

# WORKSHEET-1

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

$$m = 6 \times 6$$

$$\begin{matrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{matrix} \rightarrow ? \quad (4 \times 4)$$

$f = 3 \times 3$

Assume  $s = 1$  (stride),  $p = 0$  (padding)

$$\begin{aligned} \text{O/P size} &= \left( \frac{(n-f+2p)}{s} \right) + 1 \\ &= \left[ \frac{(6-3+2 \times 0)}{1} \right] + 1 \\ &= 4 \times 4 \end{aligned}$$

Convolution operations —

$$1) \quad 10 \times 1 + 10 \times 6 + 10 \times 1 + 10 \times 1 + \\ 10 \times 6 + 10 \times 1 + 10 \times 1 + 10 \times 0 + \\ 10 \times 1 = 0$$

10x1	10x0	10x-1
10x1	10x0	10x-1
10x1	10x0	10x-1

$$2) \quad \begin{matrix} 10x1 & 10x0 & 10x-1 \\ 10x1 & 10x0 & 10x-1 \\ 10x1 & 10x0 & 10x-1 \end{matrix} \Rightarrow 10 + 0 + 0 + 10 + 0 + 0 + 10 + 0 + 0 \\ \Rightarrow 30$$

$$3) \quad \begin{matrix} 10x1 & 0x0 & 0x-1 \\ 10x1 & 0x0 & 0x-1 \\ 10x1 & 0x0 & 0x-1 \end{matrix} = 10 + 0 + 0 + 10 + 0 + 0 + 10 + 0 + 0 = 30$$

$$4) \quad \begin{matrix} 0x1 & 0x0 & 0x-1 \\ 0x1 & 0x0 & 0x-1 \\ 0x1 & 0x0 & 0x-1 \end{matrix} = 0$$

$$5) \quad \begin{matrix} 10x1 & 10x0 & 10x-1 \\ 10x1 & 10x0 & 10x-1 \\ 10x1 & 10x0 & 10x-1 \end{matrix} = 0$$

6)  $\begin{bmatrix} 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix} = 30$

7)  $\begin{bmatrix} 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \end{bmatrix} = 30$

8)  $\begin{bmatrix} 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix} = 0$

$\Delta_{\text{O/P}}$  for O/P =

0	30	30	0
0	30	30	0
0	30	20	0
0	30	30	0

9)  $\begin{bmatrix} 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix} = 0$

10)  $\begin{bmatrix} 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} = 30$

11)  $\begin{bmatrix} 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \end{bmatrix} = 30$

12)  $\begin{bmatrix} 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix} = 0$

O/P =  $\begin{bmatrix} 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \end{bmatrix}$

13)  $\begin{bmatrix} 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix} = 0$

14)  $\begin{bmatrix} 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \end{bmatrix} = 30$

15)  $\begin{bmatrix} 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \end{bmatrix} = 30$

16)  $\begin{bmatrix} 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \end{bmatrix} = 0$

Final O/P —

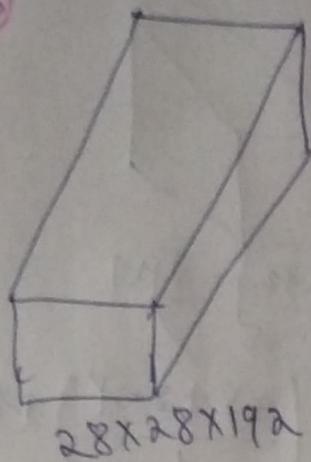
0	30	30	0
0	30	30	0
0	30	20	0
0	30	30	0

→ This filter is a vertical edge detector.

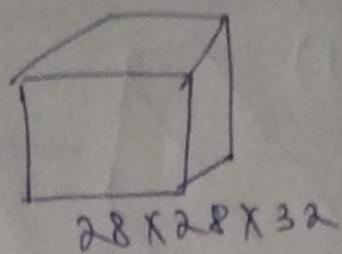
→ It detects vertical transitions from left to right.

→ Produces strong positive responses at patches that straddle the edge from 1 to 0 or 0 to 1 zero elsewhere.

(13)



$$\begin{array}{c} \xrightarrow{\text{CONV}} \\ f = 28 \times 28 \times 192 \\ \#f = 32 \\ \text{same conv} \\ (p=1) \\ s=1 \end{array}$$



$$O/P = \left( \frac{n+2p-f}{s} \right) + 1$$

$$O/P = \left( \frac{28+2-3}{1} \right) + 1 = 3 \times 3$$

which is not we're seeing in O/P dimension.  
 So, probably the given O/P dimensions are wrong.

