

CNN

=

① Image

② Filter | Kernel

③ Convolution + ReLU

④ Pooling

⑤ Flatten (2D -> 1D)

⑥ FC

Different Stages/Operations of CNN:

- Image = Collection of pixels
- Filter/Kernel
- Convolution(Matrix multiplication b/w image and filter) + ReLU(Why we add ReLU will be discussed)
- Pooling
- Flatten(2D to 1D)
- Fully connected Neural Network

Convolution means feature(pixels/kernels) extraction. For Ex: Given the image of person, we attempt to extract the features or pixels representing the hair followed by, classifying the type of the hair (Spikey, straight, curly, normal etc)

#

MATHEMATICAL
IMPLEMENTATION OF CNN
ON GRAY IMAGE(Black=0
and White=255)

256

Image

= [0-255]

R, G, B

4*

- Multiplication of image and filter is normal multiplication of corresponding pixels and filters and not matrix multiplication.
- Filter is having some window size(3X3 here). This window traverses through the whole image and perform multiplication in the respective window. During traversal filter performs jump in order to capture other pixel in next iteration which is known as STRIDE.
- Filters are square matrix(we may experiment with non square matrix) with preferably odd size?
Like: 1X1, 3X3, 5X5, 7X7 etc
- Based on how much information or pixels we want to collect in each stride we are deciding the size of the filter.

Img		
0	0	0
0	0	0
0	0	0
255	255	255
255	255	255
255	255	255

*	-1	-1	-1
*	0	0	0
*	1	1	1

3x3

$$(0 \times -1) + 0 \times (-1) + 0 \times (-1) + 0 \times 0 = 0$$

*0+0+0+0+0+0+0+0+0+0

0	0	0	0
255	255	255	255
255	255	255	255
0	0	0	0

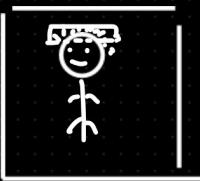
Feature map

Feature map

Here Feature Map is created by moving Stride of 1 along horizontal and vertical traversal of the image

Convolution \Rightarrow (1) extract the feature from img.

many
age and
on) as we
is
r



Feature extraction

So, Yann LeCun in his LENET architecture has created these filters randomly(not hard coded) and then applied back propagation to find the optimal filter weights in order to perform detection of the respective objective during training. This means that in LENET FILTERS become the learnable parameter. This why this become famous and treated as first official CNN architecture.

Whereas, CNN architectures before LENET like Neocognitron uses hard coded filters and there was no back propagation to embed the self learning of these filters.

0	0	0	0	0	0
0	0	0	0	0	0
0	0	0	0	0	0
0	0	0	0	0	0
0	0	0	0	0	0

This is the hard coded filter found by Scientists that can perform the Horizontal Edge Detection

$$\begin{array}{|c|c|c|} \hline 0 & 0 & 6 \\ \hline 6 & 0 & 6 \\ \hline 6 & 0 & 0 \\ \hline \end{array} \times \begin{array}{|c|c|c|} \hline -1 & -1 & -1 \\ \hline 0 & 0 & 0 \\ \hline 1 & 1 & 1 \\ \hline \end{array}$$

0	0	0	0
45	45	115	115
45	225	115	225
0	0	0	0

feature map

$$4 \times 4 =$$



0x1

Kernel | filter \rightarrow Dimension = 1×1 , 2×2 , 3×3 , 4×4 , 5×5 , 6×6 , 7×7

Initially $\Rightarrow \left(\frac{1980 - 1990}{-} \right)$



Scientist \Rightarrow Image

Hardcoded

-Odeł.

