


COURSERA COURSE

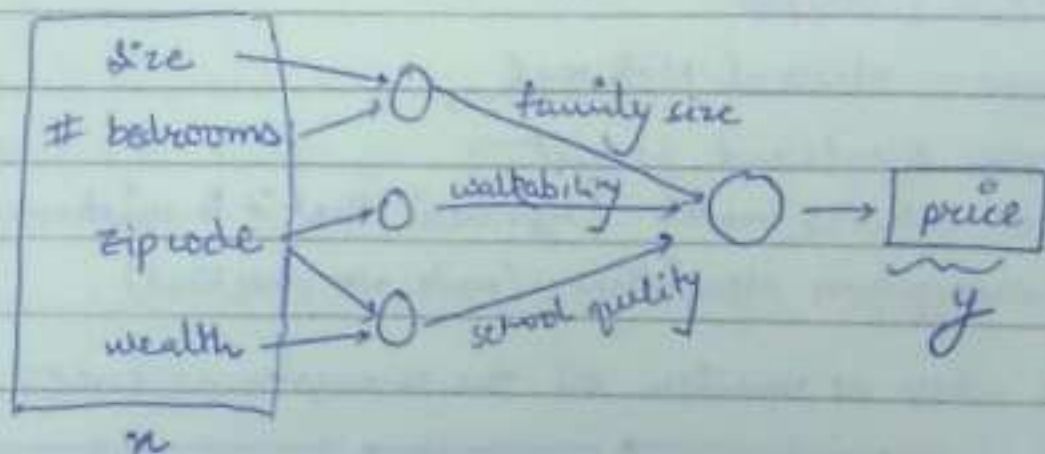
Course 1 (NN & DL)

DEEP LEARNING . ai

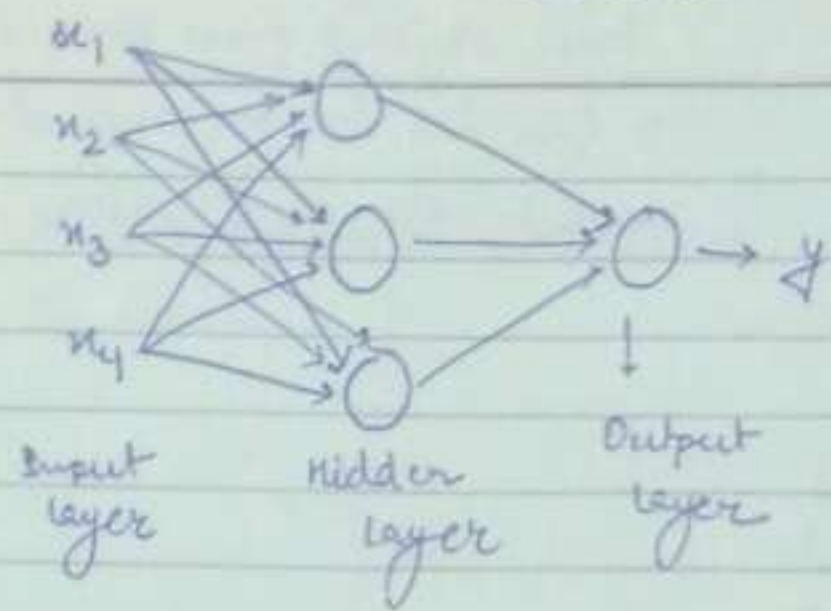
Input \rightarrow  \rightarrow Output
(Neuron)

Simplest Neural Network

All we need to feed in a neural network is input and output for entire training set. Rest everything is taken care of by itself. e.g.



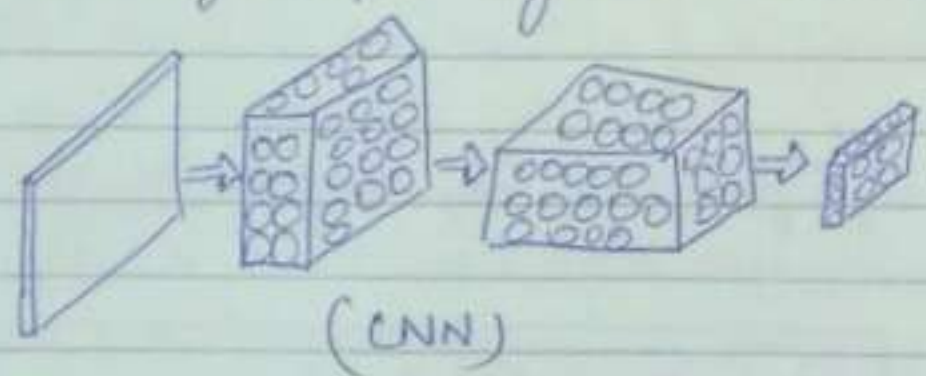
Neural network: very powerful and useful in supervised learning activities.



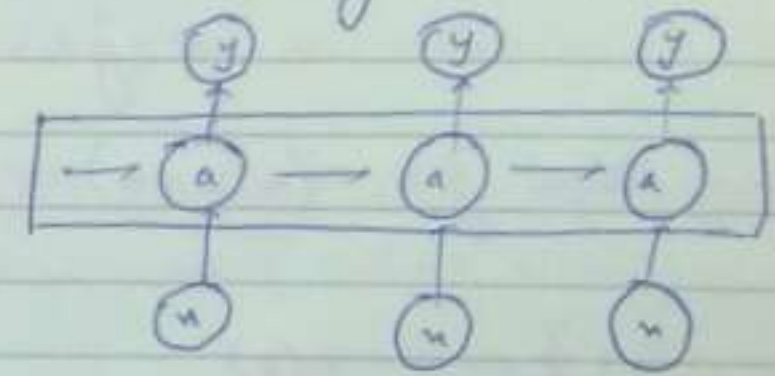
every input layer feature is connected with every hidden layer feature.

(standard NN)

for image data, we generally use CNN, and for sequence data like audio, we use RNN, for language also as the alphabets or words come one at a time and hence is sequence data. Radio or image input might need some custom or hybrid NN.



