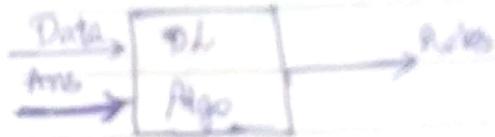
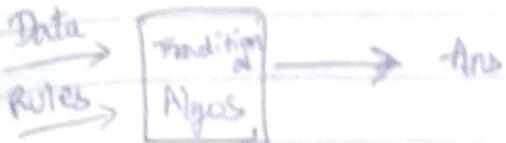


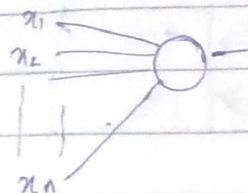
# DL



1. import tensorflow as tf
  2. import numpy as np
  3. import tensorflow import keras
  4. model = tf.keras.Sequential(Choles.layers.Dense(units=1, input\_shape=[1]))
  5. model.compile(optimizer='sgd', loss='mean\_squared\_error')
  6. model.fit(x, y, epochs=10)
  7. print(model.predict([1]))
- ↳ model - structure of entire neural network  
sequential - all neurons <sup>layers</sup> are connected in sequence.  
dense - all neurons are interconnected  
units=1 - no. of neurons  
input\_shape - type of I/p going into neuron  
6. This stmt. runs the model to find the function  
epochs - no. of iterations.

As of now we have not given I/p data & prediction  
Intended func:  $y = f(x)$ . After line 5, add:  
 $\{ x = np.array([1.0, 2.0, 3.0, 4.0, 6.0]), dtype='float'$   
 $\{ y = np.array([1.0, 2.0, 3.0, 4.0, 6.0]), dtype='float'$   
& in line 7 inside the predict give a value say 5.0

$x^w \circlearrowright \rightarrow y$  Need only 1 weight parameter.

 Need n weight parameters.

$n^{[l]}$  - no of nodes in layer l

$w^{[l]}$  - wt. parameters in layer l

$b^{[l]}$  - bias in layer l.

Gradient descent = For all  $i$ , the eqn. below is repeated until  $\frac{\delta J}{\delta w_i} = 0$

$$w_i := w_i - \alpha \frac{\delta J}{\delta w_i} \quad \begin{matrix} \text{learning rate } (\alpha > 0) \\ \text{cost} \end{matrix}$$

Repeat until we can't change  $\frac{\delta J}{\delta w_i}$  for any of them  
for  $i$  in weights  $\rightarrow$  set of all wts. irrespective of all layers

$$w_i := w_i - \alpha \frac{\delta J}{\delta w_i}$$

$$J = \frac{1}{m} \sum_{i=1}^m \delta (y^{(i)}, \hat{y}^{(i)}) \quad \begin{matrix} \text{basically avg. loss over all} \\ \text{training examples} \end{matrix}$$

Super script in parenthesis means we are talking about  $i^{th}$  sample

$$\text{loss } \leftarrow d = -(y \log \hat{y} + (1-y) \log (1-\hat{y})) \quad \begin{matrix} \text{Cost func. becomes loss} \\ \text{func. as soon as we avg.} \end{matrix}$$

We don't use methods of finding maxima / minima in calculus ( $\frac{dy}{dx}, \frac{d^2y}{dx^2}$ ) here as it will become

NP-hard in multi dimensions while gradient descent is a fast wt. algorithm.

Here we are trying to minimize cost function &

degree of  
misclassified classifications

$(\hat{y} \in (0, 1) \Rightarrow \text{If } \hat{y} \geq 0.5, \hat{y} = 1 \text{ else } 0)$

If actual  $y = 1$  & our pred. says  $\hat{y} = 0$   
 $\log 0 = -\infty$      $\log 1 = 0$

$L = -(-\infty + 0) = \infty$  i.e. we have a huge penalty.  
 $\hat{y}$  comes from  $\sigma(z) = \frac{1}{e^{-z} + 1}$  so  $\hat{y}$  can be 0 only if  $z = -\infty$ .

Now  $z$  can't be  $-\infty$ , so  $\hat{y} \neq 0$ .

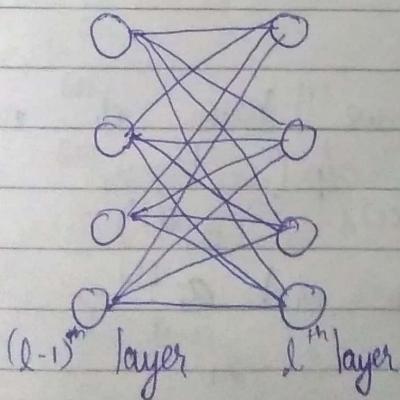
If  $\hat{y} = 0.0001$ , then loss will be very ↑.

If  $y = \hat{y}$ , then loss will be 0 i.e. no penalty.

So we want to minimise  $L$  &  $J$  as it will make our classifier better

In order to minimise  $J$ , gradient descent will help as we find  $J$ 's minimum value (using basic LA).

Dimensions of  $w^{[l]}$  in terms of  $n^{[l]} = n^{[l]} \times n^{[l-1]}$



To write previous eq. more accurately,

$$w_{ij}^{[l]} = w_{ij}^{[l]} - \alpha \frac{\delta J}{\delta w_{ij}^{[l]}}$$

By convention, we don't count I/P layer as layers!

classmate  
Date \_\_\_\_\_  
Page \_\_\_\_\_

This will contain 4 layers:

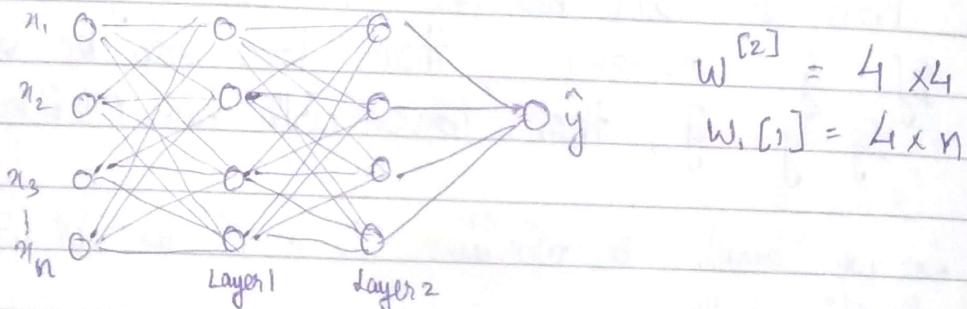
Repeat until there is no updation:

For all layers  $l = 1$  to  $l \rightarrow$  last layer

For all  $i$  this ~~weights~~  $= 1$  to  $n^{[l]}$

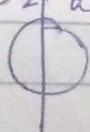
For ~~all~~  $j = 1$  to  $n^{[l-1]}$

$$w_{ij}^{[l]} = w_{ij}^{[l]} - \alpha \frac{\delta J}{\delta w_{ij}^{[l]}}$$



$\hat{y}$  is a function of  $y, \hat{y}$ ,  $\hat{y}$  is a function of  $z$   
 $z$  is a function of  $w$  which is a func. of  $a$ 's which is a func. of previous layer's  $w$ 's &  $a$ 's & so on

