

DEEP LEARNING

PAGE NO.

DATE

- Uses neural network
- * Convolutional neural network
 - good for image recognition
- * Long short - term memory network → good for speech recognition

Neural networks are inspired by brain

Neuron → a thing that holds a number

{ → One neuron corresponds to each pixel and holds grayscale value for pixel (0 - black to white) This no is called Activation

- Last layer has one neuron per each possible output
- Hidden layers b/w 1st & last layer
- Pattern is generated in each layer when image is fed to the network
- The brightest neuron in last layer is the input.

e.g - If we want to recognize a hand number

→ The second last layer

can detect loops and lines
that make up the number
→ The layer before that
can detect individual
edges of loops and lines

network shld convert :

pixels → edges → loops & lines → no

* What parameters should
network have to recognize
an edge in a region :-
→ Assign ^{weight to} connection between
neurons of 1st layer and 2nd
layer called

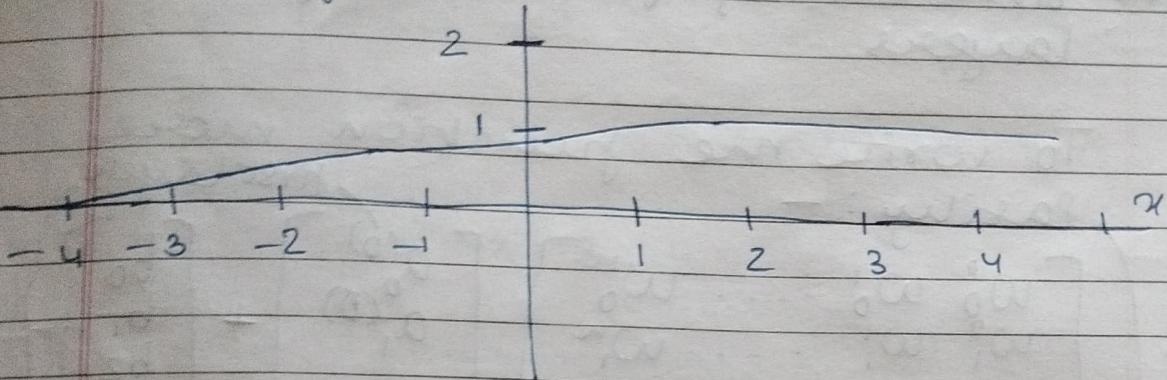
→ take all activations of
first layer and compute
weighted sum as:-
 $a_1 w_1 + \dots + a_n w_n$

→ This sum will give positive
value for region containing
a component of the loop

→ For edge detection \Rightarrow
negative weights are asso-
ciated with surrounding
pixels as they are darker

→ Since we need values
btw 0 & 1 we then pump
the weighted sum into a
function to give a value
btw 0 & 1

→ A common function that does this is the sigmoid function $\sigma(x) = \frac{1}{1+e^{-x}}$
 (logistic curve)



∴ Activation of neuron
 $= \sigma(a_1w_1 + a_2w_2 + \dots + a_nw_n)$

Suppose we don't want to activate neuron & meaning fully feel ^{weighted} sum greater than 0 but feel some other limit such as sum > 10

∴ We need bias for inactivity so we subtract some no from weighted sum before applying sigmoid func

Eg $\sigma(a_1w_1 + a_2w_2 + \dots + a_nw_n - 10)$
 ↓ Sigmoid ↑
Now positive is this? ↑ bias

→ All neurons of and layer are connected to neurons

g and layer.

