

WORKSHEET SOLUTIONS

① Let us say the output of a last FC layer is $[2.0 \ 1.0 \ 0.1]$. If we need to convert this network into predicting class probabilities for 3 classes, how will you go about doing it and what are the probability values for 3 classes?

(A) For output layer, we should Softmax function. Softmax converts raw scores (logits) into probabilities by exponentiating each score.

$$\text{softmax} = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}$$

$$S_1 = \frac{e^{x_1}}{e^{x_1} + e^{x_2} + e^{x_3}} \Rightarrow \frac{e^{2.0}}{(e^{2.0} + e^{1.0} + e^{0.1})} \Rightarrow \frac{e^{2.0}}{7.3 + 2.7 + 1.1} = \frac{e^{2.0}}{11.1} \Rightarrow \frac{e^{2.0}}{11.1} = 0.66$$

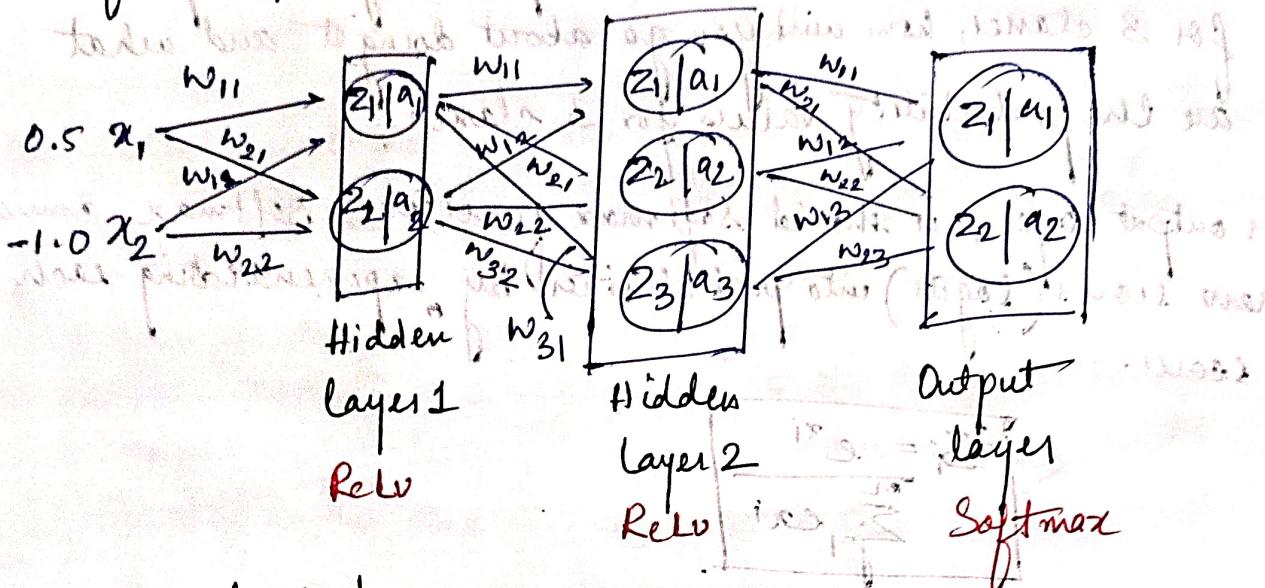
$$S_2 = \frac{e^{x_2}}{e^{x_1} + e^{x_2} + e^{x_3}} \Rightarrow \frac{e^{1.0}}{11.1} = 0.24$$

$$S_3 = \frac{e^{x_3}}{e^{x_1} + e^{x_2} + e^{x_3}} \Rightarrow \frac{e^{0.1}}{11.1} = 0.0099$$

$$S = \begin{bmatrix} S_1 \\ S_2 \\ S_3 \end{bmatrix} = \begin{bmatrix} 0.66 \\ 0.24 \\ 0.099 \end{bmatrix}$$

② Consider the following FFN. Input $x = [0.5 \ -1.0]$, two hidden layers with 2 and 3 neurons respectively. The output layer contains 2 neurons. ReLU non-linearity is applied for hidden layers and softmax operation is applied for the output layer activations. For simplicity assume that there are no bias terms connected to

neurons. Initialize the weight matrices randomly. Represent activations and weights as vectors and matrices and compute the forward pass output using a series of matrix multiplications.