Assume  $f = 3 \times 3 \times 192$ ;

$$O/P = \left( \frac{n+2p-f}{s} \right) + 1 = \left( \frac{28+2-3}{1} \right) + 1 = 28 \times 28 \times 32$$

Computational cost / no. of multiplications

$$O/P = 28 \times 28 \times 32 \xrightarrow{32 \text{ channels}}$$

Multiplications per O/P channel =  $3 \times 3 \times 192$

$$\text{So, total multiplications} = 28 \times 28 \times 32 \times 3 \times 3 \times 192$$

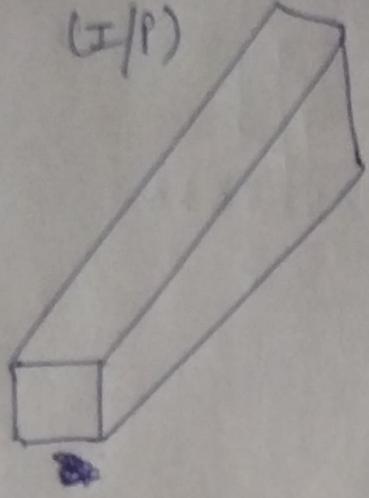
$$= 43,352,064$$

$1 \times 1$  conv —

↳ Reduces the computational complexity.

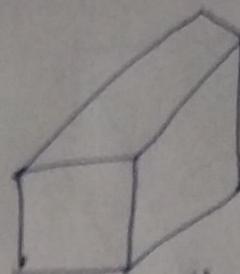
$28 \times 28 \times 192$

(I/P)



$\xrightarrow{1 \times 1 \text{ CONV}}$   
 $1 \times 1 \times 192$

$$\begin{aligned} p &= 1 \\ s &= 1 \\ 16 \text{ filters} \end{aligned}$$



$$\frac{m+2p-f+1}{s}$$

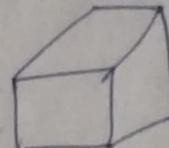
$$\Rightarrow \left( \frac{28+2-1}{1} \right) + 1$$

$$\Rightarrow 30 \times 30 \times 16$$

$$\begin{aligned} \text{Cost} &= 30 \times 30 \times 16 \times 1 \times 1 \times 192 \\ &= 2,764,800 \end{aligned}$$

$\downarrow \text{CONV}$   
 $5 \times 5 \times 16$

$$\begin{aligned} p &= 1 \\ s &= 1 \\ 32 \text{ filters} \end{aligned}$$



$$\begin{aligned} m+2p-f+1 \\ \Rightarrow 28+2-5+1 \\ \Rightarrow 28 \times 28 \end{aligned}$$

$$28 \times 28 \times 32 \text{ (O/P)}$$

~~$\text{Cost} = 28 \times 28 \times 32 \times$~~   
 ~~$30 \times 30 \times 16$~~   
 $5 \times 5 \times 16$   
 $= 10,035,200$

$$\begin{aligned} \text{Total cost} &= 2,764,800 + 10,035,200 \\ &= 12,800,000 \end{aligned}$$

cost was  
earlier  $143,352,064$   $\rightarrow$  reduced to  $12,800,000$ .  
using  $1 \times 1$  CONV.

① last  
LFC layer =  $[2.0 \quad 1.0 \quad 0.1]$   
Prob. for 3 classes?

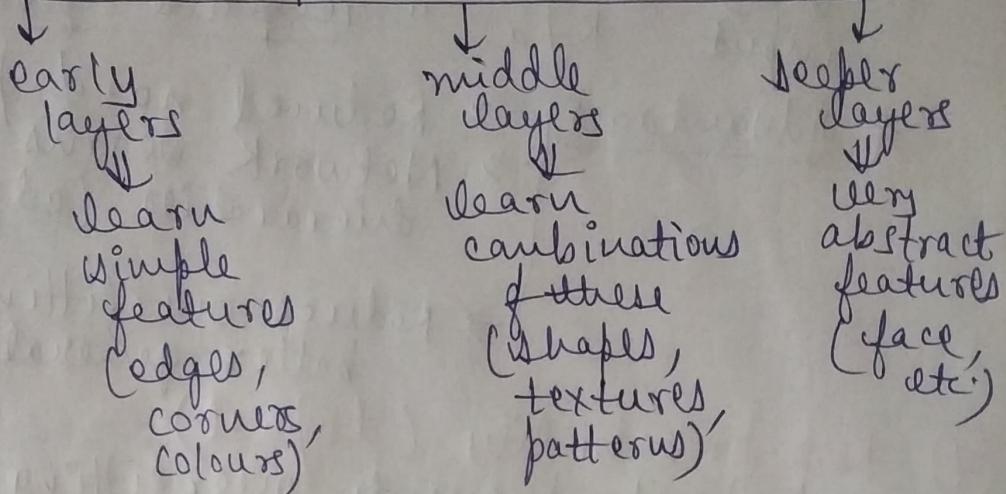
Apply softmax function  $\Rightarrow p_i = \frac{e^{x_i}}{\sum_{j=1}^3 e^{x_j}}$