(CNN)

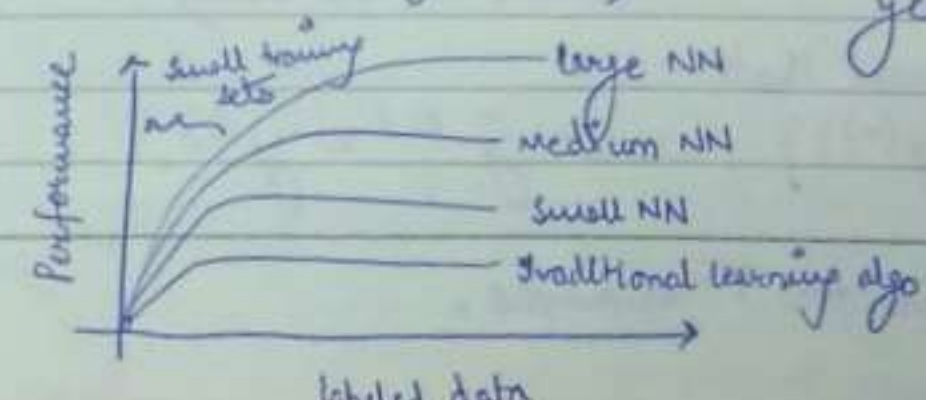


(RNN)

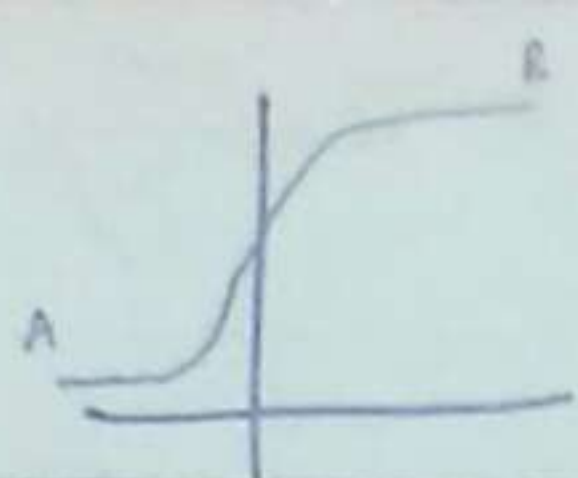
Supervised learning

- Structured data (in form of database)
- Unstructured data (in form of audio, image, text)

Deep learning and NN has helped computers better understand unstructured data in recent years.



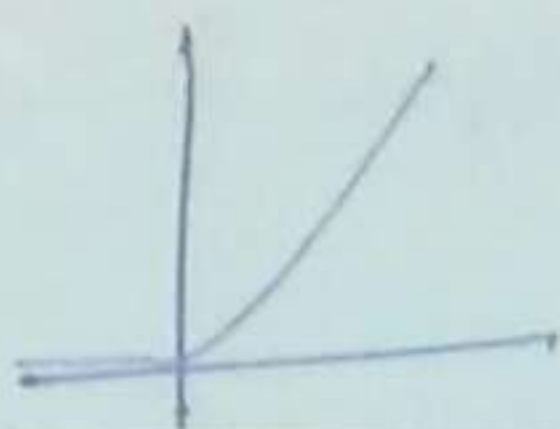
Large data and more complex NN led to the rise of deep learning.



sigmoid fn

drawbacks:

Region A and B where the gradient is nearly zero, learning becomes really slow



ReLU fn

(value fn of the rectified linear unit)

In this, gradient is 1 for all +ve values of input and is very less likely to gradually shrink to 0.

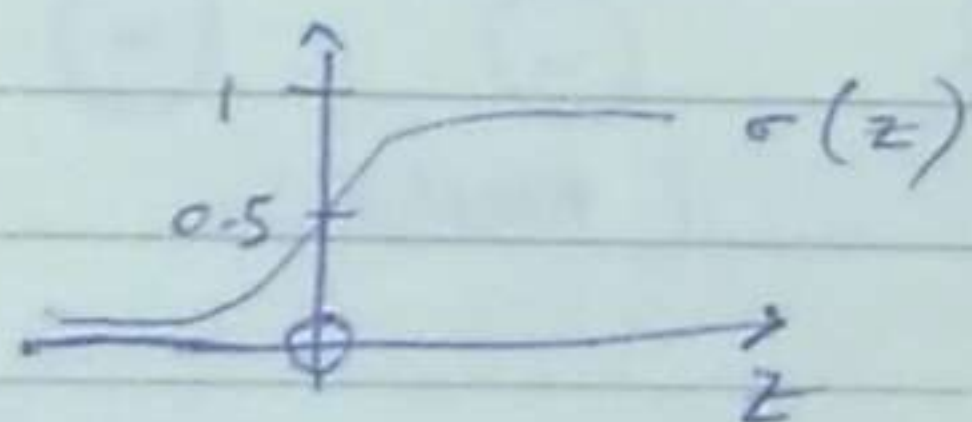
Hence, shifting from sigmoid to ReLU has led to gradient descent work much faster

→ logistic regression:

Given X , we want $\hat{y} = P(y=1|X)$

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

$$\text{Output } \hat{y} = \sigma(\underbrace{w^T x + b}_z) = \frac{1}{1 + e^{-z}} \quad (\text{So, } 0 \leq \hat{y} \leq 1)$$



$$\text{Also, } \hat{y} = \sigma(\theta^T x)$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_{n_x} \end{bmatrix} \left\{ \begin{array}{l} \theta_0 \leftarrow (\text{bias}) \\ \theta_1, \theta_2, \dots, \theta_{n_x} \leftarrow w \end{array} \right. \quad [x_0 = 1 \text{ added}]$$

(won't use this notation anywhere further in the course)

The model needs to be trained such that:

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

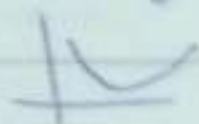
$x^{(i)}, y^{(i)}, z^{(i)}$ denote notations for i^{th} example.

