

# CNN From Scratch

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

Here five major topics of CNN is discussed

- (1) Convolution operation
- (2) Convolution Layer.
- (3) Reshape Layer
- (4) Binary Cross Entropy Loss
- (5) Sigmoid Activation.

## (1) Convolution / Correlation Operation

### Cross-Correlation

Suppose we have given input matrix, arbitrarily big and the smaller one - kernel or filter. The cross-correlation between input and kernel will create another matrix. The values contained by output matrix are obtained by sliding the kernel on the top of input and by performing a multiplication of adjacent values then summing them up.

The size of output will be size of input minus size of kernel plus one.

$$\therefore Y = I - K + 1$$

for example.

1	6	2
5	3	1
7	0	9

input  $3 \times 3$



1	2
-1	0

$2 \times 2$

a	b
c	d

$2 \times 2$

kernel

output.

where

$$a = 1 \times 1 + 6 \times 2 + 5 \times (-1) + 3 \times 0 = 8$$

$$b = 6 \times 1 + 2 \times 2 + 3 \times (-1) + 1 \times 0 = 7$$

$$c = 5 \times 1 + 3 \times 2 + 7 \times (-1) + 0 \times 0 = 4$$

$$d = 3 \times 1 + 1 \times 2 + 0 \times (-1) + 9 \times 0 = 5$$

∴ Output matrix =  $\begin{bmatrix} 8 & 7 \\ 4 & 5 \end{bmatrix}_{2 \times 2}$

Output size =  $3 - 2 + 1 = 2$ .

★ is cross-correlation operator  
 \* is convolution operator.

$\text{rot}180^\circ \equiv$  horizontal flip followed by  
vertical flip.

Date: \_\_\_\_\_

Page: \_\_\_\_\_

## Convolution Operation

It is same as cross-correlation operation.  
Only the difference is kernel is rotated by  
 $180^\circ$ .

i.e.

$$\begin{bmatrix} 1 & 6 & 2 \\ 5 & 3 & 1 \\ 7 & 0 & 4 \end{bmatrix} * \begin{bmatrix} 0 & -1 \\ 2 & 1 \end{bmatrix} = \begin{bmatrix} 7 & 5 \\ 11 & 3 \end{bmatrix}$$

Output

input ~~kernel~~

~~rot $(180^\circ)$  kernel~~

In other word

convolution between  $I$  and  $k$  is cross  
correlation between  $I$  and rotated version  
of  $k$ .

$$\text{conv}(I, k) = I \star \text{rot}180^\circ(k)$$

$$\therefore I * k = I \star \text{rot}180^\circ(k).$$

### Valid Cross-correlation

$$\begin{bmatrix} 1 & 6 & 2 \\ 5 & 3 & 1 \\ 7 & 0 & 4 \end{bmatrix} * \begin{bmatrix} 1 & 2 \\ -1 & 0 \end{bmatrix} = \begin{bmatrix} 8 & 7 \\ 4 & 5 \end{bmatrix}$$

Start computing the product by placing the kernel entirely on the top of input and stop sliding the kernel when it hits the border of input.

### Full Cross-correlation

$$\begin{bmatrix} 1 & 6 & 2 \\ 5 & 3 & 1 \\ 7 & 0 & 4 \end{bmatrix} \underset{\text{full}}{*} \begin{bmatrix} 1 & 2 \\ -1 & 0 \end{bmatrix} = \begin{bmatrix} a & b & c & d \\ e & f & g & h \\ i & j & k & l \\ m & n & o & p \end{bmatrix}$$

In this mode we start computing the product as soon as there is intersection between kernel and the input. Hence,