- Each neuron in second layer will have some bias associated with it
- Similarly for all following layers

To write the function more easily :-

$$\begin{aligned}
 & \left[\begin{array}{c} w_0^0 \quad w_1^0 \quad \dots \quad w_n^0 \\ w_0^1 \quad w_1^1 \quad \dots \quad w_n^1 \\ \vdots \\ w_0^k \quad w_1^k \quad \dots \quad w_n^k \end{array} \right] \left[\begin{array}{c} a_0^{(0)} \\ a_1^{(0)} \\ \vdots \\ a_n^{(0)} \end{array} \right] + \left[\begin{array}{c} b_0 \\ b_1 \\ b_2 \\ \vdots \\ b_m \end{array} \right] \\
 & = \left[\begin{array}{c} \sigma(a_0^{(0)}w_0^0 + \dots + a_n^{(0)}w_n^0 + b_0) \\ \vdots \end{array} \right]
 \end{aligned}$$

first layer

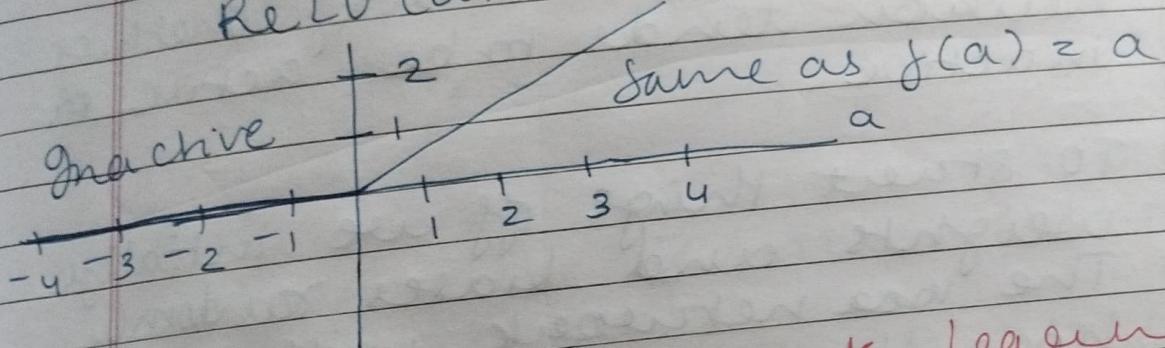
⇒ what this represents is

$$\sigma(wa^{(0)} + b) = a^{(1)}$$

↓ \ bias
 weight activation matrix of neuron
 matrix matrix in and
 layer

∴ Each neuron can be thought of as a function that takes outputs of all neurons in prev. layer and gives a no b/w 0 & 1

→ Modern networks don't use sigmoid anymore. They use ReLU as (Rectified linear unit) as it's much easier to train.
 $\text{ReLU}(a) = \max(0, a)$



How the network learns?

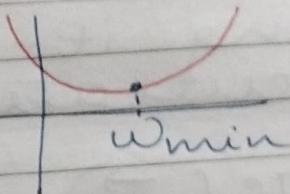
→ What we want is an algorithm where we can show this network a whole bunch of training data in form of images with labels showing what they are and the machine will adjust its weights and biases accordingly based on training data which can be generalised to images beyond training data.

3) **Machine learning** is like a calculus exercise and basically comes down to finding the minimum of a function

- Each neuron of a layer is connected to every neuron in next layer & weights are the strength of these connections
- Bias is indication of whether neuron tends to be active or inactive
- To start things off we initialize weights and biases randomly
- The ~~was~~ network will perform horribly and last neuron layer will be a mess
- So we define a cost function
To do this we add up squares of differences of the output activations and the value that we want them to have which called the cost
- Cost is small when network is more accurate
- For each image for each neuron there will be a cost
- We find the average cost for all the images (for each neuron of output layer)
 - (cost is found for each output neuron and then added up)

→ Rather than taking input of all the weights of each layer as input we use a single input: - $c(w)$ which gives a minimum output

$$\frac{dc(w)}{dw} = 0$$



w_{\min}

- However minima won't always be feasible (complicated functions)
- So instead start at a random input and which direction to step to lower the output
- However by this method the local minimum we land on may not be the lowest possible value

→ Gradient of a function tells us direction for steepest ascent whereas negative of gradient gives steepest descent