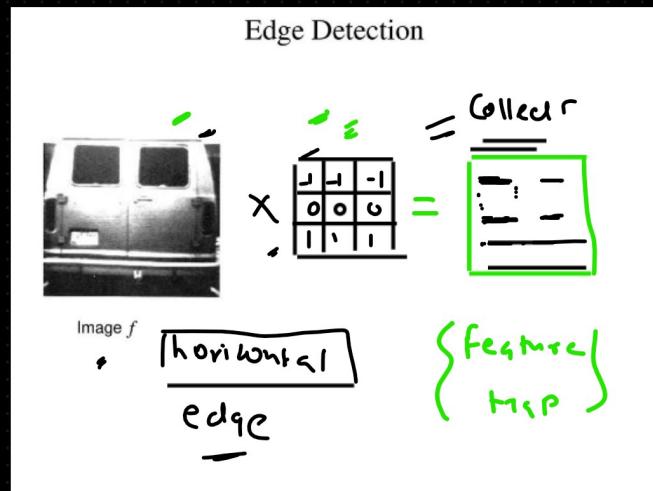
Aleksander

Wojciech Ig

Roman

Englemer

$$\begin{array}{r}
 * \overline{3} \times 3 \\
 = \\
 \overline{\overline{5} \times 5} \\
 = \\
 \overline{\overline{7} \times 7} \\
 = \\
 9 \times 9 \\
 = \\
 \overline{\overline{11} \times 11}
 \end{array}$$



$$\begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline \end{array}$$

Hard Coded Filter for Vertical Edge Detection
already discovered by scientists through extensive research

= Vertical edge detection

Viola-Jones Object Detection Framework is used for determining the hard coded filters for object detection

Filter \Rightarrow For extracting the feature from the image

$$\begin{array}{|c|c|c|} \hline -1 & -1 & -1 \\ \hline 0 & 0 & 0 \\ \hline 1 & 1 & 1 \\ \hline \end{array} \quad \begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline \end{array} \quad \begin{array}{|c|c|c|} \hline 1 & -1 & 0 \\ \hline -1 & 1 & 0 \\ \hline -1 & 0 & 0 \\ \hline \end{array}$$

CNN =
E-NET

\Rightarrow filter Architecture in CNN

$$\begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array}$$

F.P. \longrightarrow

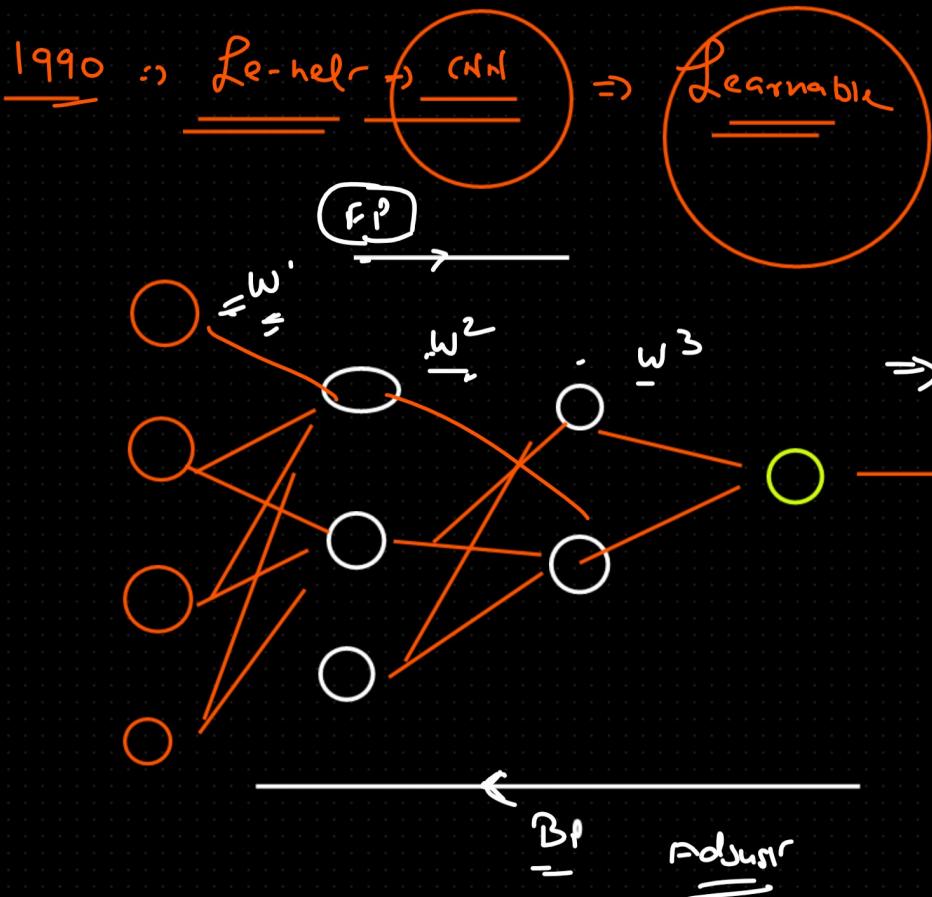
$$\begin{array}{|c|c|c|} \hline R & R & R \\ \hline L & L & L \\ \hline R & R & R \\ \hline \end{array}$$

*

B.P.

