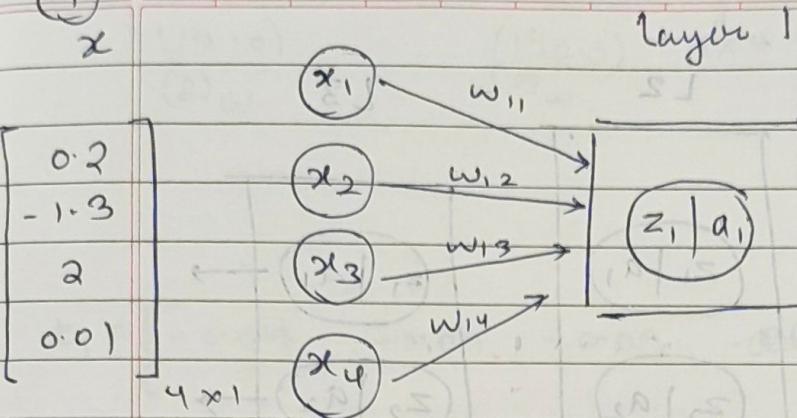


(1)



$$z_1 = w_{11}x_1 + w_{12}x_2 + w_{13}x_3 + w_{14}x_4$$

$$a_1 = \text{softmax}(z_1)$$

$$\begin{bmatrix} w_{11} \\ w_{12} \\ w_{13} \\ w_{14} \end{bmatrix} = \begin{bmatrix} 0.001 \\ 0.01 \\ -0.005 \\ -1.2 \end{bmatrix}_{4 \times 1} \quad \begin{bmatrix} a_1 = \text{softmax}(-0.0348) \\ = 0. \end{bmatrix}$$

$$1. z_1 = 0.001 \times 0.2 + 0.01 \times (-1.3) + -0.005 \times (2) + 0.01 \times (-1.2)$$