$$\begin{aligned} \text{class-1} &\Rightarrow p_1 = \frac{e^{2.0}}{11.212} = \frac{7.389}{11.212} = 0.659 \\ \text{class-2} &\Rightarrow p_2 = \frac{e^{1.0}}{11.212} = 0.242 \quad ; \quad p_3 = \frac{e^{0.1}}{11.212} = 0.099 \end{aligned}$$

$$\begin{aligned} e^{2.0} &= 7.389 \\ e^{1.0} &= 2.718 \\ e^{0.1} &= 1.105 \end{aligned}$$

⑨ Impact of increasing the no. of convolutional layers in CNN?

↳ More hierarchical feature extraction



- More layers improve performance on complex tasks (like in case of ImageNet classification).
- Risk of overfitting due to too many parameters; can memorize rather than generalizing; so to avoid this, regularization techniques such as - dropout, etc. is needed.
- Vanishing or exploding gradients problem may arise; so to avoid this ReLU is used or batch normalization is done.
- High computation and memory cost.

⑩ (a) Weight sharing in CNN —

- \* In CNN, the same set of weights (Kernel or filter) is applied across different spatial locations of the input.
- \* Advantages → huge reduction in no. of parameters.  
faster computation & easy to parallelize → Translation invariance (i.e. network learns features that are position independent).

## (b) Sparsity of connections —

- \* In CNN, each neuron (or filter) connects only to a small local region of the input K/a "receptive field".
- \* Advantages →
  - Reduced no. of parameters.
  - Network learns spatial hierarchies (edges → shapes objects)
  - Reduce overfitting better generalization.

## (7) Keep probability (dropout) —

- \* Dropout prevents overfitting by setting a fraction of neuron's activations to zero.
- \* Keep prob. ( $p$ ) is the probability that a neuron remains active (not dropped).  
$$p = 1 - \text{dropout rate}$$

e.g.: - DR = 0.2  
 $p = 0.8$  (i.e., each neuron has 80% chance to stay active).

- \* Let,  $h$  = vector of activations before dropout.  
 $r$  = random binary mask ( $1 \rightarrow p$ , and  $0 \rightarrow 1-p$ ).

Then, o/p after dropout —

$$h_{\text{drop}} = r \odot h$$

$$\mathbb{E}[h_{\text{drop}}] = \mathbb{E}[r \odot h] = p \cdot h$$

Means that on avg.  $o/p$  activations become smaller by a factor of  $p$  during training. So, when dropout is disabled at test time, activations are suddenly larger by about  $1/p$  times.

- \* Solution — Scale activations by  $1/p$  during training  $\Rightarrow h_{\text{drop}} = \frac{r \odot h}{p}$

$$\mathbb{E}[h_{\text{drop}}] = \mathbb{E}\left[\frac{r \odot h}{p}\right] = \frac{p \cdot h}{p} = h$$

*Note: o/p behaves consistent up to train. & at test.*

$\therefore$  Expected activation stays same as in non-dropout case.

(4)

I/P H.L - 1 O/P  
 3 neurons 4 neurons 2 neurons

Total no. of parameters = ?

Aus :- For layer-1 :-

$$\text{wts.} = 3 \times 4 = 12$$

$$\text{bias} = 4$$

$$\underline{\underline{16}}$$

For layer-2 :-

$$\text{wts.} = 4 \times 2 = 8$$

$$\text{bias} = \underline{\underline{2}} \\ \underline{\underline{10}}$$

Total =  $16 + 10 = 26$  parameters.