$$a = 0 \times 1 = 0$$

$$b = -1 \times 1 + 0 \times 6 = -1$$

$$c = -1 \times 6 + 0 \times 2 = -6$$

$$d = -1 \times 2 = -2$$

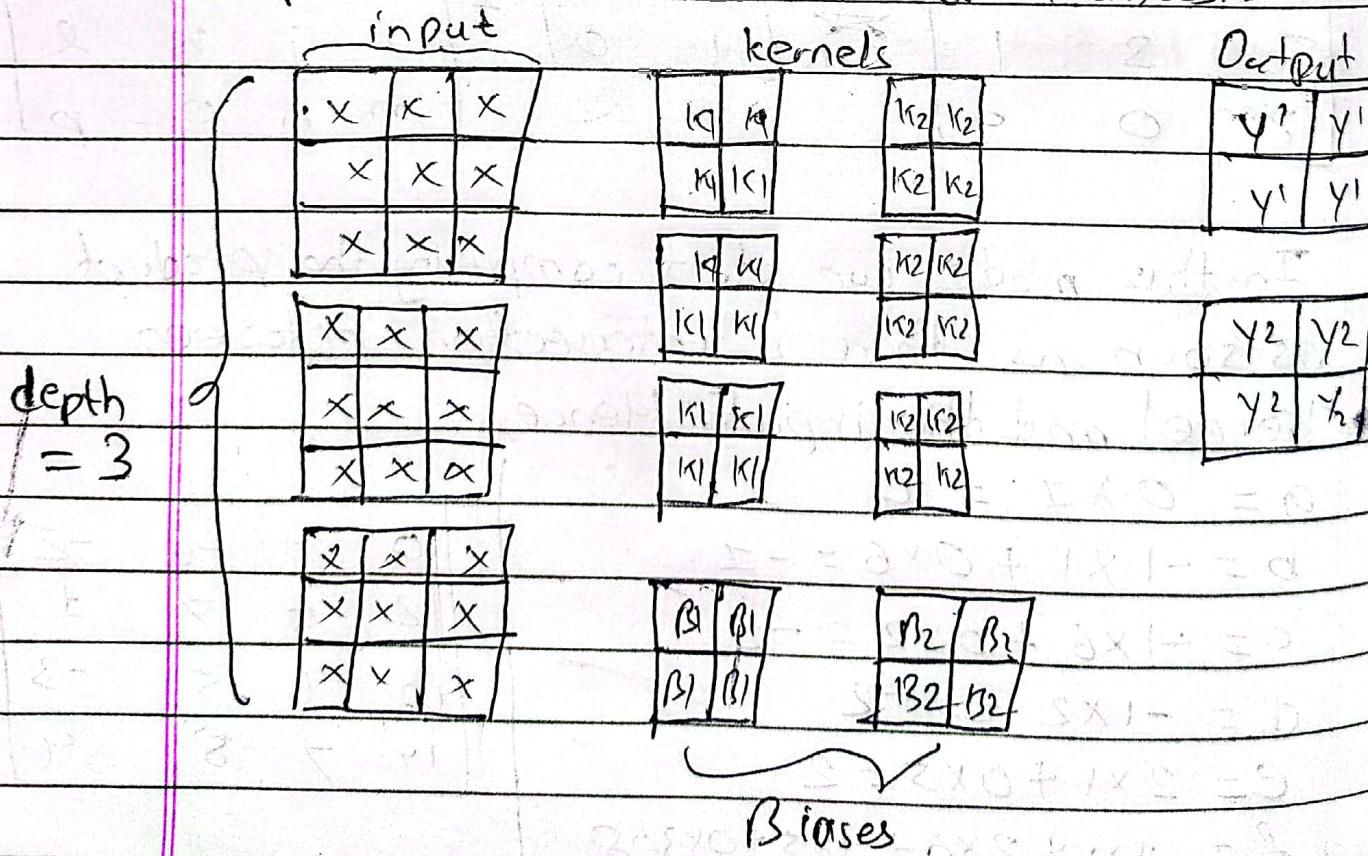
$$e = 2 \times 1 + 0 \times 5 = 2$$

$$f = 1 \times 1 + 2 \times 6 + -1 \times 5 + 0 \times 3 = 8$$

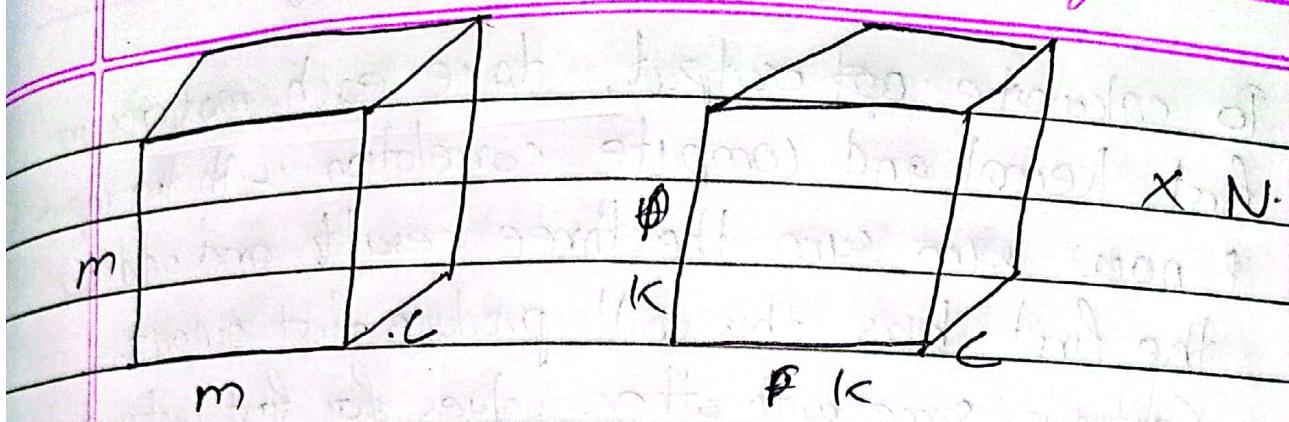
$$\begin{bmatrix} 0 & 1 & -6 & -2 \\ 2 & 8 & 7 & 1 \\ 10 & 4 & 5 & -3 \\ 14 & 7 & 8 & 4 \end{bmatrix}$$