Linear part  $\leftarrow z \leftarrow a \rightarrow$  Non Linear part o/p



Neuron

In a layer  $l$ , we'll have  $n^{[l]}$  neurons, so we will have  $n^{[l]}$  linear o/p ( $z_1^{[l]}, z_2^{[l]}, \dots, z_n^{[l]}$ ) &  $n^{[l]}$  non linear o/p ( $a_1^{[l]}, a_2^{[l]}, \dots, a_n^{[l]}$ )

How to compute each  $a_i^{[l]}$ ? It will depend on o/p from previous layer, i.e.  $a_1^{[l-1]}, \dots, a_{n^{[l-1]}}^{[l-1]}$

This whole thing is 1 row

$$z_i^{[l]} = w_{i1}^{[l]} a_1^{[l-1]} + w_{i2}^{[l]} a_2^{[l-1]} + \dots + w_{in}^{[l]} a_n^{[l-1]}$$

$$= \sum_{j=1}^{n^{[l-1]}} w_{ij} a_j^{[l-1]}$$

These are the forward passes.

classmate

Date \_\_\_\_\_  
Page \_\_\_\_\_

Every layer in a neural network has 2 steps:

$$a^{[l]} = \frac{1}{1 + e^{-z^{[l]}}}$$

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

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

any activation

Dimensions:

$$w - n^{[l]} \times n^{[l-1]}$$

$$A^{[l-1]} - n^{[l-1]} \times 1$$

$$b^{[l]} - n^{[l]} \times 1$$

$$\Rightarrow z^{[l]} - n^{[l]} \times 1$$

Capitalised notation  
for matrices

$$A^{[l]} = \dim(z^{[l]})$$

$$= n^{[l]} \times 1$$

So this will help us get rid of some of the for loops.

If so in GD, for all minima & maxima,  $\frac{\partial J}{\partial w_{ij}^{[l]}} = 0$

If minima,  $\nabla_{i,j,l} J = 0$ . (not true vice-versa)

→ GD doesn't give guaranteed minimum, it does local minimum work.

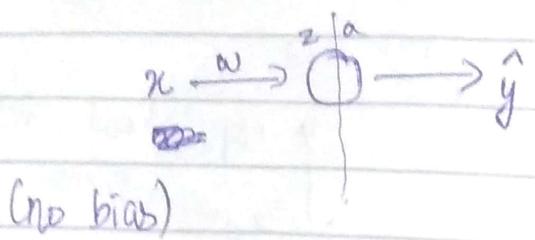
For our NN, our parameters are,  $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$   
 $w^{[2]}, b^{[2]}$

If we have 3 layers, where there are 5 I/P, 6 neurons in 1<sup>st</sup> layer, 3 neurons in 2<sup>nd</sup> layer & 1 in last layer,  
no. of parameters =  $((5 \times 6) + 6) + ((6 \times 3) + 3) + ((3 \times 1) + 1)$   
= 61. parameters

So once our GD completes, we have found best possible value of parameters for given initialisations.

How to find  $\frac{\partial J}{\partial w_{ij}}$ ?

Suppose we have only 1 neuron in 1 layer:



$$z = wx \quad \hat{y} = \frac{1}{1+e^{-z}}$$

GD eq:  $w = w - \alpha \frac{dL}{dw}$  → as there is only 1 layer,

$$\frac{dL}{dw} = \frac{dL}{d\hat{y}} \frac{d\hat{y}}{dz} \frac{dz}{dw} \quad ? \text{ Use this and it comes from dependencies \& chain rule}$$

Suppose we still have 1 neuron, but n I/p:

$$z = \sum_{i=1}^n x_i w_i + b$$

$$\hat{y} = a = \frac{1}{1+e^{-z}}$$

In prev 4 loops eq.: l & i loop disappear  
for  $j = 1$  to  $n \rightarrow$  no. of I/p  
 $w_{ij} = w_{ij} - \alpha \frac{\delta L}{\delta w_{ij}}$  → So we have  $n$  equations of the  
a one equation for bias.

Repeat

for  $l = 1$  to  $L$

for  $i = 1$  to  $n^{[l]}$

for  $j = 1$  to  $n^{[l-1]}$

$$w_{ij} = w_{ij} - \alpha \frac{\delta L}{\delta w_{ij}}$$

$$b_i^{[l]} = b_i^{[l]} - \alpha \frac{\delta L}{\delta b_i^{[l]}}$$

Just cause we are going over all  $w$ , doesn't mean it is a global minima.

$$\hat{y} = a^{[L]} = \text{as a class } g^{[L]}(z^{[L]})$$

$$z^{[L]} = w^{[L]} a^{[L-1]} + b^{[L]}$$

$$J = \frac{1}{m} \sum_{i=1}^m \alpha \delta_i (y^{[L]}, \hat{y}^{[L]})$$

So  $J$  is a func of  $\hat{y}$  which in turn is a func of  $g^{[L]}$  &  $z^{[L]}$  which is a func. of  $w^{[L]}, a^{[L-1]}, b^{[L]}$

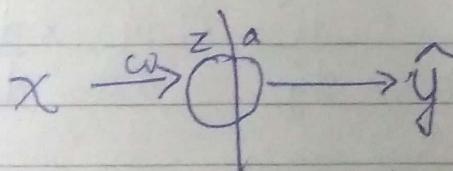
Out of these  $g^{[L]}, w^{[L]}, b^{[L]}$  are final parameters that don't depend on ~~other~~<sup>previous</sup> values.

$$a^{[L-1]} = g^{[L-1]} z^{[L-1]} \quad \text{These again don't depend on prev values!}$$

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

So in the end  $J$  depends on  $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots, w^{[L]}, b^{[L]}$

To find  $\frac{\partial J}{\partial w^{[L]}}$  we use back propagation.



$$\hat{y} = a = g(z) \quad (1)$$

$$z = w x \quad (2)$$

In prvs code, first 3 loops go away & our code becomes:

repeat

$$w = w - \frac{\delta J}{\delta w}$$

We are using 1 training sample, so  $J$  becomes  $b$  (as  $m=1$ )

How to compute  $\frac{dL}{dw}$  ?

