

CNN Architecture

3) Convolution Operation :-

assume stride = 1

1	0	0
0	1	0
0	0	1

1	0
0	1

Resultant Matrix

2	0
0	2

Resultant Height (or) width

\Rightarrow

(Image Height or width - Filter Height (or) width) / stride + 1

In this eg:-

$$\therefore \text{resultant ht or wd} = ((3-2)/1) + 1 = 2 \times 2$$

→ using im2col approach

channel 1 : $\begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \\ 13 & 14 & 15 & 16 \end{bmatrix}$ 4x4

filter1 : $\begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}$

channel 2 : $\begin{bmatrix} 17 & 18 & 19 & 20 \\ 21 & 22 & 23 & 24 \\ 25 & 26 & 27 & 28 \\ 29 & 30 & 31 & 32 \end{bmatrix}$ 4x4

filter2 : $\begin{bmatrix} 1 & -1 \\ 1 & 0 \end{bmatrix}$

filter3 : $\begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}$

im2col

im2col \Rightarrow matrix

1	2	3	5	6	7	9	10	11
2	3	4	6	7	8	10	11	12
5	6	7	9	10	11	13	14	15
6	7	8	10	11	12	14	15	16
17	18	19	21	22	23	25	26	27
18	19	20	22	23	24	26	27	28
21	22	23	25	26	27	29	30	31
22	23	24	26	27	28	30	31	32

generally :

→ img (batch-size, channels, ht, wd)

→ filter (filters, channels, ht, wd)

→ im2col size = (channels * filter(ht) * filter(wd)), (new ht * new wd * filters)

→ convolution operation resultant matrix size = ((Image ht/wd - filter ht/wd) / stride) + 1
= (new ht, new wd)

Algorithm to calculate Convolution using im2col method :-

step1: convert each stride of the convolutional filter over an image, into a column of a matrix. (X -im2col)

step2: flatten convolutional filter across the no of filters. (conv-flatten)

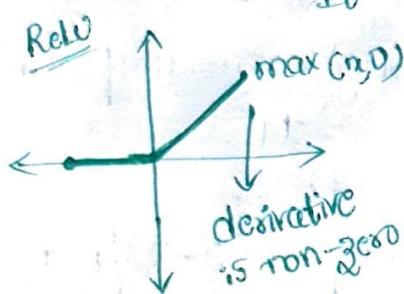
step3: calculate $C = (\text{conv-flatten}) @ (X\text{-im2col}) \Rightarrow @$ - matrix multiplication

step4: Reshape C to fit the resulting shape of matrix.

ReLU: → It is a non-linear activation function which filters out negative values.

→ It does not face problems such as vanishing gradients.

→ It result in non-zero gradients for positive values.



Max pooling: → It is a process of extract low level features in the image. This is done by picking image chunks of pre-determined sizes, and keeping the largest values from each of these chunks.

2	1	-11	0
0	10	-15	-16
23	100	2	3
-23	0	7	20

Max pooling

filter
ht/wd = 2,
stride = 2

$$\text{Resultant} = \left(\frac{(\text{image ht/wd} - \text{filter ht/wd})}{\text{stride}} + 1 \right)$$
$$\rightarrow ((4-2)/2+1) = (2, 2)$$

Soft max function: → It converts a vector of real values to a vector of values that range b/w 0 and 1. The newly transformed vector adds up to 1; the transformed vector becomes a probability distribution; A large value will be transformed to a value that is close to 1; a small value will be close to 0.

$$\text{softmax}(x_i) = \frac{e^{x_i}}{\sum_{i=1}^N e^{x_i}}$$

If $x_i = [+3, -2, -100]$, denominator = $\sum_{i=1}^N e^{x_i} = e^3 + e^{-2} + e^{-100} = 20.22$