Forward Feed

Layer 1

$$z_1 = w_{11}x_1 + w_{12}x_2 \rightarrow \text{ReLU}$$

$$= 0.4 \times (0.2 \times 0.5) + (1.0 \times -1.0) \rightarrow w_{11} = \begin{bmatrix} 0.2 & -0.4 \\ 1.0 & 0.1 \end{bmatrix}$$

$$= 0.1 + (-0.4) (-1)$$

$$= 0.1 - 0.4 - 1.0 = -0.1 - 1.0 \Rightarrow 0.2 - 0.9 \rightarrow \frac{0.2}{1.1} = \frac{-0.9}{1.1} = -0.818$$

$$z_2 = w_{21}x_1 + w_{22}x_2 \rightarrow \text{ReLU}$$

$$= (0.4 \times 0.5) + (0.1 \times -1.0) \rightarrow a_1 = [-0.9, -0.3]$$

$$= -0.2 + (-0.1)$$

$$\text{ReLU} \rightarrow -0.3, [0.1 - 2.6] = 8 \text{ (logistic function of all values)}$$

applying ReLU: $\frac{x + |x|}{2} \rightarrow a_2 = \frac{-0.9 + |-0.9|}{2}$

$$a_2 = \frac{-0.9 + 0.9}{2} \rightarrow -0.9 + 0.9 \rightarrow 0$$

$$\Rightarrow \frac{-0.3 + 0.3}{2} \rightarrow 0 \rightarrow 0$$

of both neurons

Layer 2

$$\text{taking } w_{21} = w^{[2]} = \begin{bmatrix} 0.3 & 0.8 & 0.5 \\ 0.7 & -0.6 & 0.2 \end{bmatrix}$$

$$z_1 = 0 \quad z_2 = 0 \quad z_3 = 0$$

as a_1 and a_2 are zero

0 × multiplying with any weights will give zero.

$$\text{Hence } a_1 = 0, a_2 = 0, a_3 = 0 \quad \leftarrow \text{ReLU layer}$$

Layer 3

$$\text{taking } w^{[3]} = \begin{bmatrix} 0.4 & -0.2 \\ 0.3 & 0.9 \\ -0.5 & 0.7 \end{bmatrix}$$

$$z_1 = 0 \quad z_2 = 0 \quad \cancel{z_3 = 0}$$

→ similarly everything will be zero

Applying softmax function, output probability is derived

$$a_1 = \frac{e^0}{e^0 + e^0 + e^0} \Rightarrow \frac{1}{3} \Rightarrow \cancel{0.33} \quad \cancel{0.5}$$

$$a_2 = \cancel{0.33} \quad 0.5$$

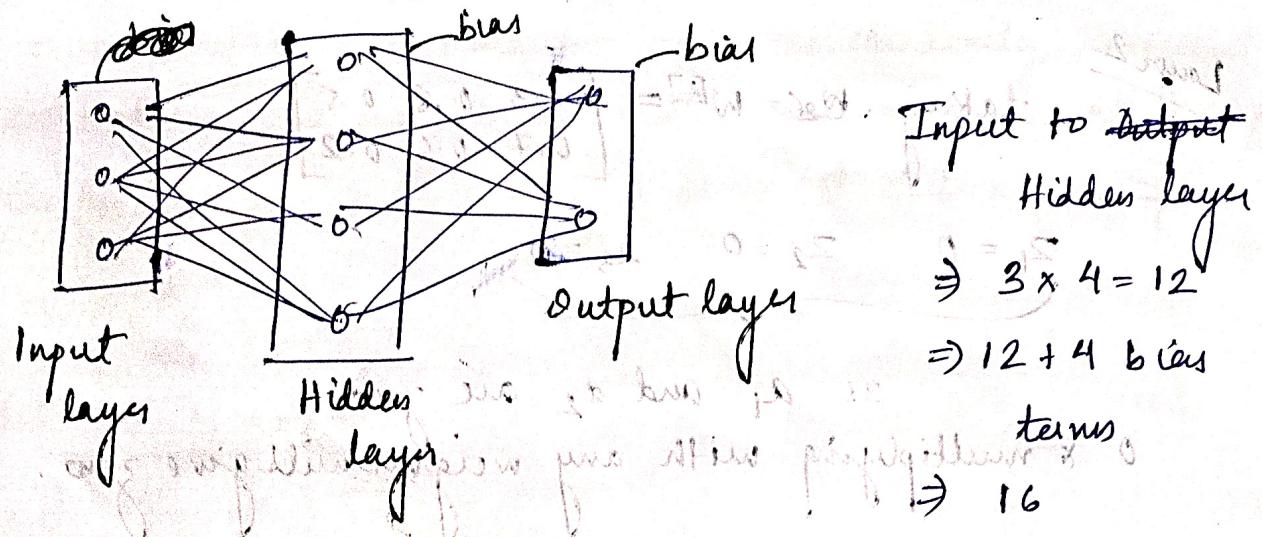
all work is zero

$$a_3 = \cancel{0.33}$$

selected axis is 0.5

$$\boxed{P = [0.5, 0.5]}$$

- ④ Consider a 2 layer fully connected feedforward network with 3 neurons in the input layer, 4 neurons in the hidden layer and 2 neurons in the output layer. Bias terms are connected to all neurons in the hidden and output layers. What will be the total no. of parameters in the network?