(6)

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

$$\text{gamma} = 2$$

$$\beta = +1$$

Batch normalization - O/P = ?

$$\text{Ans} :- \bar{x} = \frac{2+4+6+8+10}{5} = 6$$

$$\sigma^2 = \frac{1}{5} [(2-6)^2 + (4-6)^2 + (6-6)^2 + (8-6)^2 + (10-6)^2]$$

$$\sigma^2 = 8 ; \sigma = \sqrt{8} = 2.828$$

$$\hat{x}_i = \frac{x_i - \bar{x}}{\sqrt{\sigma^2 + (\theta)}} \rightarrow \text{ignore it} = \frac{x_i - \bar{x}}{\sigma}$$

$$\hat{x}_1 = \frac{2-6}{2.828} = -1.414$$

$$\hat{x}_2 = -0.707$$

$$\hat{x}_3 = 0$$

$$\hat{x}_4 = 0.707$$

$$\hat{x}_5 = 1.414$$

$$\tilde{x} = \alpha \hat{x} + \beta \quad \Rightarrow \text{Scale + shift}$$

$$\tilde{x}_1 = 2 \times (-1.414) + 1 = -1.828$$

$$\tilde{x}_2 = 2 \times (-0.707) + 1 = -0.414$$

$$\tilde{x}_3 = 2 \times (0) + 1 = 1$$

$$\tilde{x}_4 = 2 \times (0.707) + 1 = 2.414$$

$$\tilde{x}_5 = 2 \times (1.414) + 1 = 3.828$$

(14)

IP

32x32x192

3x3 CONV  
P=1  
S=1  
#128

$$\begin{aligned} & \left( \frac{n+2p-4}{s} \right) + 1 \\ & \Rightarrow \left( \frac{32+2-3}{1} \right) + 1 \\ & \Rightarrow 32 \times 32 \times 128 \end{aligned}$$

$$\begin{aligned} & 5 \times 5 \text{ CONV} \\ & P=\frac{R-1}{2}=2 \quad \Rightarrow \quad \left( \frac{32+4-5}{1} \right) + 1 \\ & S=1 \\ & \#32 \quad \Rightarrow 32 \times 32 \times 32 \end{aligned}$$

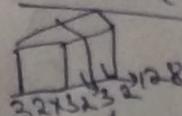
(a) multiplications

$$\Rightarrow (32 \times 32 \times 128 \times 3 \times 3 \times 192)$$

$$+ (5 \times 5 \times 192 \times 32 \times 32 \times 32)$$

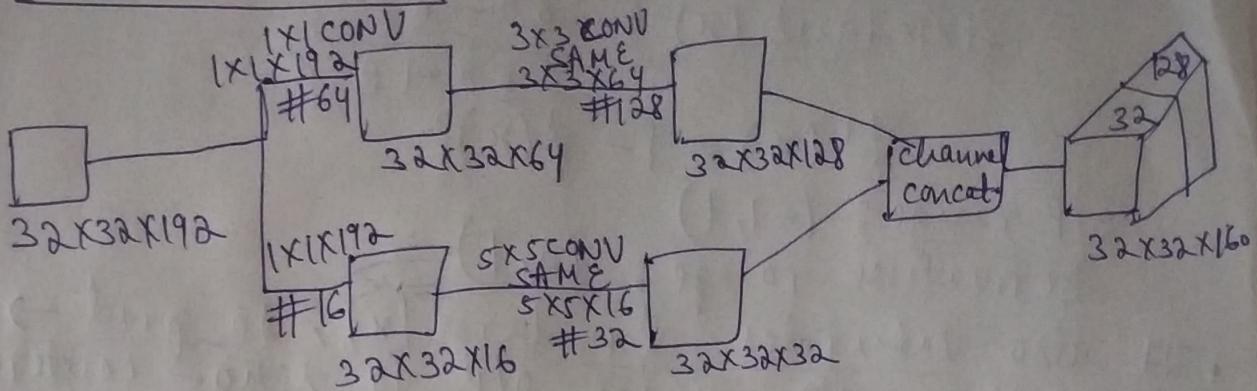
$$\Rightarrow 383,778,816 \approx 400 \text{ million}$$

channel concatenation



(b) 1x1 CONV to be applied to ↓ no. of multiplications.

## New architecture



3x3 branch :-

$$\begin{array}{r}
 32 \times 32 \times 64 \times 192 = 12582912 \\
 + \\
 32 \times 32 \times 128 \times 3 \times 3 \times 64 = 75497472 \\
 \hline
 88080384
 \end{array}$$

5x5 branch :-

$$\begin{array}{r}
 32 \times 32 \times 16 \times 192 = 3145728 \\
 + \\
 32 \times 32 \times 32 \times 5 \times 5 \times 16 = 13107200 \\
 \hline
 16252928
 \end{array}$$

$$\text{Total} = 88080384 + 16252928 = 104,333,312 \\
 \approx 100 \text{ million}$$

Hence, cost reduced with new architecture of 1x1 conv, thus reduces the computational cost to 4 folds.