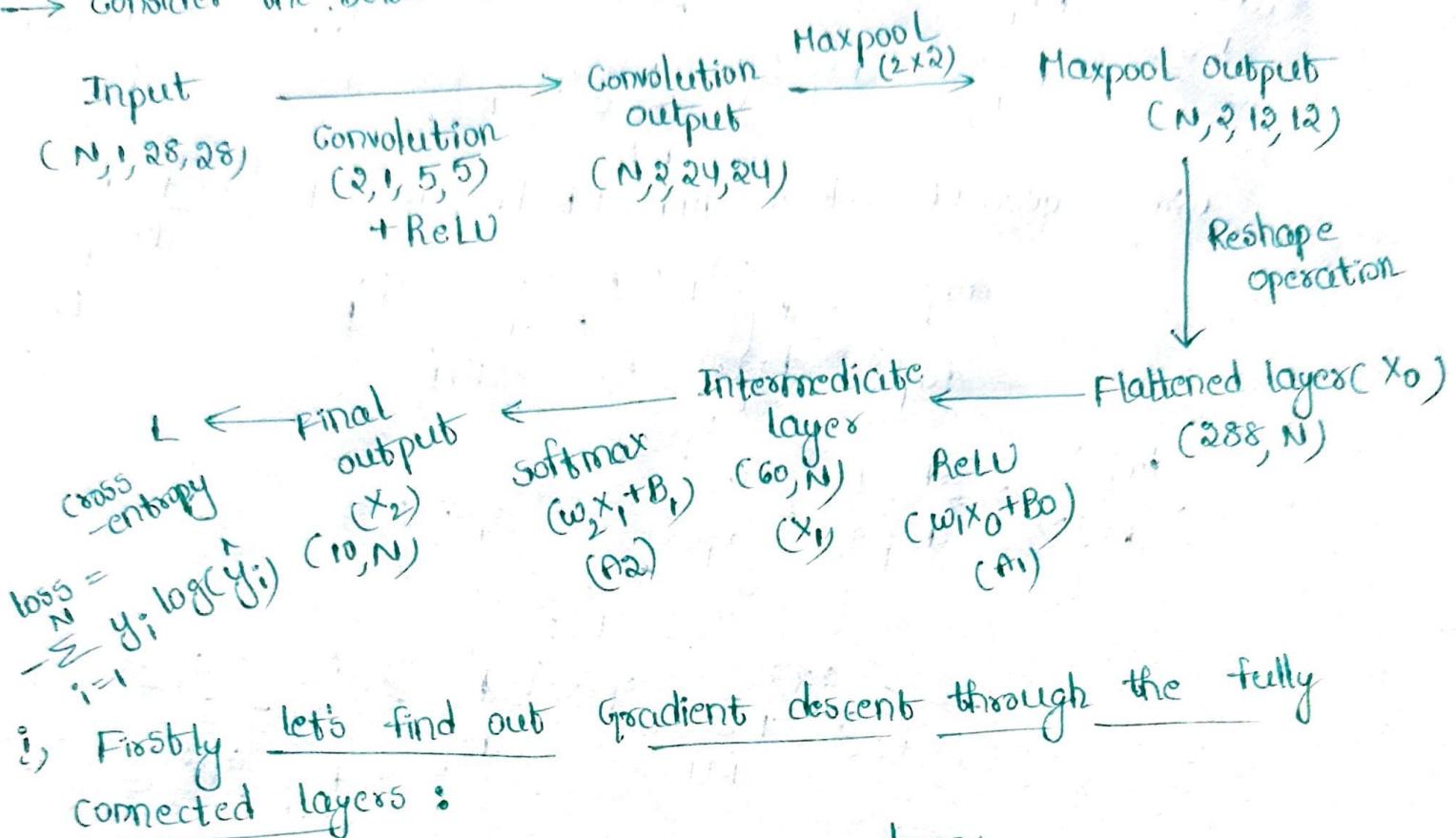
$$\rightarrow \text{softmax}(3) = \frac{e^3}{20.22} = 0.99$$

$$\rightarrow \text{softmax}(-2) = \frac{e^{-2}}{20.22} = 0.0066$$

$$\rightarrow \text{softmax}(-100) = \frac{e^{-100}}{20.22} = 1.88 \times 10^{-45} \text{ (close to zero)}$$

Back Propagation:

→ Consider the below CNN network.



i) Firstly, let's find out gradient descent through the fully connected layers :

$$\text{Gradient descent, } \theta = \theta - \alpha \frac{dF(\theta)}{d\theta}$$

α - learning rate.

ii) Calculating derivative of cross entropy with respect of w_2 and B_1 ,

$$\text{i.e., } \frac{dL}{dw_2}, \frac{dL}{dB_1} :$$

$$\text{let } a_2 = w_2 x_1 + B_1$$

$$x_2 = \text{softmax}(a_2)$$

$$\begin{aligned} \frac{dL}{dw_2} &= \frac{dL}{dx_2} \cdot \frac{dx_2}{da_2} \cdot \frac{da_2}{dw_2} = \frac{dL}{dx_2} \cdot \frac{dx_2}{da_2} \cdot \frac{d(w_2 x_1 + B_1)}{dw_2} \\ &= \frac{dL}{dx_2} \cdot \frac{dx_2}{da_2} \cdot (x_1) \end{aligned}$$

iii) Finding $\frac{dx_2}{da_2}$:

Let $a_i^{(2)}$ be the softmax output corresponding to $a_i^{(2)}$

$$\Rightarrow \frac{d a_i^{(2)}}{d a_j^{(2)}} = \frac{d}{d a_j^{(2)}} \left[\frac{e^{a_i^{(2)}}}{\sum_{j=1}^N e^{a_j^{(2)}}} \right]$$

there are two cases here, when $i=j$, $i \neq j$

$$\text{when } i=j, \frac{dx_i^{(2)}}{da_i^{(2)}} = \frac{d}{da_i^{(2)}} \left(\frac{e^{\alpha_i^{(2)}}}{\sum_{k=1}^K e^{\alpha_k^{(2)}}} \right) \quad \begin{bmatrix} \text{softmax}(x) = e^{\alpha_i^{(2)}} \\ \therefore \sum_{i=1}^K e^{\alpha_i^{(2)}} \end{bmatrix}$$

using quotient rule of differentiation, we get

$$\frac{dx_i^{(2)}}{da_i^{(2)}} = \frac{d}{da_i^{(2)}} \left[\frac{e^{\alpha_i^{(2)}}}{e^{\alpha_i^{(2)}} + \sum_{k \neq i} e^{\alpha_k^{(2)}}} \right]$$

$$\begin{aligned} & \because \frac{d}{dx} \left(\frac{u}{v} \right) \\ &= v \frac{du}{dx} - u \frac{dv}{dx} \\ & \quad v^2 \end{aligned}$$

$$\frac{dx_i^{(2)}}{da_i^{(2)}} = \frac{d}{da_i^{(2)}} \left[\frac{e^{\alpha_i^{(2)}}}{e^{\alpha_i^{(2)}} + \sum_{k \neq i} e^{\alpha_k^{(2)}}} \right]$$

$$= \left(e^{\alpha_i^{(2)}} + \sum_{k \neq i} e^{\alpha_k^{(2)}} \right) \frac{d}{da_i^{(2)}} \left(\frac{e^{\alpha_i^{(2)}}}{e^{\alpha_i^{(2)}} + \sum_{k \neq i} e^{\alpha_k^{(2)}}} \right) - e^{\alpha_i^{(2)}}$$

$$\frac{d}{da_i^{(2)}} \left(\frac{e^{\alpha_i^{(2)}}}{e^{\alpha_i^{(2)}} + \sum_{k \neq i} e^{\alpha_k^{(2)}}} \right)$$

$$\Rightarrow \left(e^{\alpha_i^{(2)}} + \sum_{k \neq i} e^{\alpha_k^{(2)}} \right)^2$$

$$= \left(e^{\alpha_i^{(2)}} + \sum_{k \neq i} e^{\alpha_k^{(2)}} \right) \left(\frac{e^{\alpha_i^{(2)}}}{e^{\alpha_i^{(2)}} + \sum_{k \neq i} e^{\alpha_k^{(2)}}} - e^{\alpha_i^{(2)}} \cdot e^{\alpha_i^{(2)}} \right)$$

$$\left(e^{\alpha_i^{(2)}} + \sum_{k \neq i} e^{\alpha_k^{(2)}} \right)^2$$

$$\frac{dx_i^{(2)}}{da_i^{(2)}} = \frac{e^{\alpha_i^{(2)}} + \sum_{k \neq i} e^{\alpha_k^{(2)}}}{\sum_{k=1}^N e^{\alpha_k^{(2)}} \cdot \sum_{k=1}^N e^{\alpha_k^{(2)}}} = \alpha_i^{(2)} (1 - \alpha_i^{(2)})$$

$$\therefore \frac{dx_2}{dA_2} = \frac{dx_1^{(2)}}{da_1^{(2)}} = \alpha_i^{(2)} (1 - \alpha_i^{(2)})$$

$$\begin{bmatrix} \alpha_i^{(2)} = \frac{\alpha_i^{(2)}}{\sum_{k=1}^N \alpha_k^{(2)}} \\ \alpha_i^{(2)} = \frac{\sum_{k=1}^N \alpha_k^{(2)}}{\sum_{k=1}^N \alpha_k^{(2)}} \end{bmatrix}$$

when
 $i \neq j$

$$\frac{d\alpha_i^{(2)}}{d\alpha_j^{(2)}} = \frac{d}{d\alpha_j^{(2)}} \left[\frac{\alpha_i^{(2)}}{\sum_{j=1}^N e^{\alpha_j^{(2)}}} \right]$$

$$= \frac{0 - e^{\alpha_i^{(2)}} \cdot e^{\alpha_j^{(2)}}}{\left(\sum_{j=1}^N e^{\alpha_j^{(2)}} \right)^2} = \frac{-e^{\alpha_i^{(2)}}}{\sum_{j=1}^N e^{\alpha_j^{(2)}}} \cdot \frac{e^{\alpha_j^{(2)}}}{\sum_{j=1}^N e^{\alpha_j^{(2)}}} = -\alpha_i^{(2)} \cdot \alpha_j^{(2)}$$