$$\text{If } y = e^{x^2+2}$$

$$\frac{dy}{dx} = e^{x^2+2}(2x)$$

Now if we write this as:

$$y = e^z$$

$$z = x^2 + 2$$

$$\frac{dy}{dx} = \frac{dy}{dz} \times \frac{dz}{dx}$$

$$= e^z \times 2x = 2x e^z$$

$$\text{So } \frac{dL}{dw} = ?$$

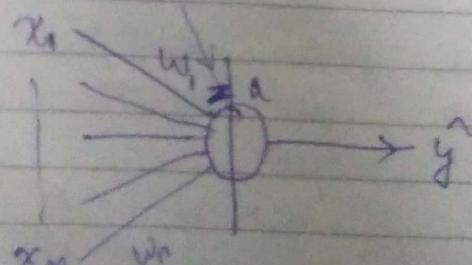
$L$  is a direct func. of  $y$  &  $\hat{y}$  but  $y$  is not a parameter (i.e. we have no control over it, it's given)

$$\frac{dL}{dw} = \frac{dL}{dy} \times \frac{dy}{dz} \quad (\text{Using ①})$$

If we had more examples to change except one  
 Next level of complication would be:

$$\frac{dL}{dw} = \frac{dL}{dy} \times \frac{dy}{dz} \times \frac{dz}{dw} \quad (\text{Using ②})$$

(logistic regression)



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

(Say  $g$  is sigmoid)

In prev code,  
 2<sup>nd</sup> for loop goes away as  $L=1$ ,  
 j loop from 1 to n goes away as  $n^{th} = 1$ , so we get

Forward pass:

$$W^{[1]} = [w_1 \quad w_2 \quad \dots \quad w_n]_{n \times n}, \quad X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

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

$$a = g(z)$$

$$\Rightarrow z = \sum_{i=1}^n x_i w_i + b.$$

$\frac{\delta L}{\delta w_1}, \frac{\delta L}{\delta w_2}, \dots, \frac{\delta L}{\delta w_n}$  &  $\frac{\delta L}{\delta b}$  will be computed  
 by gradient descent.

$$\frac{\delta L}{\delta w_1} = \frac{\delta L}{\delta \hat{y}} \times \frac{\delta \hat{y}}{\delta z} \times \frac{\delta z}{\delta w_1}$$

$$\frac{\delta L}{\delta w_2} = \frac{\delta L}{\delta \hat{y}} \times \frac{\delta \hat{y}}{\delta z} \times \frac{\delta z}{\delta w_2}$$

So in general,

$$\frac{\delta L}{\delta w_j} = \frac{\delta L}{\delta \hat{y}} \times \frac{\delta \hat{y}}{\delta z} \times \frac{\delta z}{\delta w_j} \quad \forall j \in [1, n]$$

$$\frac{\delta L}{\delta b} = \frac{\delta L}{\delta \hat{y}} \times \frac{\delta \hat{y}}{\delta z} \times \underbrace{\frac{\delta z}{\delta b}}$$

$$\frac{\delta L}{\delta b} = \frac{\delta L}{\delta \hat{y}} \times \frac{\delta \hat{y}}{\delta z} \times \frac{\delta z}{\delta b} = \frac{\delta L}{\delta \hat{y}} \times \frac{\delta \hat{y}}{\delta z}$$

This is = 1

$$\text{as } z = x_1 w_1 + x_2 w_2 + \dots + x_n w_n + b$$

Now instead of writing so many eq., say we  
 vectorise it:

$$\hat{y} = g(z) \quad z = w^{[1]} x + b$$

$$z = [w_{01} \quad w_{11} \quad \dots \quad w_{n1}]^T \begin{bmatrix} \\ \\ \vdots \\ \end{bmatrix} + b$$

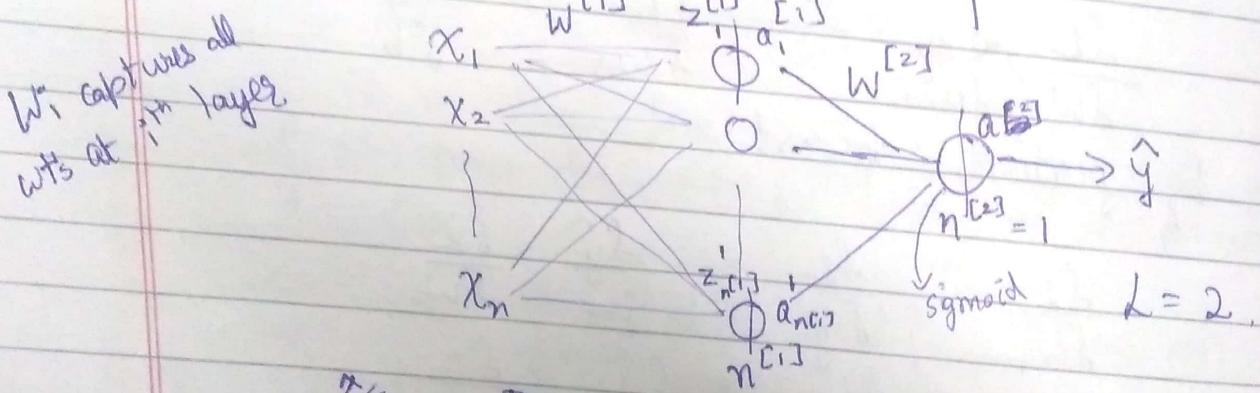
This is fine for 1 layer

Now we don't have  $x$  explicit for loops

$$w^{[1]} = w^{[1]} - \alpha \frac{\delta L}{\delta w^{[1]}}$$

This helps us get rid of j loop p.  
(using example of n S/P)

Next generalization: This part is like logistic regression



Dimension of  $w^{[1]} = n^{[1]} \times n^{[1]}$   
" "  $w^{[2]} = 1 \times n^{[1]}$

$$a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ \vdots \\ a_{n^{[1]}}^{[1]} \end{bmatrix} \quad n^{[1]} \times 1$$

Dimension of  $a$  has to be same as that of  $x$

$$z^{[1]} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ \vdots \\ z_n^{[1]} \end{bmatrix}_{n \times 1}$$

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

$$\begin{aligned} \hat{y} &= a^{[2]} = g^{[2]}(z^{[2]}) \\ z^{[2]} &= w^{[2]} a^{[1]} + b^{[2]} \end{aligned} \quad \left\{ \begin{array}{l} z^{[2]} = w^{[1]} a^{[1]} + b^{[1]} \\ \hat{y} = g^{[2]}(z^{[2]}) \end{array} \right.$$

We can't get rid of any of the for loops here.  
(without vectorisation).

$$\text{Back prop.} = \frac{\delta L}{\delta w^{[2]}}, \frac{\delta L}{\delta w^{[1]}}, \frac{\delta L}{\delta b^{[1]}}, \frac{\delta L}{\delta b^{[2]}}$$

(so with vectorisation, we don't need any of the for loops)

$$\frac{\delta L}{\delta w^{[2]}} = \frac{\delta L}{\delta \hat{y}} \times \frac{\delta \hat{y}}{\delta z^{[2]}} \times \frac{\delta z^{[2]}}{\delta w^{[2]}}$$

If we take the case of logistic regression,

$$\frac{\delta L}{\delta w_i} = \frac{\delta L}{\delta \hat{y}} \times \frac{\delta \hat{y}}{\delta z} \times \frac{\delta z}{\delta w_i} \quad \textcircled{1}$$

