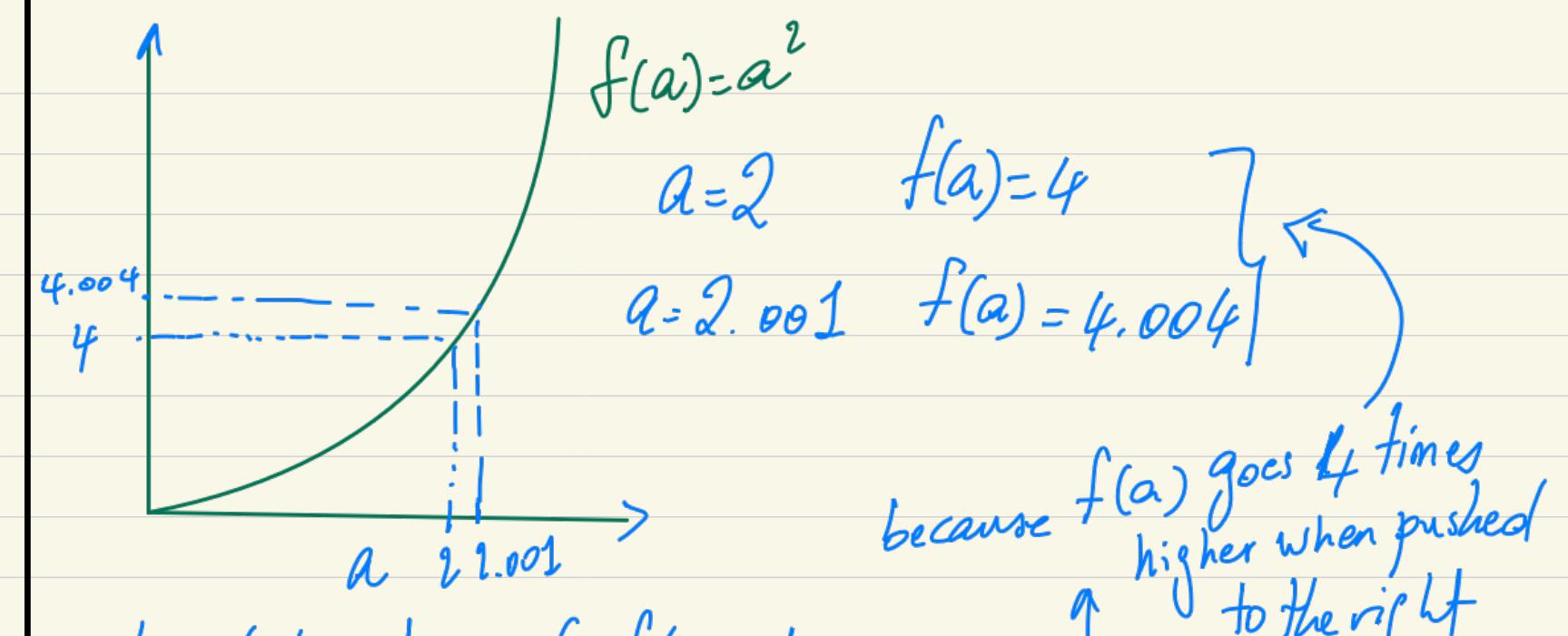
$$\text{So Slope} = \frac{\text{height}}{\text{width}} = \frac{0.003}{0.001}$$

$$\text{if } a=5 \quad f(a)=15$$

$$\text{if } a=5.001 \quad f(a) = 15.003$$

$$\text{slope of } f(a) = \frac{d f(a)}{da} = 3 = \frac{d}{da} f(a)$$

If I push a a little to the right
 $f(a)$ goes up by 3.



slope (derivative) of $f(a)$ at $a=2$ is 4

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

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

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

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

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

$$w_i = w - \alpha \frac{d J}{dw}$$

If we have a value "a", we push it slightly to the right new $f(a) = \text{old } f(a) + \frac{d}{da} f(a) + \text{distance pushed to right}$

More derivative examples:

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

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

if I knowed a to
the right by 0.001
 $f(a)$ will jump by
 $3a^2$

Let's say that $f(a) = \log(a)$ ($\ln a$ or $\log_e a$)

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

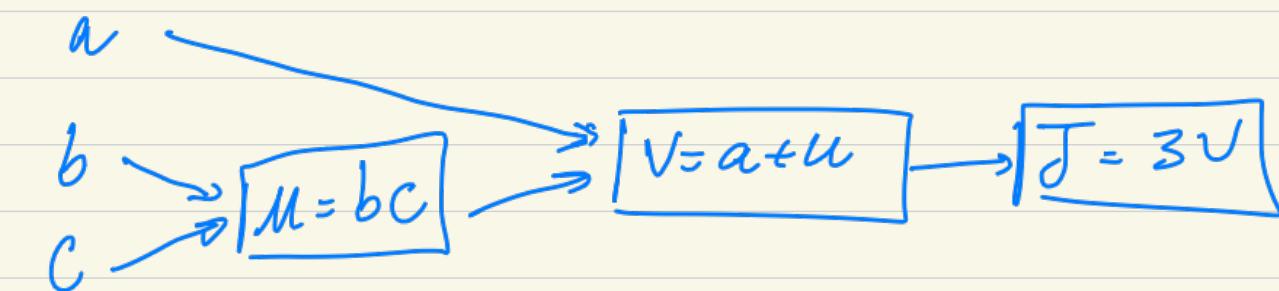
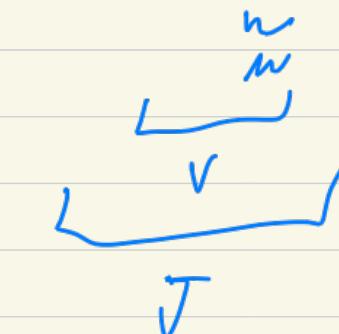
$\ln(a)$

$\boxed{\text{if } a=2 \rightarrow f(a) = \log(2)}$

$\boxed{a=2.001 \rightarrow f(a) = \log(2) + \frac{1}{a} * 0.001}$

Computation graph

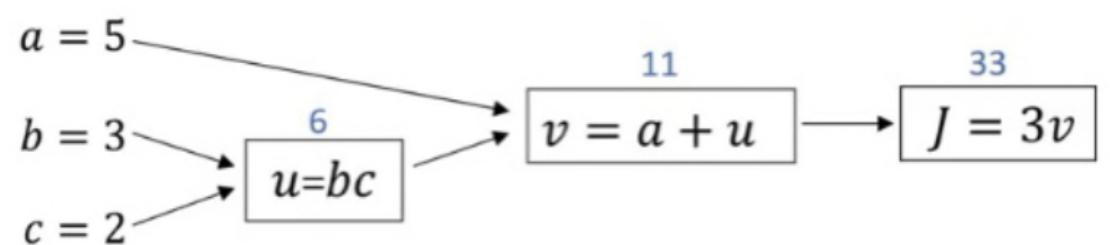
$$J(a, b, c) = 3(a+b+c)$$



In order to compute J we go from left to right
to calculate derivatives we go from right to left

This will be clear in the next couple of videos.

Derivatives with a computation graph



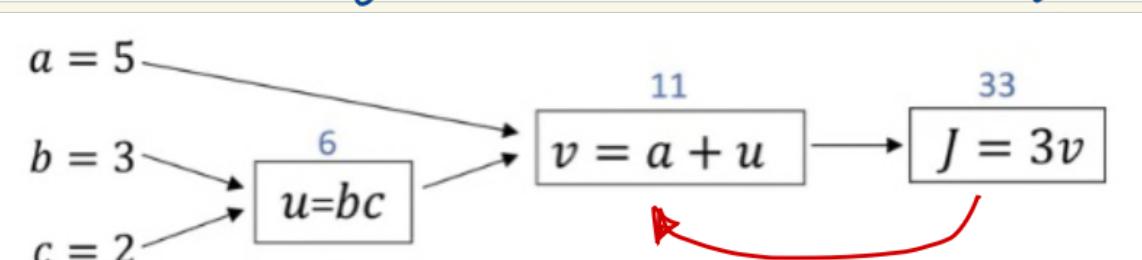
$\frac{dJ}{dv} = ?$ (If we take the value of v and change it by 0.001, how will the value of J change?)

$$J = 3v$$

$$v = 11 \rightarrow 11.001$$

$$J = 33 \rightarrow 3.003$$

In the terminology of backpropagation, to compute derivatives of J we have to go back in the computation graph



$$\frac{dJ}{da} = ? \Rightarrow a = 5 \rightarrow 5.001$$

$$v = 11 \rightarrow 11.001$$

$$so \frac{dJ}{da} = 3$$

$$J = 33 \rightarrow 33.003$$

If I change a that will change v & that will change J

This is called the chain rule

$$\frac{dJ}{da} = \frac{dJ}{dv} \frac{dv}{da}$$

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

If J change "a" by 0.001 v changes with the same amount

$$so \frac{dv}{da} = 1$$

$$\frac{dJ}{du} = \frac{dJ}{dv} \frac{dv}{du} = 3 \cdot 1 = \boxed{3}$$

$$\frac{dJ}{db} = \frac{dJ}{dv} \frac{dv}{du} \frac{du}{db} \quad \leftarrow \text{Chain Rule}$$

$$= 3 \cdot 1 \cdot 2 = \boxed{6}$$

$$b = 3 \rightarrow 3.001$$

$$u = b \cdot c = 6 \rightarrow 6.002$$

Take home message \rightarrow To take derivatives it's better to go from right to left (back propagation)

Logistic regression gradient descent

How to do gradient descent for logistic regression?

Logistic Regression Recap

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

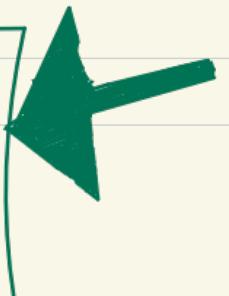
$$\hat{y} = a = \sigma(z)$$

$$\mathcal{L}(a, y) = -\left(y \log(a) + (1-y) \log(1-a)\right)$$

Let's say we have only 2 features x_1 & x_2

To compute derivatives of the loss:

$$*\frac{d\mathcal{L}(a, y)}{da} = \frac{-y}{a} + \frac{(1-y)}{1-a} *$$



$$\frac{d\mathcal{L}}{dz} = \frac{d\mathcal{L}}{da} \frac{da}{dz} = \left(\frac{-y}{a} + \frac{1-y}{1-a} \right) \frac{da}{dz}$$

Not Mentioned by Andrew Ng:

To calculate the derivative of the sigmoid function

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

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

Reminder that $\boxed{\frac{d}{dz} e^z = e^z}$

quotient rule for derivatives: $\rightarrow f(x) = \frac{g(x)}{h(x)}$

$$f'(x) = \frac{g'(x)h(x) - h'(x)g(x)}{(h(x))^2}$$

$$\frac{d}{dz} \left[\frac{1}{1+e^{-z}} \right] = \frac{(0)(1+e^{-z}) - (-e^{-z})(1)}{(1+e^{-z})^2} =$$

$$\frac{e^{-z}}{(1+e^{-z})^2} = \frac{1}{(1+e^{-z})} \cdot \frac{e^{-z}}{(1+e^{-z})} =$$

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

$$\frac{1}{1+e^{-z}} \left[\frac{1+e^{-z}}{1+e^{-z}} - \frac{1}{1+e^{-z}} \right] =$$

$$\sigma(z) [1 - \sigma(z)] =$$

$$[a [1-a]]$$

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

$\nearrow a(1-a)$

$\nearrow -y/a + \frac{1-y}{1-a}$

$= \boxed{a-y}$

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

$$\frac{dL}{dw_2} = x_2 \frac{dL}{dz}$$

$$\frac{dL}{db} = \frac{dL}{dz}$$

Gradient descent on m examples

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

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

What was shown in the previous lesson was how to calculate derivatives for 1 training example.

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

Logistic Regression on m examples:

$$J=0, dw_1=0, dw_2=0, db=0 \quad \text{init}$$

for $i=1$ to m :

$$\left. \begin{array}{l} Z^{(i)} = w^T x^{(i)} + b \\ a^{(i)} = \sigma(Z^{(i)}) \\ J += - (y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log(1-a^{(i)})) \end{array} \right\} \text{forward}$$

$$\frac{dL}{dZ^{(i)}} += a^{(i)} - y^{(i)}$$

$$\frac{dL}{dw_1} += x_1^{(i)} \frac{dL}{dZ^{(i)}}$$

$$\frac{dL}{dw_2} += x_2^{(i)} \frac{dL}{dZ^{(i)}}$$

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

$$J /= m$$

$$dw_1 /= m; dw_2 /= m; db /= m$$

After the previous algorithm, to implement one step of the gradient descent we do

$$w_1 := w_1 - \alpha \frac{dL}{dw_1}$$

$$w_2 := w_2 - \alpha \frac{dL}{dw_2}$$

$$b := b - \alpha \frac{dL}{db}$$

Vectorization

Vectorization \Rightarrow enables to get rid of for loops in your code

What is Vectorization?