∴ thus,

$$\frac{d\alpha_i^{(2)}}{d\alpha_j^{(2)}} = \begin{cases} -\alpha_i^{(2)} \alpha_j^{(2)}, & \text{for } i \neq j \\ +\alpha_i^{(2)} (1 - \alpha_i^{(2)}), & \text{for } i = j \end{cases}$$

Now,

let the estimates of the model be $\hat{\alpha}_i^{(2)}$ and
let the actual labels be $\hat{z}_i^{(2)}$.

so, cross entropy = $-\sum_{i=1}^N \hat{z}_i^{(2)} \log(\hat{\alpha}_i^{(2)})$

$$\frac{dL}{dx_2} = -\hat{z}_i^{(2)} \frac{1}{\hat{\alpha}_i^{(2)}}$$

$$\begin{aligned} \therefore \text{Now, } \frac{dL}{dx_2} \cdot \frac{dx_2}{d\alpha_2} &= -\sum_{j=1}^N \hat{z}_j^{(2)} \frac{1}{\hat{\alpha}_j^{(2)}} \alpha_j^{(2)} (1 - \alpha_j^{(2)}) \\ &\quad - \sum_{j \neq i} \hat{z}_j^{(2)} \frac{1}{\hat{\alpha}_j^{(2)}} (-\alpha_j^{(2)} \hat{x}_j^{(2)}) \\ &= \sum_j \hat{z}_j^{(2)} \alpha_j^{(2)} - \sum_{j \neq i} \hat{z}_j^{(2)} (1 - \alpha_j^{(2)}) \end{aligned}$$

combining these two

$$= \sum_j \hat{z}_j^{(2)} \alpha_j^{(2)} - \sum_{j=1}^N \hat{z}_j^{(2)}$$

one-hot encoded for single point

$$= \alpha_i^{(2)} - \hat{z}_i^{(2)}$$

$$\boxed{\therefore \frac{dL}{d\omega_2} = (x_2 - \hat{z}_2) * x_1^T}$$

$$\text{similarly, } \frac{dL}{dB_1} \rightarrow \frac{dL}{dB_1} = \frac{dL}{dx_2} \cdot \frac{dx_2}{dA_2} \cdot \frac{dA_2}{dB_1}$$

$$= (x_2 - b_2) \cdot \frac{d}{dB_1} (w_2 x_1 + b_1) = (x_2 - b_2)$$

$$\boxed{\therefore \frac{dL}{dB_1} = (x_2 - b_2)}$$

ii) Calculating derivative of cross entropy with respect to w_1 and b_0 :

i.e., $\frac{dL}{dw_1}, \frac{dL}{db_0}$: we know that, $A_1 = w_1 x_0 + b_0$
 $x_1 = \text{ReLU}(A_1)$
 $a_2 = w_2 x_1 + b_1$
 $x_2 = \text{softmax}(a_2)$

$$\begin{aligned} \therefore \frac{dL}{dw_1} &= \frac{dL}{dx_2} \cdot \frac{dx_2}{dA_2} \cdot \frac{dA_2}{dx_1} \cdot \frac{dx_1}{dA_1} \cdot \frac{dA_1}{dw_1} \\ &= (x_2 - b_2) \cdot w_2^T \cdot (\text{ReLU}'(A_1)) \cdot x_0 \end{aligned}$$

$$\boxed{\therefore \frac{dL}{dw_1} = [w_2^T * (x_2 - b_2)] \cdot [\text{ReLU}'(A_1)] * x_0}$$

where, \cdot → element wise product (dot product)
 $*$ → matrix multiplication

similarly, $\frac{dL}{db_0} = \frac{dL}{dx_2} \cdot \frac{dx_2}{dA_2} \cdot \frac{dA_2}{dx_1} \cdot \frac{dx_1}{dA_1} \cdot \frac{dA_1}{db_0}$

$$\frac{dL}{db_0} = (x_2 - b_2) \cdot w_2^T \cdot (\text{ReLU}'(A_1)) \cdot (1).$$

$$\boxed{\therefore \frac{dL}{db_0} = [w_2^T * (x_2 - b_2)] \cdot [\text{ReLU}'(A_1)]}$$

iii) Calculating the derivative of loss function with respect to x_0 :

i.e. $\frac{dL}{dx_0}$: $\frac{dL}{dx_0} = \frac{dL}{dx_2} \cdot \frac{dx_2}{dA_2} \cdot \frac{dA_2}{dx_1} \cdot \frac{dx_1}{dA_1} \cdot \frac{dA_1}{dx_0}$

$$\boxed{\therefore \frac{dL}{dx_0} = [w_1^T * [w_2^T * (x_2 - b_2)]] \cdot [\text{ReLU}'(A_1)]}$$

iii) Calculating gradients at Max pooling Layer : we know that the first fully connected layer is a reshaped version of the max pooling layer, we just need to reshape our gradient matrix at the first fully connected layer, back to the shape of the max pooling layer.