## 2) Convolution Layer

The convolution layer takes : 3-dimensional block of data as input, in that case the depth of i/p is 3. The layer has trainable parameters among them the kernels. Each kernel is also a three dimensional block that extend the full depth of input. The layer may have multiple kernels all of which must extend the full depth of i/p. To each kernel we associate a bias matrix that will have same shape as the outputs. Finally the layer will produce the three dimensional data as output. The depth of output is the same as no. of kernels.



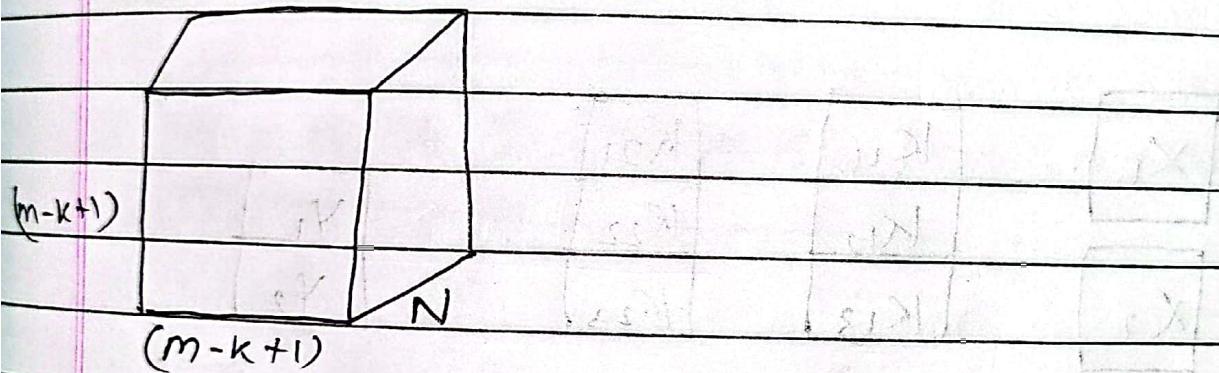
Date: \_\_\_\_\_  
Page: \_\_\_\_\_



Input of  $(m, m)$  kernel of size  $(k, k)$

height & height with depth 'c'

and 'c' depth  $N$  is number of such kernel



Output with size  $(m-k+1, m-k+1)$   
and  $N$  as depth.

To calculate opt output, take each matrix in first kernel and compute correlation with the input

now sum sum the three results and add up the first bias, this will produce first output.

Continue same with other values for first output.

Like wise second output is computed by taking the sum of correlation between input and second kernel plus second bias

$x_1$	$k_{11}$	$k_{21}$		$y_1$
$x_2$	$k_{12}$	$k_{22}$		$y_2$
$x_3$	$k_{13}$	$k_{23}$		

$x_1$	$B_1$	$B_2$
-------	-------	-------

$$y_1 = B_1 + x_1 \star k_{11} + x_2 \star k_{12} + x_3 \star k_{13}$$

$$y_2 = B_2 + x_1 \star k_{21} + x_2 \star k_{22} + x_3 \star k_{23}$$

For d number of kernels.

$$y_d = B_d + x_1 \star k_{d1} + x_2 \star k_{d2} + x_3 \star k_{d3}$$

For input of depth d.