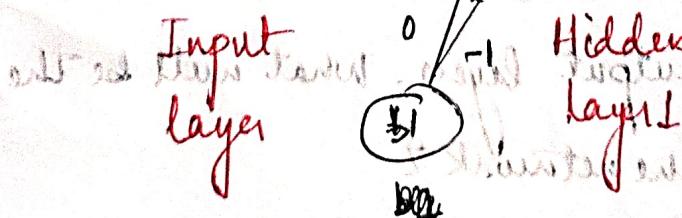
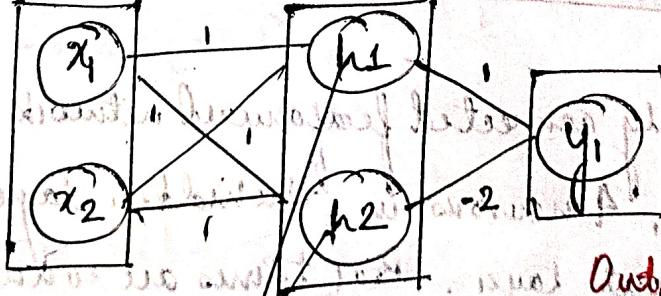


Hidden to Output $\rightarrow 4 \times 2 \rightarrow 8 + 2$ bias terms for Output layer

$$\rightarrow [10] = 10 \text{ parameters}$$

Total $\rightarrow 16 + 10 \rightarrow 26$ Parameters are there in the Network.

- ③ Construct a computational graph, compute the gradients and update the weights, for a single forward and backward pass. The activations are ReLU for hidden layers. The input X is shown in the Table below. Use all the four samples as a single batch.



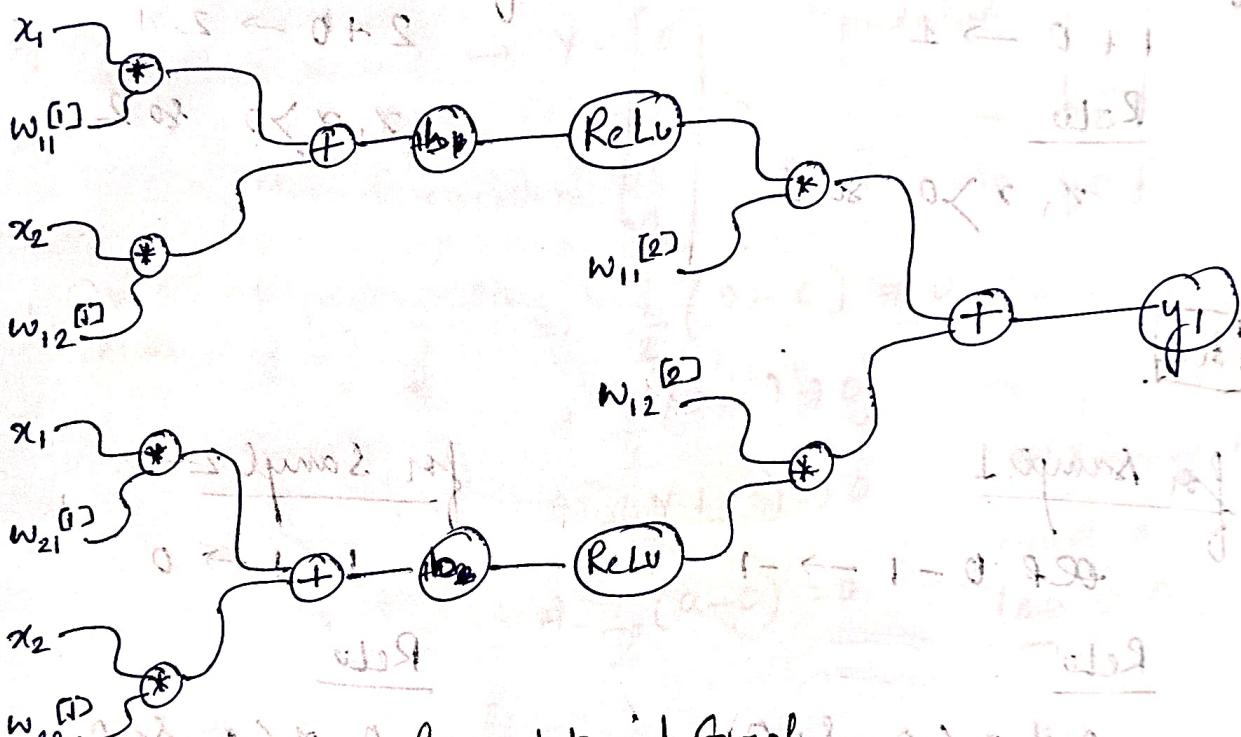
x_1	x_2	y
0	0	0
0	1	1
1	0	1
1	1	0

for the computational graph

~~Forward Pass~~

$$\text{Hidden Layer } W: \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$$

$$\text{Output layer } W = \begin{bmatrix} w_{11} & w_{12} \end{bmatrix} = \begin{bmatrix} 1 & -2 \end{bmatrix}$$