Learn

$$\begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline 4 & 5 & 6 \\ \hline 7 & 8 & 9 \\ \hline \end{array}$$



NN \Rightarrow it is learning by itself

it is learning based on data

Image Size	Filter Size	Convolved Image/Feature Map Size
<u>6×6</u>	\ast	<u>3×3</u> = <u>4×4</u>
<u>8×8</u>	\ast	<u>3×3</u> = <u>6×6</u>
<u>28×28</u>	\ast	<u>3×3</u> = <u>26×26</u>

General formula for determining the shape of convoluted image = $(n-m+1) \times (n-m+1)$ where n and m represents the dimension of Image and Filters respectively. Please note that here we have taken the consideration that all 3(image, filter and convoluted image) is in square nature (rows=columns)

$$\underline{n \times n} \ast \underline{m \times m} \Rightarrow \underline{(n-m+1) \times (n-m+1)}$$

$$6 \times 6 \ast 3 \times 3 = (6-3+1) \times (6-3+1)$$

$$= (3+1) \times (3+1)$$

$$= \boxed{4 \times 4}$$

$$\frac{\text{img_size} - \text{filter} + 1}{\text{strides}}$$

$$[64 \times 64] * 5 \times 5$$

=

$$= \frac{64 - 5 + 1}{5} = 60 \times 60$$

$$(64 - 5 + 1) \\ (59 + 1) = 60$$

$$64 \times 64 = \underline{\underline{60 \times 60}}$$

Padding

Stride x
Padding

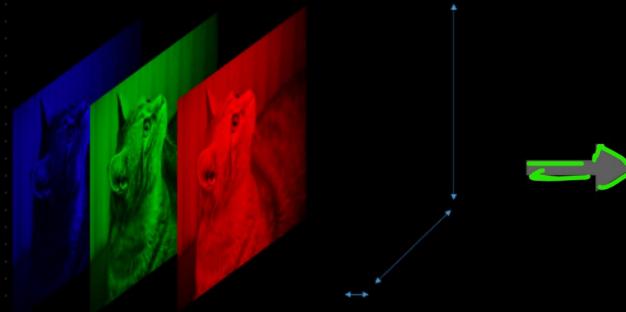
- Every feature
low level feature



MATHEMATICAL IMPLEMENTATION OF CNN ON COLORED IMAGE(3 Channels RGB):



Please note that for Gray scale image we will have image, filter and convoluted image in 2 dimensions. Whereas, for Colored Image we will be having Image and Filter in 3 dimensions and convoluted image in 2 dimensions



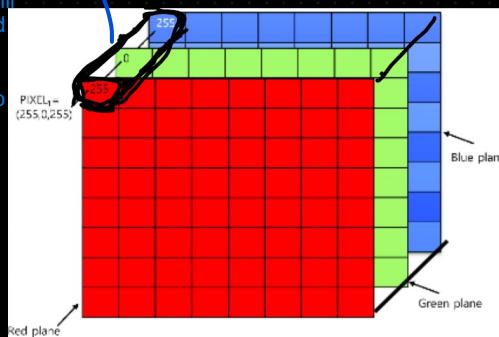
.392	.482	.576
.478	.55	.169
.580	.7	.263
.373	.60	.376
.443	.569	.674

Color \Rightarrow R, G, B

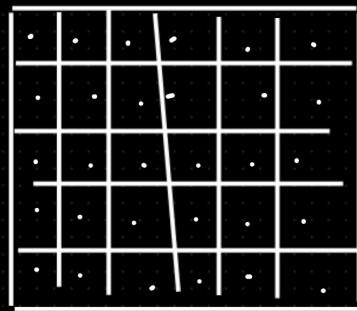


In convoluted image if we are getting pixel value more than 255 (image and filter multiplication) then round it to 255 since we have pixel value up to maximum 255 (Range is [0.255]). Same happens in the lower bound as well)

Third dimension (depth RGB) will always be fixed in each stride that's why we are having Two dimensions in the convoluted image for colored image