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

* might not usually used in logistic regression as it makes gradient

descent to not work well. So, we define

L as $-(y \log \hat{y} + (1-y) \log(1-\hat{y}))$ that gives us an optimizable problem that is convex.



[we need convex bcoz such problems only have global minima]

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

$$\text{if } y=0, L(\hat{y}, y) = -\log(1-\hat{y})$$

if $y=1$, we want $-\log \hat{y}$ to be as small as possible, i.e. \hat{y} to be as large as possible and \hat{y} is sigmoid i.e. it can be maximum 1 which we actually want.

if $y=0$, we want $-\log(1-\hat{y})$ to be as small as possible, i.e. \hat{y} to be as small as possible, as \hat{y} is sigmoid, \hat{y} can be minimum 0 which is actually our target.

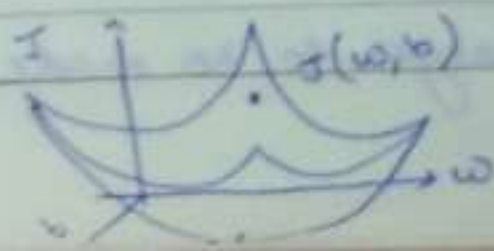
* loss fn measures the error on a single training set, but cost fn i.e. J measures how well we perform on entire training set.

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

cost of our parameters

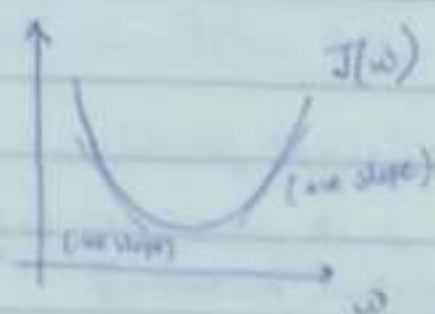
So, while training logistic regression model, we try to find w and b that minimize the overall cost fn J .

logistic regression \equiv very very small NN



convex fn, w and b are randomly initialized to some value and applying gradient descent, the lowest point is achieved

- Gradient Descent: It moves in the dirⁿ of steepest slope as quickly downhill as possible. It takes a step in 1 Iteration and another in next iterⁿ and do on.



repeat {

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

}

(Note: if slope is +ve, w will ↓ while coming down and if slope is -ve, w will ↑)

for $J(w, b)$ →

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

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

f/w propagaⁿ pass helps calculate the value of J and b/w propagaⁿ pass helps calculate the derivative of J and get final value of dp.

$$\frac{dJ}{da} = 3 \quad a = 5$$

$$\frac{dJ}{db} = 6 \quad b = 3$$

$$c = 2$$

$$u = bc$$

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

$$= \frac{dJ}{du} + \frac{dv}{du}$$

$$v = a + u$$

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

$$J = 3v$$

(backpropagaⁿ)

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

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

16/3/18

logistic regression gradient decent

$$x_1$$

$$w_1$$

$$x_2$$

$$w_2$$

$$b$$

$$z = w_1 x_1 + w_2 x_2 + b$$

$$a = \sigma(z)$$

$$L(a, y)$$

$$\frac{dL}{dz} = \frac{dL(a, y)}{dz}$$

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

$$= \frac{a-y}{a}$$

$$\frac{\partial L}{\partial w_1} = x_1 \cdot dz$$

$$dw_2 = x_2 \cdot dz$$

$$db = dz$$

$$\begin{cases} w_1 = w_1 - \alpha dw_1 \\ w_2 = w_2 - \alpha dw_2 \\ b = b - \alpha db \end{cases}$$

Through b/w propagation we get to know the changes to be made in w_1 , w_2 and b for a closer output.

Note Vectorization techniques help us get rid of explicit for loops in our code to go through the entire training set.

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
```

$w^T x$

(much faster than non-vectorized implementation)

GPU } single instruction,
CPU } multiple data.

Note Always avoid using for loops whenever possible.

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

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

```
for i in range(n):
```

```
    u[i] = math.exp(v[i])
```

vectorized:

```
import numpy as np
```

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

- Vectorizing logistic regression:

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

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

$$\text{So, } Z = w^T X + b$$

$$\text{i.e. } [z^1 \ z^2 \ \dots \ z^m] = w^T X + [b \ b \ \dots \ b]$$

$(1 \times m)$

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

$$A = [a^1 \ a^2 \ \dots \ a^m] = \sigma(Z) \quad \left\{ \begin{array}{l} \text{no looping} \\ \text{needed} \end{array} \right.$$

$$dz = A - y$$

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

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

Q Calories from carbs, proteins, fats in 100g of different foods:

$$\begin{matrix} & \text{Apples} & \text{Beef} & \text{Eggs} & \text{Potatoes} \\ \begin{matrix} \text{Carb} \\ \text{Protein} \\ \text{Fat} \end{matrix} & \begin{bmatrix} 58.0 & 0.0 & 4.4 & 68.0 \\ 1.2 & 104.0 & 52.0 & 8.0 \\ 1.8 & 135.0 & 99.0 & 0.9 \end{bmatrix} & = & A & \\ & & & & (3,4) \end{matrix}$$

Calculate % of calories from Carb, Protein, Fats without using a for loop.

Ans $\text{cal} = A \cdot \text{sum}(\text{axis}=0)$ // python will sum vertically.
 $\Rightarrow [59. \quad 239. \quad 155.4 \quad 76.9]$

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

↳ (not necessary in this case as matrix is already 11,4)

→ Broadcasting examples:

$$i) \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + 100 \equiv \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + \begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \end{bmatrix} \quad (\text{in python}) = \begin{bmatrix} 101 \\ 102 \\ 103 \\ 104 \end{bmatrix}$$

(auto expansion)

$$ii) \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + [100 \quad 200 \quad 300] = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix}$$

$$iii) \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 \\ 200 \end{bmatrix} \equiv \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 100 & 100 \\ 200 & 200 & 200 \end{bmatrix}$$

General principle: $(m,n) \xrightarrow{+} (1,n) \cong (m,n) \xrightarrow{+} (m,n)$
 $(m,n) \xrightarrow{+} (m,1) \cong (m,n) \xrightarrow{+} (m,n)$

Note $a = \text{np.random.randn}(s)$. Here 'a' is a rank 1 array in python and is either a column or a row vector.

$a = a.T$ and $\text{np.dot}(a, a.T)$ is a single number
 $a \cdot \text{shape} = (s,) \leftarrow (\text{shape})$

while, $a = \text{np.random.randn}(5, 1)$ is a $(5, 1)$ column vector and
 now a^T is a $(1, 5)$ row vector and
 $\text{np.dot}(a, a^T)$ is $(5, 5)$ vector.

Note vector is contained in two square brackets - $[[...]]$ } **
 array " " " 1 " " " [...]

to avoid making ^{rank 1} arrays, rather make vectors as their behaviour
 is easy to understand. or matrices

If sometimes, we end up with a rank 1 array, we can always
 reshape it to a vector as $a = a.\text{reshape}(5, 1)$.

$\text{assert}(a.\text{shape} == (5, 1))$ { use such statements to assert anything }

COURSE 1 (Week 2)

Programming Assignment FAQ

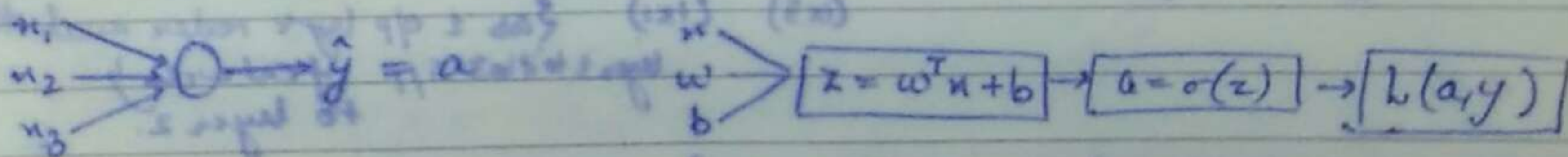
B, D, B, C, D, A, C, D, A, B