Computational Graph

Forward Pass

$$\begin{array}{l}
 \text{for } x_1=0 \quad h_1 = x_1 w_{11} + x_2 w_{12} \Rightarrow 0+0 \Rightarrow 0 \\
 \text{for } x_2=0, \quad x_1=0 \quad h_1 = x_1 w_{11} + x_2 w_{12} \Rightarrow 0+0 \Rightarrow 0 \\
 \text{for } x_1=0, \quad x_2=1 \quad h_1 = x_1 w_{11} + x_2 w_{12} \Rightarrow 0+1(1) \Rightarrow 1 \\
 \text{for } x_1=1, \quad x_2=0 \quad h_1 = x_1 w_{11} + x_2 w_{12} \Rightarrow 1+0 \Rightarrow 1 \\
 \text{for } x_1=1, \quad x_2=1 \quad h_1 = x_1 w_{11} + x_2 w_{12} \Rightarrow 1+1 \Rightarrow 2
 \end{array}$$

h₁

for sample 1

$$0 + \text{bias} \rightarrow 0$$

(Input weight)

ReLU

$$0 \text{ if } x \leq 0 \text{ so } 0$$

for sample 2

$$0 + 1 + 0 \rightarrow 1$$

ReLU

$$x, x > 0 \text{ so } 1$$

for sample 3

$$1 + 0 \rightarrow 1$$

ReLU

$$x, x > 0 \text{ so } 1$$

for sample 4

$$2 + 0 \rightarrow 2$$

h₂

for sample 1

$$0 + 0 - 1 \rightarrow -1$$

ReLU

$$0 \text{ if } x \leq 0 \text{ so } 0 \text{ if } x \leq 0 \text{ so } 0$$

for sample 2

$$1 - 1 \rightarrow 0$$

ReLU

$$0, x \leq 0 \text{ so } 0$$

for sample 3

$$0 + 1 - 1 \rightarrow 0$$

ReLU

$$0, x \leq 0 \text{ so } 0$$

for sample 4

$$2 - 1 \rightarrow 1$$

Hence $H = \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}$

(After applying
ReLU)

Input layer has 3 neurons. $(0, 1, 2 \rightarrow 0, 1, 2) = 3$ $\boxed{3}$

for Y: bias $b = 0$ and $[1 - 2]^T$ for output layer

$$H \times W \Rightarrow \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} 1 & -2 \end{bmatrix}^T$$

$$\Rightarrow \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ -2 \end{bmatrix} \Rightarrow \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix} \leftarrow \text{for Output layer.}$$

Ground Truth values $\rightarrow Y = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$

$$Y_{\text{predicted}} = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

$$MSE = \frac{1}{2}(y - y)^2 \Rightarrow \frac{1}{2}(0 - 0)^2 = 0$$

$$\Rightarrow \frac{1}{2}(1 - 1)^2 = 0$$

$$\Rightarrow \frac{1}{2}(1 - 1)^2 = 0$$

$$\Rightarrow \frac{1}{2}(0 - 0)^2 = 0$$

as $y_{\text{predicted}}$ and ground truth values are same
then $= 0$.

Then $= 0$.

with respect to output layer units

$$\frac{dl}{dw} = H^T \times \frac{dL}{dy_{\text{true}}}$$

$$= H^T \times 0 \Rightarrow \begin{bmatrix} 0 & 1 & 2 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\frac{dl}{dh} = \frac{dl}{dy_{\text{true}}} \times w^T = 0 \times [1 \ -2] = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 1 \\ -2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

and wrt hidden activations

⑥ If $x = [2 4 6 8 10]$, gamma = 2 and beta = 1 perform batch norm of all these 4 samples and write the output.

$$x = [2 4 6 8 10]$$

$$\bar{y} = \frac{2+4+6+8+10}{5} \Rightarrow 30/5 \Rightarrow 6$$

$$\sigma^2 = \frac{1}{5} \sum_{i=1}^5 (x_i - \bar{y})^2$$

$$(2-6)^2 = (-4)^2 = 16$$

$$(4-6)^2 = (-2)^2 = 4$$

$$(6-6)^2 = 0$$

$$(8-6)^2 = (2)^2 = 4$$

$$(10-6)^2 = (4)^2 = 16$$

$$\text{Adding them } 16 + 4 + 0 + 4 + 16 = 40$$

~~$\sigma^2 = \frac{1}{5} \sum_{i=1}^5 (x_i - \bar{y})^2$~~

$$\sigma^2 = \frac{1}{5} \sum_{i=1}^5 (x_i - \bar{y})^2 = 0$$

$$\sigma^2 = \frac{1}{5} \sum_{i=1}^5 (x_i - \bar{y})^2 = \frac{40}{5} = 8$$

$$\hat{z}_1 = \frac{x_1 - \bar{y}}{\sigma} = \frac{2-6}{\sqrt{8}} = \frac{-4}{\sqrt{8}} = -0.5$$

