

NEURAL NETWORK

Convolutional neural network → Good for image recognition

long short-term memory → Good for speech recognition

NEURAL NETWORK

What are neurons?

thing that holds a number (mostly 6/bt 0
8/1)

All based on its grayscale representation
ranging bt 0 (black) &
1 (white) pixels.

(x) → no. inside the neuron is called "ACTIVATION".

They have values 6/bt 0 to 1 as whatever digit is represented by the network

How are they connected?

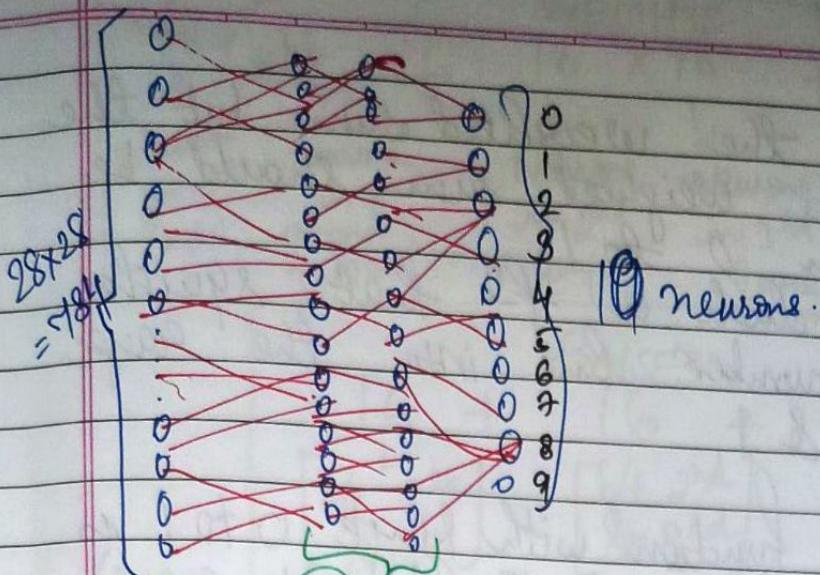
① If we have 28 by 28 pixels grid then

$28 \times 28 = 784$. becomes the 1st layer of our network.

② The last layer is the 4th layer consisting of 10 neurons each representing digits from 0 to 9.

③ The middle 2 layers in between are called as the hidden layers.

The brightest of all neurons (value closer to 1) is the value of digit that the neural network represent



assumed
can be taken as anything

2 hidden layers
selected of
(16) neurons each

Representations of Digits

$$1) [9] = [0] + [1]$$

$$2) [8] + [0] + [0]$$

$$3) [4] + [1] + [1] + [-]$$

} all this becomes the part of 3rd layer

This all loops/lines are divided into smaller sections to recognise those loops and they become the 2nd layer!

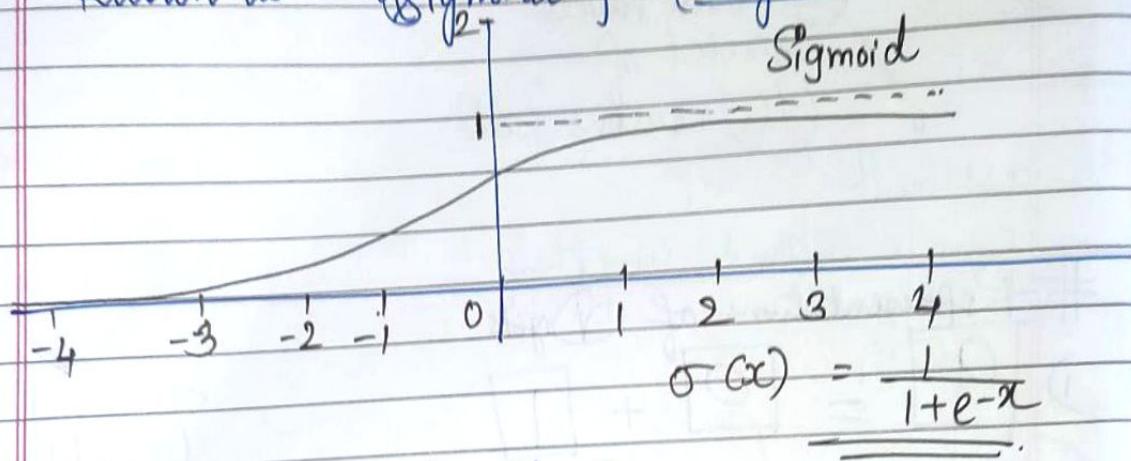
$$\text{eg: } [0] = [r] + [c] + [l] + [o]$$

$$[1] = [l] + [o]$$

We calculate the weighted sum but the range of weighted sum should be in between 0 to 1