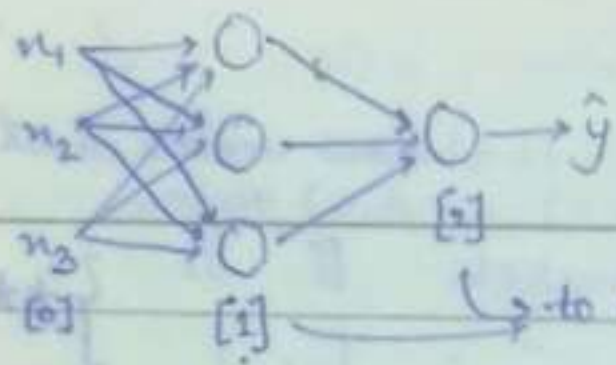
Note In ML or DL we don't use math library because it takes veel number
 inputs. In ML we generally deal with arrays and so numpy library
 is needed.

Note gradient of sigmoid f.w.r.t input $x = \sigma(x) (1 - \sigma(x))$

Note $\text{softmax}(u) = \text{softmax}([u_1, u_2, \dots, u_n]) = \left[\frac{e^{u_1}}{\sum_j e^{u_j}}, \frac{e^{u_2}}{\sum_j e^{u_j}}, \dots, \frac{e^{u_n}}{\sum_j e^{u_j}} \right]$

Neural Network:





Note $X = a^{[0]}$

Combination of many sigmoid fns give a neural network?

to represent layer 1 and 2

$$\begin{aligned}
 & \text{Forward pass: } z^{[1]} = W^{[0]}x + b^{[1]} \Rightarrow a^{[1]} = \sigma(z^{[1]}) \Rightarrow z^{[2]} = W^{[1]}a^{[1]} + b^{[2]} \Rightarrow a^{[2]} = \sigma(z^{[2]}) \\
 & \text{Loss: } L(a^{[2]}, y) \\
 & \text{Backward pass (derivatives): } \frac{\partial L}{\partial a^{[2]}} \Rightarrow \frac{\partial L}{\partial z^{[2]}} \Rightarrow \frac{\partial L}{\partial a^{[1]}} \Rightarrow \frac{\partial L}{\partial z^{[1]}} \Rightarrow \frac{\partial L}{\partial W^{[0]}} \Rightarrow \frac{\partial L}{\partial b^{[1]}}
 \end{aligned}$$

So, backward calculations are done in neural network for updating weights and biases.

[0]: Input layer (input features)

[1]: Hidden layer (things in hidden layer are not seen in the training set)

[2]: Output layer

$$a^{[1]} = \begin{bmatrix} a^{[1]}_1 \\ a^{[1]}_2 \\ a^{[1]}_3 \end{bmatrix} \quad (\text{for hidden layer in above case})$$

We don't count the i/p layer in NNs and so the above example is a 2-layered NN. The hidden layer has parameters w and b associated with it ($w^{[1]}, b^{[1]}$)

$(3,3)$ $(3,1)$ \rightarrow as we have 3 i/p feature and 3 hidden layer nodes.

Similarly o/p layer also has $w^{[2]}, b^{[2]}$ associated with it.

$(1,3)$ $(1,1)$ as 1 o/p layer nodes and the 3 features to layer 2.

Now, vectorized implementation of this can be written as:

Given x :

$$\left. \begin{aligned} z^{[1]} &= W^{[1]} x^{[0]} + b^{[1]} \\ a^{[1]} &= \sigma(z^{[1]}) \end{aligned} \right\} \text{output of layer 1}$$

$$\left. \begin{aligned} z^{[2]} &= W^{[2]} a^{[1]} + b^{[2]} \\ a^{[2]} &= \sigma(z^{[2]}) = \hat{y} \end{aligned} \right\} \text{output of layer 2.}$$

Now, we need to compute outputs on all the examples or pretty much all at the same time. If we have m training examples instead of 1,

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

So, $a^{[2]}(i)$ refers to training e.g. i
refers to the 2nd layer.

So, for $i = 1$ to m :

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

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

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

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

vectorizing
this, we get:

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ \vdots & \vdots & \ddots & \vdots \end{bmatrix} \quad (n_x, m)$$

$$\text{So, } Z^{[1]} = W^{[1]} X + b^{[1]}$$

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

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

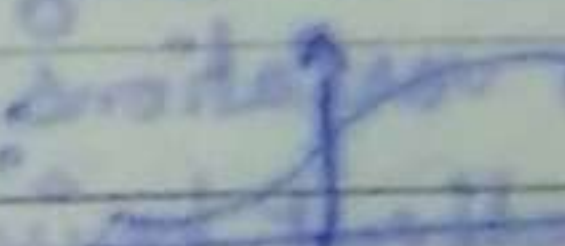
$$A^{[2]} = \sigma(Z^{[2]})$$