For input of depth  $n$ .

$$y_1 = b_1 + x_1 * k_{11} + \dots + x_n * k_{1n}$$

$$y_2 = b_2 + x_1 * k_{21} + \dots + x_n * k_{2n}$$

$$\vdots$$
$$y_d = b_d + x_1 * k_{d1} + \dots + x_n * k_{dn}.$$

This is forward propagation of ~~the~~ convolution layer.

$$\therefore y_i = b_i + \sum_{j=1}^n x_j * k_{ij}, i = 1 \dots d.$$

Here we cannot directly perform matrix operation, as variables are not scalar, they are matrices and  $*$  operator is not just multiplication operator.

let  $\cdot \star$  is an operator which can perform some task then we can write:

$$\begin{array}{c|c} y_1 \\ y_2 \\ \vdots \\ y_d \end{array} = \begin{array}{c|c} b_1 \\ b_2 \\ \vdots \\ b_d \end{array} + \begin{array}{c|c} k_{11} & k_{12} \dots k_{1n} \\ k_{21} & k_{22} \dots k_{2n} \\ \vdots & \vdots \\ k_{d1} & k_{d2} \dots k_{dn} \end{array} \cdot \star \begin{array}{c|c} x_1 \\ x_2 \\ \vdots \\ x_n \end{array}$$
$$\therefore Y = B + K \cdot \star X$$

It is the same equation for dense layer, except that instead of having matrices of scalar, we have matrices of matrices and instead of regular dot product we have cross-correlated dot product.

If each of these matrices were one dimensional (i.e. each matrix contains a single element) then the cross correlation between two matrices containing one element is equal to one-dimensional matrix containing the product of both of these elements and therefore this new operator ends of being a regular dot product. Hence convolution layer is nothing but a generalization of dense layer. (or rather the dense layer is a particular case of the convolution layer).

Backward propagation

As earlier, assume we have  $\frac{\partial E}{\partial y_i}$  then, we have

to compute derivative of  $E$  w.r.t to trainable parameter of the layer i.e. kernels and biases and also need to compute the derivative of  $E$  w.r.t to input of the layer

$$\begin{array}{c} \xrightarrow{\frac{\partial E}{\partial k_{ij}}, \frac{\partial E}{\partial B_i}} \\ \xrightarrow{\frac{\partial E}{\partial y_i}} \end{array}$$

$$\frac{\partial E}{\partial x_j}$$

① For  $\frac{\partial E}{\partial k_{ij}}$

from eq<sup>n</sup> of forward propagation

$$y_i = B_i + \sum_{j=1}^n x_j * k_{ij}$$

Expanding sum

$$y_i = B_i + x_1 * k_{i1} + \dots + x_n * k_{in}$$

Simplifying this eq<sup>n</sup>

$$y_i = B_i + x_1 * k_{i1}$$

$$\begin{bmatrix} y_{11} & y_{12} \\ y_{21} & y_{22} \end{bmatrix} = \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{bmatrix} + \begin{bmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \\ x_{31} & x_{32} & x_{33} \end{bmatrix} \star \begin{bmatrix} k_{11} & k_{12} \\ k_{21} & k_{22} \end{bmatrix}$$

$$y_{11} = b_{11} + k_{11}x_{11} + k_{12}x_{12} + k_{21}x_{21} + k_{22}x_{22}$$

$$y_{12} = b_{12} + k_{11}x_{12} + k_{12}x_{13} + k_{21}x_{22} + k_{22}x_{23}$$

$$y_{21} = b_{21} + k_{11}x_{21} + k_{12}x_{22} + k_{21}x_{31} + k_{22}x_{32}$$

$$y_{22} = b_{22} + k_{11}x_{22} + k_{12}x_{23} + k_{21}x_{32} + k_{22}x_{33}$$

Now,

$$\frac{\partial E}{\partial k_{11}} = \frac{\partial E}{\partial y_{11}} \frac{\partial y_{11}}{\partial k_{11}} + \frac{\partial E}{\partial y_{12}} \frac{\partial y_{12}}{\partial k_{11}} + \frac{\partial E}{\partial y_{21}} \frac{\partial y_{21}}{\partial k_{11}} + \frac{\partial E}{\partial y_{22}} \frac{\partial y_{22}}{\partial k_{11}}$$

$$\therefore \frac{\partial E}{\partial k_{11}} = \frac{\partial E}{\partial y_{11}} x_{11} + \frac{\partial E}{\partial y_{12}} x_{12} + \frac{\partial E}{\partial y_{21}} x_{21} + \frac{\partial E}{\partial y_{22}} x_{22}$$

Similarly

$$\frac{\partial E}{\partial k_{12}} = \frac{\partial E}{\partial y_{11}} x_{12} + \frac{\partial E}{\partial y_{12}} x_{13} + \frac{\partial E}{\partial y_{21}} x_{22} + \frac{\partial E}{\partial y_{22}} x_{23}$$

$$\frac{\partial E}{\partial k_{21}} = \frac{\partial E}{\partial y_{11}} x_{21} + \frac{\partial E}{\partial y_{12}} x_{22} + \frac{\partial E}{\partial y_{21}} x_{31} + \frac{\partial E}{\partial y_{22}} x_{32}$$

$$\frac{\partial E}{\partial k_{22}} = \frac{\partial E}{\partial y_{11}} x_{22} + \frac{\partial E}{\partial y_{12}} x_{23} + \frac{\partial E}{\partial y_{21}} x_{32} + \frac{\partial E}{\partial y_{22}} x_{33}$$