In logistic Regression, we need to compute:

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

$$w \in \mathbb{R}^{n \times 1} \quad x \in \mathbb{R}^{n \times 1}$$

Non-Vectorized implementation:

$$z = 0$$

for i in range(n):

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

$$z += b$$

Vectorized implementation

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

More Vectorization examples.

matrix \rightarrow vector
 $u = A v$

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

Non Vectorized Implementation

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

for i ...

for j ...

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

Vectorized Imp

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

Another example: Say you need to apply the exponential operation on every element of a matrix/Vector

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

Vectorized imp

$$u = np.exp(v)$$

Non Vectorized implementation

```

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

```

Logistic Regression derivatives

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

for $i = 1$ to m :

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

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

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

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

$$\begin{aligned} dw_1 &+= x_1^{(i)} dz^{(i)} && \text{if I had more than 2 features} \\ dw_2 &+= x_2^{(i)} dz^{(i)} && \text{I would do a for loop over} \\ db &+= dz^{(i)} && \text{the weights } w_1 \text{ wrt } b \end{aligned}$$

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

To eliminate this second for loop

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

for the initialization

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

will eliminate the line marked above

Vectorising logistic regression

forward propagation step in a logistic Regression:

$$\begin{aligned} Z^{(1)} &= w^T x^{(1)} + b & Z^{(2)} &= w^T x^{(2)} + b \\ a^{(1)} &= \sigma(Z^{(1)}) & a^{(2)} &= \sigma(Z^{(2)}) \end{aligned}$$

$$\begin{aligned} Z^{(3)} &= w^T x^{(3)} + b \\ a^{(3)} &= \sigma(Z^{(3)}) \end{aligned}$$

for 1st training example

for 2nd training example

for 3rd training example

The question now is how to do it without an explicit for loop

$$\text{Remember we defined a matrix } X = \begin{bmatrix} 1^{(1)} & 1^{(2)} & \dots & 1^{(m)} \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ 1 & 1 & \dots & 1 \end{bmatrix}$$

Examples are stacked by columns.

(n_x, m)

D # examples

of features

$$\text{Let's define } 1 \times m \text{ vector } Z = [Z^{(1)} Z^{(2)} \dots Z^{(m)}]$$

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

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

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

$\rightarrow 1 \times m$ vector

$$Z = [z^{(1)} \ z^{(2)} \dots z^{(m)}] = w^T X + [b \ b \ b \dots b]$$

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

In python, this implemented as:

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

added to a vector \leftarrow b is just a real number
 using "Broadcasting"

So now we implemented " Z ", how do we implement " \hat{a} ".

* We define vector $[\hat{a}^{(1)} \ \hat{a}^{(2)} \dots \hat{a}^{(m)}] = A$

$$= \sigma(Z)$$

Victor rising logistic regression gradient decent

In this lesson, we will see how to use vectorization for gradient's calculations.

When we saw the gradient's calculation:

* $\frac{dL}{dZ^{(1)}} = \hat{a}^{(1)} - y^{(1)} \quad \frac{dL}{dZ^{(2)}} = \hat{a}^{(2)} - y^{(2)} \dots$

\hookrightarrow denoted by dZ

We define $dZ = [dZ^{(1)} \ \dots \ dZ^{(m)}] \Rightarrow 1 \times m$ vector

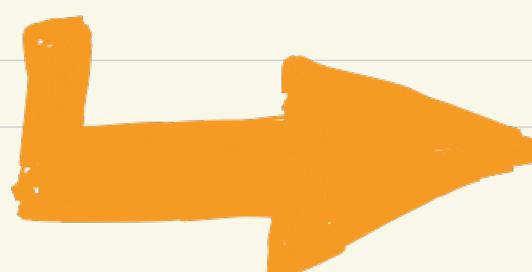
from the prev lesson: $A = [\hat{a}^{(1)} \ \hat{a}^{(2)} \dots \hat{a}^{(m)}] \Rightarrow 1 \times m$ vector

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

So $dZ = A - Y = [\hat{a}^{(1)} - y^{(1)} \ \hat{a}^{(2)} - y^{(2)} \ \dots \ \hat{a}^{(m)} - y^{(m)}]$

So with one line of code we can do this calculation

In a previous lesson, we were able to get rid of one for loop (over the set of features for one training example).
 But we still have to get rid of (a for loop over the training examples).



Reminder From a Previous Lesson

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

for i = 1 to m:

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

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

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

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

$$\begin{aligned} dw_1 &+= x_1^{(i)} dz^{(i)} \\ dw_2 &+= x_2^{(i)} dz^{(i)} \\ db &+= dz^{(i)} \end{aligned}$$

if I had more than 2 features
I would do a for loop over
the weights w_1 or b

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

To eliminate this second for loop

$$dw = np.zeros((n_x, 1)) \quad \leftarrow \text{for the initialization}$$

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

will eliminate the line marked above

What we do basically with db is

$$\begin{aligned} db &= \frac{1}{m} \sum_i^{m-1} dz^{(i)} \\ &= \frac{1}{m} np.sum(dz) \end{aligned}$$

in python

What about dw?

$$\begin{aligned} dw &= \frac{1}{m} X dz^T = \frac{1}{m} \begin{bmatrix} 1 & 1 & \dots & 1 \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix} \begin{bmatrix} dz^{(1)} \\ dz^{(2)} \\ \vdots \\ dz^{(m)} \end{bmatrix} = \\ &\frac{1}{m} \begin{bmatrix} x^{(1)} dz^{(1)} + \dots + x^{(m)} dz^{(m)} \end{bmatrix} = \end{aligned}$$

in python

so let's put everything together:

Inefficient
Python implementation

] In the previous lesson we got rid of the for loop over the features

How to implement logistic Regression Gradient Descent efficiently in Python.

$$\begin{aligned} \hookrightarrow Z &= w^T X + b \\ &= np.dot(w.T, X) + b \quad \left. \begin{array}{l} \text{forward} \\ \text{prop.} \end{array} \right\} \\ A &= \sigma(Z) \end{aligned}$$

$$\begin{aligned} dZ &= A - Y \\ dW &= \frac{1}{m} X dZ^T \quad \left. \begin{array}{l} \\ \text{Backprop} \end{array} \right\} \\ db &= \frac{1}{m} np.sum(dZ) \end{aligned}$$

$$\begin{aligned} w_i &= w - \alpha dw \\ b_i &= b - \alpha db \end{aligned}$$

This is only one iteration of Gradient Descent

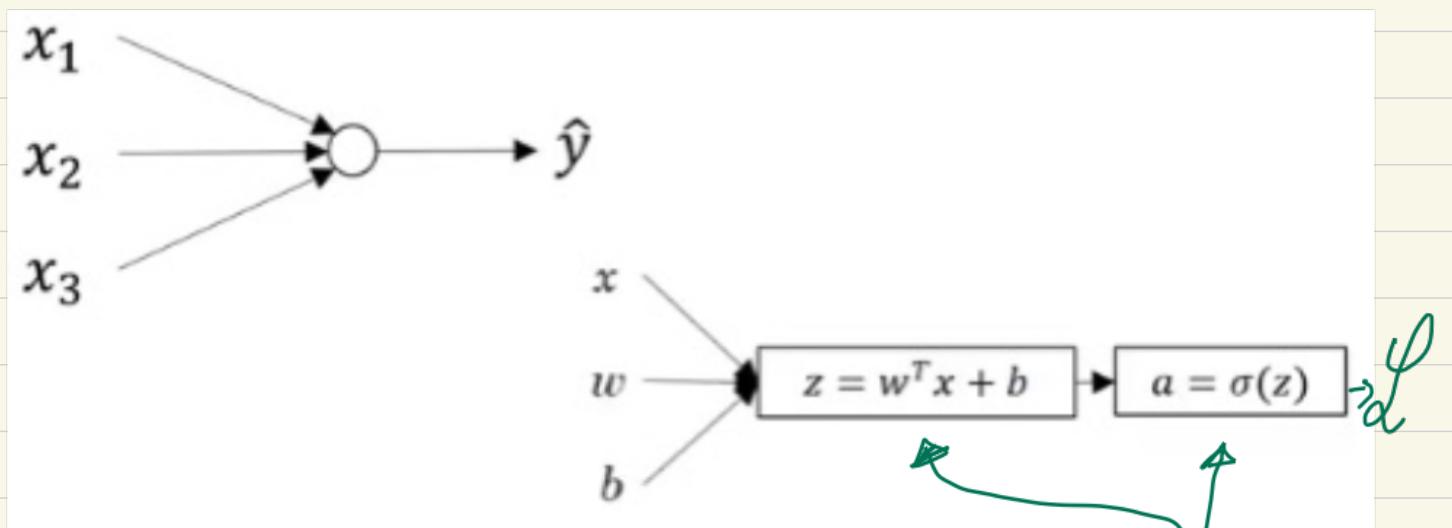
For multiple iterations \Rightarrow I will need a for loop!!!

Course 1

Week 3

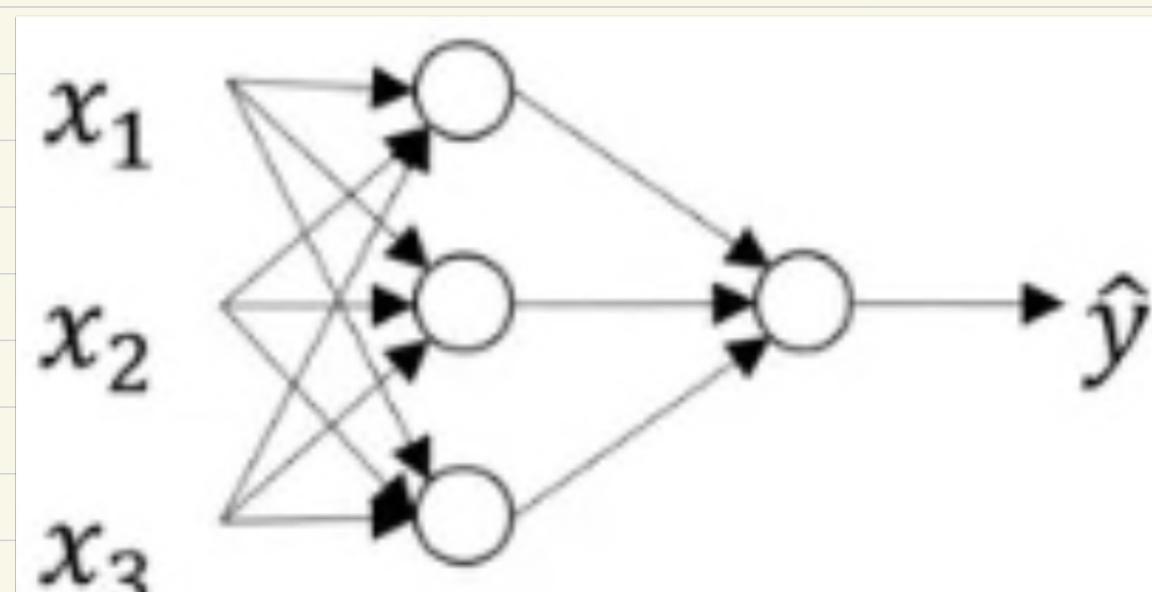
Neural networks overview

A logistic Regression is like a small Neural Netwo.