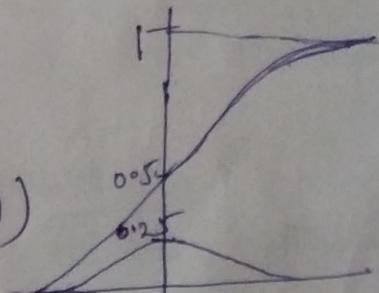
(5) Sigmoid func -

$$f(x) = \frac{1}{1+e^{-x}}$$

$$f'(x) = f(x)(1-f(x))$$

off range: (0, 1)

derivative range: (0, 0.25)



- For large +ve or -ve inputs, Sigmoid saturates near 1 or 0. In those regions,  $f'(x) \approx 0$ .
- During backprop, multiplying by many small derivatives ( $\approx 0.1$  or less), leads to vanishing grad.

## Tanh function —

$$f(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$f'(x) = 1 - f(x)^2$$

O/P range:  $(-1, 1)$

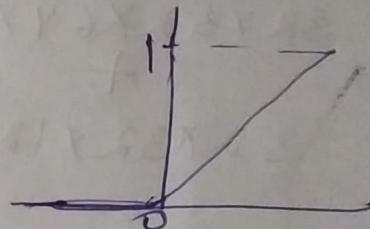
Derivative range:  $(0, 1)$

- It also saturates for large  $|x|$ , so grad.  $\rightarrow 0$ .
- However, its mean O/P is zero-centered, which helps slightly.

## ReLU func. —

$$f(x) = \max(0, x)$$

$$f'(x) = \begin{cases} 1, & x > 0 \\ 0, & x \leq 0 \end{cases}$$



O/P range:  $[0, \infty)$

Derivative range: 1 for +ve I/P, 0 for -ve I/P

- For +ve I/P, grad. = 1  $\rightarrow$  grad. neither vanish nor explode.
- For -ve I/P, grad. = 0  $\rightarrow$  neuron is inactive ("dead ReLU").
- When ~~sets~~ weights are initialized poorly, then explodes gradient.

(15)

Option - (a) is better

- ↳ Flattening destroys the 2-D spatial relationships in image.
- ↳ Nearby px are treated as unrelated features.
- ↳ weights =  $1024 \times 5 = 5120$  / or no feature extraction or hierarchical learning can't learn local patterns

Option - (b) — BETTER

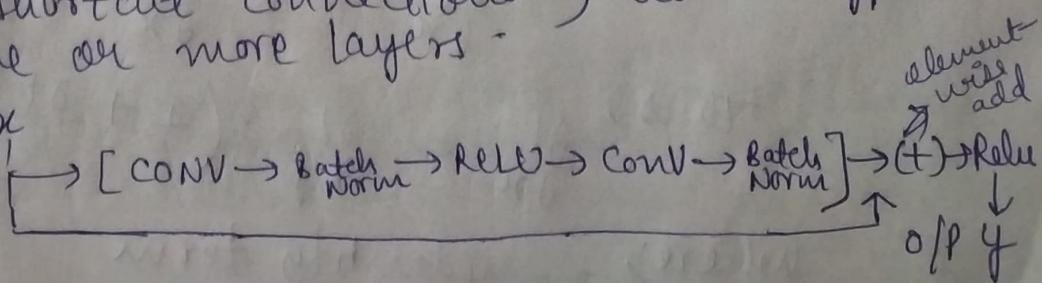
- ↳ captures local patterns
- ↳ local connectivity  $\downarrow$  es no. of parameters; and wt. sharing also. (Translation invariant pattern learning)
- ↳ Better scalability + generalization

(11)

ResNet —

- \* Earlier, researchers tried to make network deeper (20/50/100 layers) to get better accuracy, but it performed worse than shallower one. (Training error fed, vanishing gradient). So, then idea of "ResNet" came in; i.e., if some layers do nothing (act as identity function) then deeper network will act as shallow ones.
- \* It introduces "skip connections" (k/a "shortcut connections") that bypasses one or more layers.

I/P X



$$y = F(x, w_i) + x$$

where;

$F(x, w_i)$  = O/P of few layers (like two  
 CONV  $\rightarrow$  Batch Norm  $\rightarrow$  ReLU layers)  
 parameterized by  $w_i$ .

The skip connection directly adds  $x$  to  
 the output of  $F(x, w_i)$ .

\* To make residual block to behave like  
 an identity func., i.e;  $y = x$ ,

then,  $F(x, w_i) = 0$

$$y = 0 + x = \underline{\underline{x}}$$

Residual func. outputs zero  
 & entire block just pass  
 I/P through unchanged.

\* Thus, gradient can flow directly via skip  
 connections, avoiding vanishing gradients.