$$\therefore \text{delta_maxpool} = \text{reshape to maxpool o/p shape } \left(\frac{dL}{dx_0} \right)$$

iii) Calculating gradient at Convolutional layer : To calculate the gradients at the convolutional layer, we need to move each gradient element back to the positions in the convolutional layer, from where the maximum pixel values were extracted.

Forward Pass → let's say x_{conv} (before pooling) is :

$$x_{\text{conv}} = \begin{bmatrix} 1 & 3 \\ 2 & 4 \end{bmatrix}$$

If we apply 2×2 max-pooling with a stride of 2, we get,

$$\text{Max pooled output} = [4]$$

since, '4' was the selected max value, the max-indices array stores $\{\text{max-indices}\} = [(0, 1, 0, 0, 0, 0)]$

This means,

$(1, 1, 0) \rightarrow$ original location of max value '4' in x_{conv} .

$(0, 0, 0) \rightarrow$ is the corresponding location in delta_maxpool.

Backward Pass : Now, assume the gradient from max pooling (delta_maxpool) is $\delta_{\text{maxpool}} = [10]$

Using our max-indices values :

$$\text{delta_conv}[0, 1, 0] = \text{delta_maxpool}[0, 0, 0]$$

This means: $\delta_{\text{conv}}(1,1,0) = \delta_{\text{maxpool}}(0,0,0) = 10$

\therefore the updated delta-conv = $\begin{bmatrix} 0 & 0 \\ 0 & 10 \end{bmatrix}$

\rightarrow And finally we apply derivative of ReLU to delta-conv, to ensure that the gradient only flows through the active neurons (i.e., $x_{\text{Conv}} > 0$).

4) Calculating gradient of convolutional filter: Now that we have found the gradient of the convolutional layer, we need to calculate the gradient of the convolutional filter. This will be used to optimize the filter in each learning step.

Algorithm to calculate gradient of convolutional filter:-

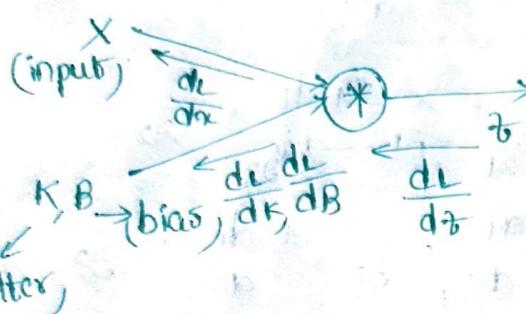
Let 'q' be the gradient matrix for the convolutional layer. This has a dimension of $(1, 3, 24, 24)$. let 'I' be the input image of shape $(1, 3, 28, 28)$.

- i) For a given channel 'c' in 'q', pick an element
- ii) use 5×5 filters with stride=1 to create 5×5 chunks of image input 'I'.
- iii) Multiply the chosen gradient element with the consecutive chunks and add them up.
- iv) The resultant matrix after the consecutive addition operation is the gradient associated with the channel 'c' in the convolutional filter.

v) Repeat steps 1 to 4 for the rest of the channels in 'q'.

Given that we have to find the gradient with respect to the convolutional filter by pixel, this process will take a lot of time and effort. One way to overcome this is to use the `im2col()` function.

Back Propagation in Convolution layer:



$$\begin{bmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \\ x_{31} & x_{32} & x_{33} \end{bmatrix} \otimes \begin{bmatrix} k_{11} & k_{12} \\ k_{21} & k_{22} \end{bmatrix} + B = \begin{bmatrix} t_{11} & t_{12} \\ t_{21} & t_{22} \end{bmatrix}$$