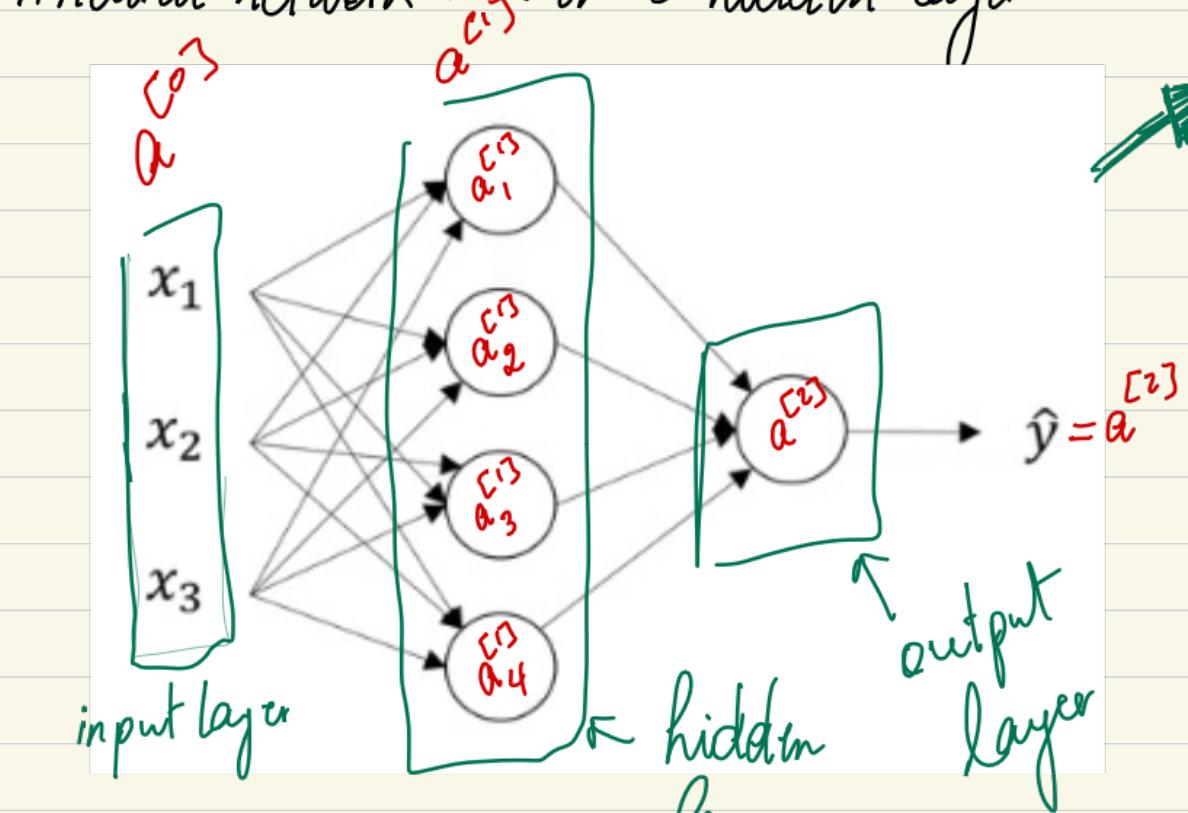
In logistic Regression we have two steps of calculations

In Neural Nets, more Calculations are needed.
Will get clearer in the next lessons



Neural network representation

A neural network with s_i be hidden layer



2 layer NN

Some new notations:

Input Layer $X = a^{[0]}$ \Rightarrow "a" for activation

Hidden Layer " $a^{[1]}$ " & each node will be called $a^{[1]}_1, a^{[1]}_2, a^{[1]}_3, a^{[1]}_4$

Output Layer " $a^{[2]}$ " \Rightarrow Which is a number

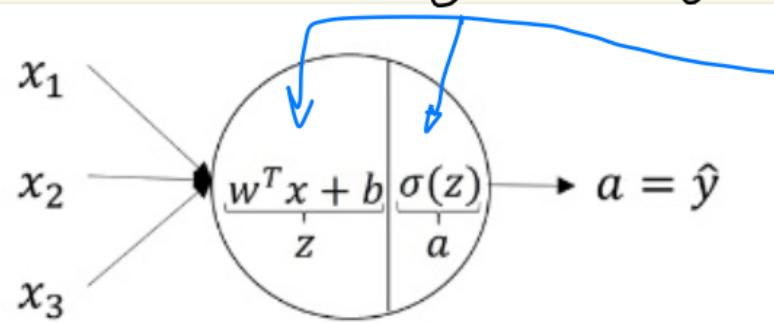
Note that: the hidden layer $a^{[1]} = \begin{bmatrix} a^{[1]}_1 \\ a^{[1]}_2 \\ a^{[1]}_3 \\ a^{[1]}_4 \end{bmatrix}$

- the hidden layer has params associated with it noted $w^{[1]}, b^{[1]}$
- $w^{[1]}$ is a $(4, 3)$ matrix \rightarrow (4, 1) vector
- $b^{[1]}$ is a $(4, 1)$ vector
- $w^{[1]}$ has 3 nodes in the i/p layer
- $b^{[1]}$ has 4 nodes in hidden layer

When we count layers in NN, we don't Count the i/o layer

Computing in neural networks output

The circle in a logistic Regression has 2 steps of computations.

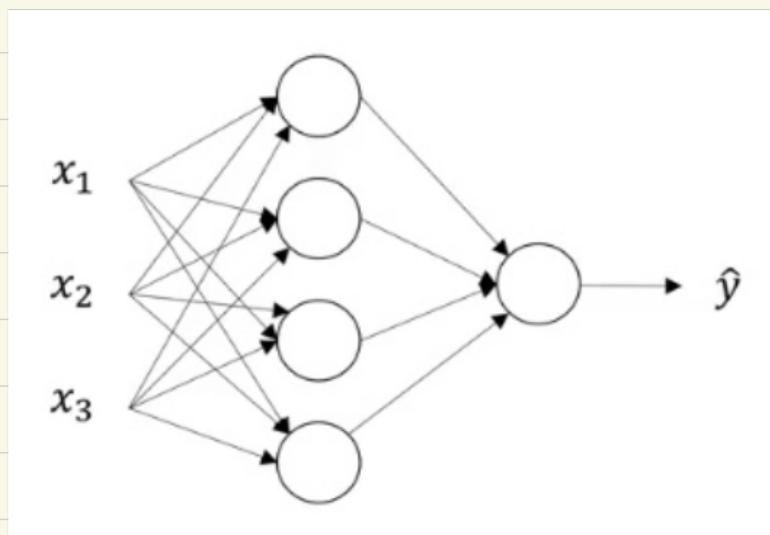


$$z = w^T x + b$$

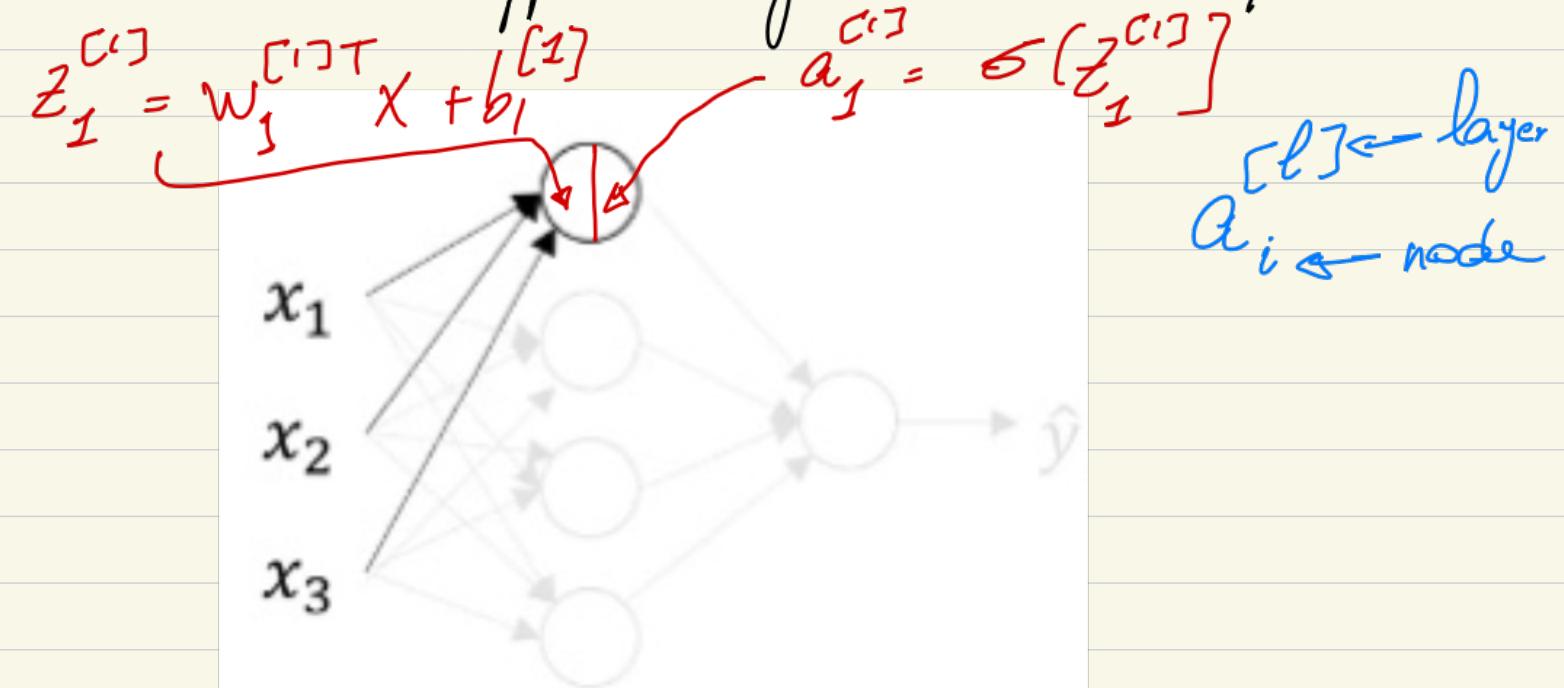
$$a = \sigma(z)$$

For Example For a neural network like the following:

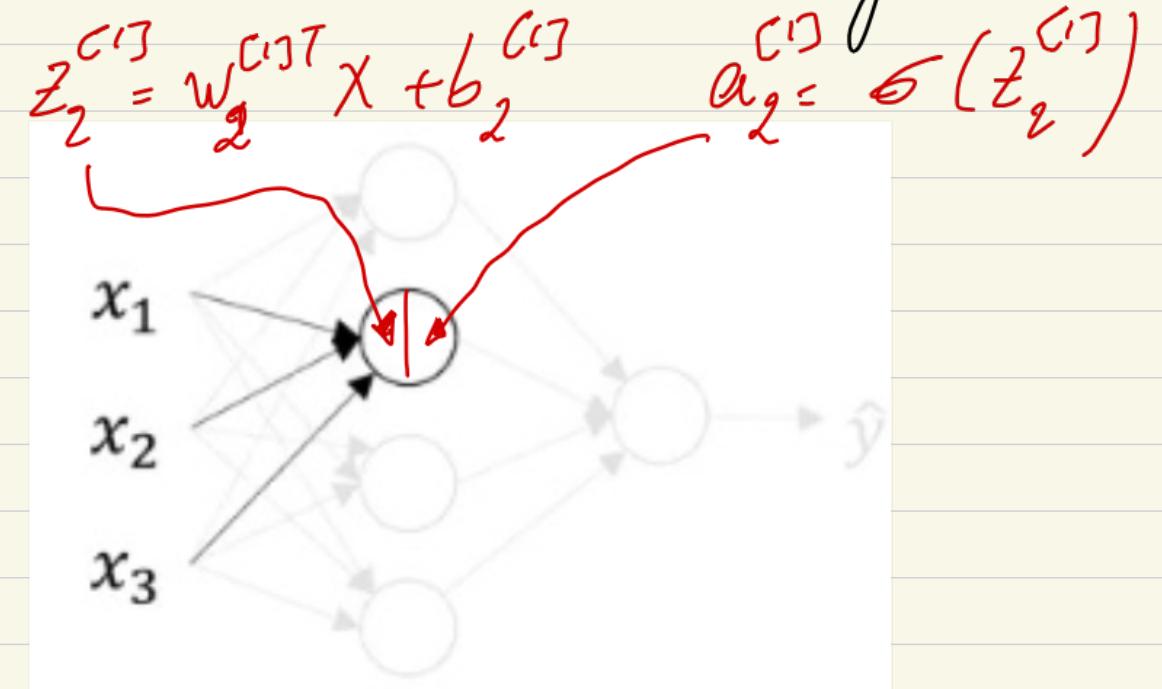
The Same is done for each node in the hidden layer



let's take a look on what happens in a single node in a hidden layer



Same for the Second in the hidden layer



Etc for all the other nodes

To do this using vectors & matrices (to avoid for loops in Python)

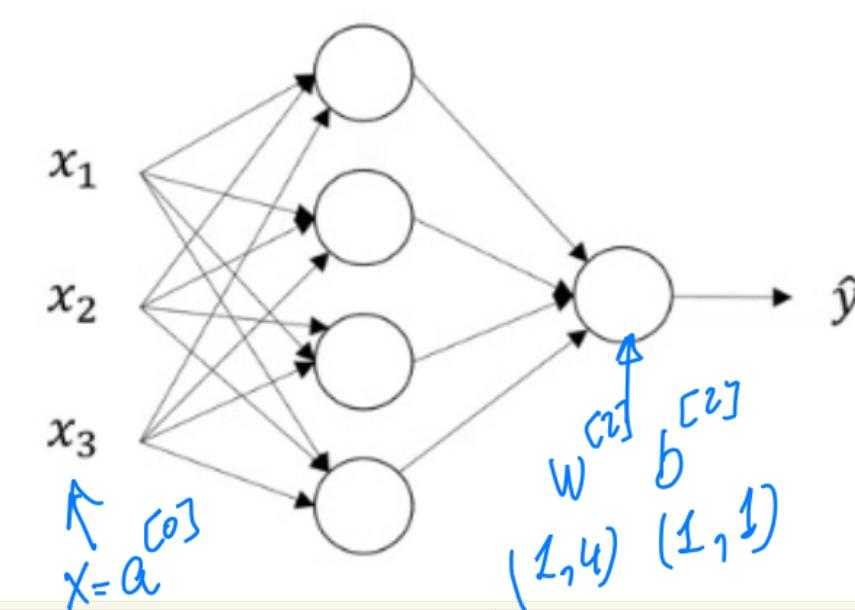
We end up with:

$$\begin{bmatrix} w^{[1]} \\ w^{[2]} \\ w^{[3]} \\ w^{[4]} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} + \begin{bmatrix} b^{[1]} \\ b^{[2]} \\ b^{[3]} \\ b^{[4]} \end{bmatrix} = \begin{bmatrix} w_1^{[1]T} x + b_1^{[1]} \\ w_2^{[2]T} x + b_2^{[2]} \\ w_3^{[3]T} x + b_3^{[3]} \\ w_4^{[4]T} x + b_4^{[4]} \end{bmatrix} = Z^{[1]}$$