This is cross-correlation between input matrix and the output gradient.

$$\begin{matrix} \frac{\partial E}{\partial K_{11}} & \frac{\partial E}{\partial K_{12}} \\ \vdots & \vdots \\ \frac{\partial E}{\partial K_{21}} & \frac{\partial E}{\partial K_{22}} \end{matrix} = \begin{bmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \\ x_{31} & x_{32} & x_{33} \end{bmatrix} \star \begin{matrix} \frac{\partial E}{\partial y_{11}} & \frac{\partial E}{\partial y_{12}} \\ \frac{\partial E}{\partial y_{21}} & \frac{\partial E}{\partial y_{22}} \end{matrix}$$

$$\therefore \frac{\partial E}{\partial K} = X \star \frac{\partial E}{\partial Y}$$

This is true only if  $Y = B + X \star K$ .

Now

$$\text{For } Y_i = B_i + \sum_{j=1}^n X_j \star K_{ij} \quad i = 1 \dots d$$

$$Y_1 = B_1 + X_1 \star K_{11} + \dots + X_n \star K_{1n}$$

$$Y_2 = B_2 + X_1 \star K_{21} + \dots + X_n \star K_{2n}$$

⋮

$$Y_d = B_d + X_1 \star K_{d1} + \dots + X_n \star K_{dn}$$

$$X \frac{\partial E}{\partial K_{21}} = \frac{\partial E}{\partial Y_1} \frac{\partial Y_1}{\partial K_{21}} + \frac{\partial E}{\partial Y_2} \frac{\partial Y_2}{\partial K_{21}} + \dots + \frac{\partial E}{\partial Y_d} \frac{\partial Y_d}{\partial K_{21}}$$

It is not valid, because  $Y$  &  $K$  are not scalar, there are matrices, so derivative of matrix is not valid.

$$\frac{d}{dx} (y_1 + y_2 + y_3) = \frac{\partial y_1}{\partial x} + \frac{\partial y_2}{\partial x} + \frac{\partial y_3}{\partial x}$$

Dots:  
Dots:

derivative of sum = sum of derivative

$$\frac{\partial E}{\partial K_{21}} = ?$$

Looking for  $K_{21}$  in forward propagation eq?

$K_{21}$  appears only in second eq<sup>n</sup> and since the derivative of sum is the sum of the derivative. As far as  $K_{21}$  is concerned  $y_{21}$  only boils down to

$$y_2 = B_2 + X_1 * K_{21}$$

So we can write

$$\frac{\partial E}{\partial K_{21}} = X_1 * \frac{\partial E}{\partial y_2}$$

In general,

$\frac{\partial E}{\partial K_{ij}} = X_j * \frac{\partial E}{\partial y_i}$
--

## ( Bias gradient)

(ii) For  $\frac{\partial E}{\partial \beta_i}$

As earlier, simplifying forward equation to

$$y_i = \beta_i + x_1 * k_{i1}$$

$$\frac{\partial E}{\partial b_i} = \frac{\partial E}{\partial y_{i1}} \frac{\partial y_{i1}}{\partial b_{i1}} + \frac{\partial E}{\partial y_{i2}} \frac{\partial y_{i2}}{\partial b_{i1}} + \frac{\partial E}{\partial y_{i1}} \frac{\partial y_{i1}}{\partial b_{i1}} + \frac{\partial E}{\partial y_{i2}} \frac{\partial y_{i2}}{\partial b_{i1}}$$

$$\frac{\partial E}{\partial b_i} = \frac{\partial E}{\partial y_{i1}}$$

Similarly,

$$\frac{\partial E}{\partial b_{i2}} = \frac{\partial E}{\partial y_{i2}}$$

$$\frac{\partial E}{\partial b_{i1}} = \frac{\partial E}{\partial y_{i1}}$$

$$\frac{\partial E}{\partial b_{i2}} = \frac{\partial E}{\partial y_{i2}}$$

$$\therefore \frac{\partial E}{\partial b_i} = \frac{\partial E}{\partial y_i}$$

This is simplified version of problem so,

$$\frac{\partial E}{\partial \beta_i} = \frac{\partial E}{\partial y_i}$$

In general

$\frac{\partial E}{\partial \beta_i} = \frac{\partial E}{\partial y_i}$
---

(en) for  $\frac{\partial E}{\partial x_j}$  (input gradient)