→ So algorithm for finding w_{\min} is

- 1) compute ∇c (gradient)
- 2) small step in $-\nabla c$ direction
- 3) repeat step 2 until convergence

$$\overrightarrow{w} = \begin{bmatrix} \text{all wts} \\ \text{& biases} \end{bmatrix}$$

$$-\nabla c(\overrightarrow{w}) = \begin{bmatrix} \text{bias} \\ \vdots \end{bmatrix}$$

PAGE NO. / / /
DATE / / /

$\vec{w} - \Delta C(\vec{w})$ — will decrease cost function

→ cost function is average over all samples so minimising it gives output more accurately in each sample

Eg: $-\nabla C(\vec{w}) = \begin{bmatrix} 0.31 \\ 0.03 \\ -1.25 \\ 0.18 \\ -0.37 \\ 0.16 \end{bmatrix}$

The sign of this func tells what shld inc/ decrease

It will take inputs equal to no of weights (ie connections b/w neurons)

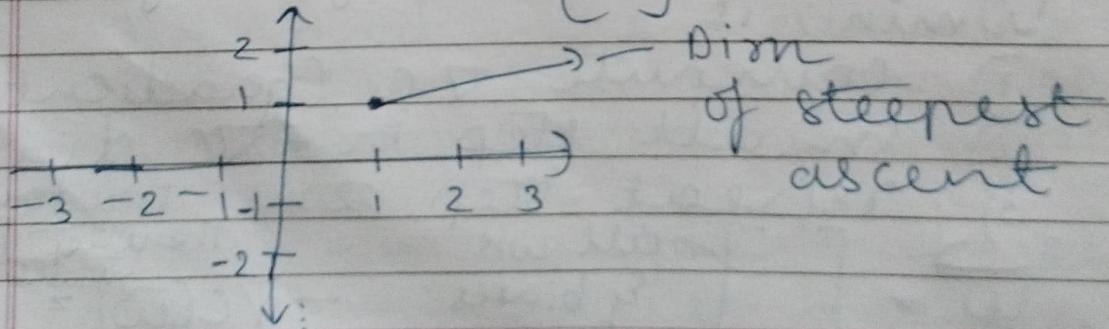
→ The value tells whether it shld decrease a lot or little or somewhere in between

→ Some adjustments matter more than others

Eg suppose we have a cost function with only 2 inputs

$$C(x, y) = \frac{3}{2} x^2 + \frac{1}{2} y^2$$

and $\nabla C(1,1) = \begin{bmatrix} 3 \\ 1 \end{bmatrix}$



since $\nabla C(1,1) = \begin{bmatrix} 3 \\ 1 \end{bmatrix}$ we can

say that change to x have 3 times as much impact as comp. to change to y

→ Even though image recognizes images it cannot generate them. So if we give a nonsense image it will confidently give a single output which ~~so~~ shouldn't happen

Backpropagation

→ Since cost function involves averaging over thousands of training examples the way ~~be~~ weights and bias per each step are adjusted depends on each training example

→ Suppose we focus on one training example only. If the network is still working badly we will get a garbage output. It is important to keep track of which values are important (eg if a wrong output is close to 0.0 we ignore it and focus on the outputs closer to 1 which need to be decreased and also focus on correct o/p that must be increased)

Suppose we need to increase activation of correct output whose activation is given by $\sigma(w_0 a_0 + \dots + w_n a_n + b)$

We can do this in 3 ways

- Increase b
- Increase w_i — In prop to a_i
- change a_i — In prop to w_i



* If we only concentrate on changing weights (w_i) we see that different weights have differing level of influence

* Connections from ~~near~~ brighter neurons of previous layer have more influence since they are multiplied by larger activation values ($w_i a_i$)

→ Hence there is a strong connection between neurons of layer that are more active and neurons of next layer that we wish to become more active

* Another way is to increase activation (a_i) of everything connected to ideal o/p neuron with +ve weight became more active and with -ve weight became less active.