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

$$\frac{\delta L}{\delta \hat{y}} = -\left(\frac{y}{\hat{y}} - \frac{1-y}{1-\hat{y}}\right) = \frac{1-y}{1-\hat{y}} - \frac{y}{\hat{y}} = \frac{\hat{y}-y-y+\hat{y}y}{(1-\hat{y})\hat{y}} \quad \textcircled{2}$$

$$\hat{y} = \frac{1}{1+e^{-z}} \Rightarrow \frac{\delta \hat{y}}{\delta z} = \frac{-1}{(1+e^{-z})^2} \times -e^{-z} = \frac{e^{-z}}{(1+e^{-z})^2}$$

This is  $g(z)$  - sigmoid

classmate  
Date \_\_\_\_\_  
Page \_\_\_\_\_

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

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

→ This is nothing but  $\hat{y}$  as there is only one place where we are using?

$$\frac{\delta z}{\delta w_i} = \frac{f_w x_i + b}{f_w} = n_i$$

∴

So ① can be written as:

$$\frac{\delta L}{\delta w_i} = \left( \frac{\delta L}{\delta z} \right) x_i.$$

→ This is same irrespective of  $i$ .

$$\begin{aligned} \frac{\delta L}{\delta w_{i,j}} &= \begin{bmatrix} \frac{\delta L}{\delta w_1} & \frac{\delta L}{\delta w_2} & \dots & \frac{\delta L}{\delta w_n} \end{bmatrix} \\ &= \begin{bmatrix} \frac{\delta L}{\delta z} x_1 & \frac{\delta L}{\delta z} x_2 & \dots & \frac{\delta L}{\delta z} x_n \end{bmatrix} = \frac{\delta L}{\delta z} \times x^T \end{aligned}$$

$\hat{y} = a^{(i)}$  for this

as  $x$  is column matrix

Upcoming diagram

④

$$\frac{\delta L}{\delta z} = \frac{\delta L}{\delta \hat{y}} \times \frac{\delta \hat{y}}{\delta z}$$

$$(From ②) = \frac{\hat{y} - y}{(1-\hat{y})\hat{y}} \times g(z)(1-g(z))$$

$$(From ③) = \frac{\hat{y} - y}{(1-\hat{y})\hat{y}} \times \hat{y}(1-\hat{y}) = \hat{y} - y. \quad ⑤$$

$$From ④ \& ⑤ : \frac{\delta L}{\delta w_{i,j}} = (\hat{y} - y) x^T \rightarrow \text{Calculate } \frac{\delta L}{\delta w} = \frac{\delta L}{\delta z} = \hat{y} - y$$

Now if we take the case of the diagram on previous page,

$z$  is func. of  $w^{[2]}, b^{[2]}, w^{[1]}, b^{[1]}$

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

$$a^{[1]} = g^{[1]}(z^{[1]}) \rightarrow \text{Superscript of } g \text{ is a general form}$$

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

$$a^{[2]} = g^{[2]}(z^{[2]}) = \hat{y}.$$

if we use more than one activation func.

$$\frac{\delta L}{\delta W^{[2]}} = \frac{\delta L}{\delta z^{[2]}} \times \frac{\delta z^{[2]}}{\delta W^{[2]}}$$

$z^{[2]}$   $\downarrow$   
as we are  
doing back prop.

$\frac{\delta L}{\delta z^{[2]}} = \hat{y} - y = a^{[2]} - y$  & similar to form,  $\frac{\delta L}{\delta w^{[2]}} = a^{[2]T}$

$$\Rightarrow \frac{\delta L}{\delta W^{[2]}} = (a^{[2]} - y) a^{[1]T}$$

$$\frac{\delta L}{\delta b^{[2]}} = a^{[2]} - y.$$

Now if consider the left part of the diagram (diff. from logistic reg.)

$\frac{\delta L}{\delta z^{[1]}} = w^{[2]T} \frac{\delta L}{\delta z^{[2]}} \times g^{[1]}(z^{[1]})$  This is scalar (S)

$\frac{\delta L}{\delta W^{[2]}} = \frac{\delta L}{\delta z^{[1]}} x^T$  Element wise multiplication  
not matrix multiplication

$\frac{\delta L}{\delta b^{[1]}} = \frac{\delta L}{\delta z^{[1]}}$  Final eq.

Now to get these:

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

$$\frac{\delta L}{\delta z^{[1]}} = \frac{\delta L}{\delta a^{[2]}} \times \frac{\delta a^{[2]}}{\delta z^{[2]}} \times \frac{\delta z^{[2]}}{\delta a^{[1]}} \times \frac{\delta a^{[1]}}{\delta z^{[1]}}$$

$$\frac{\delta L}{\delta z^{[1]}} = \frac{\delta L}{\delta z^{[2]}}$$

$$W^{[2]} = \begin{bmatrix} W_1^{[2]} & W_2^{[2]} \\ \vdots & \vdots \\ W_{n^{[1]}}^{[2]} \end{bmatrix}$$

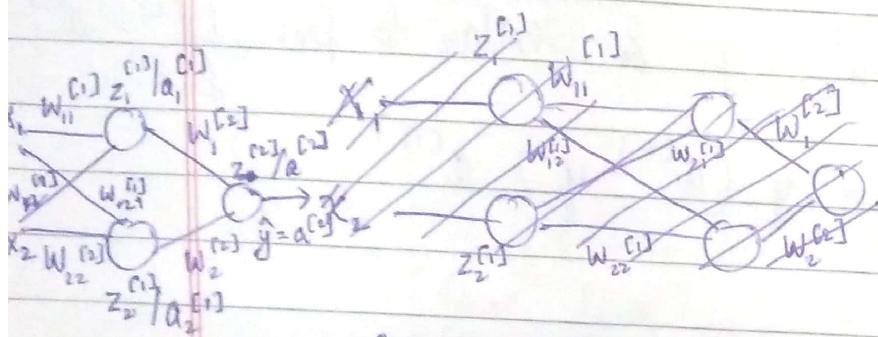
$$W^{[2]T} = \begin{bmatrix} W_1^{[2]} \\ W_2^{[2]} \\ \vdots \\ W_{n^{[1]}}^{[2]} \end{bmatrix}$$

When we multiply this by  $\frac{\delta L}{\delta z^{[2]}}$ , then we are finding

contribution of every neuron

i.e propagating the error.