$$= 0.0002 - 0.013 - 0.01 - 0.012$$

~~$= 0.0002 - 0.0025 - 0.01$~~

~~$= 0.0026 - 0.01$~~

~~$= -0.0123$~~

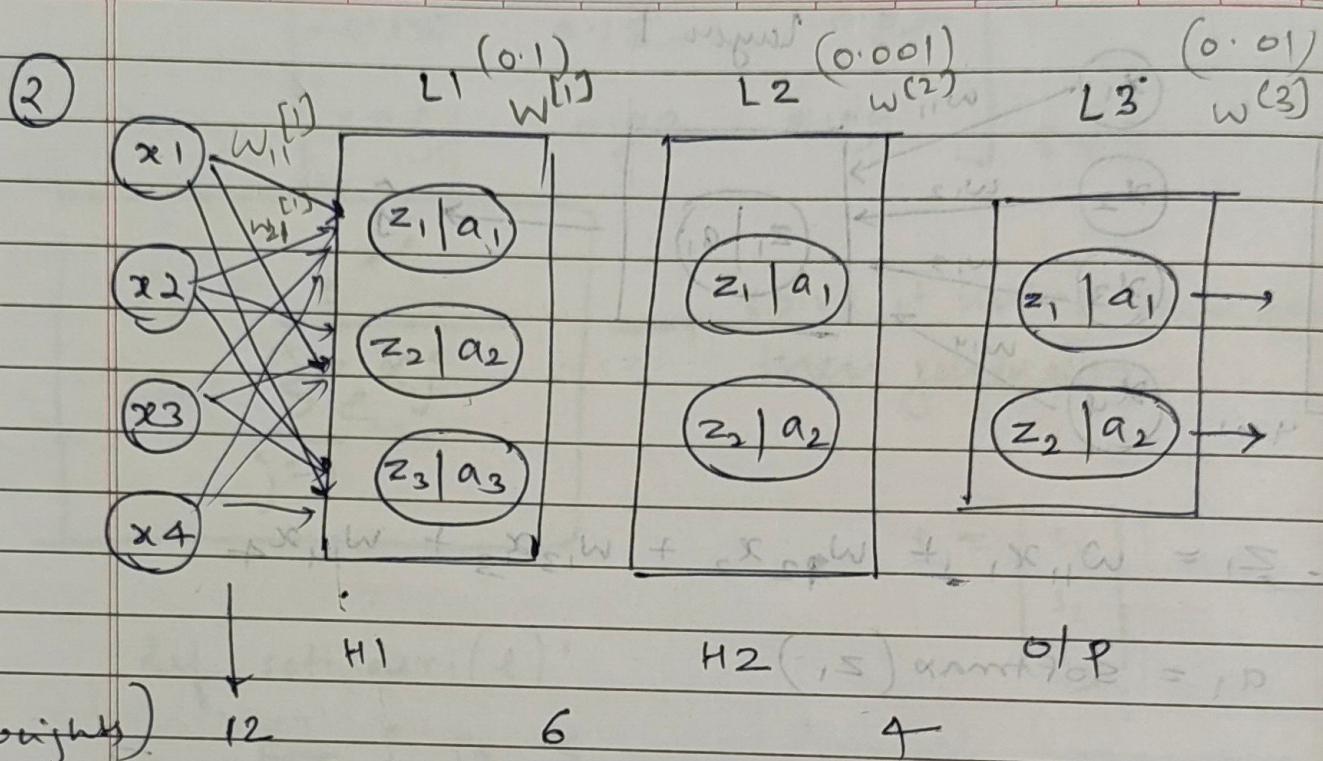
$= -0.0348$

$$2. z_1 = x^T w = \begin{bmatrix} 0.2 & -1.3 & 2 & 0.01 \end{bmatrix} \times \begin{bmatrix} 0.001 \\ 0.01 \\ -0.005 \\ -1.2 \end{bmatrix}$$

$$= 0.2 \times 0.001 - 1.3 \times 0.01 + 2 \times (-0.005) - 1.2 \times 0.01$$

$$= -0.0348$$

- * hidden layer \rightarrow ReLU
- * output layer \rightarrow Softmax (make 2 neurons)



π (weights) 12 6 4
ReLU ReLU softmax

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} -2.4 \\ 1.2 \\ -0.8 \\ 2.0 \end{bmatrix}$$

$\begin{matrix} w_{11}^{[1]} \\ w_{21}^{[1]} \\ w_{31}^{[1]} \\ w_{41}^{[1]} \end{matrix} = \begin{matrix} 0.1 \\ 0.1 \\ 0.1 \\ 0.1 \end{matrix}$

$w_{12}^{[1]} = 5.0 \times 100.0$

$w_{22}^{[1]} = 5.0 \times 100.0$

$w_{32}^{[1]} = 5.0 \times 100.0$

$w_{42}^{[1]} = 5.0 \times 100.0$

$w_{13}^{[1]} = 5.0 \times 100.0$

$w_{23}^{[1]} = 5.0 \times 100.0$

$w_{33}^{[1]} = 5.0 \times 100.0$

$w_{43}^{[1]} = 5.0 \times 100.0$

$w_{14}^{[1]} = 5.0 \times 100.0$

$w_{24}^{[1]} = 5.0 \times 100.0$

$w_{34}^{[1]} = 5.0 \times 100.0$

$$b_1 = \begin{bmatrix} x_1 & x_2 & x_3 & x_4 \end{bmatrix}^T \begin{bmatrix} w_{11} & w_{21} & w_{31} \\ w_{12} & w_{22} & w_{23} \\ w_{13} & w_{23} & w_{33} \\ w_{14} & w_{24} & w_{34} \end{bmatrix}$$

row or weights from the activation function is n (input).

$$= \begin{bmatrix} -2.4 & 1.2 & -0.8 & 1.1 \end{bmatrix} \begin{bmatrix} 0.1 & 0.1 \\ 0.1 & 0.1 \\ 0.1 & 0.1 \\ 0.1 & 0.1 \end{bmatrix}$$

Date _____
Page _____

4×3

$$t_1 = \begin{bmatrix} -0.09 & -0.09 & -0.09 & -0.09 \end{bmatrix} - 1.5 \rightarrow \text{target}$$

$$t_2 = \begin{bmatrix} -0.09 & -0.09 & 0.09 & 0.09 \end{bmatrix}$$

$$\delta P = [a_1(-0.09) \quad a_2(-0.09) \quad a_3(-0.09)]$$

$$\delta P = [0 \quad 0 \quad 0]$$

with other two part below of #

$$L_2 = \begin{bmatrix} 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0.001 & 0.001 & 0.001 \\ 0.001 & 0.001 & 0.001 \\ 0.001 & 0.001 & 0.001 \end{bmatrix}$$

using 3x3 to make 3x2
so 3x3 to 3x2
 $(!) =$ reverse two rows to make #

$$= [0 \quad 0]$$

$$L_2 = [a_1(0) \quad a_2(0)] \quad \text{and} \quad [0 \leftarrow 0] \quad \begin{matrix} 0.6 & \vdots & x \\ 0.1 & \vdots & x \\ 1.0 & \vdots & x \end{matrix}$$

$$L_3 = \begin{bmatrix} 0 & 0 \end{bmatrix}_{1 \times 2} \begin{bmatrix} 0.01 & 0.01 \\ 0.01 & 0.01 \end{bmatrix}$$