(Q) \* Problem :- without skip connections, backprop:-

$$\text{chain rule} : - \frac{\partial L}{\partial x_0} = \frac{\partial L}{\partial x_L} \cdot \frac{\partial x_L}{\partial x_{L-1}} \cdot \frac{\partial x_{L-1}}{\partial x_{L-2}} \cdots \frac{\partial x_1}{\partial x_0}$$

If each derivative has values smaller than 1,  
 their product gets smaller & smaller as we move backward (i.e; Vanishes gradient).

\* Soln. :- Gradient flow via a residual block

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial x}$$

Since,  $y = F(x, w) + x$ , so

$$\frac{\partial y}{\partial x} = \frac{\partial F(x, w)}{\partial x} + I$$

Thus,

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial y} \cdot \left( I + \frac{\partial F(x, w)}{\partial x} \right)$$

where,  $I$  = Identity term

Even if  $\frac{\partial F(n, w)}{\partial x}$  becomes very small, it ensures gradients flow directly backward via skip connection. Thus avoids the above issue of exploding grad.

## (16) CNN Visualization methods →

### 1) Filter (Kernel) Visualization

- \* Each filter in CNN learns to detect a particular feature or pattern (layer)
  - edges / colour gradients, deeper layers
  - more abstract features).
- \* Take learned filter weights (from conv layers) after training.

For the 1st conv layer, each filter directly operates on image pixels, so you can visualize it as a small image patch.

Can display these filters as grids of small images.

Thus, can detect if the network has learnt meaningful patterns or noise.

### 2) Activation (feature map) visualization

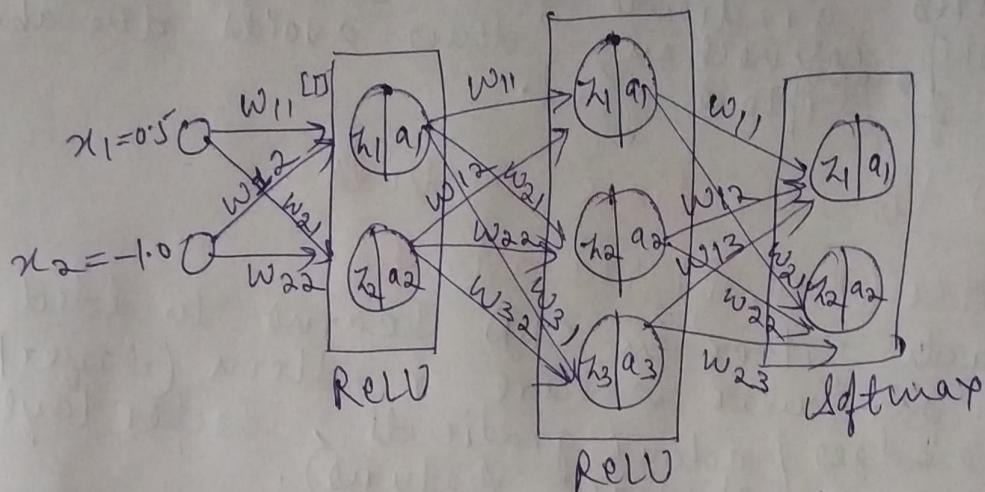
- \* Input an image to trained CNN.
    - Extract o/p of chosen layer (activations)
    - Visualise each as a heatmap.
- Deep layers act.  
Early layer act.      Mid-level act.      Whole face  
act.                    (edges)      (ears/legs)      body

### 3) Maximally activating image patches

- \* To look at which parts of real images cause strongest activations in neuron.
- \* Pass many images via CNN → Record the patches that display those patches. ← patches that maximise act.

(2)

Given :-  $X = \begin{bmatrix} 0.5 \\ -1.0 \end{bmatrix}$ , no bias term



Initialize weights randomly;

$$W_1 = \begin{bmatrix} w_{11} & w_{12} \\ 0.2 & -0.1 \\ 0.4 & 0.3 \end{bmatrix}_{2 \times 2}$$

$$W_2 = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ 0.2 & 0.3 & -0.1 \\ -0.3 & 0.4 & 0.2 \\ 0.2 & -0.2 & -0.1 \end{bmatrix}_{3 \times 2}$$

$$W_3 = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ 0.3 & -0.2 & 0.1 \\ -0.1 & 0.2 & 0.4 \end{bmatrix}_{2 \times 3}$$