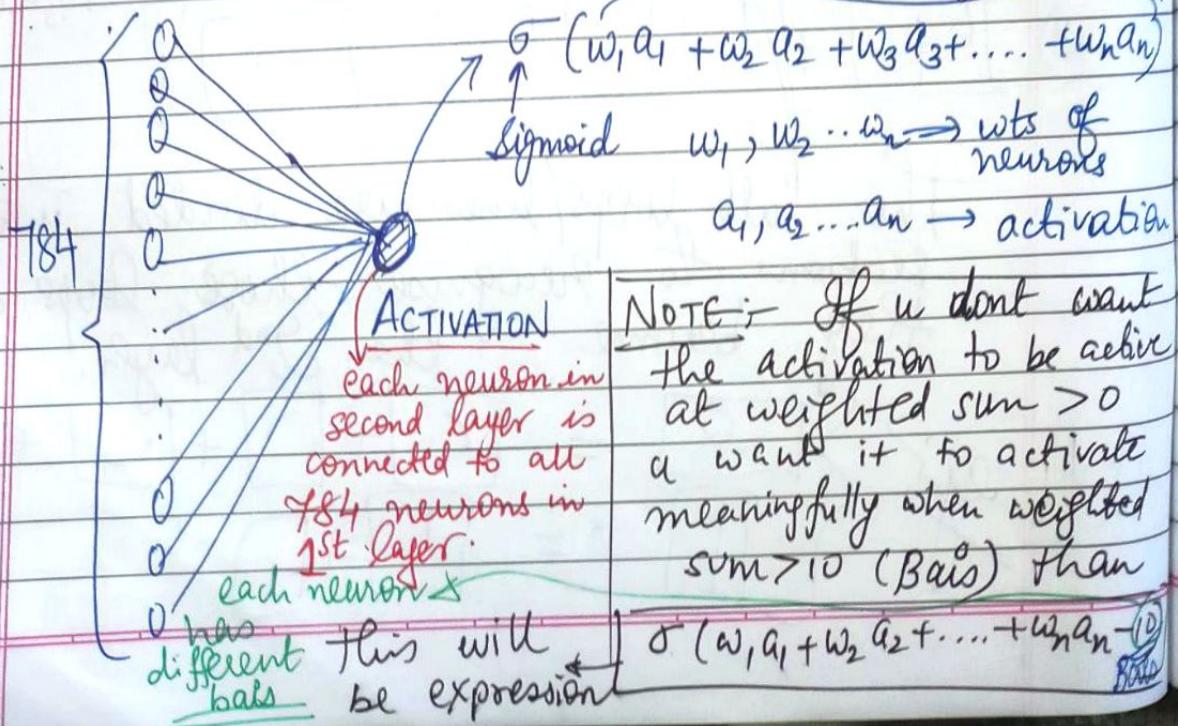
so we include a fn that squishes the real number line into the range between 0 & 1

The Common function with range 0 to 1 is known as Sigmoid fn (Logistic curve).



-ve inputs \rightarrow close to 0
+ve inputs \rightarrow close to 1

how positive the weighted sum is



WEIGHTS.

$$784 \times 16 + 16 \times 16 + 16 \times 10$$

↓ ↓ ↓

neurons in 1st layer. neurons in second layer. Neurons in 2nd layer. Neurons in 3rd layer. Neurons in 3rd layer. Neurons in 4th layer.