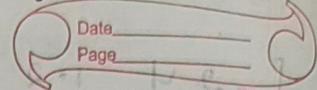
deriving
 $i = j = k$

$$L_3 = \begin{bmatrix} 0 & 0 \end{bmatrix} = [\text{softmax}(0) \quad \text{softmax}(0)]$$

$$\delta P = [0.5 \quad 0.5] = [(x = 1) \quad (x = 2)]$$

$(1.0)^3 + (0.1)^3 + (0.01)^3$

Assignment 1 - Deep learning - Saarthak



① Output of last FC-layer = $\begin{bmatrix} 2.0 \\ 1.0 \\ 0.1 \end{bmatrix}$

input $\rightarrow L_1 \rightarrow \dots \rightarrow$ last layer $\rightarrow [p_{0.0} \dots p_{0.0}] = 1.0$
 $\circ \rightarrow 2.0$

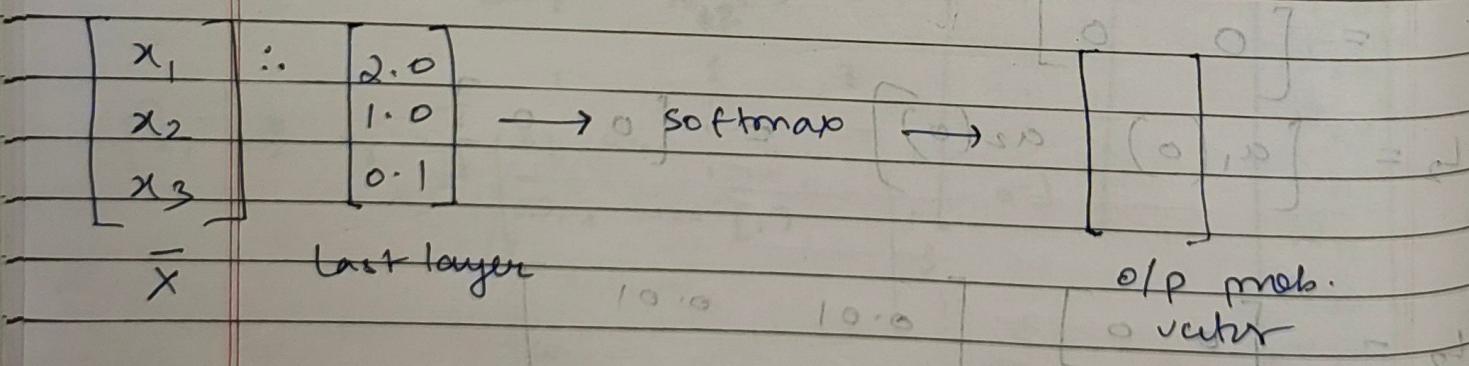
$$[p_{0.0} \dots p_{0.0}] \rightarrow 1.0 \rightarrow p_{0.0}$$

$$[(p_{0.0})^{10^3} (p_{0.0})^{10^3} \dots (p_{0.0})^{10^3}] = 96$$

To convert this net into predicting class prob. for 3 classes, we can use a softmax classifier.

Softmax classifier takes an input vector i.e. logits & returns an output vector of prob. for each class in the input vector.

Sum of this output vector = 1



$$\text{softmax} = P(y=i) = \frac{e(x_i)}{\sum e(x_j)}$$

prob. of i^{th} class
in vector (input)