Similarly,

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

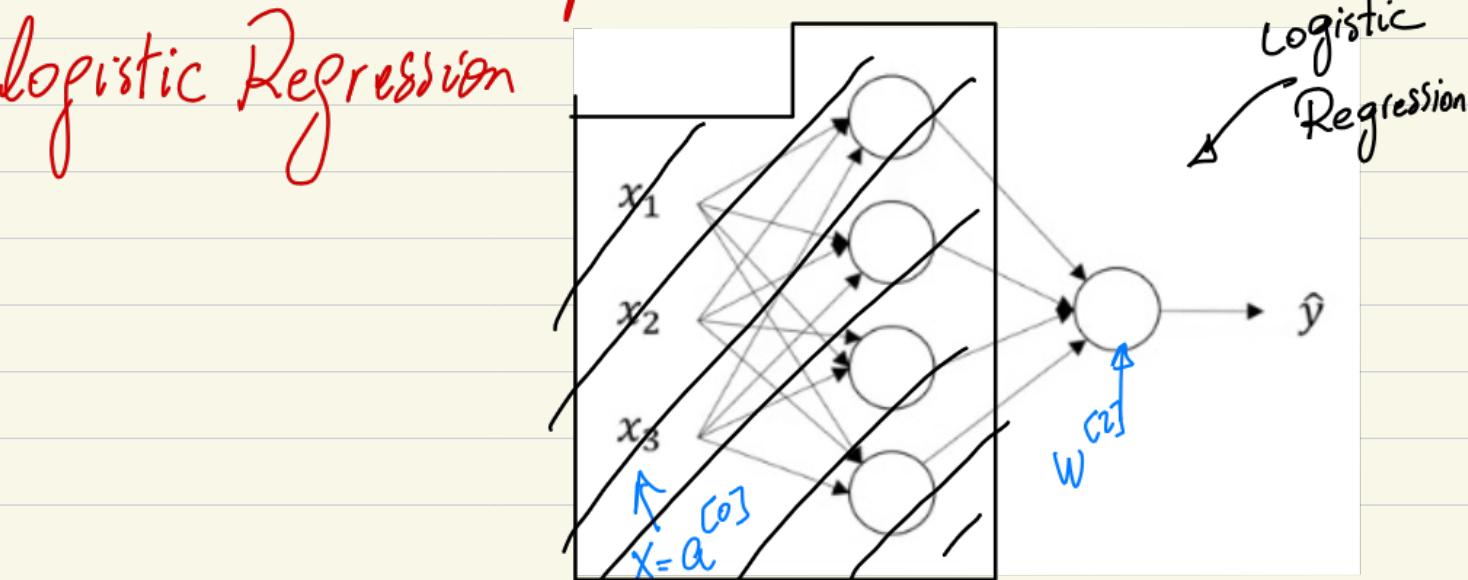
So given an i/p X :



$$\begin{aligned} Z^{[1]} &= W^{[1]} X + b^{[1]} \\ a^{[1]} &= \sigma(Z^{[1]}) \\ Z^{[2]} &= W^{[2]} a^{[1]} + b^{[2]} \rightarrow \text{Value} \\ a^{[2]} &= \sigma(Z^{[2]}) \end{aligned}$$

Remember that we said the input layer X is $a^{[0]}$
so in the equation we can replace X by $a^{[0]}$

If we take the figure above & cover up the left part
like this, we end up with a structure like a
logistic Regression



Vectorising across multiple examples

One hidden layer

These are our four equations from the previous lesson.

They tell given an input vector x how to get a prediction a^2

If I have m feature example we should repeat these equations

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

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

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

$$\vdots$$

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

Round bracket

(i) means the i^{th} example

$a^{[2](i)}$ i^{th} example layer 2

$$Z = W X + b$$

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

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

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

A non vectorized implementation will consist of a for loop to loop on each example.

$$\text{for } i = 1 \text{ to } m,$$

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

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

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

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

Recall We define matrix $\hat{X} = \begin{bmatrix} 1 & 1 & 1 \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ 1 & 1 & \dots & 1 \end{bmatrix}_{n \times m}$ # of features # Examples

$$Z = W X + b$$

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

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

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

Basically what we have done is to use Capital \hat{X} matrix instead of lower case $X^{(i)}$ by stacking $x^{(i)}$ in columns. Same thing can be done with Z

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

Same thing can be done with the A matrix

$$A = \begin{bmatrix} | & | & | \\ c_{1,1}^{(1)} & c_{1,2}^{(1)} & c_{1,m}^{(1)} \\ a & a & \dots \\ | & | & | \end{bmatrix}$$

Remember from previous lessons that
 $Z^{(1)(1)}$ is $(4, 1)$ matrix where 4 was the
number of nodes. So Capital Case Z matrix

Contains $\begin{bmatrix} Z^{c_{1,1}^{(1)}} & Z^{c_{1,2}^{(1)}} & \dots & Z^{c_{1,m}^{(1)}} \end{bmatrix}$ ↑ hidden units

* a number of columns = number of examples
& # of rows = number of nodes in the layer 1 *

Remember the node number was node $Z_1^{(1)(1)}$
first node →

Explanation for vectorised implementation

Justification for vectorized implementation

$$Z^{c_{1,1}^{(1)}} = W^{c_1} X^{(1)} + b^{c_1} \quad \leftarrow \text{First training example}$$

$$Z^{c_{1,2}^{(1)}} = W^{c_2} X^{(2)} + b^{c_2} \quad \leftarrow \text{Second training example}$$

$$Z^{c_{1,m}^{(1)}} = W^{c_m} X^{(m)} + b^{c_m} \quad \leftarrow \text{Third example}$$

For the moment we assume $b^{c_i} = 0$

↳ for Simplification

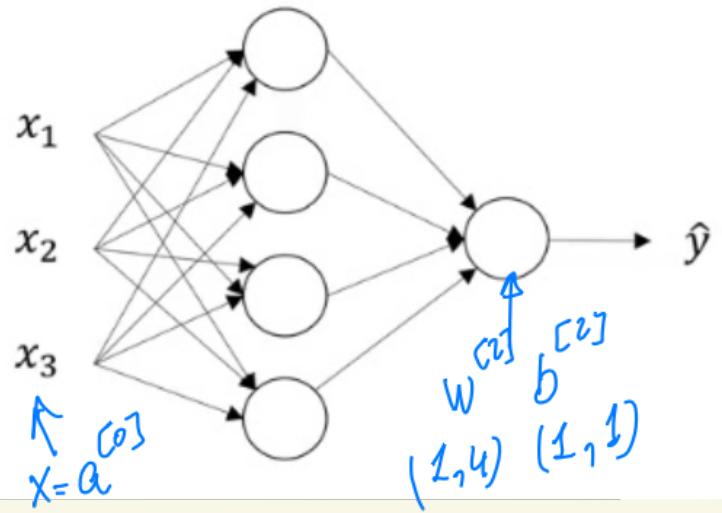
$$W^{c_i} = \begin{bmatrix} \rightarrow \\ \rightarrow \\ \rightarrow \\ \rightarrow \end{bmatrix}$$

$$\text{So } W^{c_i} X^{(i)} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \Rightarrow \text{Column Vector}$$

To know why it's a column vector ⇒

Check next slide

So given an i/p X :



(The example from a previous lesson)

$W^{(1)}$ here is $(4, 3)$ matrix 3 because of the i/p features
4 because of the # nodes

$x^{(1)}$ is a vector $(3, 1)$

Multiplying the will
result in a column vector.

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

$$\begin{aligned} a^{(2)} &= \sigma(z^{(1)}) \\ z^{(2)} &= W^{(2)} a^{(1)} + b^{(2)} \rightarrow \text{Value} \\ a^{(2)} &= \sigma(z^{(2)}) \end{aligned}$$

$$W^{(1)} X^{(1)} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ - \end{bmatrix}$$

$$W^{(1)} X^{(1)} = \begin{bmatrix} 0 \\ 1 \\ 1 \\ - \end{bmatrix}$$

$$W^{(1)} X^{(1)} = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

So now if we consider the training set Capital \vec{X}

$$\vec{X} = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & | \end{bmatrix} \Rightarrow \text{multiplying it by } W^{(1)}$$

$$W^{(1)} \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & x^{(3)} \\ | & | & | \end{bmatrix} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} =$$

$$\begin{aligned} \text{Assuming } b^{(1)} &= 0 \\ \begin{bmatrix} | & | & | \\ z^{(1)(1)} & z^{(1)(2)} & z^{(1)(3)} \\ | & | & | \end{bmatrix} &= \begin{bmatrix} | & | & | \\ z^{(1)(1)} & z^{(1)(2)} & z^{(1)(3)} \\ | & | & | \end{bmatrix} \\ &= Z^{(1)} \end{aligned}$$

We can also add $b^{(1)}$
won't change anything

Activation function

One of the choices when designing a NN is to choose the activation function. So far we have been using the Sigmoid function.

$$a = \frac{1}{1 + e^{-z}} \Rightarrow \text{Sigmoid function}$$

↳ the Sigmoid goes between 0 & 1

↳ Another activation function that works better most of the time

is the tanh function.

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

* The tanh function is a shifted version of the sigmoid function

It works better than the sigmoid because it has a mean of zero so the mean of the activations coming out of the hidden layer is close to zero.

⇒ We will see more in the second course

In the hidden layer tanh is almost always superior than the Sigmoid fn.

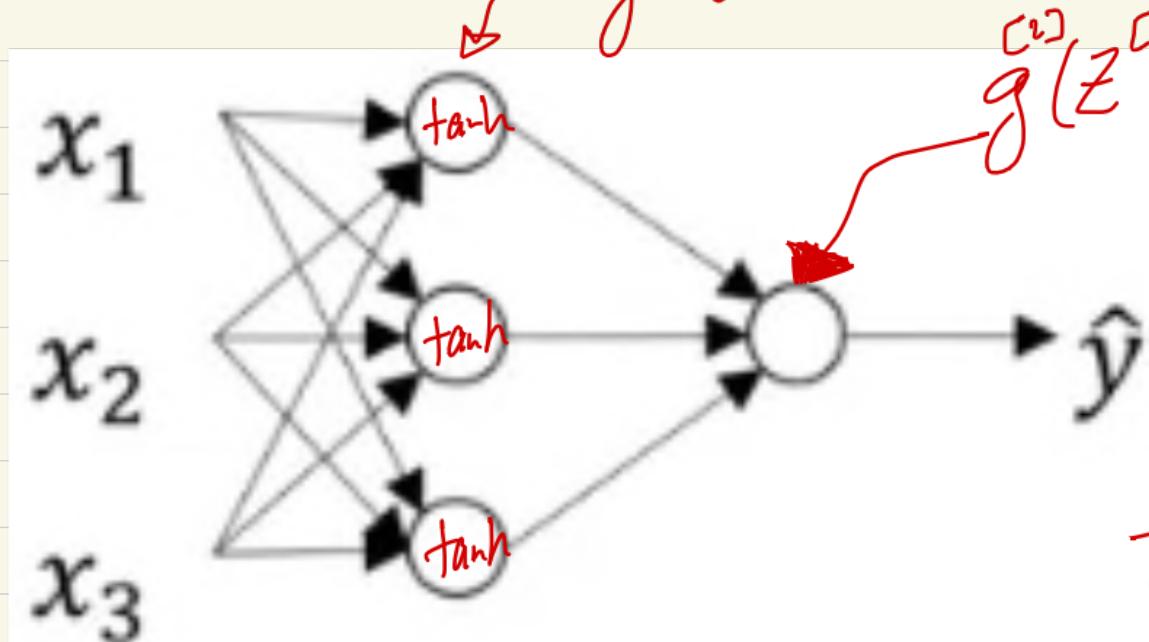
→ The only exception is for the output layer

↳ if y is 0 or 1 \Rightarrow So the Sigmoid is used or between 0 & 1 \Rightarrow

So we can use tanh for hidden layers & Sigmoid for o/p layer.

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

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



Problem with
Tanh & Sigmoid

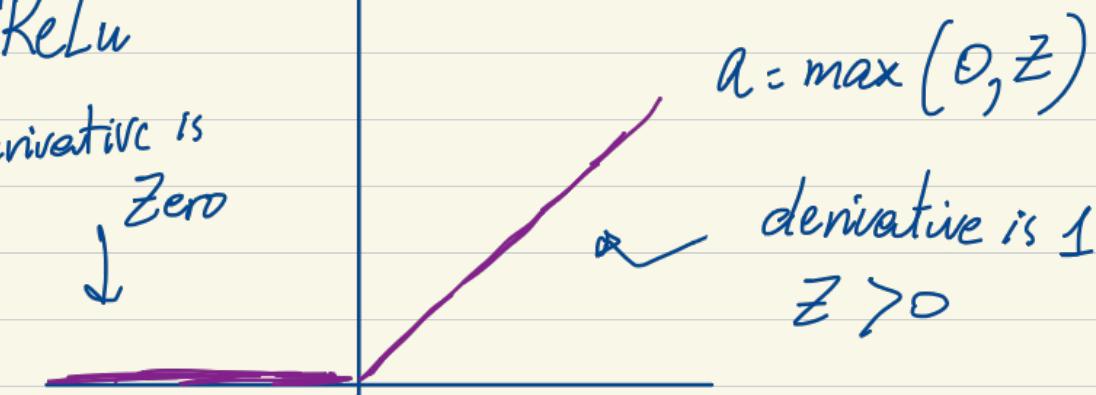
Again the problem of tanh & Sigmoid is that when z is very large or very small, the slope (gradient) is close to zero.