$$\hat{z}_2 = \frac{4-6}{\sqrt{8}} = \frac{-2}{\sqrt{8}} = -0.25$$

$$\hat{z}_3 = \frac{6-6}{\sqrt{8}} = \frac{0}{\sqrt{8}} = 0$$

$$\hat{z}_4 = \frac{8-6}{\sqrt{8}} = \frac{2}{\sqrt{8}} = 0.25$$

$$\hat{z}_5 = \frac{10-6}{\sqrt{8}} = \frac{4}{\sqrt{8}} = 0.5$$

calculating

2 score

for each

sample

$$\text{Bco } y = 2 \quad \beta = -1 \quad \boxed{\tilde{z} = y \tilde{z}' + \beta}$$

$$\tilde{z}_1 = -2(-0.5) + (-1) = 1 + 0 + 0 + 0 + 0 = 1 \quad (\text{Ans})$$

$$= -1 \Rightarrow -1 - 1 \Rightarrow -\frac{2}{2} = 0 + 0 + 0 - 0 + 0 \quad (\text{Ans})$$

$$\tilde{z}_2 = 2(-0.25) - 1 \Rightarrow -0.5 - 1 \Rightarrow -\frac{1.5}{2} = 0 + 0 + 0 - 0 + 0 \quad (\text{Ans})$$

$$\tilde{z}_3 = 2(0) - 1 \Rightarrow 0 - 1 \Rightarrow -1$$

$$\tilde{z}_4 = 2(0.25) - 1 \Rightarrow 0.5 - 1 \Rightarrow -0.5$$

$$\tilde{z}_5 = 2(0.5) - 1 \Rightarrow 1 - 1 \Rightarrow 0$$

8. For the following image and a filter, apply the convolution operation and write the output image. What did the convolution operation achieve in this example?

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

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

Input image

$$6 \times 6 = 36$$

Filter

$$3 \times 3 = 9$$

$$= \begin{array}{|c|c|c|c|} \hline 0 & 30 & 30 & 0 \\ \hline 0 & 30 & 30 & 0 \\ \hline 0 & 30 & 30 & 0 \\ \hline 0 & 30 & 30 & 0 \\ \hline \end{array}$$

$$\Rightarrow 10 + 0 - 10 + 10 + 0 - 10 + 10 + 0 - 10 \Rightarrow 0$$

$$\Rightarrow 10 + 0 - 0 + 10 + 0 - 0 + 10 + 0 - 0 \Rightarrow 30$$

$$\Rightarrow 10 + 0 - 0 + 10 + 0 - 0 + 10 + 0 - 0 \Rightarrow 30$$

$$\Rightarrow 0 + 0 - 0 + 0 + 0 - 0 + 0 + 0 - 0 \Rightarrow 0$$

$$\Rightarrow 10 + 0 - 10 + 10 + 0 - 10 + 10 + 0 - 10 \Rightarrow 0$$

$$\Rightarrow 10 + 0 - 0 + 10 + 0 - 0 + 10 + 0 - 0 \Rightarrow 30$$

$$\Rightarrow 10 + 0 - 0 + 10 + 0 - 0 + 10 + 0 - 0 \Rightarrow 30$$

$$\Rightarrow 0 + 0 - 0 + 0 + 0 - 0 + 0 + 0 - 0 \Rightarrow 0$$

$$\Rightarrow 10 + 0 - 10 + 10 + 0 - 10 + 10 + 0 - 10 \Rightarrow 0$$

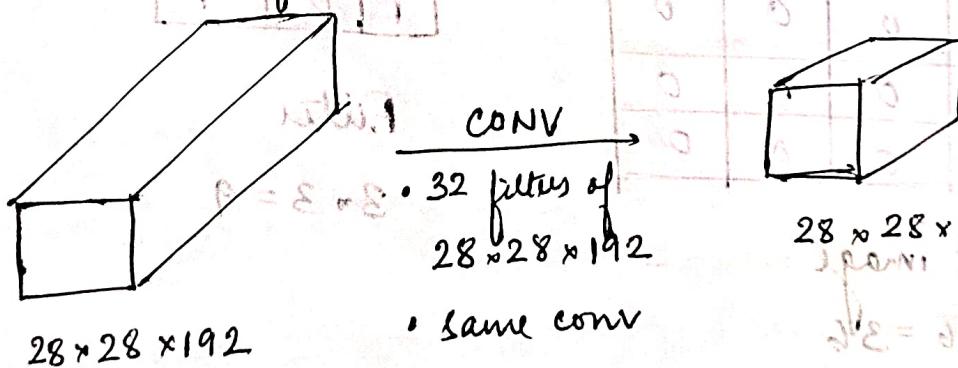
$$\Rightarrow 10 + 0 - 0 + 10 + 0 - 0 + 10 + 0 - 0 \Rightarrow 30$$

$$\Rightarrow 10 + 0 - 0 + 10 + 0 - 0 + 10 + 0 - 0 \Rightarrow 30$$

$$\Rightarrow 0 + 0 - 0 + 0 + 0 - 0 + 0 + 0 - 0 \Rightarrow 0.$$

(13)

For the following CONV operation, compute the computational cost (number of multiplication for example)



Design a new network architecture using 1×1 conv layer to reduce the computational cost. Show the savings from your architecture. Please note that your new proposed architecture should take the same sized input image ($28 \times 28 \times 192$) and output image of the same size shown in the above diagram ($28 \times 28 \times 32$)

Input - $28 \times 28 \times 192$

Filter - 32 each of size of $28 \times 28 \times 192$

Output - $28 \times 28 \times 32$

Each filter is applied across all 192 channels.