$$\therefore P(y=x_1) = \frac{e(2.0)}{e(2.0) + e(1.0) + e(0.1)} = 0.658$$

$$P(y=x_2) = \frac{e(1.0)}{11.204} = \frac{2.71}{11.204} = 0.2418$$

Date _____
Page _____

$$P(y=x_3) = \frac{e(0.1)}{11.204} = \frac{1.105}{11.204} = 0.0986$$

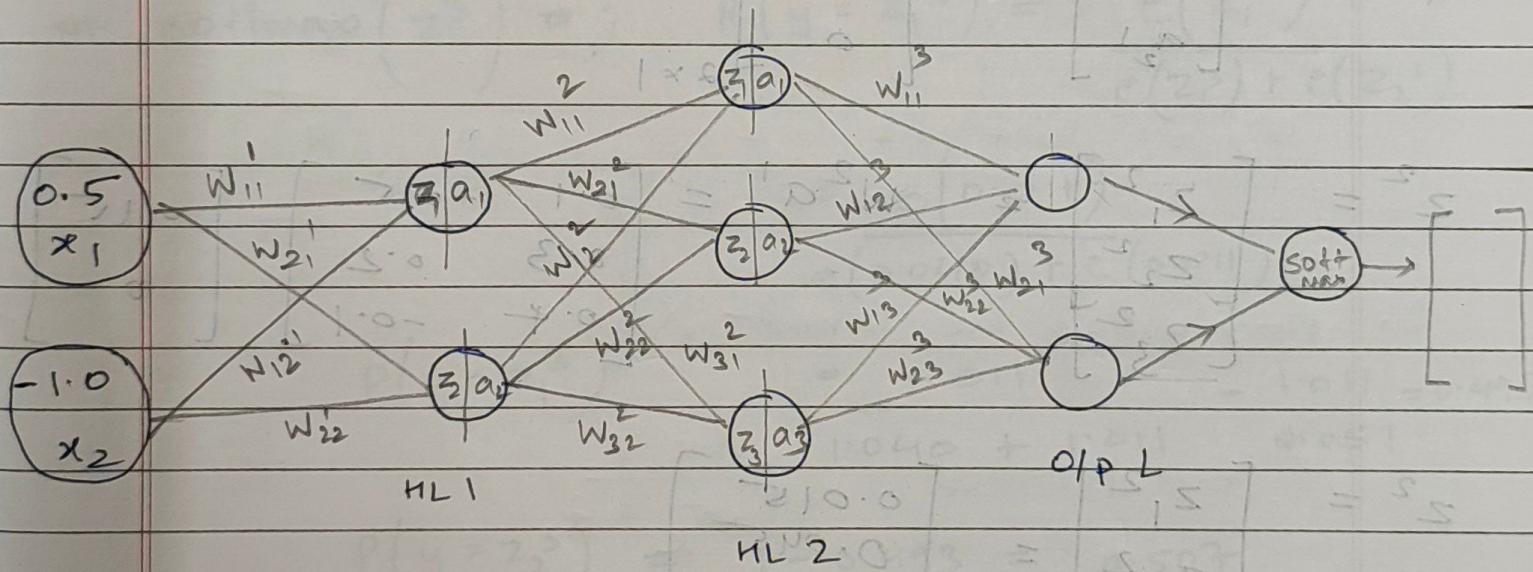
∴ Output vector i.e. probability vector for 3 classes is

$$\begin{bmatrix} 0.658 & 0.2418 & 0.0986 \end{bmatrix}$$

Q. Given: $x = [0.5 \ -1.0]$, 2 hidden layers with 2 & 3 neurons

O/P layer = 2 neurons

ReLU in hidden layers / softmax at O/P layer



$$\text{Input } x = \begin{bmatrix} 0.5 \\ -1.0 \end{bmatrix} \quad \text{weight } w^1 = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix}$$

Using uniform distribution for random initialization

$$\therefore \begin{bmatrix} 0.1 & -0.1 \\ 0.2 & 0.3 \end{bmatrix}$$

$$\text{Now, } z^1 = w^1 x = \begin{bmatrix} 0.1 & -0.1 \\ 0.2 & 0.3 \end{bmatrix} \begin{bmatrix} 0.5 \\ -1.0 \end{bmatrix} = \begin{bmatrix} 0.15 \\ -0.2 \end{bmatrix}$$

$$z^1 = \begin{bmatrix} z_1^1 \\ z_2^1 \end{bmatrix} = \begin{bmatrix} 0.15 \\ -0.2 \end{bmatrix}$$

Q Data
Page

$$a^1 = \text{ReLU}(z^1) = \begin{bmatrix} \text{ReLU}(z_1^1) \\ \text{ReLU}(z_2^1) \end{bmatrix} = \begin{bmatrix} 0.15 \\ 0 \end{bmatrix} = \begin{bmatrix} a_1^1 \\ a_2^1 \end{bmatrix}$$

layer 2 initialize weights

$$W^2 = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \\ w_{31} & w_{32} \end{bmatrix} = \begin{bmatrix} 0.1 & -0.5 \\ 0.3 & 0.2 \\ 0.4 & -0.1 \end{bmatrix}$$

$$a^1 = \begin{bmatrix} a_1^1 \\ a_2^1 \end{bmatrix} = \begin{bmatrix} 0.15 \\ 0 \end{bmatrix}$$

$$z^2 = \begin{bmatrix} z_1^2 \\ z_2^2 \\ z_3^2 \end{bmatrix} = W^2 a^1 = \begin{bmatrix} 0.1 & -0.5 \\ 0.3 & 0.2 \\ 0.4 & -0.1 \end{bmatrix} \begin{bmatrix} 0.15 \\ 0 \end{bmatrix}$$

$$z^2 = \begin{bmatrix} z_1^2 \\ z_2^2 \\ z_3^2 \end{bmatrix} = \begin{bmatrix} 0.015 \\ 0.045 \\ 0.006 \end{bmatrix}$$