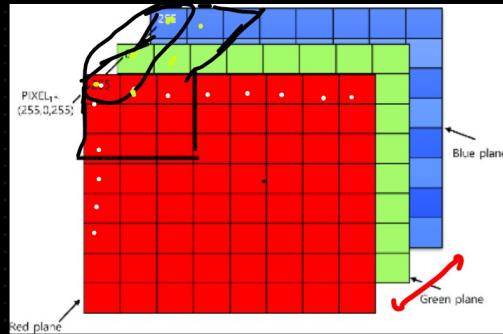
gray Scale image



Filter size is also the hyperparameter



In convolution we are loosing the data. So in case we do not want to lose the original size of the data then we can use PADDING



$(\underline{8 \times 8 \times 3})$

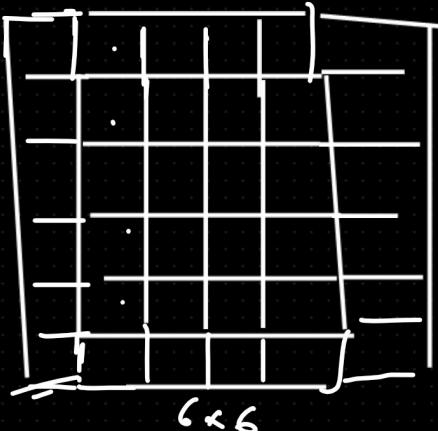
color

*

$(\underline{3 \times 3 \times 3})$

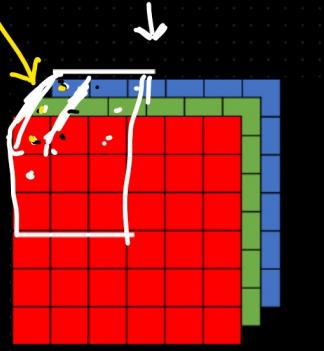
$(\underline{n - m + 1})$

=



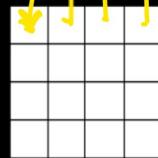
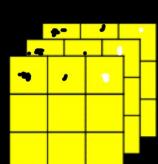
$\underline{f \times m}$

Single Channel



$6 \times 6 \times 3$

$$P_1 + P_2 + P_3 + P_4 + P_5 + P_6 + P_7 + P_8 + P_9$$



Here also for each stride multiply each corresponding image pixel with corresponding filter weight to get respective pixel for convoluted image

$3 \times 3 \times 3$



4×4

Practical

We can apply as much as filter as we want (where each feature is performing different detection or extraction) on top of image as number of filter applied is a hyperparameter and at the end we are basically stacking up all the convoluted images formed with respective filters which is called the Feature Map.



image \times filters

"

"

"

"

"

"

"

"

filter 1 \rightarrow vertical edge

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

X

Image \times filters

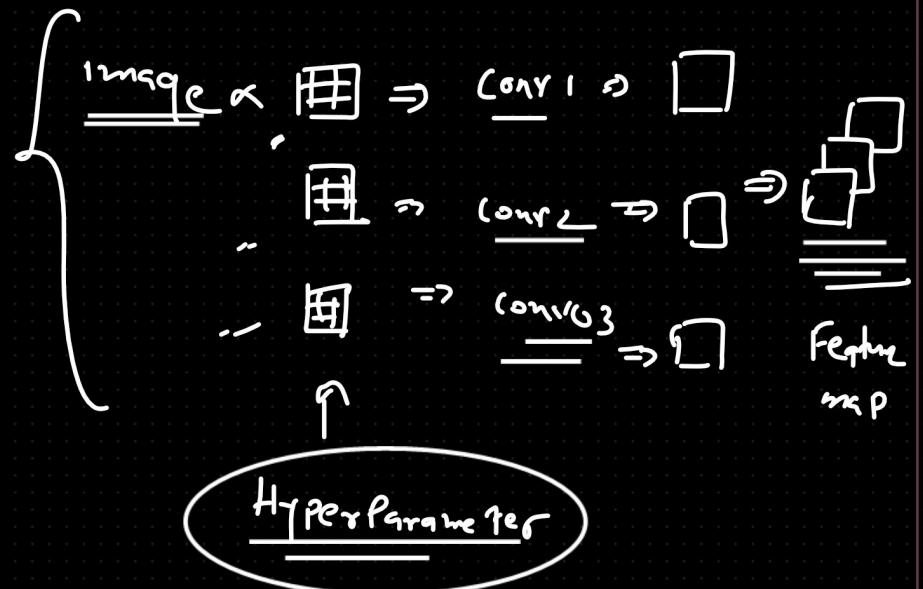
↓

Extract information from img

Image \Rightarrow

{
 side edge
 horizontal edge
 clock
}

filter2 Horizontal edges



Padding | Stride | Pooling

Padding

Convert a original image of 5×5 image into 7×7 .