One other choice is

Rectified Linear Unit ReLU

ReLU

derivative is zero



Technically when $z=0$, the slope is **NOT DEFINED**

but we never get exactly $z=0$

A Rule of Thumb \Rightarrow If the output is 0 or 1

\downarrow \hookrightarrow Binary classif.
The output layer \Rightarrow Sigmoid

for all other units ReLU is the

default choice of activation function.

Disadv. of ReLU is the slope of 0 when z is -ve

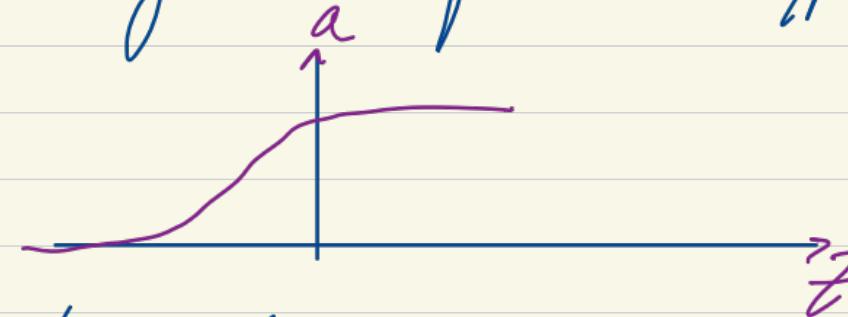
So there is another version of the ReLU called

Leaky ReLU

faster than
ReLU (or Leaky ReLU) makes the Algo learns a lot
the sigmoid or tanh

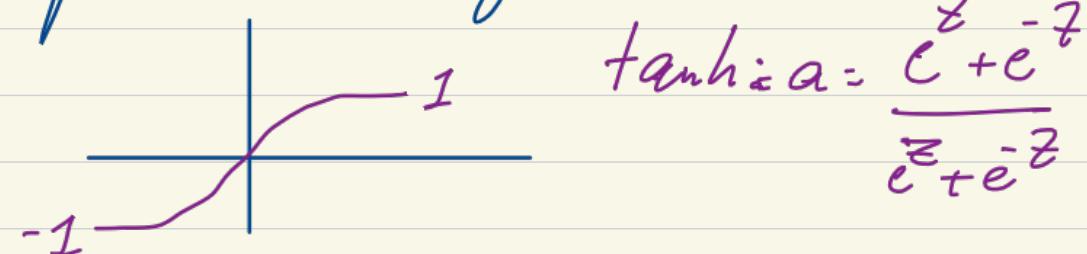
Pros & Cons of activation functions

- Never use Sigmoid except for the output layer



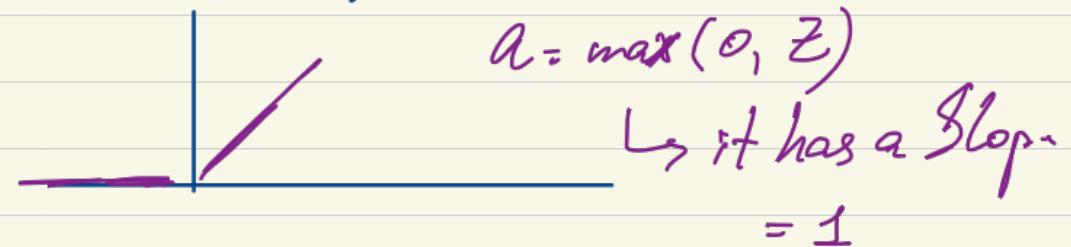
$$\text{Sigmoid} = \frac{1}{1 + e^{-z}}$$

- tanh is strictly superior than the sigmoid



$$\tanh: a = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

- We can always use tanh by default

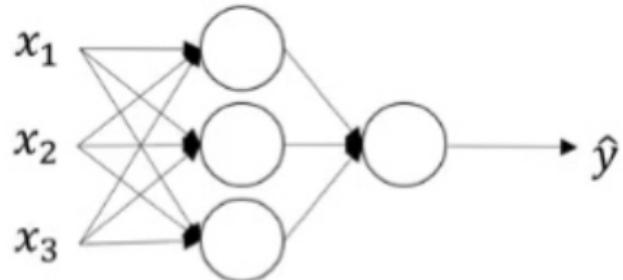


- We can also use Leaky ReLU

Why do you need nonlinear activation functions

*

Activation function



Given x :

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

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

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

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

let's try to get rid
of the function $g()$
& make it a linear fn

or $g(z) = z$ ↳ linear activation
function

It turns out that in this case

In this case my NN
is learning a linear function
of my l/p. NO matter how
many layers I have !!

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

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

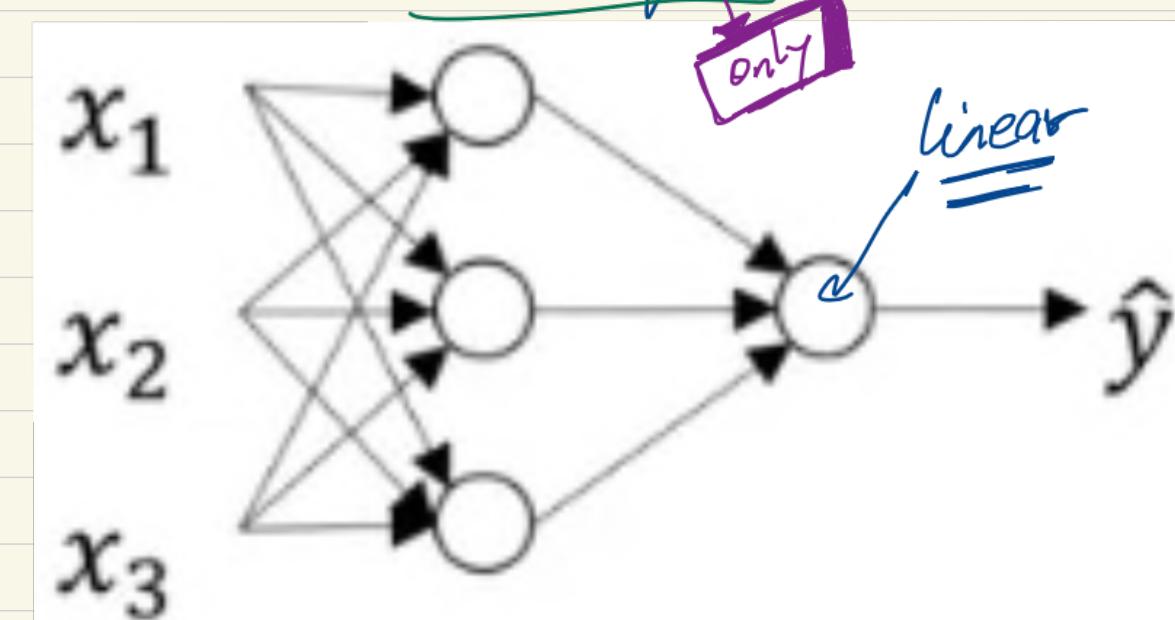
$$= W^{[2]}(W^{[1]}x + b) + b =$$

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

* So a Linear hidden layer is useless *

One reason to use where I might use a linear activation fn is in Regression problems. → where y takes a continuous value \Rightarrow a real number

↳ So I use it in [my o/p layer] → so my output can $\in]-\infty, +\infty[$



Derivatives of activation functions

When implement Backprop \Rightarrow I will need to calc the slope of activation function

Sigmoid function

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

So if $z=10$

$$g(z) \approx 1$$

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

\hookrightarrow Slope close to Zero **

I proved it before in
my notes.

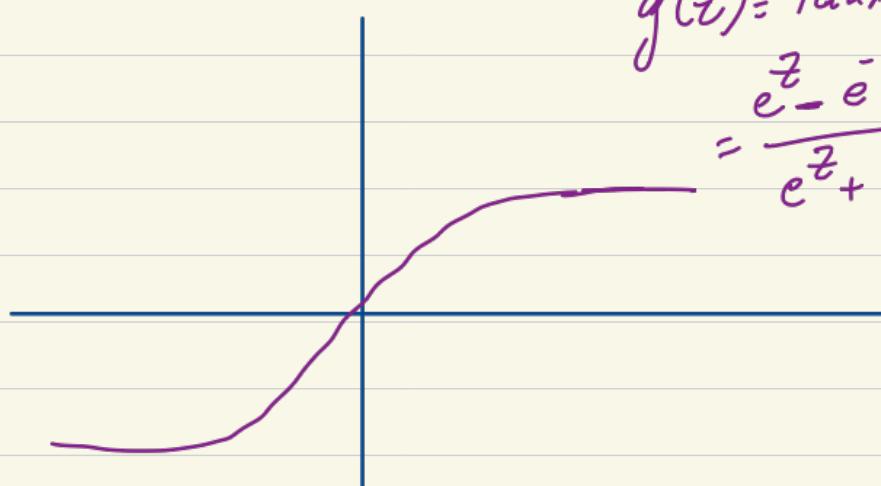
$$\text{if } z=0 \quad g(z) = \frac{1}{2}$$

$$\frac{d}{dz} g(z) = \boxed{\frac{1}{4}} *$$

$\hookrightarrow g'(z)$

Tanh activation function

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



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

$$\text{if } z=10 \quad \tanh(z) \approx 1 \quad g'(z) \approx 0$$

$$\text{if } z=-10 \quad \tanh(z) \approx -1 \quad g'(z) \approx 0$$

$$\text{if } z=0 \quad \tanh(z)=0 \quad g'(z) \approx 1$$

Relu activation function

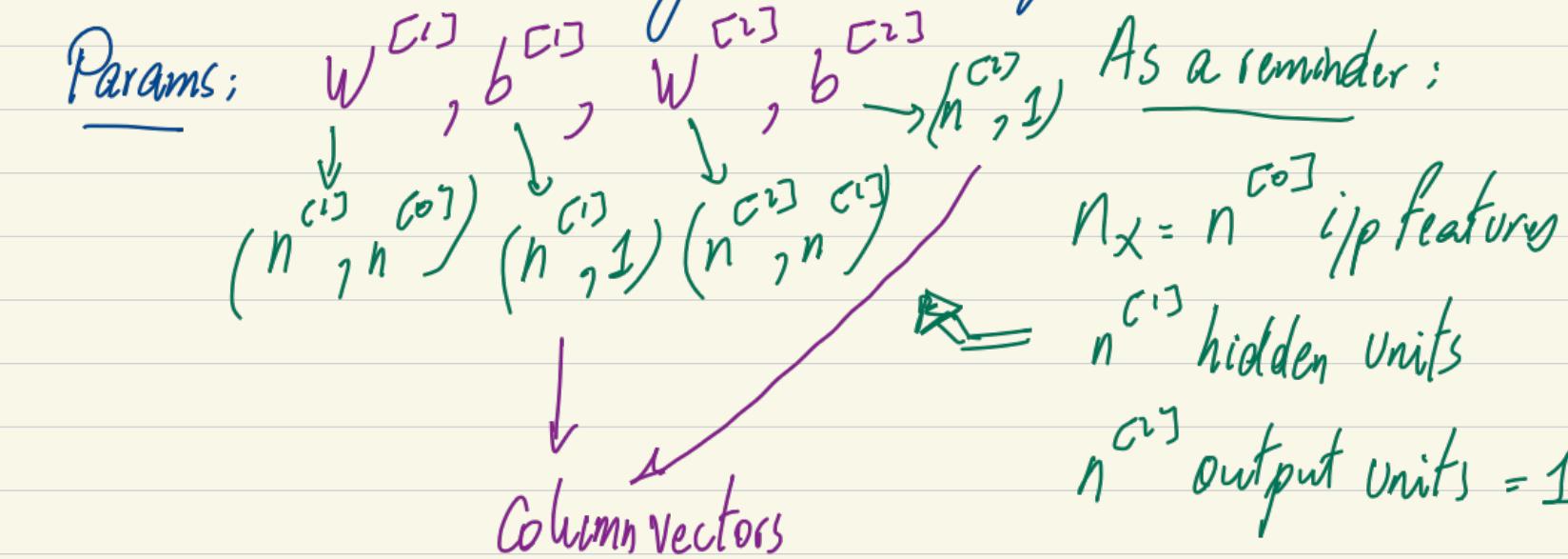
$$g(z) = \max(0, z)$$

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

↳ technically it will be $z=0.....$.
 (It won't be undefined) !!