Most effective when seeking changes proportional to

size of corresponding weights (w_i)
We cannot change activation
we can only change weights
& biases (since activation
depends on brightness of
pixels) It's helpful to keep a
note of what desired changes
are

- This is only for the output neuron and what it wants
- We also have to decrease activation of other neurons which have their own desires as per what should happen per their activation to decrease
- All these desires of each output neuron as to what should happen to and last layer is added together in proportion to w_i and in proportion to a_i
- By adding together all these desired changes we get a list of nudges that you want to happen to and last layer
- This is backpropagation
- It shows how each training example wishes to nudge weights & biases
- If we do this for only one image sample the network would only be able to identify that category of images

→ So we do it per all images
 → And average it out
 Hence the collection obtained of all the nudges to weights & bias ultimately gives the -ve gradient ie

$$-\nabla C(w_1, w_2, \dots, w_{13}, b) \\ = \begin{bmatrix} -0.08 \\ +0.12 \\ -0.06 \\ \vdots \\ 0.04 \end{bmatrix}$$

→ It takes computers extremely long time to add up nudges per all examples and gradient descent steps. So instead:

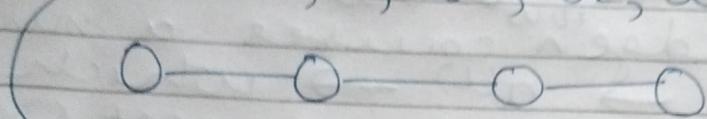
- * Training data is shuffled into mini batches
- * Then compute cost of a step per mini batch which won't be the actual gradient of cost func.

Show that
 Each mini batch gives a good approximation as well as a computational speed up

↓
 use one
 mini for
 each step
 ↓
 rather than
 all data
 all for step

Backpropagation Calculus

→ If we start with a very simple layer with 3 weights and 3 biases
 $\{w_1, b_1, w_2, b_2, w_3, b_3\}$



→ Our goal is to understand how sensitive the cost function is to these variables

→ Let activation of last neuron be $a^{(L)}$ & second last as $a^{(L-1)}$

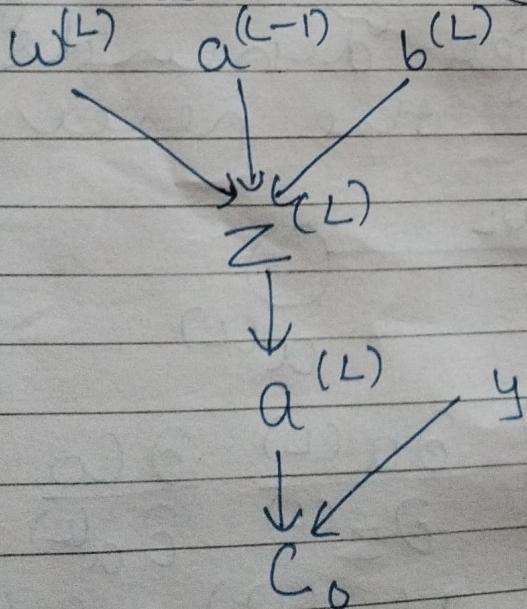
→ Suppose desired value of $a^{(L)}$ is y (eg 1.0)

→ ∴ Cost of this network for single training eg is

$$C_0 = \dots = (a^{(L)} - y)^2$$

$$z^{(L)} = w^{(L)} \cdot a^{(L-1)} + b^{(L)}$$

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



We can imagine a separate no. line for w, z, a, c_o

→ The goal is to find how sensitive the cost function c_o is to w
 \therefore we want $\frac{\partial c_o}{\partial w^{(L)}}$

→ Tiny nudge $\Delta w^{(L)}$ causes a nudge to $z^{(L)} \rightarrow a^{(L)} \rightarrow$ which directly influences cost (Δc_o)
 $\therefore \boxed{\frac{\partial c_o}{\partial w^{(L)}} = \frac{\partial z^{(L)}}{\partial w^{(L)}} * \frac{\partial a^{(L)}}{\partial z^{(L)}} * \frac{\partial c_o}{\partial a^{(L)}}}$