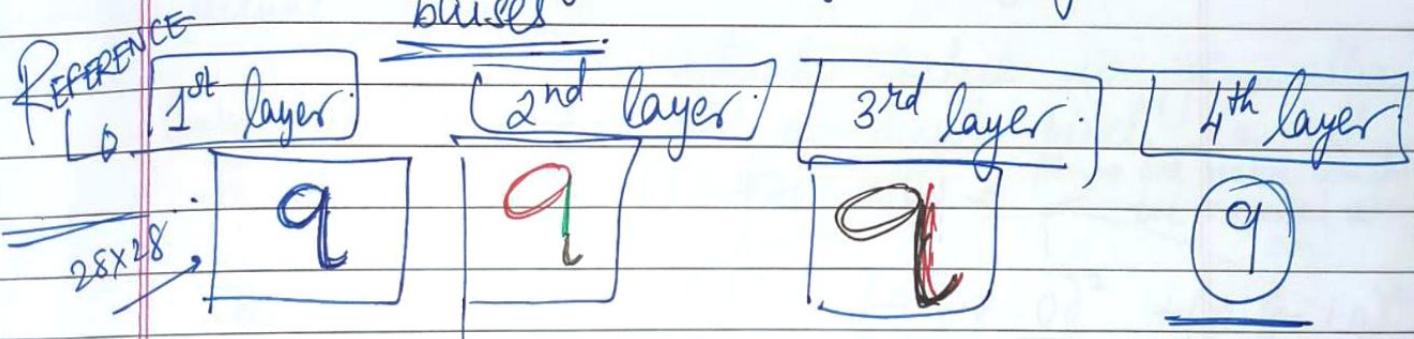
BAISES.

$$16 + 16 + 10$$

[b/t 1st & 2nd] [b/t 2nd & 3rd] [b/t 3rd and 4th]

Total 13002 bases & weights

Learning :- Finding the right weights and biases



$$a_0^{(1)} = \sigma(w_{0,0}a_0^{(0)} + w_{0,1}a_1^{(0)} + w_{0,2}a_2^{(0)} + \dots + w_{0,n}a_n^{(0)} + b_0)$$

Basis

$$\begin{bmatrix} w_{0,0} & w_{0,1} & \dots & w_{0,n} \\ w_{1,0} & w_{1,1} & \dots & w_{1,n} \\ \vdots & \vdots & & \vdots \\ w_{k,0} & w_{k,1} & \dots & w_{k,n} \end{bmatrix} \begin{bmatrix} a_0^{(0)} \\ a_1^{(0)} \\ \vdots \\ a_n^{(0)} \end{bmatrix} = \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_n \end{bmatrix}$$

Sigmoid

① ORGANISE ALL THE ACTIVATIONS FROM ONE LAYER INTO A COLUMN AS A VECTOR.

② ORGANISE ALL THE WEIGHTS AS MATRIX where each row of the matrix corresponds to the connections b/t one layer and a particular neuron in the

③ Bias Vector to add all bias to the previous vector product of weighted sum.
This can be represented as

$$a_j^{(1)} = \sigma(wa^{(0)} + b) \quad \begin{matrix} \text{bias vector} \\ \text{Initial layer} \end{matrix}$$

1st layer

→ Nowadays, RELU function is used instead of sigmoid fn

We define a cost function which is the difference between 2 values which we get and the actual value

Consider avg cost of all data. 1.0 for a single pixel & 0 for others → one value which we expected as ans.