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

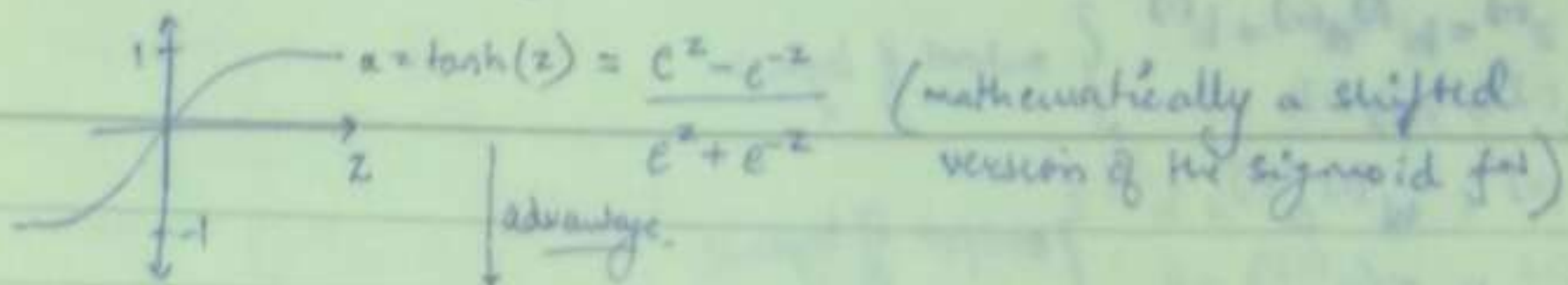
vertically, we have nodes in the layer (i.e. the different hidden units)
horizontally is the no. of training examples.

$$\text{Now } X = A^{[0]}$$

\Rightarrow Activation fn (which sigmoid is just an example)

Sigmoid fn  activation fn almost always work better than sigmoid fn and can be non-linear.

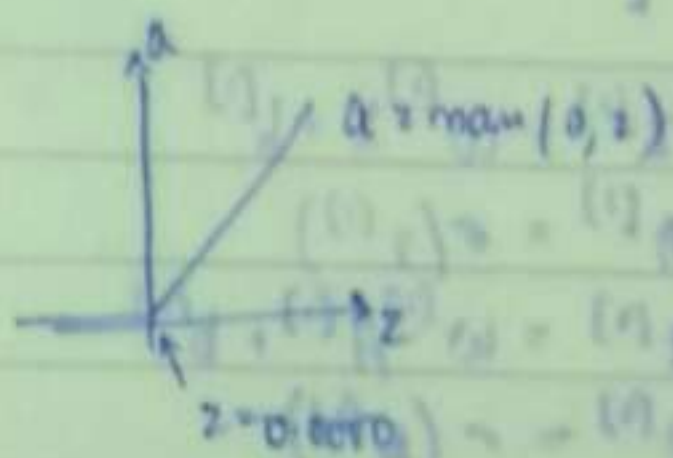
$a = \tanh(z)$ instead of $\sigma(z)$ is superior activⁿ fⁿ for sure



(mean of data comes closer to zero and hence easier to analyze)
while in case of σ (sigmoid fⁿ) data is centered about around 0.5 which makes learning for the next layer a little bit hard.

Note activⁿ fⁿ can be different for different layers. Sigmoid fⁿ used in o/p layer when classification is to be done, as the value varies b/w 0 & 1 and not -1 and 1 like the tanh fⁿ.

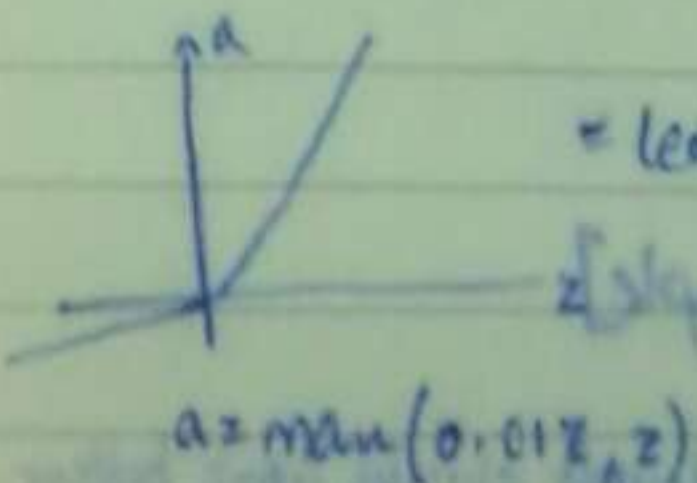
Drawback of these 2 activⁿ fⁿ is that when z is very large or very small, the slope of the fⁿ ends up being close to zero, so this can slow down gradient descent.



= ReLU (Rectified Linear Unit)

[Network learns much faster & ReLU as compared to Sigmoid or tanh]

for o/p layer, sigmoid is used by default in case of binary classification only. In all other cases generally ReLU is used now a days, or tanh is also sometimes used.



= Leaky ReLU. to avoid problems of dying ReLU.

beoz in dying ReLU, those neurons which entered the 0 activⁿ state of ReLU, will stop responding to variations in error / input beoz of gradient being 0 in that state.

⇒ Need for non-linear activaⁿ fn:

$a = g(z) = z$ (sometimes called linear activaⁿ fn^s)
or 'identity activaⁿ fn^s.'

↓

This means our model is computing \hat{y} has as a linear fn of the input features i.e. $a = z = wX + b$. So, no matter how many hidden layers are there in the NN, always a linear fn of the i/p is calculated hence a hidden layer (linear) is more or less useless, bcoz even combinaⁿ of linear fn is a linear fn, leading to no new fn being made in the hidden layers to understand the i/p-o/p relaⁿ.

It might be okay to have a linear activaⁿ fn if y is a regression variable, so one can use linear activaⁿ in the o/p layer if y is to be obtained in b/w $-\infty$ to $+\infty$ (i.e. regression). e.g. in housing prices example, Relu can also be used bcoz all prices are +ve and lie b/w 0 to ∞ .

⇒ Derivatives of activaⁿ fn^s:

① Sigmoid: $g(z) = \frac{1}{1+e^{-z}}$, So, $\frac{d}{dz} g(z) = \frac{1}{1+e^{-z}} \left(\frac{e^{-z}}{1+e^{-z}} \right) = g(z)(1-g(z))$
Slope of $g(z)$ or activaⁿ fn.
i) $z=10 \Rightarrow \frac{dg(z)}{dz} \approx 0$
ii) $z=-10 \Rightarrow \frac{dg(z)}{dz} \approx 0$
iii) $z=0 \Rightarrow \frac{dg(z)}{dz} \approx \frac{1}{4}$

② tanh: $g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$, So, $g'(z) = 1 - (\tanh(z))^2$
i) $z=10 \Rightarrow \tanh(z) \approx 1 \Rightarrow g'(z) \approx 0$
ii) $z=-10 \Rightarrow g'(z) \approx 0$
iii) $z=0 \Rightarrow g'(z) \approx 1$

③ Relu: $g(z) = \max(0, z)$, So, $g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z > 0 \end{cases}$