Forward pass :-

layer-1 —

$$z^{[1]} = w_{11}x_1 + w_{12}x_2 \quad \text{and} \quad w_{21}^{[1]}x_1 + w_{22}^{[1]}x_2$$

$$z^{[1]} = \begin{bmatrix} 0.2 & -0.1 \\ 0.4 & 0.3 \end{bmatrix} \begin{bmatrix} 0.5 \\ -1.0 \end{bmatrix} = \begin{bmatrix} 0.2 \\ -0.1 \end{bmatrix} \rightarrow z_1^{[1]}$$

$$a^{[1]} = \text{ReLU}(z^{[1]}) = \begin{bmatrix} \max(0, 0.2) \\ \max(0, -0.1) \end{bmatrix}$$

$$a^{[1]} = \begin{bmatrix} 0.2 \\ 0.0 \end{bmatrix} \rightarrow a_1^{[1]} \rightarrow a_2^{[1]}$$

layer-2 —

$$z^{[2]} = \begin{bmatrix} 0.1 & 0.2 \\ -0.3 & 0.4 \\ 0.2 & -0.2 \end{bmatrix} \begin{bmatrix} 0.2 \\ 0.0 \end{bmatrix} = \begin{bmatrix} 0.02 \\ -0.06 \\ 0.04 \end{bmatrix}$$

$$a^{[2]} = \begin{bmatrix} 0.02 \\ 0.0 \\ 0.04 \end{bmatrix}$$

Layer-3 —

$$z^{[3]} = \begin{bmatrix} 0.3 & -0.2 & 0.1 \\ -0.1 & 0.2 & 0.4 \end{bmatrix} \begin{bmatrix} 0.02 \\ 0.0 \\ 0.04 \end{bmatrix}$$

→ 0.006 + 0  
 + 0.004  
 = 0.01  
 → -0.002 + 0  
 + 0.016  
 = 0.014

$$z^{[3]} = \begin{bmatrix} 0.01 \\ 0.014 \end{bmatrix}$$

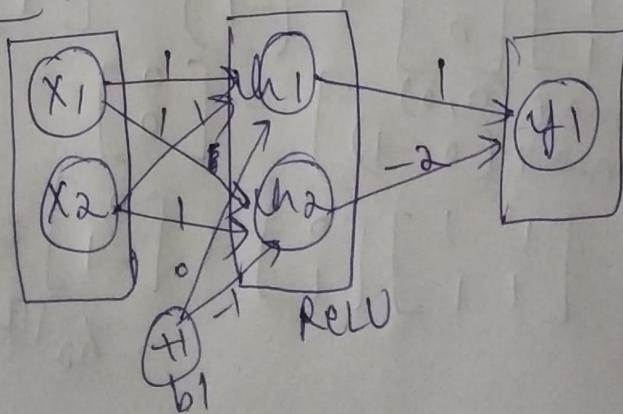
$$a^{[3]} = \begin{bmatrix} \max\{0, 0.01\} \\ \max\{0, 0.014\} \end{bmatrix} = \begin{bmatrix} 0.01 \\ 0.014 \end{bmatrix} \rightarrow \begin{array}{l} \text{raw} \\ \text{logits} \end{array}$$