$$\text{Cost} = (0.43 - 0.0)^2 + (0.28 - 0.0)^2 + (0.8 - 1.0)^2 + \dots$$

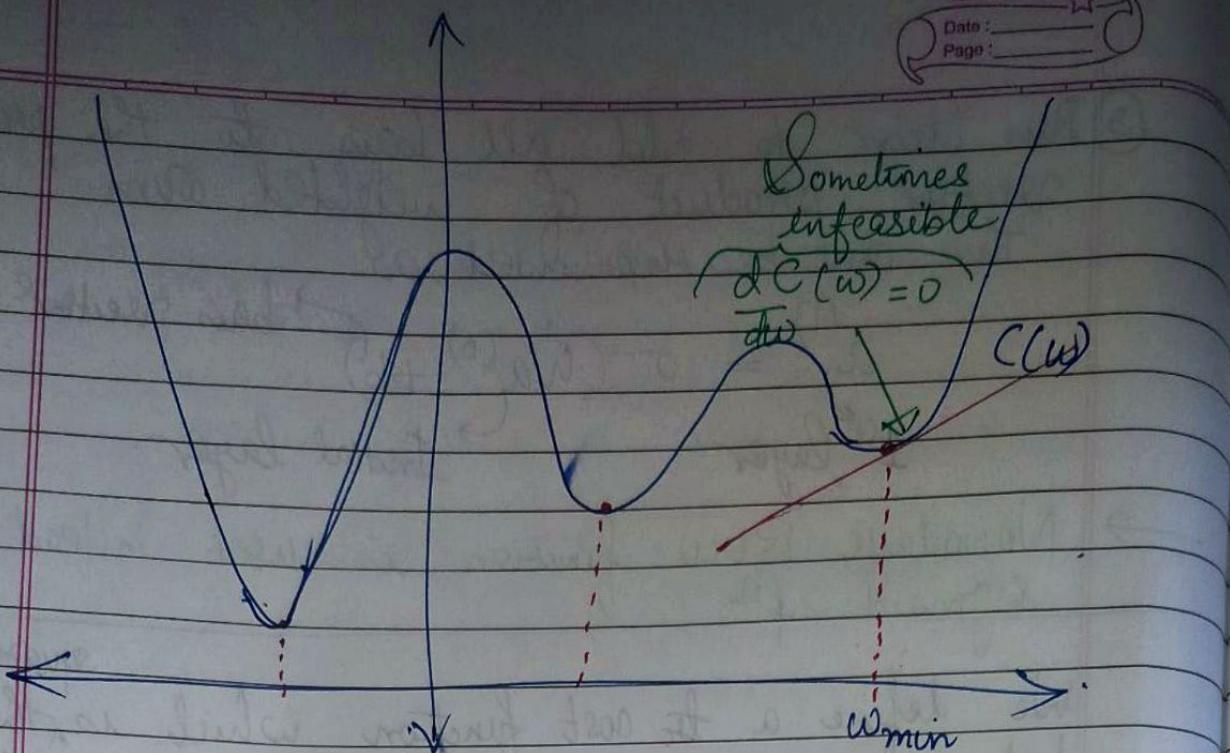
Neural Network function	Cost function
-------------------------	---------------

Input: 784 numbers (pixels) Input: 13,002 weights/baises

Output: 10 numbers.

Parameters: 13,002 weights/baises

Parameters: Many, many training examples



There are many minima for a complicated f^n but ~~not~~ all minima don't have the minimum value of the f^n , only one minima has the minimum value of the f^n

The minima is determined by calculating the slope of f^n at that point

If Slope is +ve \rightarrow shift to left.
If slope is -ve \rightarrow shift to right.

by this way we get a point where slope is zero and we find a minimum.

GRADIENT

Gradient of a function gives you the direction of steepest ascent, basically which direction should you step to increase the function most quickly. Taking the negative of that gradient gives you direction to step that decreases the function most quickly. Length of the gradient vector is actually an indication for just how steep that steepest slope is.

"Gradient" → the direction of steepest increase
 $\nabla C(x, y)$.

$$\vec{w} = \begin{bmatrix} 2.93 \\ -1.57 \\ 1.98 \\ \vdots \\ -1.16 \\ 3.82 \\ 1.21 \end{bmatrix}$$

13,002 weights & biases.

negative gradient of cost f₂

$$-\nabla C(\vec{w}) = \begin{bmatrix} 0.18 \\ 0.45 \\ -0.51 \\ \vdots \\ 0.40 \\ -0.32 \\ 0.82 \end{bmatrix}$$

This tells which changes to which weight matters most

How to nudge all weights & biases

The negative gradient of cost function is just a vector in some direction inside this huge input space that tells you which nudges to all those nos is going to cause the most rapid decrease to the cost function

We need to minimise cost funct² as small as possible.

The process of repeatedly nudging an input of a function by some multiple of the negative gradient is called gradient descent. It's a way to converge towards some local minimum of cost function.

each component of the negative gradient ($-\nabla C(\vec{w})$) tells us 2 things.