Gradient descent for neural networks

Here we will see how to implement Gradient Descent for a NN with a single hidden layer



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

(If I am doing binary classification)

I can use the same loss function



$$\mathcal{L}(a_i, y) = -\left(y \log(a_i) + (1-y) \log(1-a_i)\right)$$

To train the algo I will need to perform gradient descent:

first I have to initialize the parameters (Randomly)

Repeat {

Compute preds \hat{y}_i $i: 1 \dots m$ for all examples

$$dW^{[1]} = \frac{dJ}{dW^{[1]}} \quad db^{[1]} = \frac{dJ}{db^{[1]}}$$

$$W^{[1]} = W^{[1]} - \alpha dW^{[1]} \quad \begin{matrix} \text{Compute} \\ \text{derivatives} \end{matrix}$$

$$b^{[1]} = b^{[1]} - \alpha db^{[1]}$$

$$W^{[2]} = W^{[2]} - \alpha dW^{[2]}$$

:

}

Formulas for computing derivatives:

• Forward propagation

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

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

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

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

Sigmoid
if I am
doing binary
classification

Back propagation

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

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

These equations
was seen in the
LR part

$$Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}] \quad \hookrightarrow \text{1x}m \text{ vector}$$

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

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

$(m \times 1^m)$ element wise product \hookrightarrow the derivative of the activation function in the hidden layer.

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

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

Course 1

Week 4

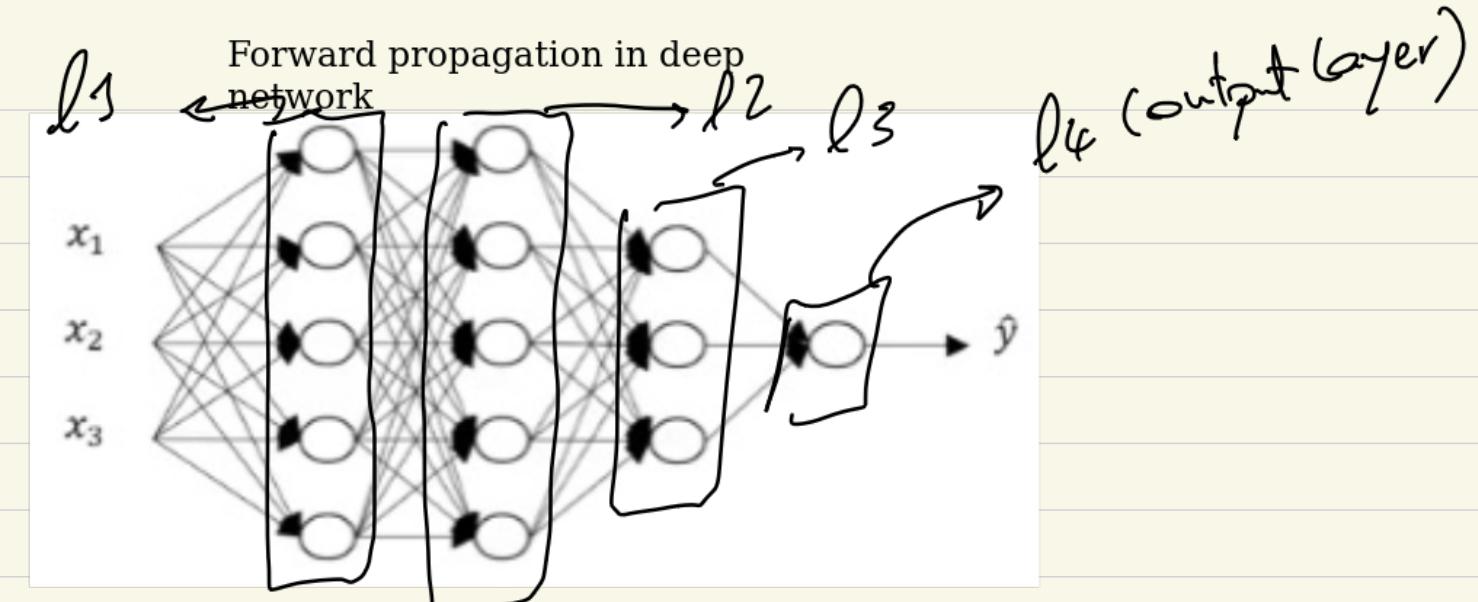
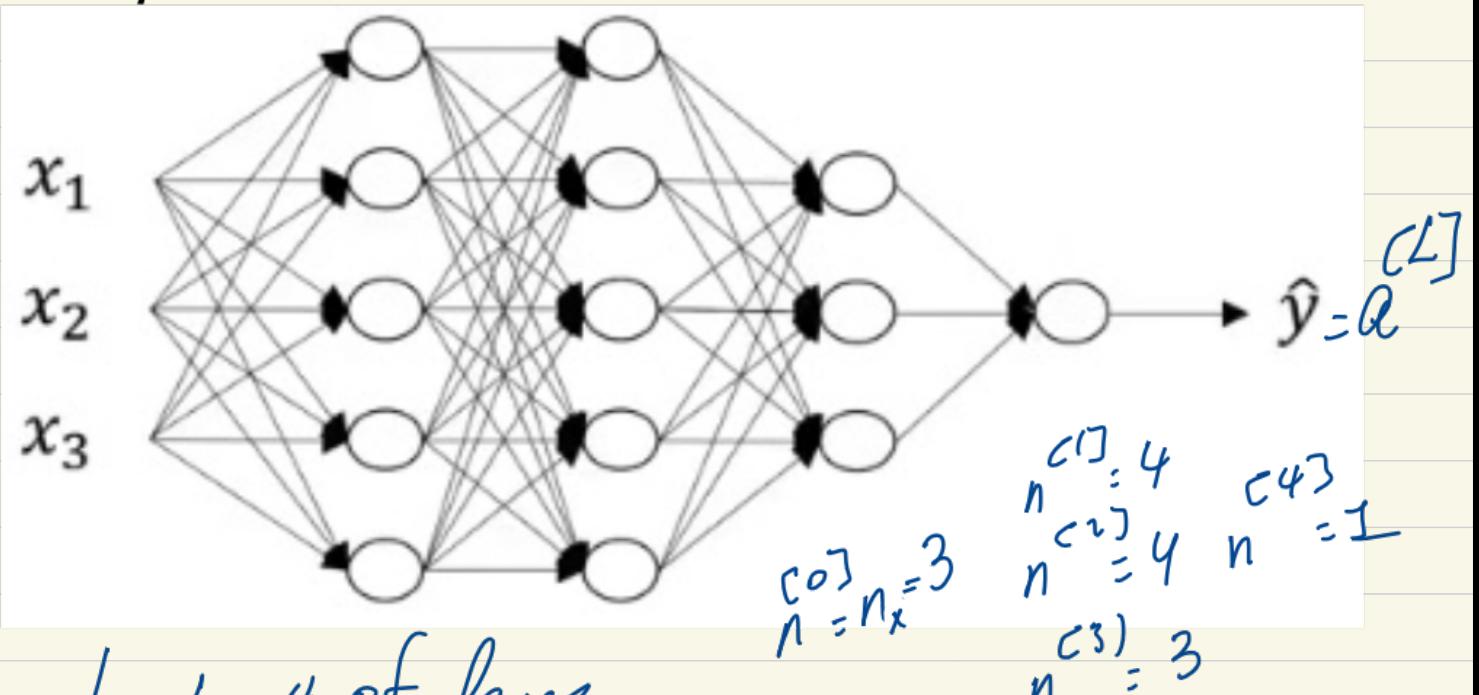
Deep L layer neural network

What is a deep neural network?

↳ We have more than 1 hidden layer

* A Logistic Regression is a shallow NN!! *

Deep Neural Network notation



* How to do forward propagation for a single example?

$$x: z^{[0]} = W^{[0]} X + b^{[0]} \rightarrow X = a^{[0]} \rightarrow \text{the activation of layer 0}$$

$$a^{[1]} = g^{[1]}(z^{[1]}) \Rightarrow \text{activation function}$$

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

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

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

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

Vectorization: (here we will use capital letters)

→ A reminder for a previous lesson
is in the next page.

So given an i/p X :

$W^{(1)} X^{(1)} = \begin{bmatrix} \cdot & \cdot \\ \cdot & \cdot \end{bmatrix}$

$W^{(1)} X^{(1)} + b^{(1)} = \begin{bmatrix} \cdot & \cdot \\ \cdot & \cdot \end{bmatrix}$

$Z^{(1)} = \sigma(W^{(1)} X^{(1)} + b^{(1)})$

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

$W^{(2)} = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix}$

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

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

$Y = g(Z^{(3)}) = A^{(3)}$

$Y = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix}$

$\text{Now if we consider the training set Capital } X$

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

$W^{(1)} X^{(1)} = \begin{bmatrix} \cdot & \cdot \\ \cdot & \cdot \end{bmatrix}$

$W^{(1)} X^{(1)} + b^{(1)} = \begin{bmatrix} \cdot & \cdot \\ \cdot & \cdot \end{bmatrix}$

$\Rightarrow \text{multiplying it by } W^{(1)}$

$W^{(1)} \left[\begin{array}{c|ccccc} 1 & x^{(1)} & x^{(2)} & \dots & x^{(n)} \\ \hline 1 & 1 & 1 & \dots & 1 \end{array} \right] = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix}$

$\text{Assuming } b^{(1),0} \Leftrightarrow$

$\text{We can also add } b^{(1)} \text{ won't change anything}$

When doing Vectorization:

$$Z^{(1)} = W^{(1)} A^{(0)} + b^{(1)}$$

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

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

$$A^{(2)} = g^{(2)}(Z^{(2)})$$

$$Y = g(Z^{(3)}) = A^{(3)}$$

we have

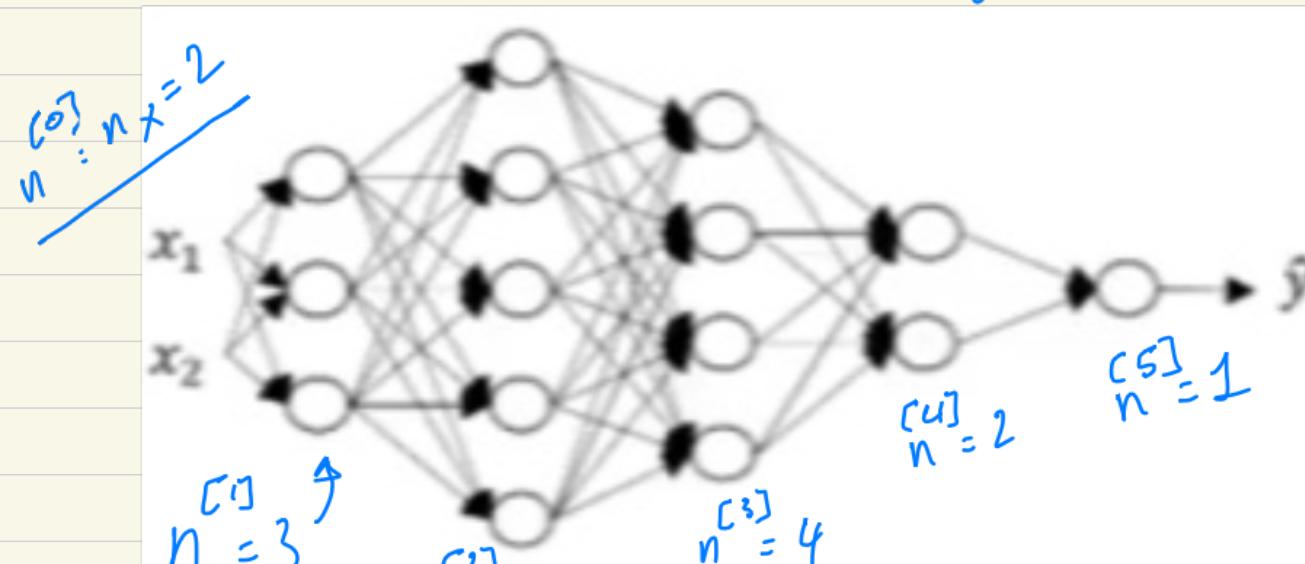
a for loop on the layers !! for $l=1\dots 4$

When doing forward prop
it's okay to use an explicit
for loop

To have a bug free implementation we have to
think somehow systematically about the
dimension of matrices !!