Working on simplified problem

$$\frac{\partial E}{\partial x_{11}} = \frac{\partial E}{\partial y_{11}} k_{11} + \frac{\partial E}{\partial y_{12}} k_{12} + \frac{\partial E}{\partial y_{21}} k_{21} + \frac{\partial E}{\partial y_{22}} k_{22}$$

$$\text{or } \frac{\partial E}{\partial x_{11}} = \frac{\partial E}{\partial y_{11}} k_{11}$$

$$\frac{\partial E}{\partial x_{12}} = \frac{\partial E}{\partial y_{11}} k_{12} + \frac{\partial E}{\partial y_{12}} k_{11}$$

$$\frac{\partial E}{\partial x_{21}} = \frac{\partial E}{\partial y_{12}} k_{12}$$

$$\frac{\partial E}{\partial x_{21}} = \frac{\partial E}{\partial y_{11}} k_{21} + \frac{\partial E}{\partial y_{21}} k_{11}$$

$$\frac{\partial E}{\partial x_{22}} = \frac{\partial E}{\partial y_{11}} k_{22} + \frac{\partial E}{\partial y_{12}} k_{12} + \frac{\partial E}{\partial y_{21}} k_{11} + \frac{\partial E}{\partial y_{22}} k_{21}$$

$$\frac{\partial E}{\partial x_{23}} = \frac{\partial E}{\partial y_{12}} k_{22} + \frac{\partial E}{\partial y_{22}} k_{12}$$

$$\frac{\partial E}{\partial x_{31}} = \frac{\partial E}{\partial y_{11}} k_{21}$$

$$\frac{\partial E}{\partial x_{32}} = \frac{\partial E}{\partial y_{12}} k_{22} + \frac{\partial E}{\partial y_{22}} k_{21}$$

$$\frac{\partial E}{\partial x_{33}} = \frac{\partial E}{\partial y_{22}} k_{22}$$

Page \_\_\_\_\_

It is the full cross-correlation between output gradient and rotated kernel

$$\begin{array}{|c|c|c|} \hline
 \frac{\partial E}{\partial x_{11}} & \frac{\partial E}{\partial x_{12}} & \frac{\partial E}{\partial x_{13}} \\ \hline
 \end{array}
 \quad
 \begin{array}{|c|c|c|} \hline
 \frac{\partial E}{\partial y_{11}} & \frac{\partial E}{\partial y_{12}} \\ \hline
 \end{array}
 \quad
 \begin{array}{|c|c|} \hline
 K_{22} & K_{21} \\ \hline
 K_{12} & K_{11} \\ \hline
 \end{array}$$
  

$$\begin{array}{|c|c|c|} \hline
 \frac{\partial E}{\partial x_{21}} & \frac{\partial E}{\partial x_{22}} & \frac{\partial E}{\partial x_{23}} \\ \hline
 \end{array}
 \quad
 \begin{array}{|c|c|c|} \hline
 \frac{\partial E}{\partial y_{21}} & \frac{\partial E}{\partial y_{22}} \\ \hline
 \end{array}
 \quad
 \begin{array}{|c|c|} \hline
 K_{11} \\ \hline
 \end{array}$$
  

$$\begin{array}{|c|c|c|} \hline
 \frac{\partial E}{\partial x_{31}} & \frac{\partial E}{\partial x_{32}} & \frac{\partial E}{\partial x_{33}} \\ \hline
 \end{array}
 \quad
 \begin{array}{|c|c|c|} \hline
 \frac{\partial E}{\partial y_{31}} & \frac{\partial E}{\partial y_{32}} \\ \hline
 \end{array}
 \quad
 \begin{array}{|c|c|} \hline
 K_{21} \\ \hline
 \end{array}$$

$$\frac{\partial E}{\partial x} = \frac{\partial E}{\partial y} \star_{full} \text{rot 180}(k)$$

Also we can say

$$\frac{\partial E}{\partial x} = \frac{\partial E}{\partial y} \star_{full} K$$

This is simplified version.