↓

the sign tells us whether the corresponding component of input vector should be nudged up or down
 (relative magnitudes of all comp)

↓

the relative magnitudes of components states that which changes matter more.

eg

$$-\nabla C(\vec{w}) = \begin{bmatrix} 0.31 \\ 0.03 \\ -1.25 \\ \vdots \\ 0.48 \\ -0.37 \\ 0.16 \end{bmatrix}$$

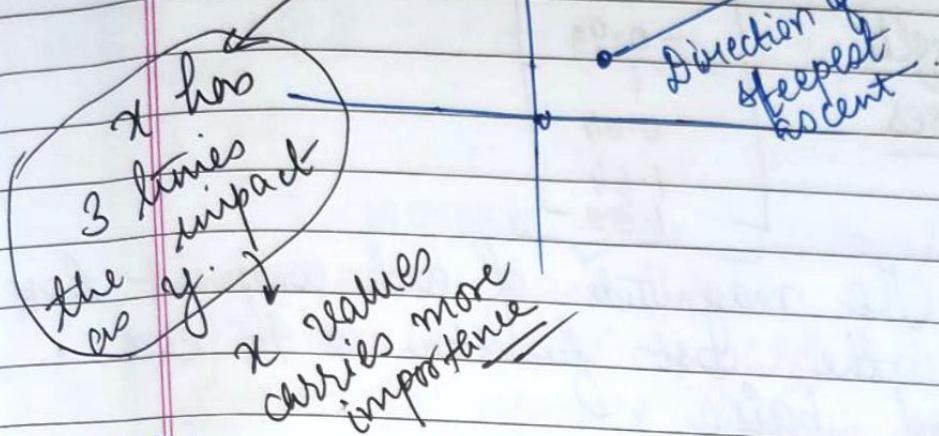
1300 ~

w_0 should increase some what.
 w_1 ~ a little.
 w_2 should decrease a lot.
 w_{1300} should increase a lot.
 w_{13001} should decrease somewhat.
 w_{13002} should increase a little.

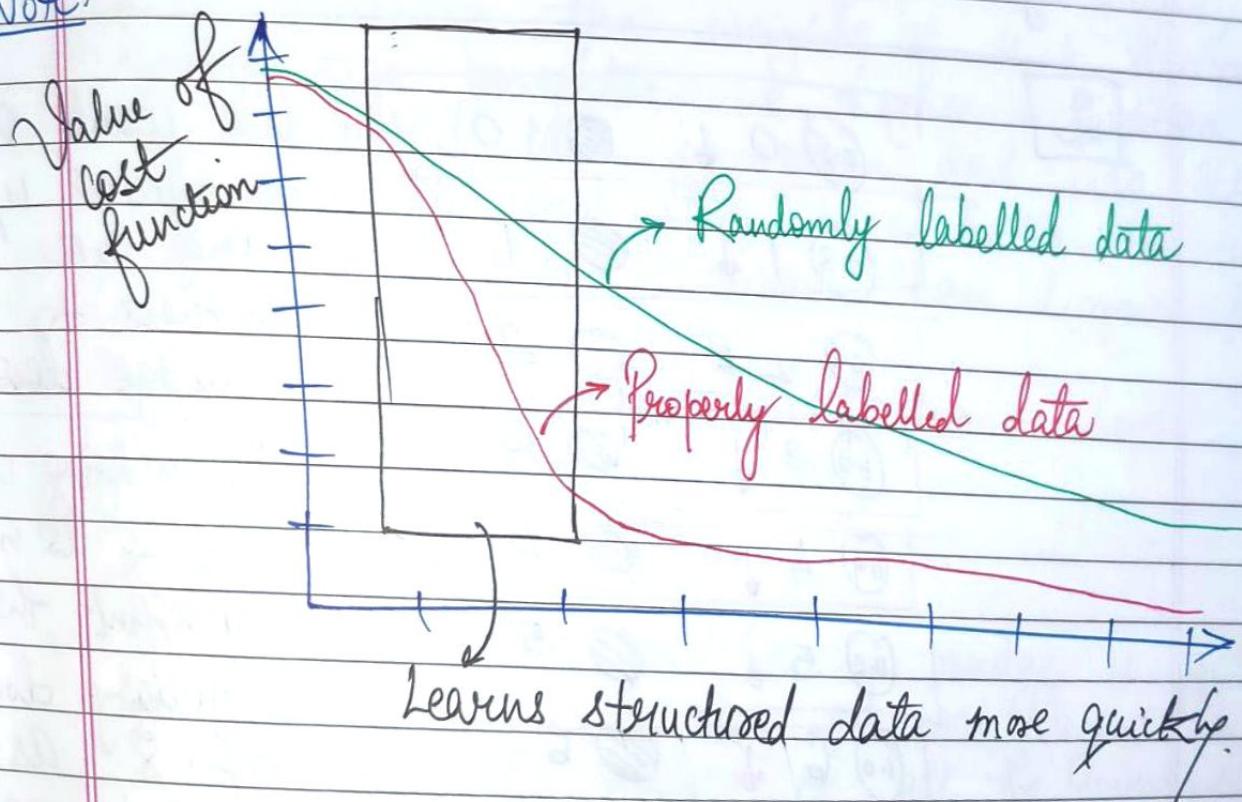
If a function has 2 variables as an input, and we compute that its gradient at some particular point comes out as $(3, 1)$