Getting your matrix dimensions right

$L=5$ # of layers (not counting the output layer as usual)



$$Z^{[l]} = W^{[l]} X + b^{[l]} \quad \left(\begin{array}{l} \text{(let's ignore the bias term)} \\ \text{for now} \end{array} \right)$$

$$\left. \begin{array}{c} \downarrow \\ \begin{matrix} (3,1) \\ (n^{[0]}, 1) \end{matrix} \end{array} \right\} \text{More generally: } W^{[l]} = \left(\begin{matrix} n^{[l]} \\ n^{[0]} \end{matrix} \right)$$

$$\left. \begin{array}{c} \downarrow \\ \begin{matrix} (3,2) \\ (n^{[0]}, 2) \end{matrix} \end{array} \right\} \left. \begin{array}{c} \downarrow \\ \begin{matrix} (2,2) \\ (n^{[0]}, 2) \end{matrix} \end{array} \right\} \left. \begin{array}{c} \downarrow \\ \begin{matrix} (2,1) \\ (n^{[0]}, 1) \end{matrix} \end{array} \right\}$$

$$W^{[l]} = \left(\begin{matrix} n^{[l]} & n^{[l-1]} \end{matrix} \right)$$

$$\hookrightarrow \text{Accordingly: } Z^{[l]} = \left(\begin{matrix} n^{[l]} \\ 1 \end{matrix} \right) *$$

What are the dimensions of the bias term $b^{[l]}$?

$$b^{[l]} = \left(\begin{matrix} n^{[l]} \\ 1 \end{matrix} \right)$$

Note that when doing backpropagation

$dW^{[l]}$ have the same dimensions of $W^{[l]}$

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

Let's continue with the dimensions of other terms:

$$a^{[l]} = g^{[l]}(Z^{[l]}) \quad a^{[l]} \text{ have the same dimensions of } Z^{[l]}$$

Let's look at vectorized implementation to consider several examples at the same time:

$$\text{For only one example: } Z^{[l]} = \frac{\left(\begin{matrix} n^{[l]} \\ 1 \end{matrix} \right)}{\left(\begin{matrix} n^{[l]} \\ n^{[0]} \end{matrix} \right)} X + \frac{\left(\begin{matrix} n^{[l]} \\ 1 \end{matrix} \right)}{\left(\begin{matrix} n^{[0]} \\ 1 \end{matrix} \right)} \rightarrow \left(\begin{matrix} n^{[l]} \\ 1 \end{matrix} \right)$$

for multiple examples we consider terms in Capital letters

\hookrightarrow In next Page

Let's look at vectorized implementation to consider

Several examples at the same time:

For only one example: $Z^{[1]} = W^{[1]} X + b^{[1]}$

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

for multiple examples (we consider terms in Capital letters):

* We stack the terms of all examples together

$$\hookrightarrow Z^{[1]} = \left[\begin{matrix} Z^{1} & Z^{[1](2)} & \dots & Z^{[1](m)} \end{matrix} \right]_{(n^{[1]}, m)}$$

$$Z^{[1]} = W^{[1]} X + b^{[1]} \rightarrow \text{the same}$$
$$\begin{matrix} \downarrow & \downarrow & \downarrow \\ (n^{[1]}, m) & (n^{[1]}, n^{[0]}) & (n^{[0]}, m) \end{matrix}$$

as for only 1 example

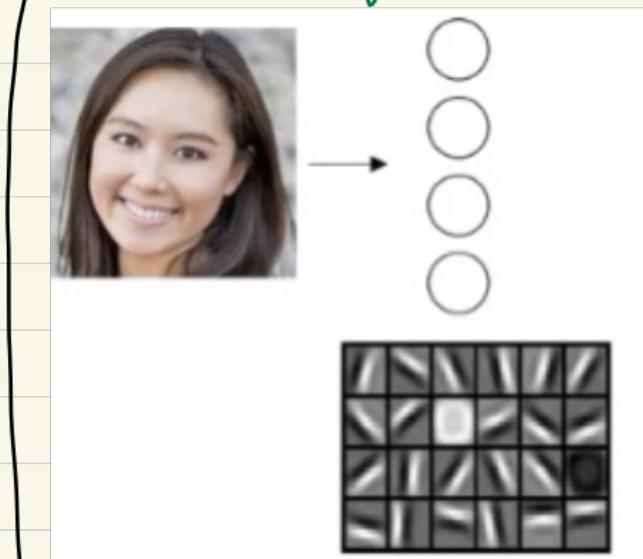
↓
It changes
↓
all the examples
↓
doesn't change
↓
change. It is added through Python broadcasting.

Why deep representations?

Why deep networks are better than shallow nets?

Intuition about deep representation

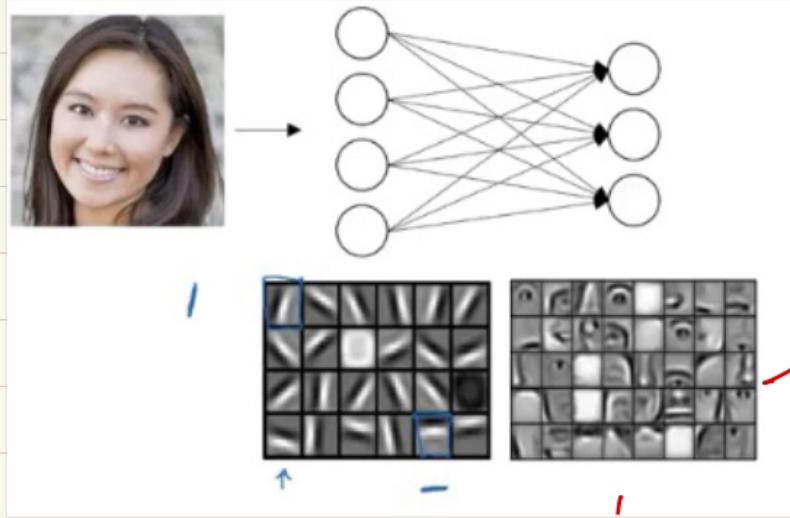
↪ if I am building a system for face detection, the first layer can be an edge detector



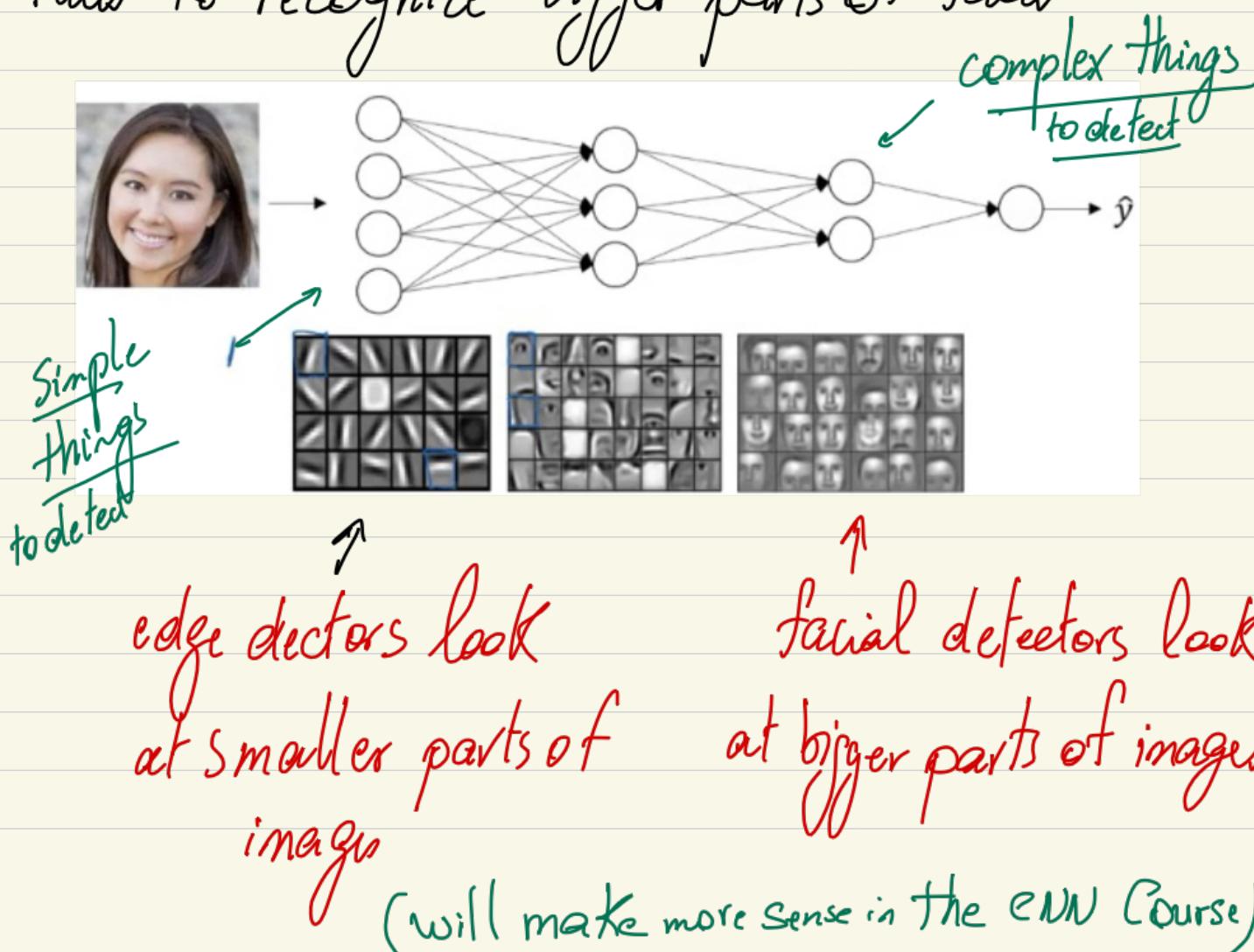
↓
edge detection by 24 hidden units

↪ A second layer may be used to group detected edges from previous layer to find parts of faces

↪ example in next page



Finally, last layer may be useful to group parts of faces to recognize bigger parts of faces



group detected edges to form parts of faces (eyes, nose, etc--)

This hierarchical schema can also be applied on different types of data!!
(Audio for example)

Circuit Theory & 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.

of hidden layers is small!!

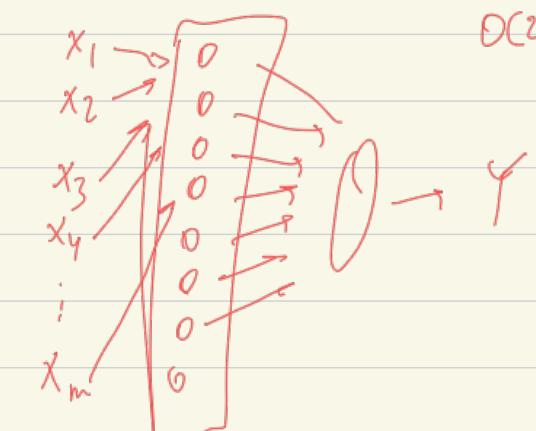
Let's say we want to compute XOR of all my i/p features:

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

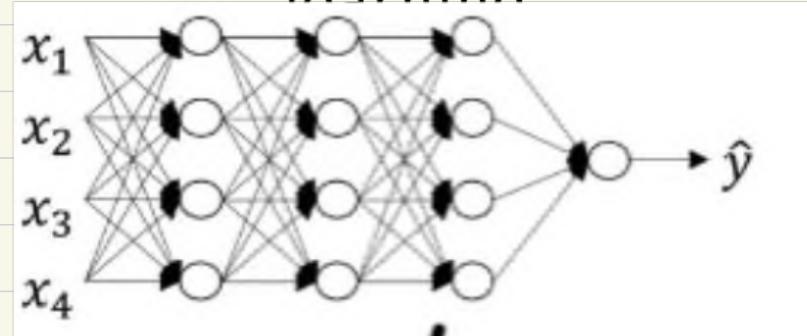
$x_1 \rightarrow (\text{XOR})$ the parity of all i/p
 $x_2 \rightarrow (\text{XOR})$
 $x_3 \rightarrow (\text{XOR})$

$x_4 \rightarrow (\text{XOR})$
 $x_5 \rightarrow (\text{XOR})$
 $x_6 \rightarrow (\text{XOR})$
So here the depth of the network will be of order $\log N$

If I am forced to compute the same function with only one layer!!
The # of neurons will be 2^{n-1}
 $O(2^n)$



Building blocks of deep learning



Layer l ; we have parameters: $W^{[l]}, b^{[l]}$

Forward: i/p: $a^{[l-1]}$, output: $a^{[l]}$

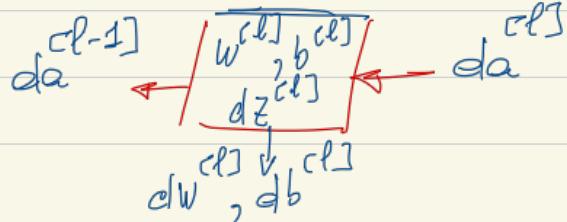
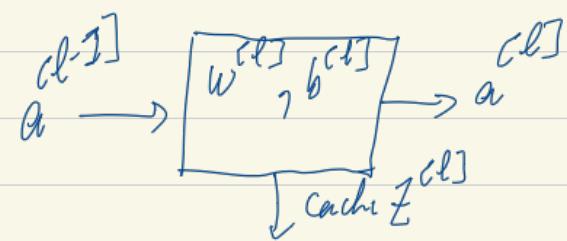
$$z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}$$