If we fill the extended pixels using Padding with 0 then it is called Zero padding.

Motive if to prevent the data loss by retaining the original image size in convoluted image formed by multiplying respective filter with the padded image

image \times filter \Rightarrow convolving

$$6 \times 6 \quad 3 \times 3 \Rightarrow 4 \times 4$$

$$n-p+1 \Rightarrow 6-3+1$$

$$3+1 = 4$$

Here we have done zero padding on 5×5 image ad converted it into 7×7 image

5x5

padding

0	0	0	0	0	0	0
0	+	2	3	3	8	0
0	1	5	3	8	9	0
0	3	3	2	8	1	6
0	2	8	7	2	7	0
0	1	9	4	5	4	0
0	0	0	0	0	0	0

*

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

=

6	.	.
.	.	.
.	.	.

3×3

5×3 ?

$$7 \times 1 + 4 \times 1 + 3 \times 1 + 2 \times 0 + 5 \times 0 + \\ 3 \times 0 + 3 \times -1 + 3 \times -1 + 2 \times -1 = 6$$

Convolved image will be 3×3 if filter multiplied on image without padding but if multiplied with padded image then will get convoluted image of shape 5×5 which is similar to the original shape of the image.

$$n \times n = \boxed{5 \times 5}$$

$$\begin{array}{l} f \times f \\ 3 \times 3 \end{array} = \frac{n-f+1}{5-3+1} = \boxed{3 \times 3}$$

$$\underline{5 \times 5 \Rightarrow 5 \times 5 ?}$$

$$\underline{\underline{n-f+1 = 5}}$$

$$n-f+1$$

$$7-3+1$$

$$4+1 = \boxed{5}$$

$$h = 5-1+f$$

$$= 5-1+3$$

$$h = 7$$

$$\underline{\underline{n=7}} \Rightarrow \boxed{5}$$



Generalizing this formula further using padding

$$\boxed{5 \times 5}$$

$$\underline{\underline{n-f+1}}$$

$$5-3+1 = 3 \Rightarrow \boxed{3 \times 3}$$

$$\underline{\underline{P=0}}$$

$$\begin{array}{l} n+2 \times 0 - f + 1 \\ (n-f+1) \end{array}$$

$$\underline{\underline{n+2 \times 1 - f + 1}}$$

$$\underline{\underline{P=1}}$$

$$5+2(1)-3+1$$

$$= 5+2-3+1 = \boxed{5}$$

Here P represents the number of pixels padded on each corner of the original image. If P=0 means no padding hence dimension of convoluted image will be n-f+1. Whereas, if P=1 means we will pad the original image with 1 pixel on all 4 corners. In such a case dimension of convoluted image will be given as n+2P-f+1.

Jump from one pixel to another pixel ↪

~~#~~ Stride ⇒ Jump is called Stride

1	2	3	4	5	6	7
11	12	13	14	15	16	17
21	22	23	24	25	26	27
31	32	33	34	35	36	37
41	42	43	44	45	46	47
51	52	53	54	55	56	57
61	62	63	64	65	66	67
71	72	73	74	75	76	77

Convolve with 3×3 filters filled with ones

108	126
288	306

1 Jump ⇒ Stride

while we are iterating over the ring

Stride = 1

$$\frac{n-f+1}{s} \Rightarrow \left(\frac{n-f+1}{s} \right) \underset{s=1}{=} \Rightarrow \left(\frac{n-f+1}{1} \right) \circled{n-f+1}$$

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

$$6 \times 6 \quad 3 \times 3 \Rightarrow \left(\frac{6-3+1}{1} \right) \Rightarrow 6-3+1 = 4$$

4×4

Same as above this is the generalized formula for the dimension convoluted image including both Stride and Padding

Ques: When to use padding?