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

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



Note:-

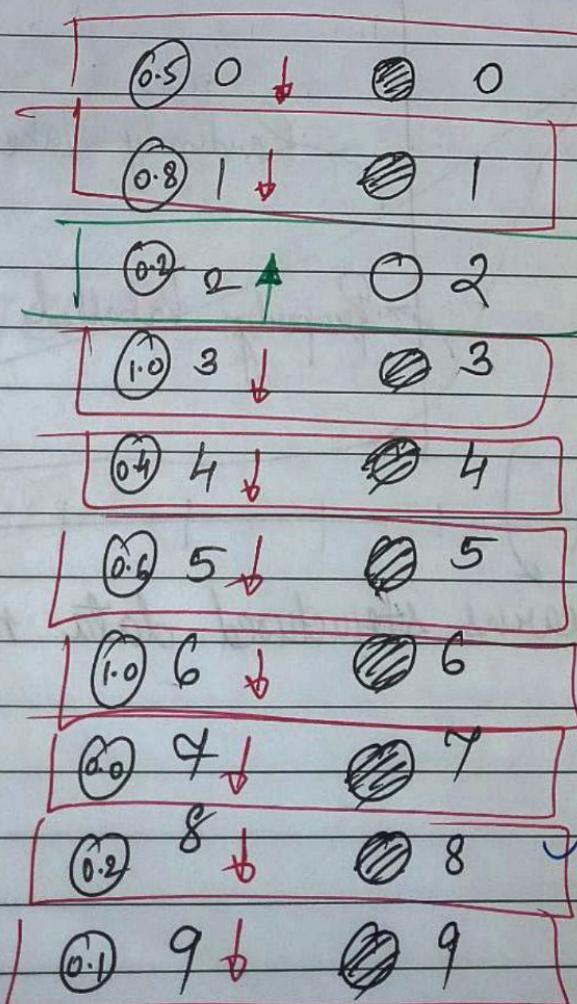


BACK PROPAGATION

↓
 Algorithm for computing the complicated gradient. An algorithm for determining how a single training example would like to nudge the weights and biases,
 $\Delta C(\dots) = \begin{bmatrix} 0.16 \\ 0.72 \\ -0.93 \\ \vdots \\ 0.04 \\ 1.64 \\ 1.52 \end{bmatrix}$
 All weights and biases

how sensitive the cost function is to each weight and bias.

[2]



We want 2 to nudge up and all others to nudge down

The nudging up of 2 is more important than nudging down of 8 as 0.2 is pretty close to 1.0 but not to

focusing on neuron whose activation is to be increased
 ②.

Activation \rightarrow weighted sum of all activation in previous layer + bias in the sigmoid squification fn.

$\rightarrow \text{② neuron of no 2.}$

$$= \sigma (w_0 a_0 + w_1 a_1 + \dots + w_n a_n + b)$$

To increase the activation, 3 methods.

1] Increase bias (b).

2] Increase w_i (weights) in proportion to ap

Note:— the weighted sum is of previous layer, so if u increase the weights of that neurons which are brighter i.e higher activation than the weighted sum increases, and has stronger influence to cost fn; ∴ it is better to increase weights of brighter neurons than dimmer neurons.

3] Change a_i (activation) of previous layer in proportion to w_i

Back Propagation

Whatever upper or lower nudge u get for the neuron in the 2nd last layer by considering the nudges of all the neurons in the last layer, we can also find the nudge of its previous (i.e 3rd last layer). This is Back Propagation

→ For all data, we do the following and take the average to generalise the negative gradient of cost fn.

We adjust the weights so that the output of last layer is w_0 by making 2 by making the weights w_1 of brighter neurons more similarly for all.