here,

$$t_{11} = x_{11}k_{11} + x_{12}k_{12} + x_{21}k_{21} + x_{22}k_{22} + B$$

$$t_{12} = x_{12}k_{11} + x_{13}k_{12} + x_{22}k_{21} + x_{23}k_{22} + B$$

$$t_{21} = x_{21}k_{11} + x_{22}k_{12} + x_{31}k_{21} + x_{32}k_{22} + B$$

$$t_{22} = x_{22}k_{11} + x_{23}k_{12} + x_{32}k_{21} + x_{33}k_{22} + B$$

$$t_{ij} \rightarrow t_{11}, t_{12}, t_{21}, t_{22}$$

$$K_{mn} \rightarrow k_{11}, k_{12}, k_{21}, k_{22}$$

From Einstein notation of chain rule multiplication

$$\frac{dL}{dK_{mn}} = \sum \frac{dL}{d\tau_{ij}} * \frac{d\tau_{ij}}{dK_{mn}}$$

$$\therefore \frac{dL}{dK_{11}} = \left(\frac{dL}{d\tau_{11}} * \frac{d\tau_{11}}{dK_{11}} \right) + \left(\frac{dL}{d\tau_{12}} * \frac{d\tau_{12}}{dK_{11}} \right) + \left(\frac{dL}{d\tau_{21}} * \frac{d\tau_{21}}{dK_{11}} \right) + \left(\frac{dL}{d\tau_{22}} * \frac{d\tau_{22}}{dK_{11}} \right)$$

$$\frac{dL}{dK_{12}} = \left(\frac{dL}{d\tau_{11}} * \frac{d\tau_{11}}{dK_{12}} \right) + \left(\frac{dL}{d\tau_{12}} * \frac{d\tau_{12}}{dK_{12}} \right) + \left(\frac{dL}{d\tau_{21}} * \frac{d\tau_{21}}{dK_{12}} \right) + \left(\frac{dL}{d\tau_{22}} * \frac{d\tau_{22}}{dK_{12}} \right)$$

similarly $\frac{dL}{dK_{21}}$ and $\frac{dL}{dK_{22}}$.

Transforming, $\frac{dL}{dK_{11}}$ by substituting the derivative values,

$$\frac{dL}{dK_{11}} = \frac{dL}{d\tau_{11}} * x_{11} + \frac{dL}{d\tau_{12}} * x_{12} + \frac{dL}{d\tau_{21}} * x_{21} + \frac{dL}{d\tau_{22}} * x_{22}$$

Similarly, substitute for $\frac{dL}{dK_{12}}, \frac{dL}{dK_{21}}$ and $\frac{dL}{dK_{22}}$.

$$\frac{dL}{dK_{12}} = \frac{dL}{dt_{11}} * x_{12} + \frac{dL}{dt_{12}} * x_{13} + \frac{dL}{dt_{21}} * x_{22} + \frac{dL}{dt_{22}} * x_{23}$$

$$\frac{dL}{dK_{21}} = \frac{dL}{dt_{11}} * x_{21} + \frac{dL}{dt_{12}} * x_{22} + \frac{dL}{dt_{21}} * x_{31} + \frac{dL}{dt_{22}} * x_{32}$$

$$\frac{dL}{dK_{22}} = \frac{dL}{dt_{11}} * x_{22} + \frac{dL}{dt_{12}} * x_{23} + \frac{dL}{dt_{21}} * x_{32} + \frac{dL}{dt_{22}} * x_{33}$$

$$\begin{bmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \\ x_{31} & x_{32} & x_{33} \end{bmatrix} \otimes \begin{bmatrix} \frac{dL}{dt_{11}} & \frac{dL}{dt_{12}} \\ \frac{dL}{dt_{21}} & \frac{dL}{dt_{22}} \end{bmatrix} = \begin{bmatrix} \frac{dL}{dK_{11}} & \frac{dL}{dK_{21}} \\ \frac{dL}{dK_{21}} & \frac{dL}{dK_{22}} \end{bmatrix} = \frac{dL}{dK}$$

$$X \otimes \frac{dL}{dt} = \frac{dL}{dK}$$

$$\therefore \frac{dL}{dK} = \text{cov}(X, \frac{dL}{dt})$$

2)

$$\boxed{\frac{dL}{dB}}$$

$$= \frac{dL}{dt_{11}} * \frac{dt_{11}}{dB}$$

$$\therefore \frac{dL}{dB} = \sum \frac{dL}{dt_{ij}} * \frac{dt_{ij}}{dB}$$

$$= \frac{dL}{dt_{11}} * \frac{dt_{11}}{dB} + \frac{dL}{dt_{12}} * \frac{dt_{12}}{dB} + \frac{dL}{dt_{21}} * \frac{dt_{21}}{dB} + \frac{dL}{dt_{22}} * \frac{dt_{22}}{dB}$$

$$\therefore \frac{dL}{dB} = \frac{dL}{dt_{11}} + \frac{dL}{dt_{12}} + \frac{dL}{dt_{21}} + \frac{dL}{dt_{22}}$$

$$\therefore \frac{dL}{dB} = \text{sum}(\frac{dL}{dt_{ij}})$$

3)