For actual forward equation.

$$\frac{\partial E}{\partial x_1} = \frac{\partial E}{\partial y_1} \star_{full} K_{11} + \dots + \frac{\partial E}{\partial y_n} \star_{full} K_{n1}$$

$$\therefore \frac{\partial E}{\partial x_j} = \frac{\partial E}{\partial y_1} \star_{full} K_{1j} + \dots + \frac{\partial E}{\partial y_n} \star_{full} K_{nj}$$

$$\boxed{\frac{\partial E}{\partial x_j} = \sum_{i=1}^n \frac{\partial E}{\partial y_i} \star_{full} K_{ij} \quad j = 1 \dots n}$$

from scipy import signal.

class Convolution(Layer):

def \_\_init\_\_(self, input\_shape, kernel\_size, kernel\_number):

input\_depth, input\_height, input\_width = input\_shape

self.depth

self.kernel\_number = kernel\_number

self.input\_shape = input\_shape

self.output\_shape = (kernel\_number, input\_height - kernel\_size +  
input\_width - kernel\_size + 1)

self.kernels\_shape = (kernel\_number, input\_depth,  
kernel\_size, kernel\_size)

self.kernels = np.random.randn(\*self.kernels\_shape)

self.biases = np.random.randn(\*self.output\_shape)

def forward(self, input):

self.input = input

self.output = np.copy(self.biases)

for i in range(self.kernel\_number):

for j in range(self.input\_depth):

self.output[i] += signal.correlate2d(self.input[i],

self.kernels[i, j], "valid")

return self.output.

`def backward(self, output-gradient, learning-rate):`

`kernels-gradient = np.zeros(self.kernels-shape)`

`input-gradient = np.zeros(self.input-shape)`

`for i in range(self.depth):`

`for j in range(self`

`for i in range(self.kernel-number):`

`for j in range(self.input-depth):`

`kernels-gradient[i, j] = signal.correlate2d(`

`self.input[j], output-gradient[i], "valid")`

`input-gradient[j] += signal.convolve2d(`

`output-gradient[i], self.kernels[i, j], "full")`

# bias gradient should not be calculated as it is equal to

# output gradient given.

`self.kernels -= learning-rate * kernels-gradient`

`self.biases -= learning-rate * output-gradient`

`return input-gradient.`

### 3) Reshape Layer

It is needed because the output of the convolution layer is 3d block and at the end of network we use dense layers which takes in a column vector as input. Therefore we need to reshape data. Hence the reshape layer

```
import numpy as np
from layer import Layer
```

```
class Reshape(Layer):
    def __init__(self, input_shape, output_shape):
        self.input_shape = input_shape
        self.output_shape = output_shape
```

```
    def forward(self, input):
        return np.reshape(input, self.output_shape)
```

```
    def backward(self, output_gradient, learning_rate):
        return np.reshape(output_gradient, self.input_shape)
```

## 1) Binary Cross Entropy

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_i \end{bmatrix} \rightarrow \text{Predicted}$$

$y_i \in \{0, 1\}$

$$\hat{y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_i \end{bmatrix} \rightarrow \text{Actual}$$

$$E = -\frac{1}{n} \sum_{i=1}^n \hat{y}_i \log(\hat{y}_i) + (1-\hat{y}_i) \log(1-\hat{y}_i)$$

$$\frac{\partial E}{\partial y_i} = \frac{\partial}{\partial y_i} \left( -\frac{1}{n} \sum_{i=1}^n \hat{y}_i \log(\hat{y}_i) + (1-\hat{y}_i) \log(1-\hat{y}_i) \right)$$

$$= \frac{\partial}{\partial y_i} \left( -\frac{1}{n} \left( \hat{y}_i \log(\hat{y}_i) + (1-\hat{y}_i) \log(1-\hat{y}_i) \right) \right) - \dots$$

$$\dots - \frac{1}{n} \left( \hat{y}_n \log(\hat{y}_n) + (1-\hat{y}_n) \log(1-\hat{y}_n) \right)$$

$$= \frac{\partial}{\partial y_i} \left( -\frac{1}{n} \left( \hat{y}_i \log(\hat{y}_i) + (1-\hat{y}_i) \log(1-\hat{y}_i) \right) \right)$$

$$= -\frac{1}{n} \left( \frac{\hat{y}_i}{\hat{y}_i} - \frac{1-\hat{y}_i}{1-\hat{y}_i} \right)$$

$$= \frac{1}{n} \left( \frac{1-\hat{y}_i}{1-\hat{y}_i} - \frac{\hat{y}_i}{\hat{y}_i} \right)$$

$$\frac{\partial E}{\partial y_i} = \frac{1}{n} \left( \frac{1-y_i}{1-\hat{y}_i} - \frac{y_i}{\hat{y}_i} \right)$$

import numpy as np.

```
def binary_crossentropy(y_true, y_pred):  
    return -np.mean(y_true * np.log(y_pred) +  
                    (1-y_true) * np.log(1-y_pred))
```

```
def binary_crossentropy_prime(y_true, y_pred):  
    return ((1-y_true)/(1-y_pred) - y_true/y_pred) / np.size(y_true)
```