Softmax applied on raw logits —

$$\text{softmax}(z)_i = \frac{e^{z_i}}{\sum_j e^{z_j}}$$

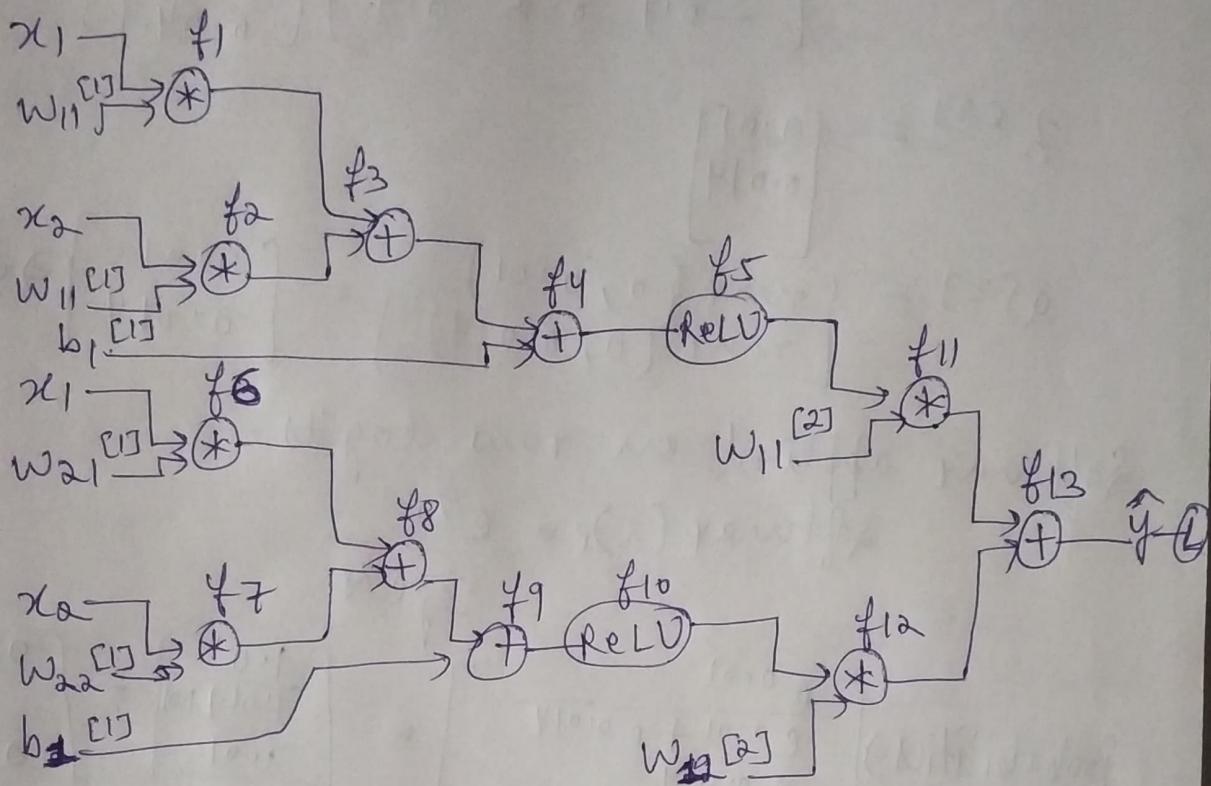
$$\text{Probabilities} \Rightarrow \begin{bmatrix} \frac{e^{0.01}}{e^{0.01} + e^{0.014}} \\ \frac{e^{0.014}}{e^{0.01} + e^{0.014}} \end{bmatrix} = \begin{bmatrix} \frac{1.01}{1.01 + 1.01} \\ \frac{1.01}{1.01 + 1.01} \end{bmatrix} = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} \approx \hat{y}$$

③ Given :-



x1	x2	y
0	0	0
0	1	1
1	0	1
1	1	0

## Computational graph :-



Forward pass :-

$$\text{output: } z_i^{[1]} = w_1 x_i + b_1 ; \quad i = \text{sample}$$

$$\text{sample-1} \Rightarrow z_1^{[1]} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ -1 \end{bmatrix} = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$$

$$a_1^{[1]} = \begin{bmatrix} \max(0, 0) \\ \max(0, -1) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\text{sample-2} \Rightarrow z_2^{[1]} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} + \begin{bmatrix} 0 \\ -1 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ -1 \end{bmatrix}$$

$$z_2^{[1]} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$a_2^{[1]} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$\text{sample-3} \Rightarrow z_3^{[1]} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} + \begin{bmatrix} 0 \\ -1 \end{bmatrix} = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$$

$$a_3^{[1]} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$\text{sample-4} \Rightarrow z_4^{[1]} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} + \begin{bmatrix} 0 \\ -1 \end{bmatrix} = \begin{bmatrix} 2 \\ 2 \end{bmatrix} + \begin{bmatrix} 0 \\ -1 \end{bmatrix}$$

$$z_4^{[1]} = \begin{bmatrix} 2 \\ 1 \end{bmatrix} = a_4^{[1]}$$

O/P layer :-

$$\text{sample-1} \Rightarrow h = \begin{bmatrix} 1 \\ -2 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Rightarrow a = \begin{bmatrix} 0 \\ 0 \end{bmatrix} = 0$$

$$\text{sample-2} \Rightarrow h = \begin{bmatrix} 1 \\ -2 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \Rightarrow a = \begin{bmatrix} 1 \\ 0 \end{bmatrix} = 1$$

$$\text{sample-3} \Rightarrow h = \begin{bmatrix} 1 \\ -2 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \Rightarrow a = \begin{bmatrix} 1 \\ 0 \end{bmatrix} = 1$$

$$\text{sample-4} \Rightarrow h = \begin{bmatrix} 1 \\ -2 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \Rightarrow a = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$



$$O/P \Rightarrow \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix} = \hat{y} = y$$

$$\text{loss} = \frac{1}{2} (\hat{y} - y)^2 = \frac{1}{2} (0 + 0 + 0 + 0) = 0$$

Backward pass :-

Since loss ( $L$ ) is exactly zero, all the gradients of loss w.r.t.  $x$ ,  $w$  and  $b$  will also be zero.

Update weights :-

$\mathbf{w}_i := \mathbf{w}_i - \eta \nabla_{\mathbf{w}_i} L$  with  $\nabla_{\mathbf{w}_i} L = 0$ , every parameter (weight) remains unchanged.