Ans: When in convolved image data loss as compared to the dimensions of the original image is very huge

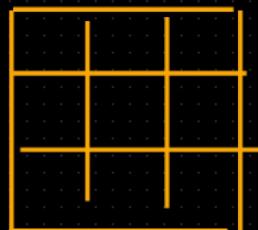
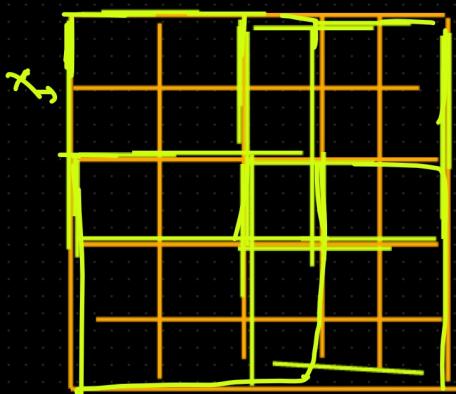
$$\underline{\underline{6 \times 6}} \quad \underline{\underline{3 \times 3}} \Rightarrow \frac{6 + 2 \times 0 - 3}{2} + 1$$

$$\left\{ \begin{array}{l} P=0 \\ S=2 \end{array} \right\} \Rightarrow \frac{6 + (-3)}{2} + 1 \Rightarrow \frac{6-3}{2} + 1 = \frac{3}{2} + 1$$

$$= 1.5 + 1 = 2.5$$

$$\downarrow \\ \underline{\underline{2 \times 2}}$$

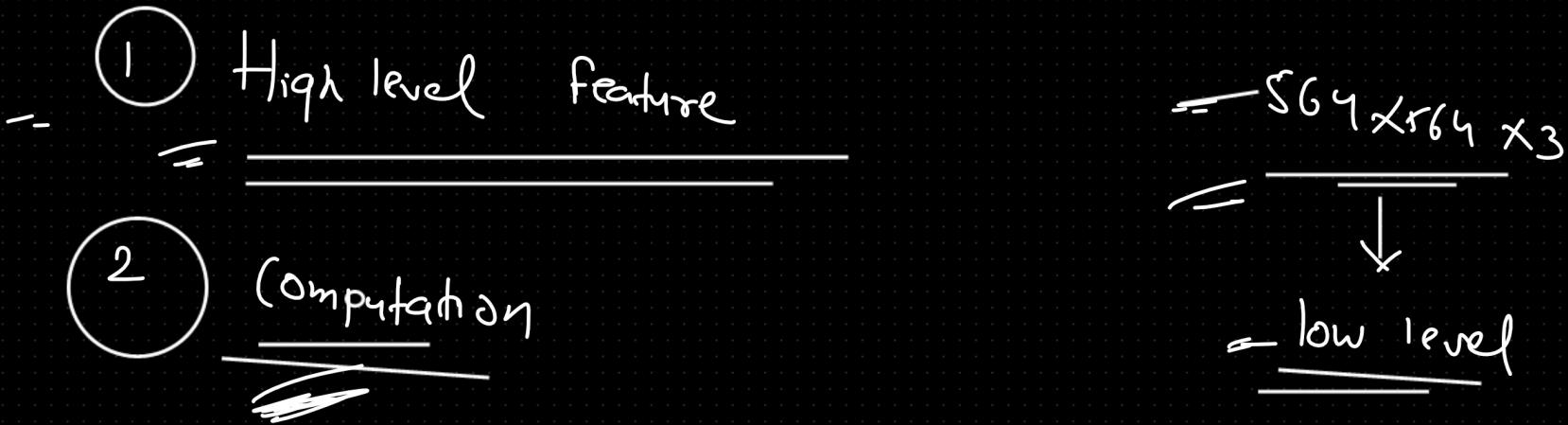
$$6 \times 6 \rightarrow 2 \times 2$$



$$\Rightarrow \underline{\underline{2 \times 2}}$$

$$\underline{\underline{5 \times 5}} \quad \underline{\underline{3 \times 3}} \Rightarrow$$

$$\left(\frac{n+2P-f}{S} + 1 \right) \Rightarrow \left(\frac{5+2 \times 0 - 3}{2} \right) + 1 \Rightarrow \left(\frac{5-3}{2} \right) + 1 = \underline{\underline{2}}$$



$$\left(\frac{n+2p-f}{s} + 1 \right) \Rightarrow \underline{\text{Residual-image}}$$

h = Size of image

~~P~~ = Padding -

\equiv $f =$ filter size

S = static value

$S = 2, 3, 4, 5$

Ques: In CNN why we need ReLU?