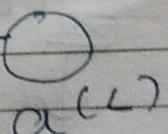
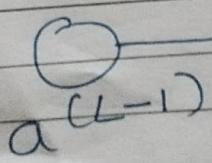
(chain rule)

$$\frac{\partial c_o}{\partial w^{(L)}} = \frac{\partial (a^{(L)} - y)}{\partial w^{(L)}}$$

$$\frac{\partial a^{(L)}}{\partial z^{(L)}} = \sigma'(z^{(L)})$$

$$\left[\frac{\partial z^{(L)}}{\partial w^{(L)}} = a^{(L-1)} \right]$$

The amount that the small nudge of w_L influenced the last layer depends on how strong the pre neuron is



$$\begin{aligned} \frac{\partial c_o}{\partial w^{(L)}} &= \frac{\partial z^{(L)}}{\partial w^{(L)}} \frac{\partial a^{(L)}}{\partial z^{(L)}} \frac{\partial c_o}{\partial a^{(L)}} \\ &= a^{(L-1)} \sigma'(z^{(L)}) 2(a^{(L)} - y) \end{aligned}$$

$$\frac{\partial C}{\partial w^{(L)}} = \frac{1}{n} \sum_{k=0}^{n-1} \frac{\partial C_k}{\partial w^{(L)}}$$

Since C is found by averaging cost over all examples, the derivative of C is by averaging derivative of cost over all examples

This is just one component of gradient vector

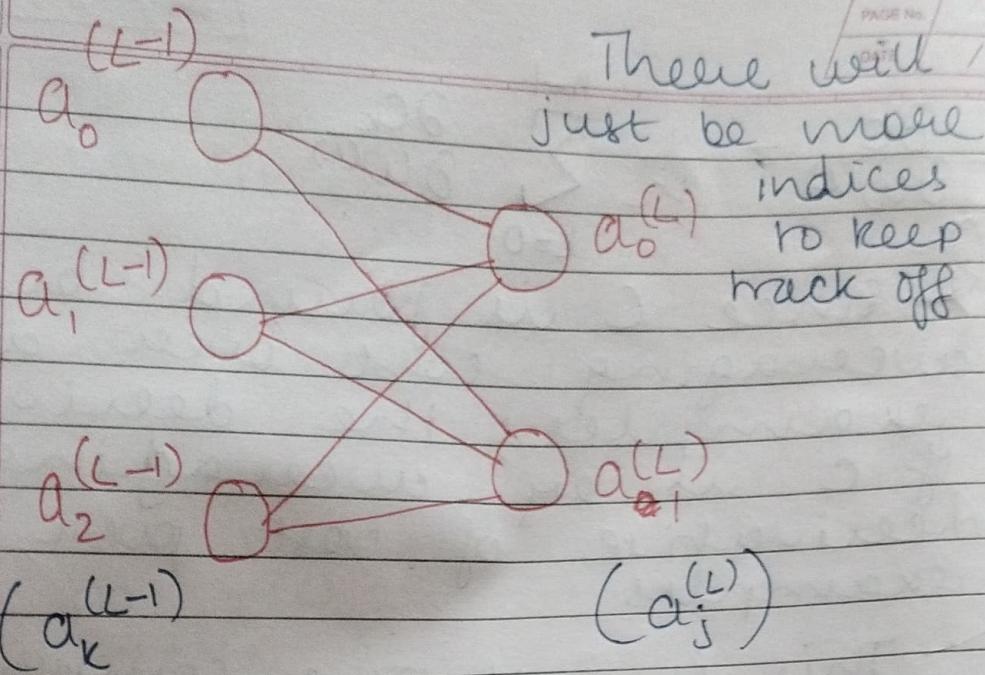
$$\nabla C = \left[\begin{array}{c} \frac{\partial C}{\partial w^{(1)}} \\ \frac{\partial C}{\partial w^{(2)}} \\ \vdots \\ \frac{\partial C}{\partial w^{(L)}} \end{array} \right]$$