$$a^2 = \text{ReLU}(z^2) = \begin{bmatrix} 0.015 \\ 0.045 \\ 0.006 \end{bmatrix}$$

output layer

Date _____
Page _____

$$w^3 = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix} \quad o_1 = \begin{bmatrix} 0.1 \\ 0.2 \\ 0.1 \end{bmatrix} \quad \begin{bmatrix} 0.1 & 0.2 & -0.1 \\ 0.4 & 0.7 & 0.3 \end{bmatrix}$$

$$z^3 = \begin{bmatrix} z_1^3 \\ z_2^3 \end{bmatrix} = w^3 \cdot a^2 = \begin{bmatrix} 0.1 & 0.2 & 0.1 \\ 0.4 & 0.7 & 0.3 \end{bmatrix} \begin{bmatrix} 0.015 \\ 0.045 \\ 0.005 \end{bmatrix}$$

$$z^3 = \begin{bmatrix} 0.011 \\ 0.040 \end{bmatrix}$$

$$\# \text{ softmax}(z^3) \Rightarrow P(y = z_1^3) = \frac{e^{z_1^3}}{e^{z_1^3} + e^{z_2^3}}$$

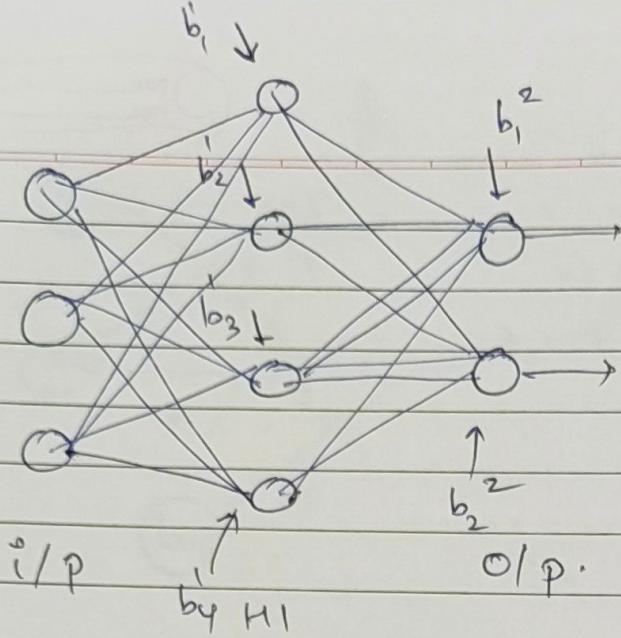
$$\textcircled{8} \quad \frac{\mu - x}{2} = \frac{0.1 + 0.2 + 0.1}{2} = 0.2 \\ = \frac{e^{0.011}}{e^{0.040} + e^{0.011}}$$

$$P(y = z_1^3) = \frac{0.011}{0.040 + 0.011} = \frac{0.011}{0.051} = 0.493$$

$$P(y = z_2^3) = 1 - 0.493 = 0.507$$

$$\therefore \text{output vector} = \begin{bmatrix} 0.493 \\ 0.507 \end{bmatrix}$$

(4)



Total no. of parameters.

$$\left\{ \begin{array}{l} W^1 = 4 \times 3 = 12 \\ b^1 = 4 \\ W^2 = 4 \times 2 = 8 \\ b^2 = 2 \end{array} \right.$$

$$\therefore \text{Total} = 26$$

$$(6) \quad x = [2, 4, 6, 8, 10] \quad \gamma = 2 \quad | \quad \beta = -1$$

Batch Normalization.

1. calculate mean & variance for sample set.

$$\mu = \frac{2+4+6+8+10}{5} = \frac{30}{5} = 6.$$

$$\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2}$$

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

$$\sigma = \sqrt{\frac{(2-6)^2 + (4-6)^2 + (6-6)^2 + (8-6)^2 + (10-6)^2}{5}}$$

$$\sigma = \sqrt{\frac{16+4+4+16}{5}} = \sqrt{8} = 2\sqrt{2} = 3.464$$

Date _____
Page _____

$$\therefore \bar{x} = 6 \quad \sigma = 3.464.$$

$$\begin{array}{r} 11.73^2 \\ - 108 \\ \hline 3 \\ \times 6^4 \\ \hline 3 \end{array}$$