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

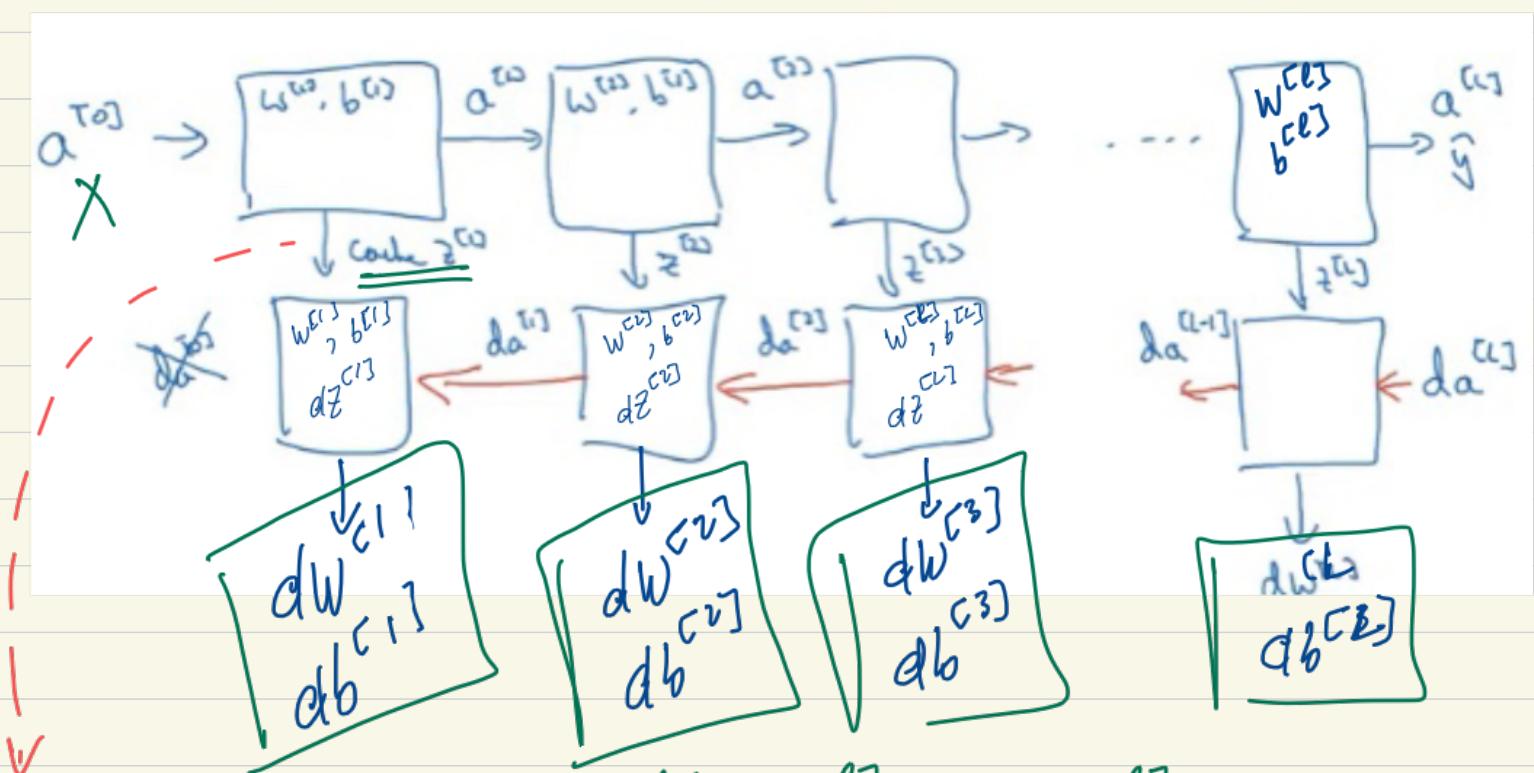
Backprop: Input: $da^{[l]}$
 $+ z^{[l]}$ output: $da^{[l-1]}$
 $dW^{[l]}$
 $db^{[l]}$

So This is the basic structure for implementing forward & backward steps

Layer l :



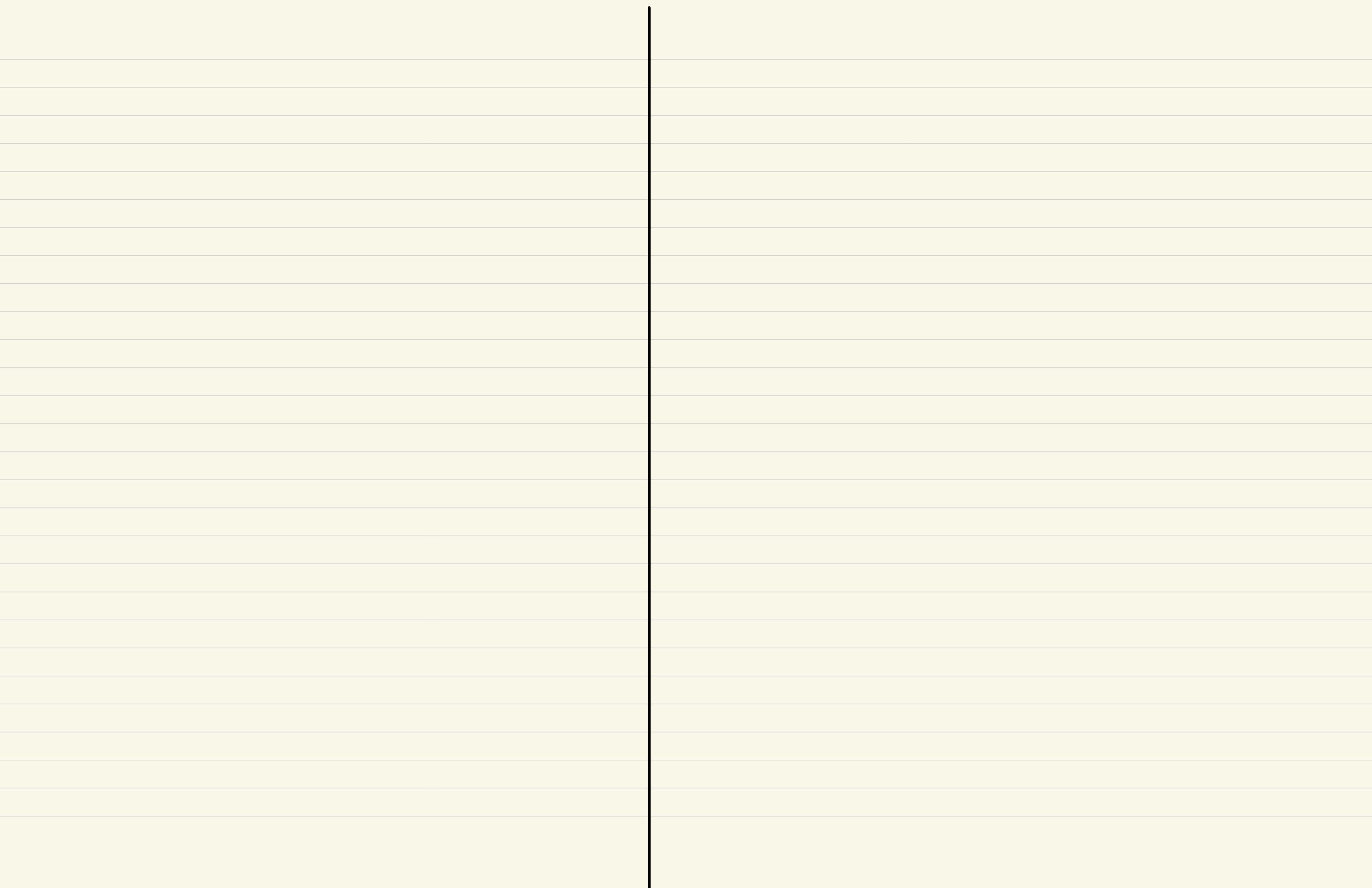
To have a fuller overview:



We store
in the Cache
 $Z^{[l]}, W^{[l]}, b^{[l]}$

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

$$b^{[l]} = b^{[l]} - \alpha db^{[l]}$$



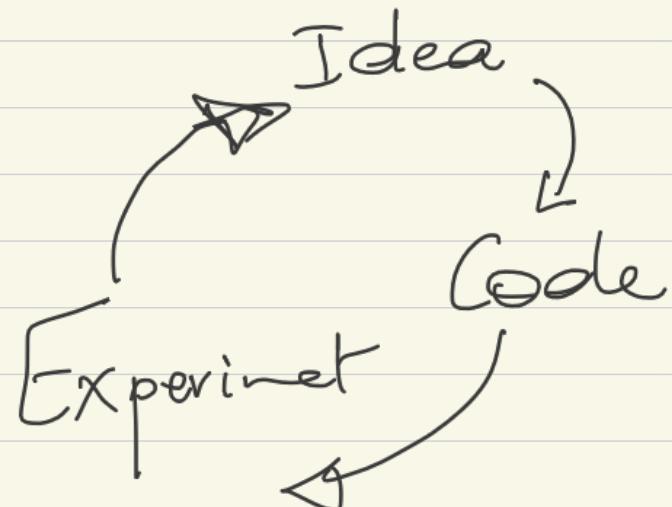
Train / Dev / Test sets



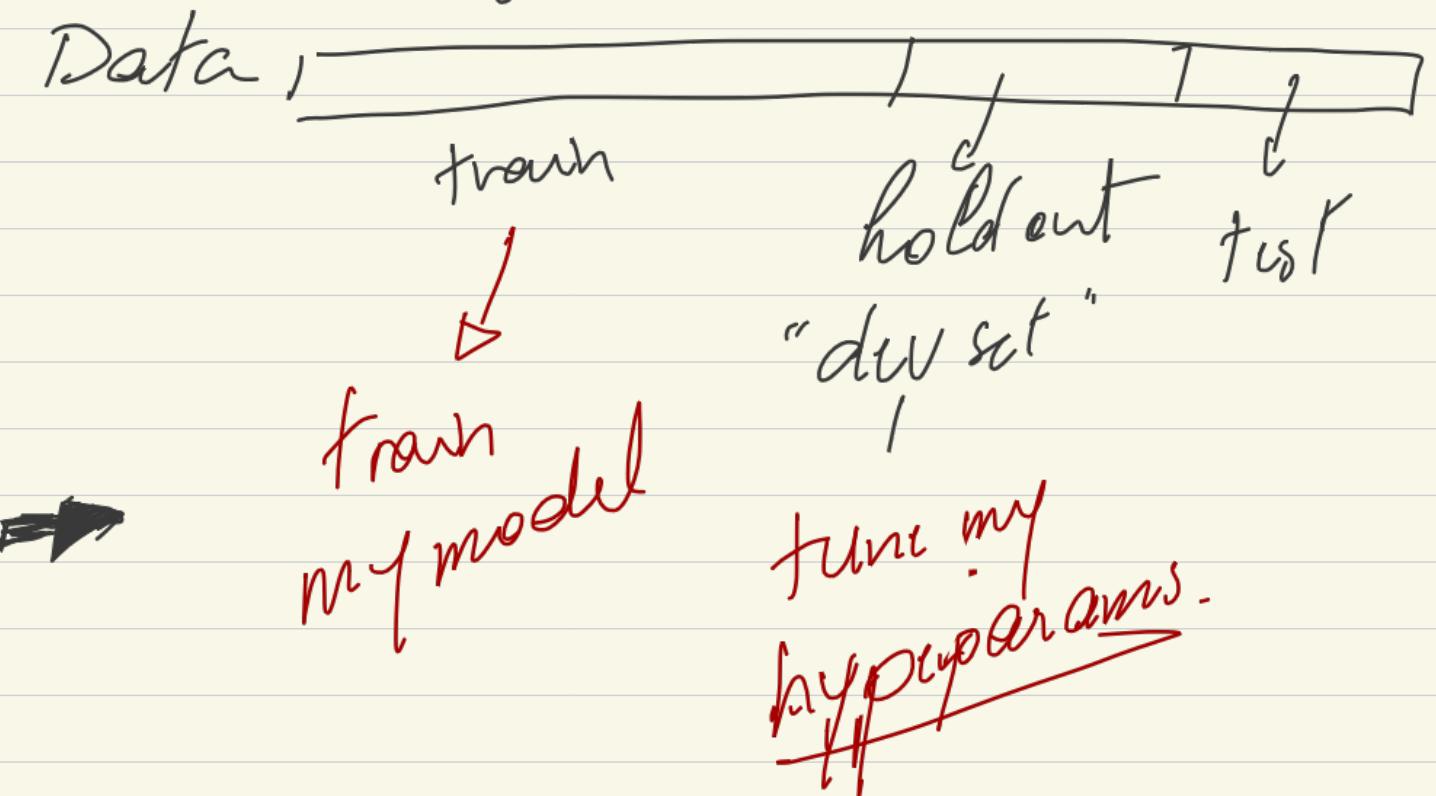
Applied ML is highly iterative Process

layers
hidden units
learning rates
activation fn

} choices of
Several parameters.



Traditionally:



Ratios can be 60%, 20%, 20%

Big Data → We can increase the ratio
of training data 98% / 1% / 1%

* Mismatched train / test distro is another

problem * Make sure we don't have
a deployment bias!!

↳ the Rule of thumb that dev & test sets
come from the same distro.

* We can also not have a test set (& it's okay :)