	2	5	0	4	1	9	Avg
w_0	-0.08	+0.02	-0.02	+0.11	-0.05	-0.14	-0.08
w_1	-0.11	+0.11	+0.07	+0.02	+0.09	+0.05	+0.12
w_2	-0.07	-0.04	-0.01	+0.02	+0.13	-0.15	-0.06
$w_{1,3002}$	+0.13	+0.08	-0.06	-0.09	-0.02	+0.04	+0.01

This is the negative gradient of cost function:

$$-\nabla C(w_1, w_2, \dots, w_{1,3002})$$

Average over all training data

The above process is little slow, so we randomly divide the input data of nos. into some minibatches, and repeatedly going through all of the mini batches and making these adjustments, we get a local minimum of cost fn.

Stochastic gradient descent

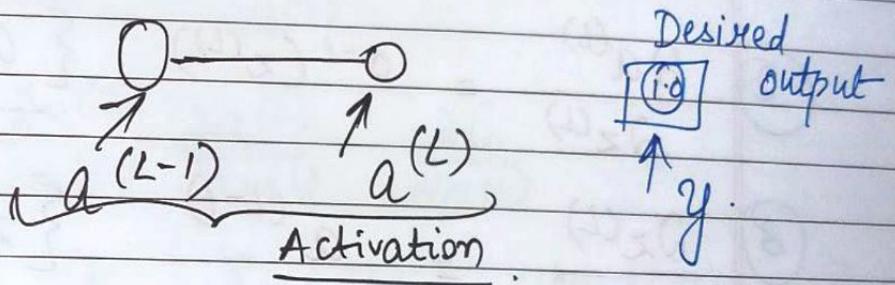
Consider a simple network, where each layer has one neuron each.



Focus on connection b/t last 2 neurons

This network has 3 weights & three bias (w_1, w_2, w_3) & (b_1, b_2, b_3)

Cost fn. $\rightarrow C(w_1, b_1, w_2, b_2, w_3, b_3)$
 Our goal is to understand how sensitive is the cost function to these variables.
 Also, we need to know which adjustments to these terms, cause most efficient decrease in cost function.



Cost of one training example $\rightarrow C_0(\dots) = (a^{(L)} - y)^2$

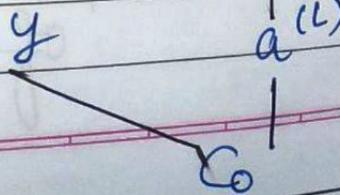
$$a^{(L)} = \sigma(w^{(L)} a^{(L-1)} + b^{(L)})$$

$z^{(L)}$ $w^{(L-1)}$ $a^{(L-2)}$ $b^{(L-1)}$

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

$w^{(L-1)}$ $a^{(L-1)}$ $b^{(L)}$

FLOWCHART



The flowchart continues.

A small change in weight ($w^{(l)}$) causes the small change in $z^{(l)}$ and that further cause change to $a^{(l)}$ and by this way C_0 is changed. Our goal is to find how cost function changes by a small change in weight so by Chain RULE

$$\frac{\partial C_0}{\partial w^{(l)}} = \frac{\partial z^{(l)}}{\partial w^{(l)}} \cdot \frac{\partial a^{(l)}}{\partial z^{(l)}} \cdot \frac{\partial C_0}{\partial a^{(l)}}$$

$$\textcircled{1} \quad \frac{\partial C_0}{\partial a^{(l)}} = 2(a^{(l)} - y) \quad \left\{ C_0 = (a^{(l)} - y)^2 \right\}$$

$$\textcircled{2} \quad \frac{\partial a^{(l)}}{\partial z^{(l)}} = \sigma'(z^{(l)}) \quad \left\{ a^{(l)} = \sigma(z^{(l)}) \right\}$$

$$\textcircled{3} \quad \frac{\partial z^{(l)}}{\partial w^{(l)}} = a^{(l-1)} \quad \left\{ z^{(l)} = w^{(l)}a^{(l-1)} + b^{(l)} \right\}$$

So,

$$\frac{\partial C_0}{\partial w^{(l)}} = a^{(l-1)} \sigma'(z^{(l)}) 2(a^{(l)} - y)$$

→ Average of all training examples is

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

Derivative of full cost function

→ This is one component of gradient vector

Gradient vector is formed by the partial derivatives of the cost fn w.r.t all those weights and biases.