$$(\hat{y} - y) \log(\hat{y}) + (1 - \hat{y}) \log(1 - \hat{y}) \quad b = 36$$

$$(\hat{y} - y) \log(\hat{y}) + (1 - \hat{y}) \log(1 - \hat{y}) \quad b = 36$$

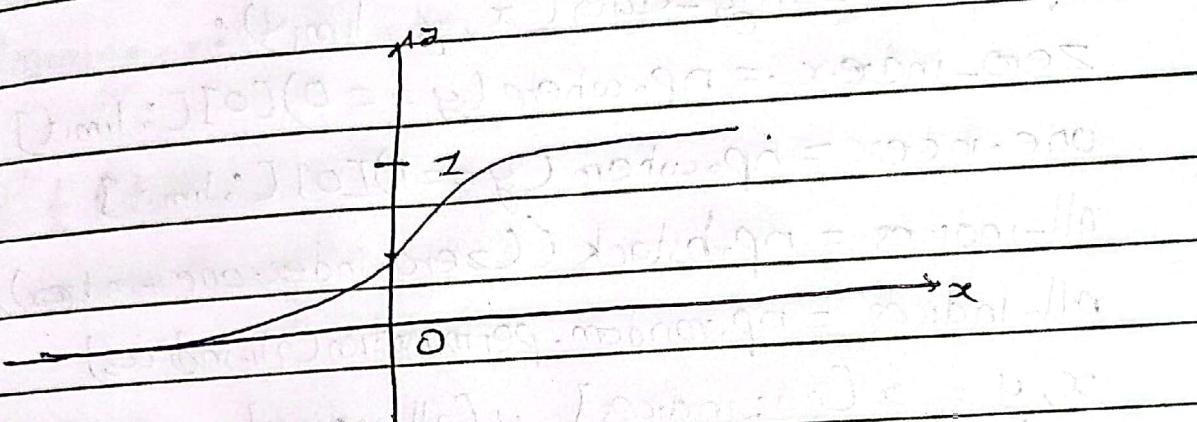
$$((\hat{y} - y) \log(\hat{y}) + (1 - \hat{y}) \log(1 - \hat{y})) / n \quad b = 36$$

$$((\hat{y} - y) \log(\hat{y}) + (1 - \hat{y}) \log(1 - \hat{y})) / n \quad b = 36$$

$$((\hat{y} - y) \log(\hat{y}) + (1 - \hat{y}) \log(1 - \hat{y})) / n \quad b = 36$$

## (5) Sigmoid Activation

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



$$\sigma'(x) = \frac{e^{-x}}{(1 + e^{-x})^2}$$

$$\sigma'(x) = \sigma(x)(1 - \sigma(x))$$

```
import numpy as np
```

```
from activation import Activation
```

```
class Sigmoid(Activation):
```

```
    def __init__(self):
```

```
        def sigmoid(x):
```

```
            return 1 / (1 + np.exp(-x))
```

```
    def sigmoid_prime(self):
```

```
        s = sigmoid(x)
```

```
        return s * (1 - s)
```

```
    super().__init__(sigmoid, sigmoid_prime)
```

Applying CNN in mnist. for for only two classes

```
# import all necessary modules made earlier
from keras.datasets import mnist
from keras.utils import np_utils

def preprocessing_data(x, y, limit):
    zero_index = np.where(y == 0)[0][:limit]
    one_index = np.where(y == 1)[0][:limit]
    all_indices = np.hstack((zero_index, one_index))
    all_indices = np.random.permutation(all_indices)
    x, y = x[all_indices], y[all_indices]
    x = x.reshape(len(y), 1, 28, 28)
    x = x.astype("float32") / 255
    y = np_utils.to_categorical(y)
    y = y.reshape(len(y), 2, 1)
    return x, y
```

(x-train, y-train), (x-test, y-test) = mnist.load\_data()  
 x-train, y-train = preprocessing\_data(x-train, y-train, 100)  
 x-test, y-test = preprocessing\_data(x-test, y-test, 100)

network = [  
 Convolutional((1, 28, 28), 3, 5),  
 Sigmoid(),  
 Reshape((5, 26, 26), (5 \* 26 \* 26, 1)),  
 Dense(5 \* 26 \* 26, 100),  
 Sigmoid(),

Dense(100, 2),  
Sigmoid()

epochs = 20

learning-rate = 0.1

fit(network, binary-cross-entropy, binary-cross-entropy-pr,  
x-train, y-train, epochs, learning-rate)

# test

for x, y in zip(x-test, y-test):

~~output = x~~

for layer in network:

~~output = layer.forward(output)~~

print

for x, y in zip(x-test, y-test):

~~output = predict(network, x)~~

print(f "pred: {np.argmax(output)},  
true: {np.argmax(y)}")