$$\frac{dL}{dx} = \frac{dL}{dt} * \frac{dt}{dx} = \sum \frac{dL}{dt_{ij}} * \frac{dt_{ij}}{dx}$$

$$\frac{dL}{dx_{11}} = \frac{dL}{dt_{11}} * \frac{dt_{11}}{dx_{11}} = \frac{dL}{dt_{11}} * K_{11}$$

$$\frac{dL}{dx_{12}} = \frac{dL}{dt_{11}} * \frac{dt_{11}}{dx_{12}} + \frac{dL}{dt_{12}} * \frac{dt_{12}}{dx_{12}}$$

$$= (\frac{dL}{dt_{11}} * K_{12}) + (\frac{dL}{dt_{12}} * K_{11})$$

$\frac{dt_{11}}{dx_{11}}, \frac{dt_{21}}{dx_{11}}, \frac{dt_{22}}{dx_{11}} = 0$ as any change in x_{11} does not effect t_{12}, t_{21}, t_{22}

$$\text{Therefore, } \frac{dL}{dx_{11}} = \frac{dL}{d\omega_{11}} * K_{11}$$

$$\frac{dL}{dx_{12}} = \frac{dL}{d\omega_{11}} * K_{12} + \frac{dL}{d\omega_{12}} * K_{21}$$

$$\frac{dL}{dx_{13}} = \frac{dL}{d\omega_{12}} * K_{12}$$

$$\frac{dL}{dx_{21}} = \frac{dL}{d\omega_{11}} * K_{21} + \frac{dL}{d\omega_{21}} * K_{11}$$

$$\frac{dL}{dx_{22}} = \frac{dL}{d\omega_{11}} * K_{22} + \frac{dL}{d\omega_{12}} * K_{21} + \frac{dL}{d\omega_{21}} * K_{12} + \frac{dL}{d\omega_{22}} * K_{11}$$

$$\frac{dL}{dx_{23}} = \frac{dL}{d\omega_{12}} * K_{22} + \frac{dL}{d\omega_{22}} * K_{12}$$

$$\frac{dL}{dx_{31}} = \frac{dL}{d\omega_{21}} * K_{21}$$

$$\frac{dL}{dx_{32}} = \frac{dL}{d\omega_{21}} * K_{22} + \frac{dL}{d\omega_{22}} * K_{21}$$

$$\frac{dL}{dx_{33}} = \frac{dL}{d\omega_{22}} * K_{22}$$

By
Simplifying,

$$\begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & \frac{dL}{d\omega_{11}} & \frac{dL}{d\omega_{12}} & 0 \\ 0 & \frac{dL}{d\omega_{21}} & \frac{dL}{d\omega_{22}} & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \otimes \begin{bmatrix} K_{22} & K_{21} \\ K_{12} & K_{11} \end{bmatrix} = \begin{bmatrix} \frac{dL}{dx_{11}} & \frac{dL}{dx_{12}} & \frac{dL}{dx_{13}} \\ \frac{dL}{dx_{21}} & \frac{dL}{dx_{22}} & \frac{dL}{dx_{23}} \\ \frac{dL}{dx_{31}} & \frac{dL}{dx_{32}} & \frac{dL}{dx_{33}} \end{bmatrix}$$

(180° inverted
'K' matrix)

$$\therefore \frac{dL}{dx} = \text{conv}(\text{padded}(\frac{dL}{d\omega}), \text{180° rotated filter } K)$$

$$\text{Therefore, } \frac{dL}{dK} = \text{conv}(x, \frac{dL}{d\omega})$$

$$\frac{dL}{dB} = \text{sum}(\frac{dL}{d\omega})$$

$$\frac{dL}{dx} = \text{conv}(\text{padded}(\frac{dL}{d\omega}), \text{180° rotated filter } K)$$

Forward Pass

let's assume, Input image $I = \begin{bmatrix} 1 & 2 & 3 & 4 & 5 \\ 6 & 7 & 8 & 9 & 10 \\ 11 & 12 & 13 & 14 & 15 \\ 16 & 17 & 18 & 19 & 20 \\ 21 & 22 & 23 & 24 & 25 \end{bmatrix}$ $(1, 1, 5, 5)$

And 3×3 filter F is:

$$F = \begin{bmatrix} 1 & 0 & +1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix}$$

o/p feature map after convolution = $I @ F = \begin{bmatrix} -6 & -6 & -6 \\ -6 & -6 & -6 \\ -6 & -6 & -6 \end{bmatrix}$ $(1, 1, 3, 3)$

Backward Propagation

Finding $\frac{dL}{dF}$ or $\frac{dL}{dk}$:

let ' δ ' be the gradient from the next layer (same shape as 'O')

$$\delta = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

$$\frac{dL}{dF} = \sum (\delta \times \text{input patch})$$

for example, the top-left value in the gradient for $F(0,0)$

$\sum_{\text{all positions}} (\text{patch}_0 \times \delta_{\text{corresponding position}})$ this is done

for all filter positions.