2. Normalize the sample values x

$$\bar{x} = \frac{x - \bar{x}}{\sigma} \quad \therefore \tilde{x}_1 = \frac{x_1 - \bar{x}}{\sigma} = \frac{2 - 6}{3.464} \\ \tilde{x}_1 = -1.154$$

$$x_2 = \frac{4 - 6}{3.464} = \frac{-2}{3.464} = -0.577$$

$$x_3 = 0, x_4 = +0.577, x_5 = +1.154$$

$$\bar{x} = [-1.154, -0.577, 0, +0.577, +1.154]$$

Q: 8

3. Scale & shift the output.

$$y = \gamma \bar{x} + \beta$$

$$\gamma = 2$$

$$\beta = -1$$

$$y = 2\bar{x} - 1$$

$$y = 2[-1.154 \ -0.577 \ 0 \ +0.577 \ 1.154] - 1$$

$$= [-2.308 \ -1.154 \ 0 \ 1.154 \ 2.308] - 1$$

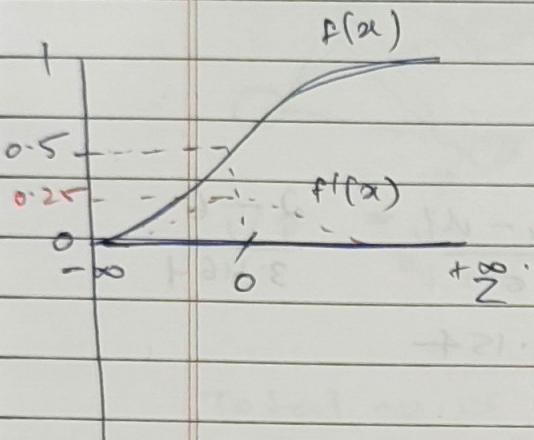
$$y = [-3.308 \ -2.154 \ -1 \ 0.154 \ 1.308]$$

⑤

Gradient Vanishing problem.

i) Sigmoid:

- output value $R \in [0, 1]$
- derivative $f'(x) = f(x)(1-f(x))$



- As the input becomes large $\rightarrow \infty$,
the output value saturates at 1,
→ As the outputs are saturated at 1,
, the derivative becomes 0
- During backprop this leads to
vanishing gradient problem.
- Gradient shrinks as the backprop
operation involves chain rule
multiplication
- Training stayrates, no learning
occurs.

Q.7.

Role of Dropout

Role of 'keep probability' 'p' in dropout.

- Dropout randomly removes / deactivates bunch of neurons during training to prevent overfitting.
- keep probability (p) is a hyperparameter which decides the number of neurons to be kept active during training.
- Dropout using keep-prob is a regularization technique
 - i) low ' p ' value signifies higher regularization & vice-versa.
- it causes randomization, so no memorization of the training data.

To keep the activation same during training & testing stage, we scale the activation by a factor of keep probability.

It is done in two ways \rightarrow 1)

1) standard dropout:

- In inference, there is no dropout occurring. \therefore , to ensure constant magnitude activation at test time, the activations are scaled \downarrow by multiplying with keep prob.
- This scales down the activation at testing

2) Inverted dropout

\hookrightarrow here, the activations are scaled at training time itself by dividing them with $(\frac{1}{p})$

\hookrightarrow no scaling at test time, as it's already accounted for at training time.

(8)

input image -

10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0

Filter -

1	0	-1
1	0	-1
1	0	-1

6×6

303

output

0	30	30	0
0	30	30	0
0	30	30	0
0	30	30	0

→ 4×4 , marks the vertical

edge between

10's region & 0's

region.

4×4 .

#

using the given convolution filter helps in detecting vertical edges present in the sample image between 10 & 0 pixel value.

- (9) The impact of increasing the no. of convolutional layers in a CNN include:-
- 1) vanishing gradient problem.
 ↗ As the number of layers increasing ↑'s q if the sigmoid is used as an activation function; then during the backpropagation due to chain multiplication of derivative the gradient gets pushed to very small values.
 ↗ This affects the training process.
 - 2) With deeper layers, the number of parameters ↑ significantly, leading to high computational cost & memory.
 - 3) With increasing no. of layers in the CNN, the original input image is lost due to series of conv. operation. At deeper layers, the net starts to learn from the noise & representation.
 - 4) Higher risk of overfitting on a small dataset.

(10)
(10)

Weight sharing

- same filter i.e. set of weights is applied across the entire input image.
- It reduces the no. of parameters as compared to the fully connected nets.
- makes the model translation - invariant
 ↗ feature can be detected anywhere in the sample input.