und. fixed if $z=0$

④ Leaky Relu: $g(z) = \max(0.01z, z)$, So $g'(z) = \begin{cases} 0.01 & \text{if } z < 0 \\ 1 & \text{if } z > 0 \end{cases}$

⇒ Gradient decent for NN:

We will know equations in order to get back propagation of the gradient decent working.

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

i/p features hidden layer units o/p units

Parameters: $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$
 $(n^{[0]}, n^{[1]}) \quad (n^{[1]}, n^{[2]})$
 $(n^{[0]}, 1) \quad (n^{[1]}, 1)$

(for a single hidden layer)

Cost fn: $J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}) = \frac{1}{n} \sum_{i=1}^n L(\hat{y}_i, y_i)$ (for binary classification)

Gradient decent: (to learn parameters for NN) less fn
 (to train the parameters)

L is same as that used in logistic regression.

Repeat { compute predic' ($\hat{y}^{(i)}, \tilde{y}^{(i)}$ for $i=1$ to m)

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

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

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

}

So,

Feed Propaga:

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

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

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

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

Back propaga:

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

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

$db^{[2]} = \frac{1}{n} \sum dz^{[2]}, \text{axis}=1, \text{keepdims}=\text{True}$

for summing matrix horizontally

(to ensure that o/p vectors are of form $(n, 1)$ & not $(1, n)$)

Back propaga' contd...

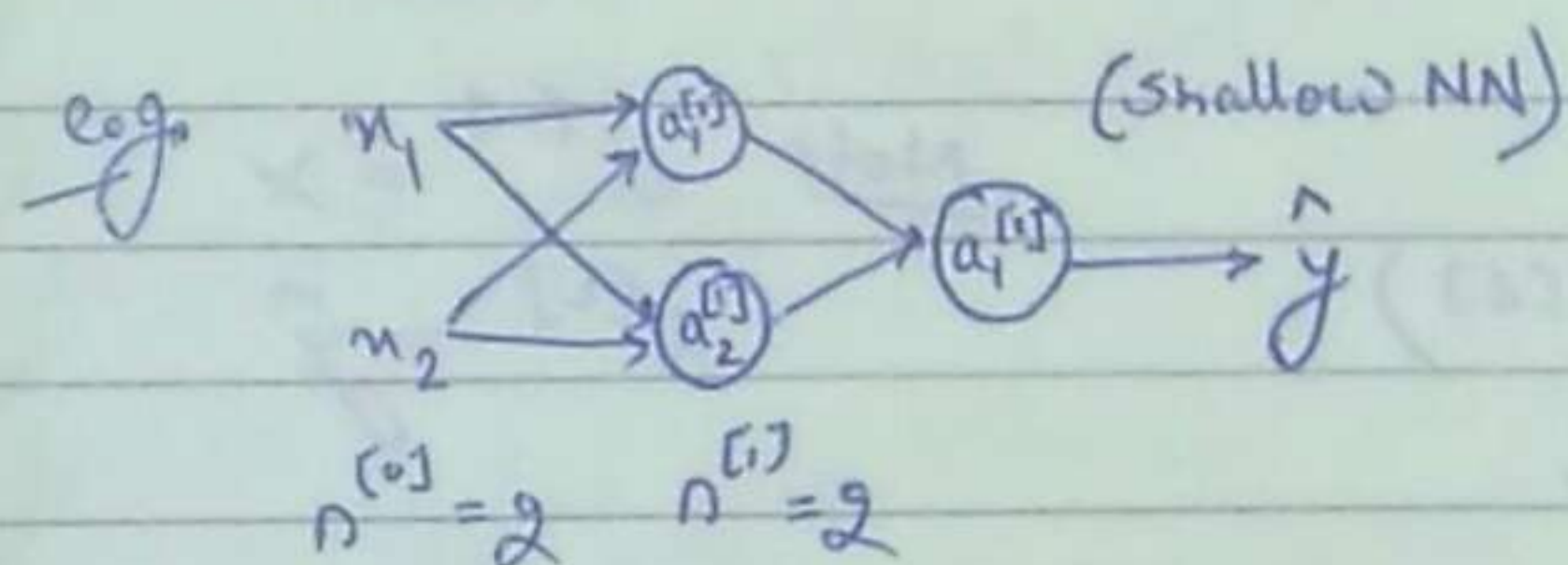
$$dz^{[l]} = \underbrace{w^{[l+1]T}}_{(n^{[l+1]}, m)} dz^{[l+1]} \underset{\substack{\text{elementwise} \\ \text{product}}}{*} \underbrace{g^{[l+1]'}(z^{[l]})}_{(n^{[l]}, m)} = da^{[l+1]} \underset{\substack{\text{elementwise} \\ \text{product}}}{*} g^{[l+1]'}(z^{[l]})$$

$$da^{[l-1]} = w^{[l]T} \cdot dz^{[l]}$$

$$dw^{[l]} = \frac{1}{m} dz^{[l]} X^T, \quad dw^{[l+1]} = dz^{[l+1]} \cdot a^{[l]}$$

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

Note Initializing parameters not to zero but randomly turns out to be very important for training neural network. Initializing to zero would work for logistic regression but not for NN, while applying Gradient Descent becoz:



$$w^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad \& \quad b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

↓
If this is the case, $a_1^{[1]} = a_2^{[1]}$

becoz both the hidden units are calculating exactly the same fn.

Also, $dz_1^{[1]} = dz_2^{[1]}$ while

back propaga'. As $a_1^{[1]} = a_2^{[1]}$,

output weights also seem to be equal.

i.e. $w_2 = [0 \ 0]$. So, after every 'time'

we will see, both the hidden units are

computing exactly the same fn,

and are symmetric. Hence, there

is no point having more than 1

hidden unit with weights initialized

to zero.

Soluⁿ is to initialize the parameters randomly.

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

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