\Rightarrow so, for a 3×3 filter of 5×5 input image, we will extract $3 \times 3 = 9$ patches, and each patch will be a flattened vector of size 9 (since $3 \times 3 = 9$).

$$\therefore \text{img2col} = \begin{bmatrix} 1 & 2 & 3 & 6 & 7 & 8 & 11 & 12 & 13 \\ 2 & 3 & 4 & 7 & 8 & 9 & 12 & 13 & 14 \\ \vdots & \vdots \\ 13 & 14 & 15 & 18 & 19 & 20 & 23 & 24 & 25 \end{bmatrix}_{9 \times 9}$$

Reshape the gradient matrix (G) or (δ)

$$G = [1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1]_{1 \times 9}$$

gradient is calculated as

$$\frac{dL}{dF} = \text{im2col}(I) \times G^T \\ (9 \times 9), (9 \times 1)$$

$$\therefore \frac{dL}{dF} = (\text{im2col}(I) \times G^T) \cdot \text{reshape}(3, 3)$$

→ reshape to filter size.

Algorithm for faster computation of gradient of convolutional filter:

- (i) Convert the ip image to im2col format; im2col matrix is a 2D-matrix where each column is a flattened vector of elements covered in a single stride of filter.
- (ii) Reshape the error matrix ' δ ' to a 2D-matrix; each row is a flattened vector of the error in each channel).
- (iii) Multiply these two matrices and reshape the results.

Adam Optimizer: Adam (Adaptive Moment Estimation) is an optimization algorithm that

- i) uses Momentum (exponential moving average of gradients) to smooth updates
- ii) uses velocity (exponential moving average of squared gradients) to adjust the learning rate.

The update rule is:

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$$

$$\theta = \theta - \frac{\alpha \cdot m_t}{\sqrt{v_t} + \epsilon}$$

where, $m_t \rightarrow$ Momentum (1st moment estimate)
 $v_t \rightarrow$ Velocity (2nd moment estimate)
 $g_t \rightarrow$ gradient of weight
 $\beta_1 \rightarrow$ Controls momentum decay (e.g., 0.9)
 $\beta_2 \rightarrow$ Controls velocity decay (e.g., 0.995)
 $d \rightarrow$ learning rate (e.g., 0.001)
 $\epsilon \rightarrow$ small number to avoid division by zero (e.g., 1e-7)

Example

Consider a small neural network with 2 weights (w_1, w_2), 2 biases (b_0, b_1), 1 convolutional filter (conv1) & let the gradients be dw_1, dw_2, db_0, db_1 and conv1-delta.

- Initial weights $w_1 = 0.5, w_2 = 0.3, b_0 = 0.2, b_1 = -0.1, \text{conv1} = 0.8$
- Gradients $dw_1 = 0.02, dw_2 = 0.01, db_0 = 0.01, db_1 = -0.03, \text{conv1-delta} = 0.05$
- Hyperparameters $\beta_1 = 0.9, \beta_2 = 0.995, \text{learning-rate} = 0.01$

calculating First Moment:

$$\begin{aligned}\text{momentum_}w_1 &= \beta_1 * \text{momentum_}w_1 + (1 - \beta_1) * dw_1 \\ &= 0.9 * 0 - 0.02 + (1 - 0.9) * 0.02 = 0.002\end{aligned}$$

similarly, $\text{momentum_}w_2 = -0.004, \text{momentum_}b_0 = 0.001,$

$\text{momentum_}b_1 = -0.003, \text{momentum_conv1} = 0.005$

calculating Second Moment:

$$\begin{aligned}\text{velocity_}w_1 &= \beta_2 * \text{velocity_}w_1 + (1 - \beta_2) * (dw_1)^2 \\ &= 0.995 * 0 + (1 - 0.995) * (0.02)^2 = 0.000002\end{aligned}$$

similarly, $\text{velocity_}w_2 = 0.000008, \text{velocity_}b_0 = 0.0000005,$

$\text{velocity_}b_1 = 0.995 * 0 + (1 - 0.995) * (-0.03)^2 = 0.0000045,$

$\text{velocity_conv1} = 0.0000125.$

Finally update weights using Adam formula:

$$w_1 = w_1 - \text{learning_rate} * \left(\frac{\text{momentum_}w_1}{\sqrt{\text{velocity_}w_1 + \epsilon}} \right)$$