Sparsity of Connection

- Each neuron in the conv. layer is connected only to a small local region of the input (receptive field), not to every input pixel.
- leads to fewer computations & parameters, making CNNs efficient

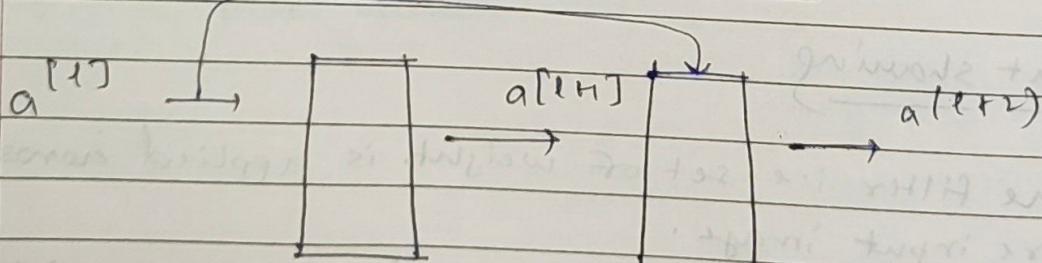
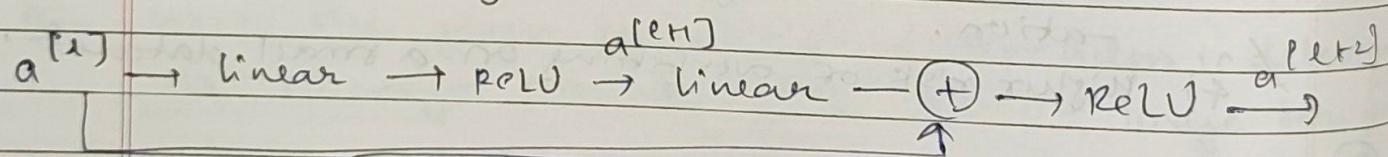
(11)

ResNet

- it introduces skip (residual) connections → instead of each layer learning $h(x)$, it learns a residual function $F(x) = h(x) - x$.
- output of residual block: $y = f(x) + x$.

$$h(x) = f(x) + x$$

letting me in part 2
 if the function is just the identity mapping
 i.e. ($h(x) = x$), the residual becomes $F(x) = 0$.
 which is easy to learn.



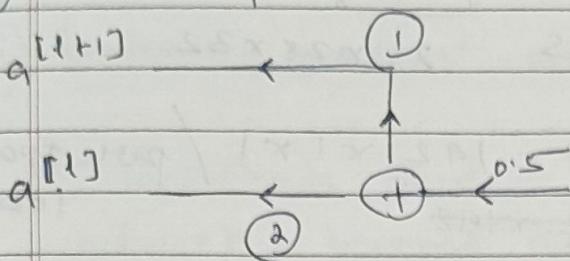
when $w^{(t)}$ becomes very small then

$$a^{(t+2)} = g[w^{(t+2)} a^{(t+1)} + a^{(t)}]$$

becomes small

↳ becomes an identity map, as it learns from the previous representation.

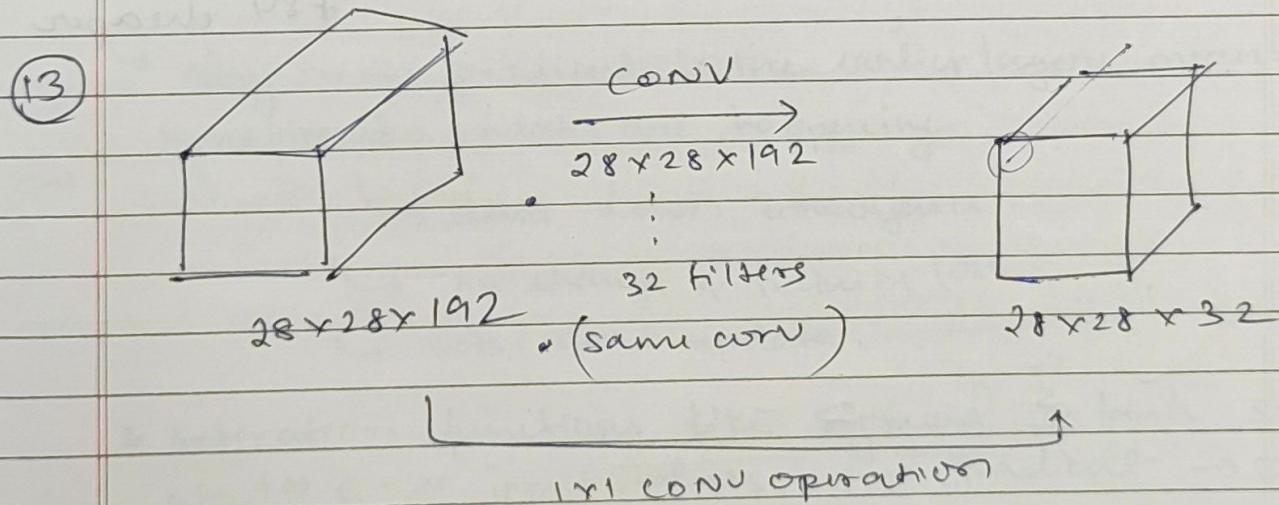
(12) How skip connection improve gradient flow in ResNet