$\Rightarrow$  Thus we get why  $w^{[2]T} \frac{\delta L}{\delta z^{[2]}}$



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

$$\frac{\delta L}{\delta w_1^{[2]}} = \frac{\delta L}{\delta \hat{y}} \frac{\delta \hat{y}}{\delta z^{[2]}} \frac{\delta z^{[2]}}{\delta w_1^{[2]}} \quad \frac{\delta L}{\delta w_2^{[2]}} = \frac{\delta L}{\delta \hat{y}} \frac{\delta \hat{y}}{\delta z^{[2]}} \frac{\delta z^{[2]}}{\delta w_2^{[2]}}$$

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

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

$$\textcircled{1} \quad \frac{\delta L}{\delta w_{11}^{[1]}} = \frac{\delta L}{\delta \hat{y}} \times \frac{\delta \hat{y}}{\delta z^{[2]}} \times \frac{\delta z^{[2]}}{\delta a_1^{[1]}} \times \frac{\delta a_1^{[1]}}{\delta z_1^{[1]}} \times \frac{\delta z_1^{[1]}}{\delta w_{11}^{[1]}}$$

$$\textcircled{2} \quad \frac{\delta L}{\delta w_{12}^{[1]}} = \frac{\delta L}{\delta \hat{y}} \times \frac{\delta \hat{y}}{\delta z^{[2]}} \times \frac{\delta z^{[2]}}{\delta a_1^{[1]}} \times \frac{\delta a_1^{[1]}}{\delta z_2^{[1]}} \times \frac{\delta z_2^{[1]}}{\delta w_{12}^{[1]}}$$

$$\textcircled{5} \quad dz^{[1]} = W^{[2]T} dz^{[2]} * g'(z^{[1]})$$

$$\textcircled{3} \quad \frac{\delta L}{\delta w_{21}^{[1]}} = \frac{\delta L}{\delta \hat{y}} \times \frac{\delta \hat{y}}{\delta z^{[2]}} \times \frac{\delta z^{[2]}}{\delta a_2^{[1]}} \times \frac{\delta a_2^{[1]}}{\delta z_1^{[1]}} \times \frac{\delta z_1^{[1]}}{\delta w_{21}^{[1]}}$$

$$\textcircled{4} \quad \frac{\delta L}{\delta w_{22}^{[1]}} = \dots \times \dots \times \frac{\delta z^{[2]}}{\delta a_2^{[1]}} \times \frac{\delta a_2^{[1]}}{\delta z_2^{[1]}} \times \frac{\delta z_2^{[1]}}{\delta w_{22}^{[1]}}$$

Now let us compare  $\textcircled{5}$  with  $\textcircled{1}, \textcircled{2}, \textcircled{3}, \textcircled{4}$ :

We need only till  $\delta z^{[1]}$ , so we can ignore last terms of  $\textcircled{1}, \textcircled{2}, \textcircled{3}, \textcircled{4}$ .

First 2 terms of  $\textcircled{1}, \textcircled{2}, \textcircled{3}, \textcircled{4}$ ,

$$\frac{\delta L}{\delta \hat{y}} \times \frac{\delta \hat{y}}{\delta z^{[2]}} = \frac{\delta L}{\delta z^{[2]}} = dz^{[2]}$$

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

$$= \begin{bmatrix} w_1^{[2]} & w_2^{[2]} \end{bmatrix} \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \end{bmatrix} + b^{[2]}$$

$$\frac{\partial z^{[2]}}{\partial a_1^{[1]}} = w_1^{[2]}$$

$$\frac{\partial z^{[2]}}{\partial a_2^{[1]}} = w_2^{[2]}$$

⑥

$$a_1^{[1]} = g_1(z_1^{[1]})$$

$$\frac{\partial a_1^{[1]}}{\partial z_1^{[1]}} = g_1'(z_1^{[1]}) - ⑦$$

So from eq. ⑤:

$$\frac{\partial L}{\partial z_1^{[1]}} = \underbrace{\frac{\partial L}{\partial y} \times \frac{\partial y}{\partial z^{[2]}}}_{\downarrow} \times \underbrace{\frac{\partial z^{[2]}}{\partial a_1^{[1]}} \times \frac{\partial a_1^{[1]}}{\partial z_1^{[1]}}}_{\downarrow}$$

This gives a term in terms of wt(⑥)

This gives ⑦

- Similar to ⑤.

$$\Rightarrow \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \end{bmatrix} = \begin{bmatrix} w_1^{[1]} x_1 + b_1^{[1]} \\ w_{12}^{[1]} x_1 + w_{21}^{[1]} x_2 + b_2^{[1]} \end{bmatrix}$$

$$\begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \end{bmatrix}_{2 \times 1} = \begin{bmatrix} w_{11}^{[1]} x_1 + w_{12}^{[1]} x_2 + b_1^{[1]} \\ w_{21}^{[1]} x_1 + w_{22}^{[1]} x_2 + b_2^{[1]} \end{bmatrix}$$

So even  $\partial z^{[1]}$  must be of dim.  $2 \times 1$ .

Now  $\partial z^{[2]}$  is a scalar =  $a^{[2]} - y$  (calculated previously)

$$\frac{\partial L}{\partial z_1^{[1]}} = \frac{\partial z_1^{[1]}}{\partial z_1^{[1]}} = (a^{[2]} - y) w_1^{[2]} \cdot g_1'(z_1^{[1]}) \rightarrow \text{not matrix mult, element wise mult}$$



$$\begin{bmatrix} \frac{\partial L}{\partial z_1^{[1]}} \\ \frac{\partial L}{\partial z_2^{[1]}} \end{bmatrix} = \begin{bmatrix} (a^{[2]} - y) & w_1^{[2]} g_1^{[1]}(z_1^{[1]}) \\ (a^{[2]} - y) & w_2^{[2]} g_2^{[1]}(z_2^{[1]}) \end{bmatrix}$$

$$= (a^{[2]} - y) \begin{bmatrix} w_1^{[2]} & g_1^{[1]}(z_1^{[1]}) \\ w_2^{[2]} & g_2^{[1]}(z_2^{[1]}) \end{bmatrix} \quad (8)$$

$$W^{[2]} = [w_1^{[2]} \quad w_2^{[2]}]_{1 \times 2} \Rightarrow W^{[2]T} = \begin{bmatrix} w_1^{[2]} \\ w_2^{[2]} \end{bmatrix}_{2 \times 1}$$

$$(8) = (a^{[2]} - y) \underbrace{W^{[2]T}}_{\substack{\text{wise mult} \\ \text{element}}} \underbrace{g^{[1]}(z^{[1]})}_{\substack{\text{of} \\ \text{size}}} \quad (5)$$

If we have  $L$  layers, to find  $\frac{\partial L}{\partial w^{[i]}}$ , we'll have approx. 2 terms per layer i.e.  $\frac{\partial L}{\partial w^{[i]}}$  we'll have no. of terms of order  $2L$ .

The problem with this is say if  $L = 100$  & we have around 200 partial. der., & say about 100 are  $\approx 1$  & 100 are  $\approx 0.1$ , then our

$$\frac{\partial L}{\partial w_{ij}^{[i]}} \text{ will become very close to } 0$$

$$\Rightarrow w_{ij}^{[i]} = w_{ij}^{[i]} - \alpha \frac{\partial L}{\partial w_{ij}^{[i]}} \rightarrow 0$$

So our G. descent algo will stop updating parameters, so algo will stop as it will think it found optimal values.