$$\nabla C = \begin{bmatrix} \frac{\partial C}{\partial w^{(1)}} \\ \frac{\partial C}{\partial b^{(1)}} \\ \vdots \\ \frac{\partial C}{\partial w^{(L)}} \\ \frac{\partial C}{\partial b^{(L)}} \end{bmatrix}$$

for bias term $\frac{\partial C}{\partial b^{(1)}}$

$$\frac{\partial C}{\partial b^{(1)}} = \frac{\partial z^{(1)}}{\partial b^{(1)}} \cdot \frac{\partial a^{(1)}}{\partial z^{(1)}} \cdot \frac{\partial C}{\partial a^{(1)}}$$

$$(4) \quad \frac{\partial z^{(1)}}{\partial b^{(1)}} = 1 \quad \left\{ z^{(1)} = w^{(1)} a^{(1)} + b^{(1)} \right\}$$

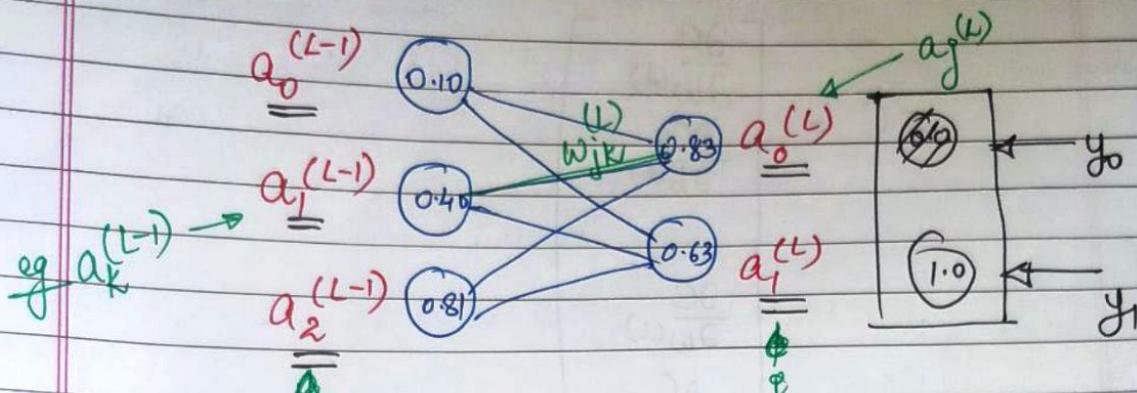
$$\therefore \frac{\partial C}{\partial b^{(1)}} = 1 \cdot \sigma'(z^{(1)}) 2(a^{(1)} - y)$$

Now, we can use the same chain rule idea backwards, this is called back propagation

$$\frac{\partial C}{\partial a^{(L-1)}} = \frac{\partial z^{(L)}}{\partial a^{(L-1)}} \cdot \frac{\partial a^{(L)}}{\partial z^{(L)}} \cdot \frac{\partial C}{\partial a^{(L)}}$$

w_L

For more than one neuron in a layer there will also be subscript to denote those neurons



Let's use k to denote the index of second last layer.

Let's use j to denote index of last layer.

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

$$z_j^{(L)} = w_{j0}^{(L)} a_0^{(L-1)} + w_{j1}^{(L)} a_1^{(L-1)} + w_{j2}^{(L)} a_2^{(L-1)} + b_j^{(L)}$$

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

$$\frac{\partial C_0}{\partial a_j^{(L)}} = \sum_{j=0}^{n_L-1} \frac{\partial z_j^{(L)}}{\partial a_j^{(L)}} \cdot \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \cdot \frac{\partial C_0}{\partial a_j^{(L)}}$$

Sum over layer L

$$\frac{\partial C_0}{\partial a_k^{(L-1)}}$$

$$\frac{\partial C_0}{\partial w_{jk}^{(L)}} = \frac{\partial z_j^{(L)}}{\partial w_{jk}^{(L)}} \cdot \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \cdot \frac{\partial C_0}{\partial a_j^{(L)}}$$

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

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

OR

$$2(a_j^{(l+1)} - y_j)$$

this chain rule expression gives you the derivatives that determine each component in the gradient that minimize the cost of the network by repeatedly stepping downhill.