Ans: To remove the -ve value or in other words to capture the non-linear relationship as well.

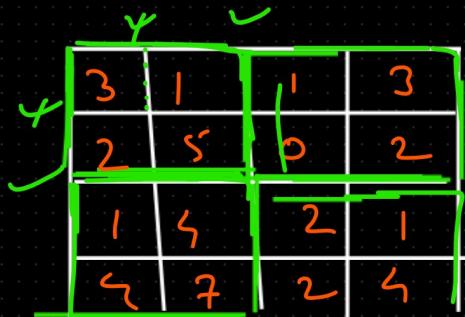
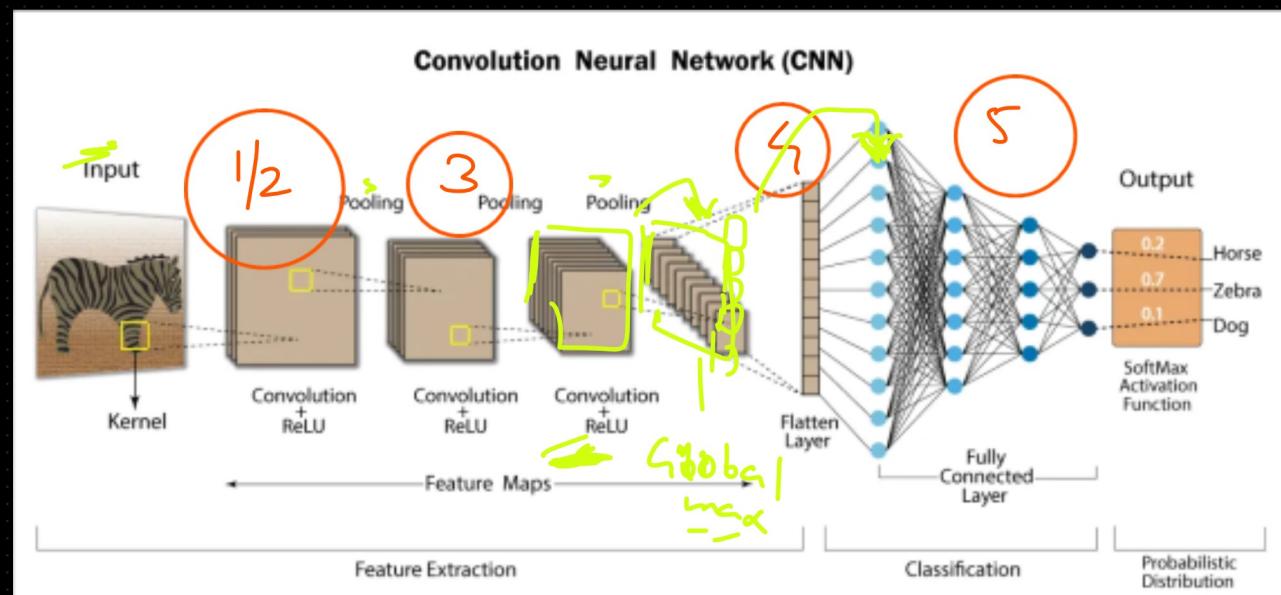
1 Convolution

2 ReLU (to remove -ve value)

3 Pooling

4 Flatten

5 FC



Feature map

Convolved image

⇒ Pooling (I want to Pool the imp feature)

1 max Pooling

2 min Pooling

3 average Pooling

4 global Pooling

Size = 2×2

Stride = 2

Type = max

For performing pooling also we take pooling window and stride

↓
image × filter

5	3
7	4

2×2

specific feature

intensity of image

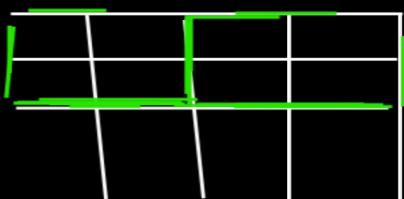
This is the Feature Extracted for the above Feature map using Max Pooling where in each widow of respective pool size we are picking up the one having larger value. Similarly, in the case of Min. Pooling we will pick up the smallest value in each window whereas, in the case of Average Pooling we will simply take the average of all the pixels in each window