There is no solution to this problem (Vanishing gradient problem).

classmate  
Date \_\_\_\_\_  
Page \_\_\_\_\_

similarly if our derivatives  $> 1$ , then our  $\frac{dJ}{dW_{ij}}$  becomes very large & we will have to subtract a huge number.

gradient problem

To this problem, we have a solution (gradient clipping), so we fix a threshold & everytime our term crosses the threshold, we make the  $\frac{dJ}{dW_{ij}} = \text{threshold}$

This happens only with deep NN

Try this with 3 nodes in 1<sup>st</sup> layer & 1 node in the 2<sup>nd</sup> layer

for 1 training example

$$\begin{cases} dZ^{[2]} = A^{[2]} - y \\ dW^{[2]} = dZ^{[2]} a^{[1]T} \\ db^{[2]} = dZ^{[2]} \end{cases}$$

$$\begin{aligned} dZ^{[1]} &= W^{[2]T} dZ^{[2]} * g^{[1]'}(z^{[1]}) & (dZ^{[2]} = W^{[2]T} dZ^{[1]} * g^{[2]'}(z^{[2]})) \\ dW^{[1]} &= dZ^{[1]} X^T * g^{[1]'}(z^{[1]}) & (dW^{[2]} = dZ^{[2]} a^{[1]T}) \\ db^{[1]} &= dZ^{[1]} \end{aligned}$$

Vectorised form (to get rid of some for loops):

$$\begin{aligned} dZ^{[2]} &= A^{[2]} - y \\ dW^{[2]} &= \frac{1}{m} dZ^{[2]} A^{[1]T} \end{aligned}$$

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

$$dZ^{[1]} = W^{[2]T} dZ^{[2]} * g^{[1]'}(z^{[1]})$$

Vectorise last 2 eqs. on your own

1<sup>st</sup> training example

$$y = [y^{(1)} \quad y^{(2)} \quad \dots \quad y^{(m)}]_{1 \times m}$$

$$A^{[2]} = [a^{[1](1)} \quad a^{[1](2)} \quad \dots \quad a^{[1](m)} \quad a^{[2](1)} \quad a^{[2](2)} \quad \dots \quad a^{[2](m)}]$$

Here we can write  $\hat{y} = \hat{y}^{(1)} \quad \hat{y}^{(2)} \quad \dots \quad \hat{y}^{(m)}$  or  $a^{[1](1)} \quad a^{[1](2)} \quad \dots \quad a^{[1](m)}$

$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \Rightarrow X = \begin{bmatrix} | & | & | & | \\ x^{(1)} & x^{(2)} & \cdots & x^{(m)} \\ | & | & | & | \end{bmatrix}_{n \times m}$$

no. of features

$$\dim(Z^{[2]}) = \dim(dZ^{[2]}) = [x]_{[2](m)'}^T$$

$$Z^{[2]} = \begin{bmatrix} Z^{[2](1)} \\ Z^{[2](2)} \\ \vdots \\ Z^{[2](m)} \end{bmatrix}_{1 \times m}$$

$$dZ^{[2]} = \begin{bmatrix} (a^{[2](1)} - y^{(1)}) \\ a^{[2](2)} \\ \vdots \\ a^{[2](m)} \end{bmatrix}_{1 \times m}$$

$$A^{[1]} = \begin{bmatrix} a^{[1](1)} \\ a^{[1](2)} \\ \vdots \\ a^{[1](m)} \end{bmatrix}_{n \times m}$$

$$A^{[1]T} = \begin{bmatrix} a^{[1](1)} \\ a^{[1](2)} \\ \vdots \\ a^{[1](m)} \end{bmatrix}_{m \times n^{[1]}}$$

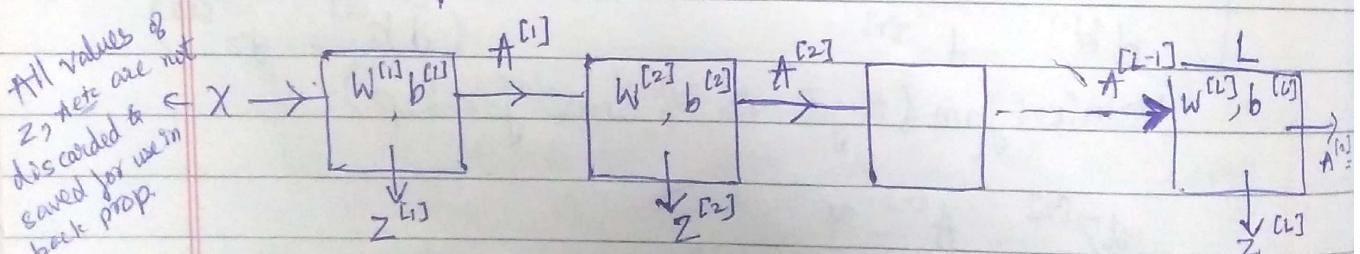
↓ when  $a^{[1](1)}$  is also written  
in complete form as it is also  
a vector

$$\dim(W^{[2]}) = \dim(dW^{[2]}) = 1 \times n^{[1]}$$

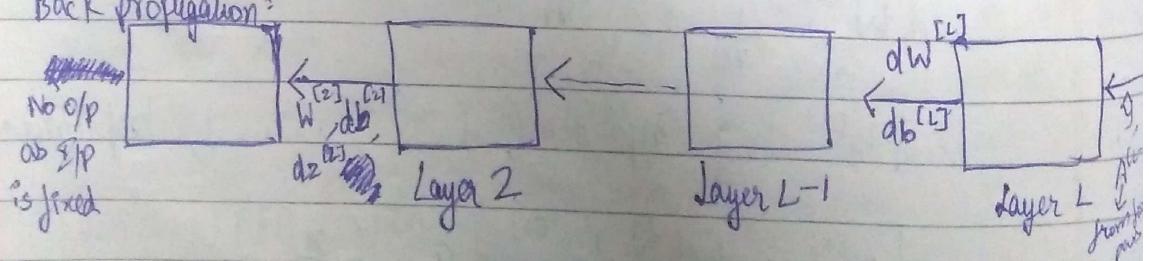
$$dW^{[2]} = \begin{bmatrix} dW_1^{[2]} & dW_2^{[2]} & \cdots & dW_{n^{[1]}}^{[2]} \end{bmatrix}$$

The  $1/m$  factor comes in as we have to sum  
& avg. over  $m$  examples

Forward pass:



Back propagation:

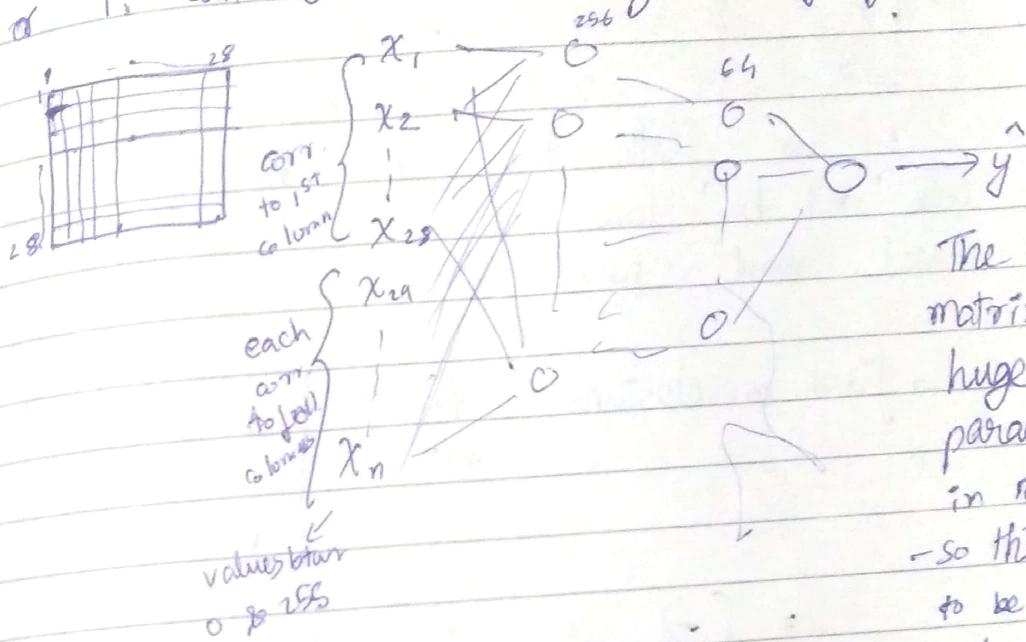


Initialisation of weights: don't give all 0's as then all values of  $x_0/x_1/x_2$  will be zero, don't give all 0's as then all values will be same, so we give random values

We want wts to be small, so if we multiply by a small no like 0.001 if range of rand no is large

See notes of in btm class.

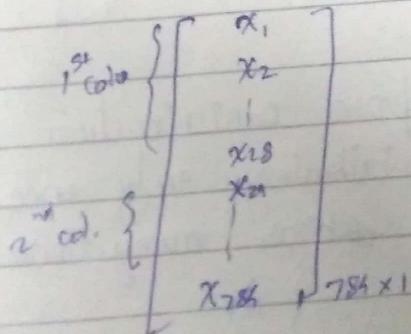
Say we have  $28 \times 28$  images which have either 0 or 1. Can we write a binary classifier for this?



The wt. matrix will be huge. & many parameters (about in millions)

- so this was thought to be inconvenient.

In tensorflow we have a layer flatten which will put the image in one column matrix.



Using this loses some properties of the image.

CNN's came into use in late 80's when banks needed automatic reading of docs. We can't use above arch as it whom we have huge I/p's, so we use CNN's which give an edge above this.

We want a NN which keeps / preserves local features in an image & uses less no of parameters

ReLU puts -ve  
values as 0

classmate

Date \_\_\_\_\_  
Page \_\_\_\_\_

$$A = \begin{bmatrix} 5 & 3 & 6 & 8 & 9 & 1 \\ 1 & 6 & 8 & 3 & 2 & 1 \\ 9 & 9 & 3 & 2 & 1 & -1 \\ 1 & 1 & 8 & 4 & 6 & 5 \\ 9 & 3 & 4 & 0 & 0 & 1 \\ 0 & 1 & 2 & 3 & 4 & 3 \end{bmatrix}_{6 \times 6}$$

row elements  
 $x^1 + x^0 + x^{-1}$

If we flatten this, our no. of parameters  
will not be low.  
Read about CNN's

To find convolutions:  $A \in \mathbb{R}^{m \times n}$

$$\begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix}_{f \times f}$$

filter: helps in  
differentiation  
of right place with  
vertical edge  
prob see fig

$$= \begin{bmatrix} -2 & 5 & 5 & 10 \\ -8 & 7 & - & - \\ -4 & - & - & - \end{bmatrix}_{4 \times 4}$$

$$\rightarrow (m-f+1) \times (n-f+1)$$

$\rightarrow$  image size  
is reducing  
which is good

Now if we compare contribution of circled 5 & 3, 5 is contributing only once to o/p while 3 is contributing ~~as~~ as many times as no. of elements in filter i.e. 9.

Now if our image is such if our corner is also imp, above problem will cause an issue, so the easiest way to do is to pad:

$$\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & \text{elements of } A & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}_{8 \times 8}$$

$$n-f+1$$

$$= 8-3+1$$

$$= 6$$

$$\text{padding} = p$$

$$n' = n + 2p$$

$$n_o = n' - f + 1 = n + 2p - f + 1$$

$$\text{Dimension: } (n+2p-f+1) \times (n+2p-f+1)$$

$$\begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix} : \text{vertical edge filter.}$$

$$\begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{bmatrix} : \text{horizontal edge filter.}$$

Strides:

$$\begin{array}{c}
 \xrightarrow{2} \\
 \downarrow \\
 \begin{bmatrix} 0 & 150 \\ 0 & 150 \\ 0 & 150 \\ 0 & 150 \end{bmatrix} \times \begin{bmatrix} 1 & -1 \\ 1 & -1 \end{bmatrix} = \begin{bmatrix} -300 & 300 \\ -300 & 300 \end{bmatrix}
 \end{array}$$

from 6x6 matrix convolution  
 stride 2  
 6x6 input  
 3x3 kernel  
 2x2 output

(Stride of 2)

If stride = s

$$n_o = \left\lfloor \frac{n + 2p - f}{s} \right\rfloor + 1 : \text{Proof?} \quad \star \star$$

Andrew's course gives this formula ab:

$$n_o = \left\lfloor \frac{n + 2p - f + 1}{s} \right\rfloor$$

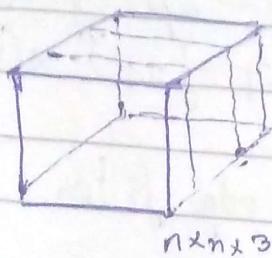
But if we take  $p=0$ ,  $6 \times 2$   $f=4$  ( $s>2$ ),  
 $n=44$ ,  $f=3$ .

$$\textcircled{1} \quad n_o = \left\lfloor \frac{44-3}{4} \right\rfloor + 1 = 11 \rightarrow \text{This was found to be better}$$

$$\textcircled{2} \quad n_o = \left\lfloor \frac{44-3+1}{4} \right\rfloor = 10 \quad (\text{no fundamental reason})$$

Say we have a  $n \times n \times 1$  image

Coloured image:  $(R, G, B)$ : 24 bits



$0-255$

3 channels

i.e. 3 2-D matrices put one over other

How to convolve?

Filter:  $3 \times 3 \times 1 \rightarrow O/P$  dimension:  $4 \times 4 \times 3$   
There is nothing fundamentally wrong with this but we  
don't use this, as below one performs better.

If our filter was  $3 \times 3 \times 2$ ,  
then o/p dim. will  
be  $4 \times 4 \times 2$ .