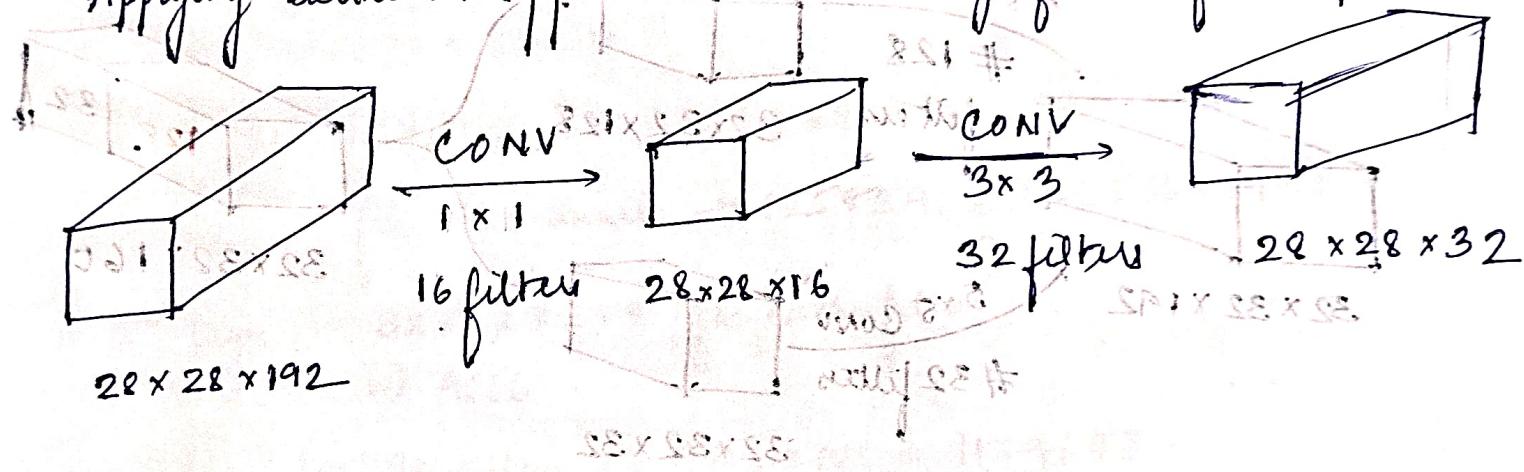
Each filter on one ^{channel} surface of the input image = $28 \times 28 \times 28 \times 28$

$$\text{For 32 filters} = 614656 \times 32 \Rightarrow 19668992 = 614656$$

$$\text{for 192 channels} = 19668992 \times 192 = 3776446464 \approx 3 \text{ billion}$$

Architecture with reduced multiplication:

Applying Bottleneck approach: 1×1 conv layer for the first layer.



Calculations:-

$$\text{FLOPs} = \frac{1}{2} \times (\text{input size})^2 \times \text{output size}^2 \times \text{number of filters}$$

$$\text{First layer} \rightarrow 1 \times 1 \text{ filter} \times 28 \times 28 \text{ input image} = 784$$