Note becoz b doesn't have symmetry problem, as long as z is initialized randomly.

we multiplied $w^{[1]}$ by 0.01 becoz

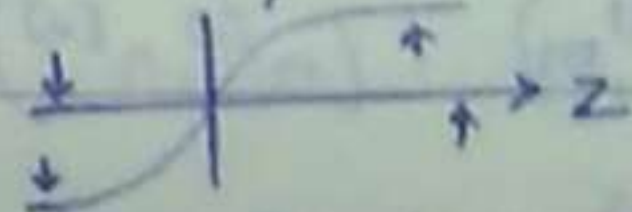
we initialize weights with very

very small values (random)

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

will have very small gradient

or slope.



as $z^{[1]}$ is very large or very small

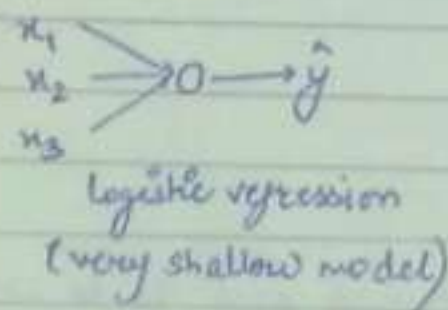
for large $w^{[1]}$, meaning learning

will be very slow for less $da^{[1]}$,

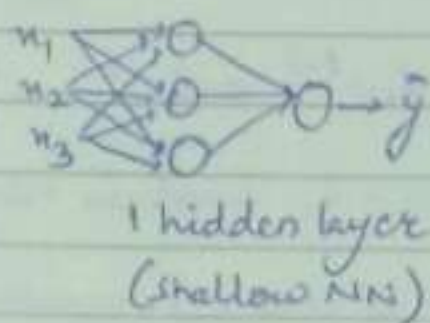
as long as tanh or sigmoid fn are

used in the NN.

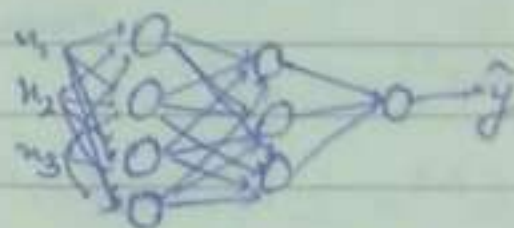
* Deep NN: (to solve percepⁿ problems in real world).



1 layer NN



2-layer model



3-layer model

(but we don't include i/p layer as a layer)

No. of layers = L

No. of units/nodes in layer $l = n^{[l]}$

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

$w^{[l]}$ = weights for computing $z^{[l]}$.

Note $a^{[0]} = x$
 $a^{[L]} = \hat{y}$

general fwd propagaⁿ algo for DNN: (considering vectorial representaⁿ)

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

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

* computing activaⁿ for all layers. $l = 1$ to L . { for loop needed to calculate all activaⁿ }

Note $z^{[1]} = w^{[1]}x + b^{[1]}$
 $(n^{[0]}, m) \quad (n^{[0]}, n^{[1]}) \quad (n^{[1]}, 1)$

(and similarly for next layers)

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

$db^{[1]} = (n^{[1]}, m)$

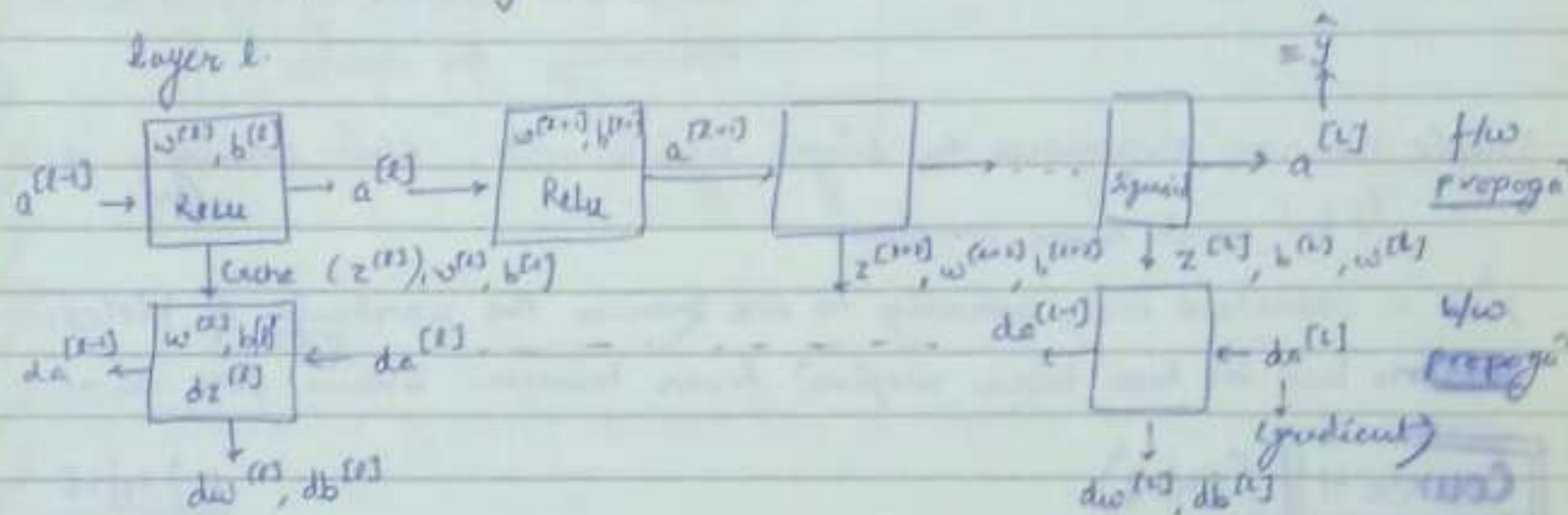
broadcasting will lead this to $(n^{[0]}, m)$

for m input training examples.

→ Why deep representations?

Ans Earlier layers of NN detect simple fns and they are composed together in the later layers of NN, so as to learn more & more complex fns. This type of simple to complex hierarchical representation is used in images, speech recognition.

Q If we try to compute a fn with lesser hidden layers then we will need to have exponentially more hidden units in it.



through this b/w step, we get all derivative values and w can get updated as $w^{(l)} = w^{(l)} - \alpha \cdot dw^{(l)}$ and $b^{(l)} = b^{(l)} - \alpha \cdot db^{(l)}$. This comprises of 1 "iter" of gradient descent for NN.

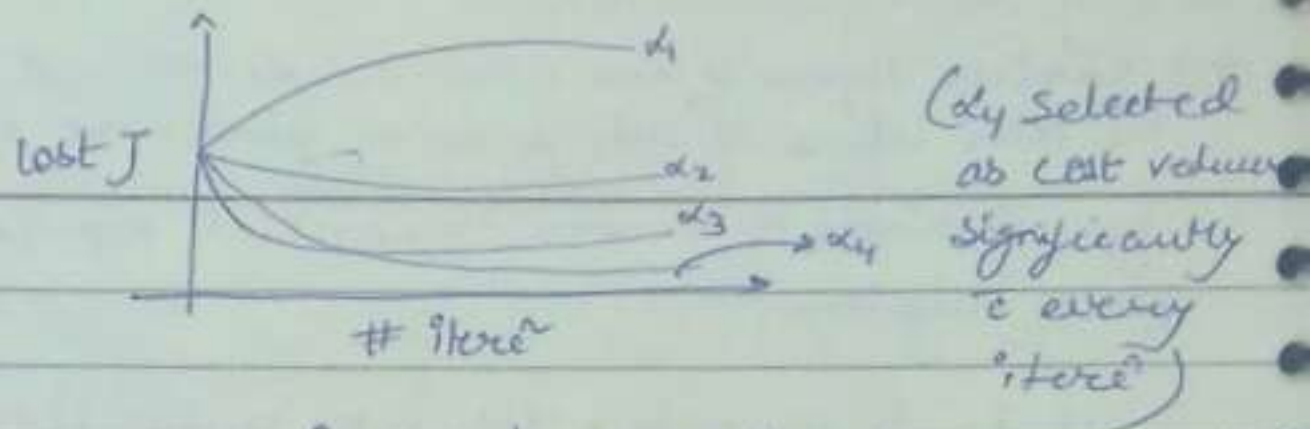
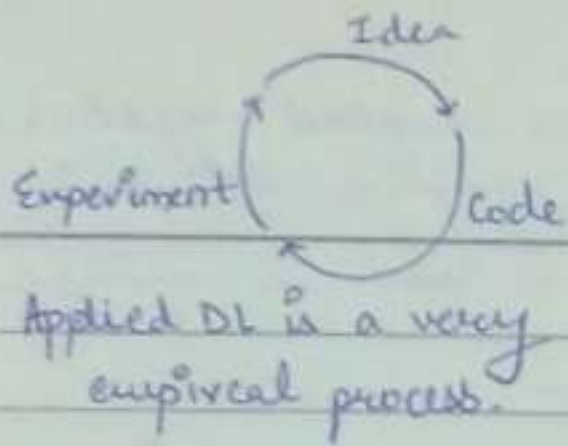
cache: pass info from one to the other so that it can be used during backpropag.

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

Hyperparameters: Learning rate α , # hidden units ($n^{(1)}, n^{(2)}, n^{(3)}, \dots, n^{(L)}$), # iterations, # hidden layers L , choice of activation fns.

Hyperparameters need to be given to the learning algorithm as they finally control the parameters w and b .

More hyperparameters: momentum term, mini batch size, regularization parameter.



In similar way other hyperparameter values are chosen, by trying and evaluating the results.

Single neuron analogous to single logistic unit. This analogy is
↳ but much more complex.

just a seductive one, actually no one knows the working of biological neurons, but DL has taken inspiraⁿ from human brain for sure.