- Skip connections creates shortcuts for gradient during backprop.
- In path ①, the standard flow or gradient from the i th layer to initial layers, the gradient might diminish when it reaches initial layers.
- In path ②, the skip connection allows the flow of gradient directly to the initial layers.

$$\therefore \frac{\partial L}{\partial w} = \frac{\partial L}{\partial w} + \frac{\partial L}{\partial w}$$

diminishes



$$\text{Size of filter} = 28 \times 28 \times 192 = 150,528 \text{ wts.}$$

$$\text{Total points} = 28 \times 28 = 784.$$

$$1 \text{ filter} = 150,528 \times 784$$

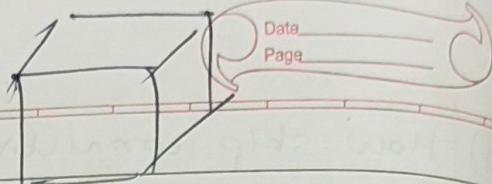
$$\therefore 32 \text{ filters} = 32 \times 150,528 \times 784$$

$\approx 3.78 \times 10^9$ multiplications.



$28 \times 28 \times 192$

$$\begin{array}{c} 1 \times 1 \text{ CONV} \\ \hline 1 \times 192 \\ \hline 32 \text{ filters} \end{array}$$



$28 \times 28 \times 32$

no. of weights in 1 filter = $192 \times 1 \times 1$ / per spatial location
~~32 filters = 32×192~~

$$\therefore \text{Total locations} = 28 \times 28 = 784$$

$$\therefore \text{Total multiplication / filter} = 784 \times 192 \\ = 150,528$$

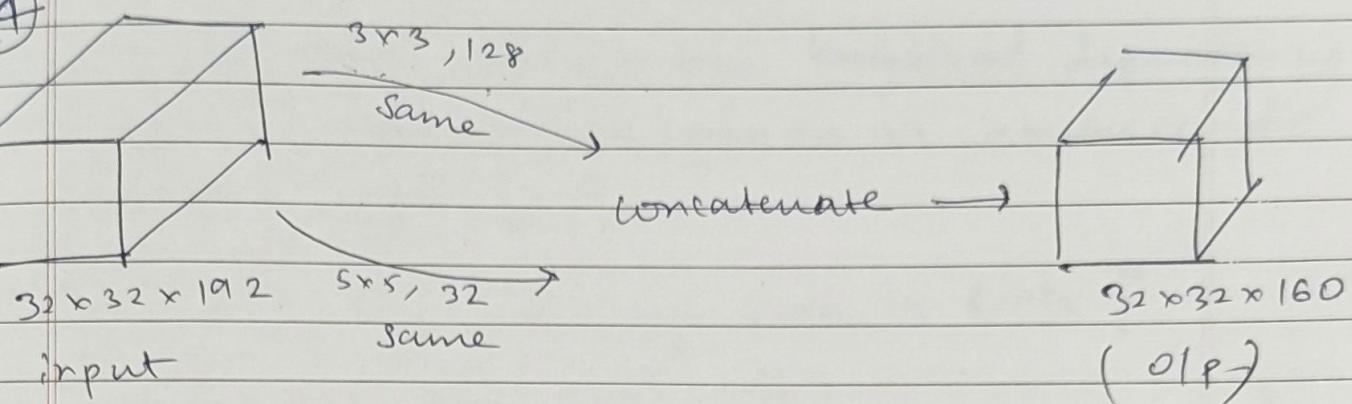
$$32 \text{ filters} = 32 \times 150,528 = 4,816,896.$$

$\approx 4.8 \text{ mil}$ parameter multiplication

$$\# \text{ Saving factor} = \frac{3.78 \times 10^9}{4.82 \times 10^6} \approx 784$$

$1 \times 1 \text{ conv}'s$
 784 cheaper

(14)



a) input = $32 \times 32 \times 192$

output = $32 \times 32 \times 160$

128 filters of 3×3

32 filters of 5×5

Multiplications for 3×3 CONV filters.

Filter size = $3 \times 3 \times 192 = 1728$ mults / position.

Spatial positions = $32 \times 32 = 1024$

Filters = 128.

Total = $128 \times 1024 \times 1728 = 226,492,416$.

multiplications for 5×5

Total = 157,286,400.

Total of 3×3 & 5×5 = 383,778,816.

b)