$$\text{for 192 channels} = 784 \times 192 = 150528$$

$$\text{for 16 filters} = 2408448$$

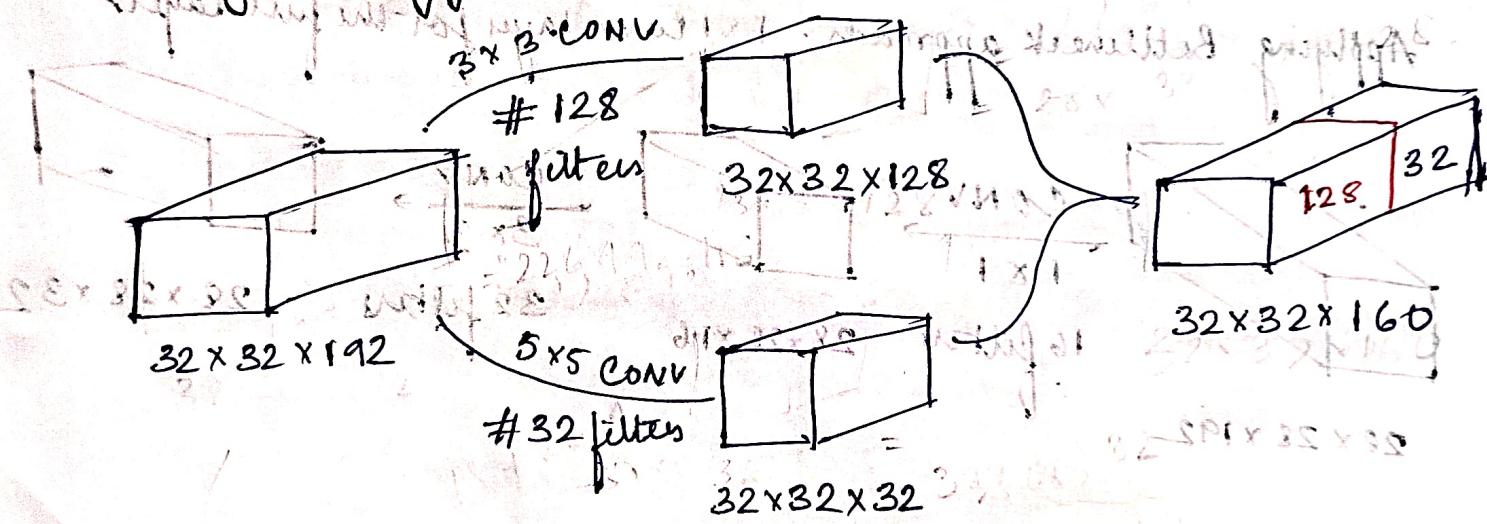
$$2^{\text{nd}} \text{ layer} \rightarrow 28 \times 28 \text{ input image} \times 3 \times 3 \text{ filter size} = 7056$$

$$\text{for 16 channels} = 112896$$

$$\text{for 32 filters} = 3612672$$

$$\text{Adding the outputs} = 2408448 + 3612672 \\ (2408448 + 3612672) = 6021120 \approx 6 \text{ million}$$

- (14) I would like to construct an inception module with $128 \times 3 \times 3$ conv filters, $32 \times 5 \times 5$ conv filter, same convolution stride = 1. The input to the inception module (coming from previous layer output activation) size is $32 \times 32 \times 192$ and the output image is $32 \times 32 \times 160$.
- What will be the number of multiplications in the above inception module?
 - Provide a strategy to reduce the number of multiplications and write down the no. of multiplications after adopting your strategy.



First branch \rightarrow filter dimensions $= 32 \times 32 \times 3 \times 3$
 $= 9216$

$$\text{for } 128 \text{ filters} = 9216 \times 128$$

$$= 1179648$$

$$\text{for } 192 \text{ channels} = 1179648 \times 192$$

$$= 226,492,416$$

Second branch \rightarrow filter dimensions $= 32 \times 32 \times 5 \times 5$
 $= 25,600$

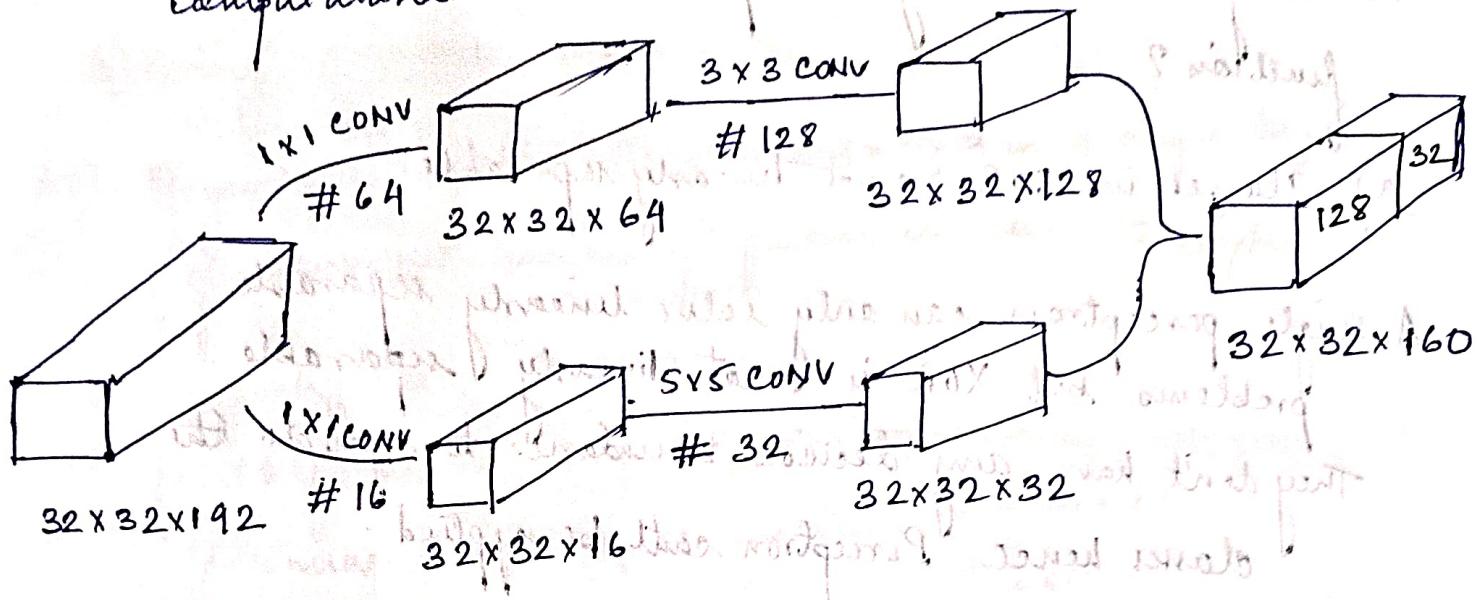
$$\text{for } 32 \text{ filters} = 25600 \times 32 = 819200$$