If our filter:  $3 \times 3 \times 3 \rightarrow O/P$  dimension:  $4 \times 4 \times 1$   
To convolve we multiply  $R^G B$  channel of I/P &  
filter correspondingly  
To get one element of the O/P we ~~will~~ do 27  
multiplications.

Notations:

Generalised dimension of volume:  $n_w^{[l]} \times n_h^{[l]} \times n_c^{[l]}$   
width height channel,

$p^{[l]}$  - padding in  $l^{th}$  layer

$s^{[l]}$  - strides in " "

$f^{[l]}$  - dimension of filter " "

Even filter size used in pooling.

Do in :

$$n_w^{[0]} \times n_h^{[0]} \times n_c^{[0]} \times f^{[1]} \times n_w^{[1]} \times n_h^{[1]} \times n_c^{[1]} = 4 \times 4 \times 1$$

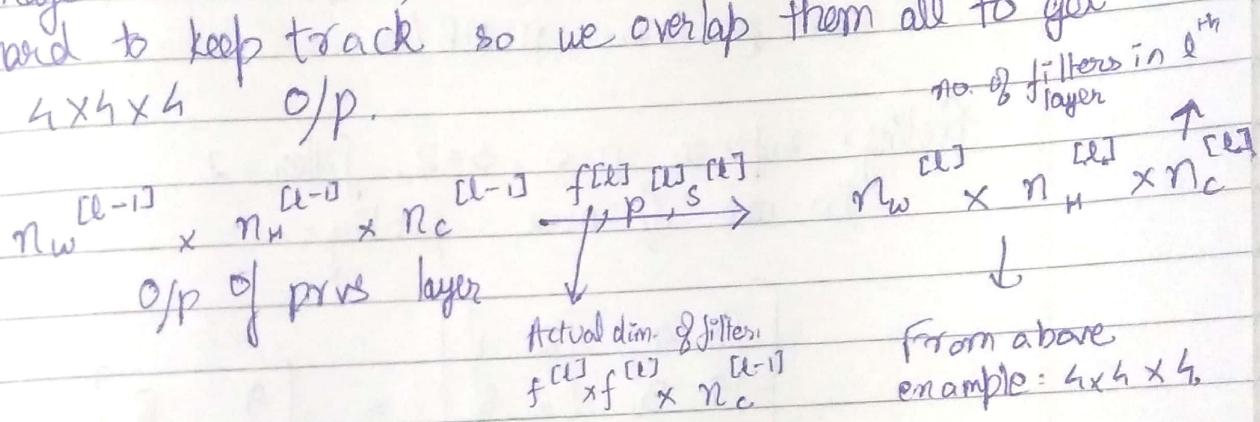
is  $p^{[0]}$  doesn't make sense, padding done in I/P

$$p^{[1]} = 0, s^{[1]} = 1$$

If we use 4 diff. filters of same dimensions:

$$\begin{array}{l}
 \text{* } f_1 (3 \times 3 \times 3) \\
 \text{* } f_2 \\
 \text{* } f_3 \\
 \text{* } f_4 \\
 \hline
 6 \times 6 \times 3
 \end{array}
 \quad =
 \quad
 \begin{array}{l}
 4 \times 4 \times 1 \\
 4 \times 4 \times 1 \\
 4 \times 4 \times 1 \\
 4 \times 4 \times 1
 \end{array}$$

Every filter will highlight a diff. aspect of the image (horizontal lines, vertical lines etc), but it will be hard to keep track so we overlap them all to get a  $4 \times 4 \times 4$  o/p.



Convolution is a linear operation (it is equivalent to  $Z^{[L]} = W^{[L]} A^{[L-1]} + b^{[L]}$  - we are yet to add bias)

$$\begin{array}{ccc}
 n_w^{L-1} \times n_h^{L-1} \times n_c^{L-1} & \xrightarrow{\quad} & n_w^L \times n_h^L \times n_c^L \\
 \begin{matrix} \boxed{\phantom{00}} \\ \boxed{\phantom{00}} \\ \boxed{\phantom{00}} \end{matrix} & \xrightarrow{\quad} & \begin{matrix} \boxed{\phantom{00}} \\ \boxed{\phantom{00}} \\ \boxed{\phantom{00}} \end{matrix} \\
 \text{f}^{[L]} & & 
 \end{array}$$

To add bias  
we add  $b^{[L]}$  to all  
elements in the o/p  
matrix.

$b^{[L]} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \vdots \\ \cdot \end{bmatrix} \quad \# \text{filter} \times 1$

So we have completed the  $Z$  part i.e linear part  
Our parameters are incorporated in the filters.  
They are learnt:  $\begin{bmatrix} w_1 & w_2 & \dots \end{bmatrix}$

If we don't have non-linearity in every layer, then we can just remove layers & compress to one as linear func. Non-linear func. is a linear func.

$$\text{To find } A^{[l]} = g^{[l]}(z^{[l]}) \\ \text{ReLU } (n_w^{[l]} \times n_h^{[l]} \times n_c^{[l]}) = n_w^{[l]} \times n_h^{[l]} \times n_c^{[l]}$$

In practice

we normalise the values so  $> 255$  won't arise.

Except case of ReLU, all values  $< 0 = 0$ , otherwise same. In case of image value  $> 255 = 255$

CNN's are closer to the vision part of our networks. Our conscious mind can process 50 b/s, our sensor can perceive 11 MB/b.

Pooling: Suppose  $p=0$ ,  $s=2$ , filter = 2.

In case of max pooling:

$$\begin{bmatrix} [8 \ 4] & [2 \ 1] \\ [9 \ 3] & [2 \ 8] \\ [1 \ 5] & [7 \ 1] \\ [2 \ 0] & [3 \ 7] \end{bmatrix} = \begin{bmatrix} 9 & 8 \\ 5 & 7 \end{bmatrix}_{2 \times 2}$$

We can take any other func other than max, chosen empirically.

Pooling is done only after convolution by convention.

Translation is invariant in pooling (to some extent).

Rotation is not " " " but you can make it learn that invariance. There are no learning parameters in pooling.

Say we have a  $100 \times 100 \times 3$  & we use 10 filters of  $3 \times 3 \times 3$  to do 1 convolution + pooling. What is the no. of learnable parameters?

I/p has no parameters.

Filter has  $(27 \times 10) + 10$  parameters  
 $= 280$

Computations cost to ~~Op~~ size  
will size

classmate

Date \_\_\_\_\_  
Page \_\_\_\_\_

If we change image size to  $200 \times 200 \times 3$  i.e. 4 times the size, using CNN the no. of parameters will still remain same, i.e. it is invariant to pic size while if we were using our normal DNN, we would have a massive increase in no. of parameters. That's why in images we use CNN as I/p size is huge.

Suppose we have a  $6 \times 6$  b/w image which we flatten & say we have 5 neurons in the 1<sup>st</sup> layer i.e.  $W^{(1)} = 5 \times 36$