→ sensitive similarly for bias

$$\frac{\partial C}{\partial b^{(L)}} = \frac{\partial z^{(L)}}{\partial b^{(L)}} \frac{\partial a^{(L)}}{\partial z^{(L)}} \frac{\partial C}{\partial a^{(L)}} = \frac{\partial z^{(L)}}{\partial b^{(L)}} \frac{\partial a^{(L)}}{\partial z^{(L)}} \times \frac{\partial C}{\partial a^{(L)}} = \frac{\partial z^{(L)}}{\partial b^{(L)}} \times 2(a^{(L)} - y)$$

→ change per bias remains same almost always

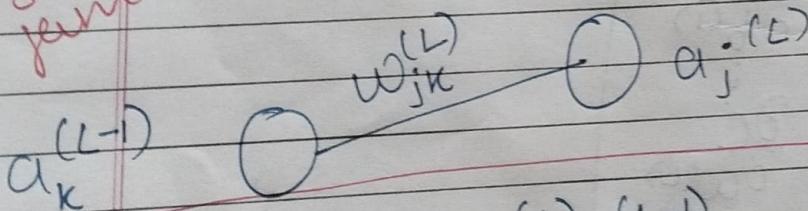
if we do this same thing for a more complicated layer



Now it will be

$$C_0 = \sum_{j=0}^{n_L-1} (a_j^{(L)} - y_j)^2$$

since j starts from 0



$$z_j^{(L)} = w_{j0}^{(L)} a_0^{(L-1)} + w_{j1}^{(L)} a_1^{(L-1)} + \dots + w_{jk}^{(L)} a_k^{(L-1)} + b_j^{(L)}$$

$$a_j^{(L)} = \sigma(z_j^{(L)})$$

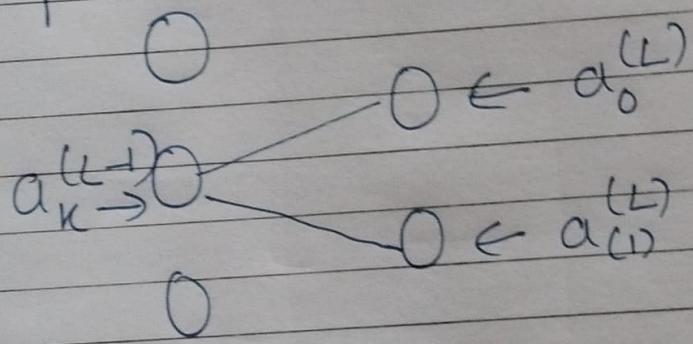
$$C_0 = \sum_{j=0}^{n_L-1} (a_j^{(L)} - y_j)^2$$

$$\frac{\partial C_o}{\partial w_{jk}^{(L)}} = \frac{\partial z_j^{(L)}}{\partial w_{jk}^{(L)}} \times \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \times \frac{\partial C_o}{\partial a_j^{(L)}}$$

$$\frac{\partial C_o}{\partial a_k^{(L-1)}} = \sum_{j=0}^{n_L-1} \frac{\partial z_j^{(L)}}{\partial a_k^{(L-1)}} \times \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \times \frac{\partial C_o}{\partial a_j^{(L)}}$$

sum over L

The difference is that the neuron k influences cost function through multiple different paths



$\nabla C \leftarrow$

$$\frac{\partial C}{\partial w_{jk}^{(L)}} = a_k^{(L-1)} \sigma'(z_j^{(L)}) \frac{\partial C}{\partial a_j^{(L)}}$$

$$\sum_{j=0}^{n_{L+1}-1} w_{jk}^{(L+1)} \sigma'(z_j^{(L+1)}) \frac{\partial C}{\partial a_j^{(L+1)}}$$

or

$$\lambda(a_j^{(L)} - y_j)$$