$$\text{for } 192 \text{ channels} = 157,286,400$$

$$\text{Concatenating the two} = 226,492,416 + 157,286,400$$

$$= 383,778,816 \quad (\underline{\text{Am}})$$

Constructing a structure to reduce the multiplications and computational cost.



Calculations :- 1st branch 1st layer :-

$$\text{each filter} \rightarrow 32 \times 32 \times 1 \times 1 = 1024$$

$$\text{for 64 filters} \rightarrow 1024 \times 64 \Rightarrow 65,536$$

$$\text{for 192 channels} \rightarrow 12582912 \quad (65536 \times 192)$$

2nd layer :- $32 \times 32 \times 3 \times 3$ (for each filter)
= 9216

$$\text{for 128 filters} \rightarrow 9216 \times 128 \Rightarrow 1179648$$

$$\text{for 64 channels} \rightarrow \begin{array}{r} 1179648 \\ \times 64 \\ \hline \end{array} = \underline{\underline{75497472}}$$

2nd branch

$$\text{1st layer :- } 32 \times 32 \times 1 \times 1 \Rightarrow 1024$$

$$(1024 \times 16) \rightarrow \text{for 16 filters} \Rightarrow 16384 \rightarrow 16384 \times 192 \text{ channels} \Rightarrow 3145728$$

$$\text{2nd layer :- } 5 \times 5 \times 32 \times 32 \Rightarrow 25600$$

$$\text{for 32 filters} \Rightarrow 25600 \times 32 \Rightarrow 819200$$

$$\text{for 16 channels} \Rightarrow 819200 \times 16 \Rightarrow 13107200$$

$$\text{Concatenate} \Rightarrow 75497472 + 13107200 \Rightarrow \underline{\underline{88,604,672}} \text{ (fm)}$$

WORKSHEET CONTINUED:-

- Q7) What is the role of Keep Probabilty in Dropout. the context of dropout & and why scaling the activations in "keep probability" is necessary?

Ans:- Keep Probabilities are used in "Dropouts" which is a type of regularization technique where we randomly drop neurons from a layer.

The fraction of neurons kept active \leftarrow Keep probability (p)
for example :- if $p=0.8$ then for each sample 80% of the neurons are kept and the rest 20% are dropped.

Draw the example from B2.

- Q8) What is the impact of increasing the number of convolutional layers in a CNN?

Ans:- Advantages:

- ① Hierarchical Feature Learning
 - Shallow CNN - detects very simple features (edges, corners, color blobs)
 - Deep CNN - detects semantic structures
 - Very Deep CNN (eg ResNet)
 - can learn highly abstract task specific features

- ② Best Accuracy - More the layers better the performance of the model as it will capture more details.

More CNN layers increase the model's accuracy capacity to fit complex data distributions, potentially leading to higher accuracy.

Disadvantages

- ① Vanishing/Exploding Gradients - ~~Q.~~ With many layers, gradients ~~we~~ might shrink/explode while backpropagating. The gradients will eventually stop updating and the initial layers ~~won't~~ will stop learning.
(For this problem we use ResNet, skip connections)
- ② Overfitting - If the dataset is small, adding many layers makes the model learn only the ~~so~~ training data causing overfitting (can also happen if ~~so~~ proper regularization is not done).
- 15) You have to develop a 5-way classifier to detect objects in an image. You are given a gray scale input image of size 32×32 . Which of the following approaches is better and why?
- Flatten the image (32×32) to form a 1024 vector and pass this to a fully connected layer with 5 neurons.
 - Directly pass the image to a CNN layer having five 32×32 filters.
1. Fully Connected Layer (Flattening)
The 32×32 grayscale image is flattened into a 1024-dimensional vector and passed to a ~~one~~ layer with 5 neurons.
This treats all pixels equally without exploiting their spatial arrangement.
Results in large no. of parameters and higher risk of overfitting - (especially with small datasets or images).

2. CNN layers

- CNNs use parameter sharing and local connectivity, allowing filters to learn local patterns (edges, textures, shapes) and build hierarchical feature representations.
- This gives better performance on images and fewer parameters compared to a fully connected layer for the same problem.

CNN is better:-

Precise spatial information

Parameter efficiency

Feature hierarchy