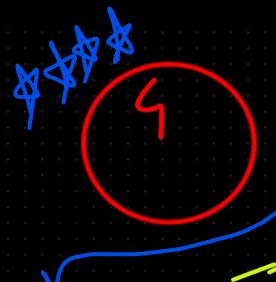
filter

$$\text{Input: } 228 \times 228 \times 3 \quad * \quad \frac{3 \times 3 \times 3}{(100)} = \frac{228 \times 228 \times 100}{(100)} \Rightarrow (2 \times 2) \rightarrow 2 \Rightarrow 113 \times 113 \times 100$$

Below are the reasons behind performing Pooling:

- ① reduce the size
- ② computational cost
- ③ focus on the specific feature

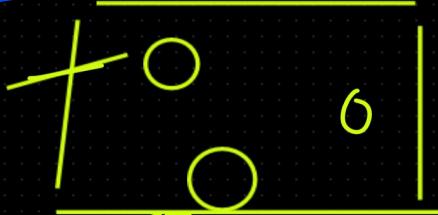




Translation invariance:
Dog is in left during training whereas, when testing on image where Dog is in left then without pooling (selecting important feature) then we will get the misleading results that's we are performing pooling

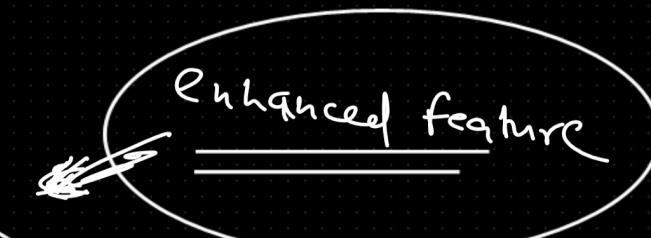
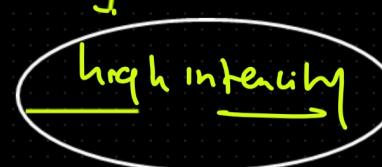
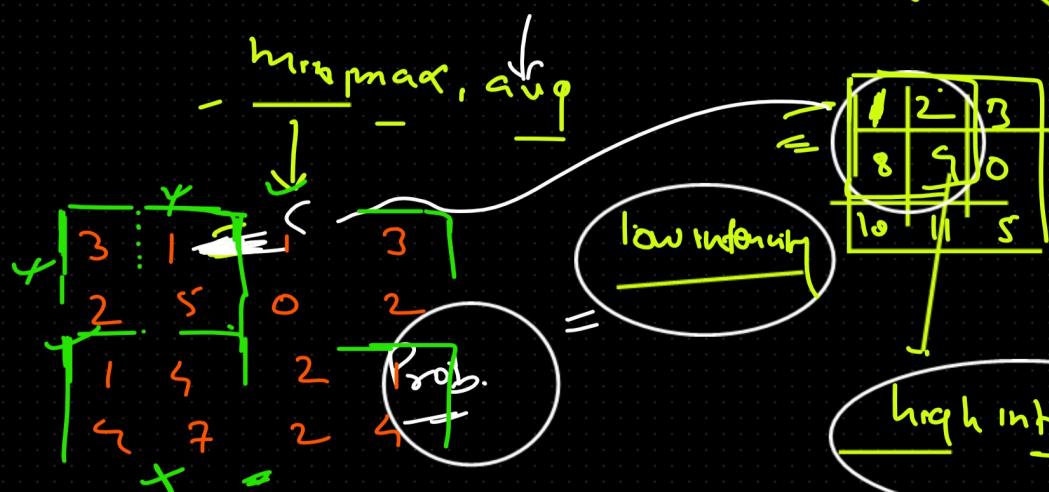
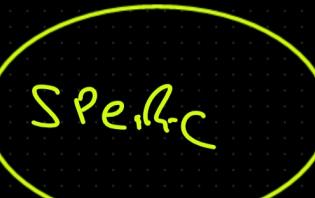
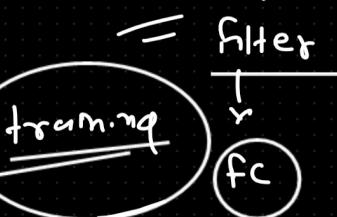


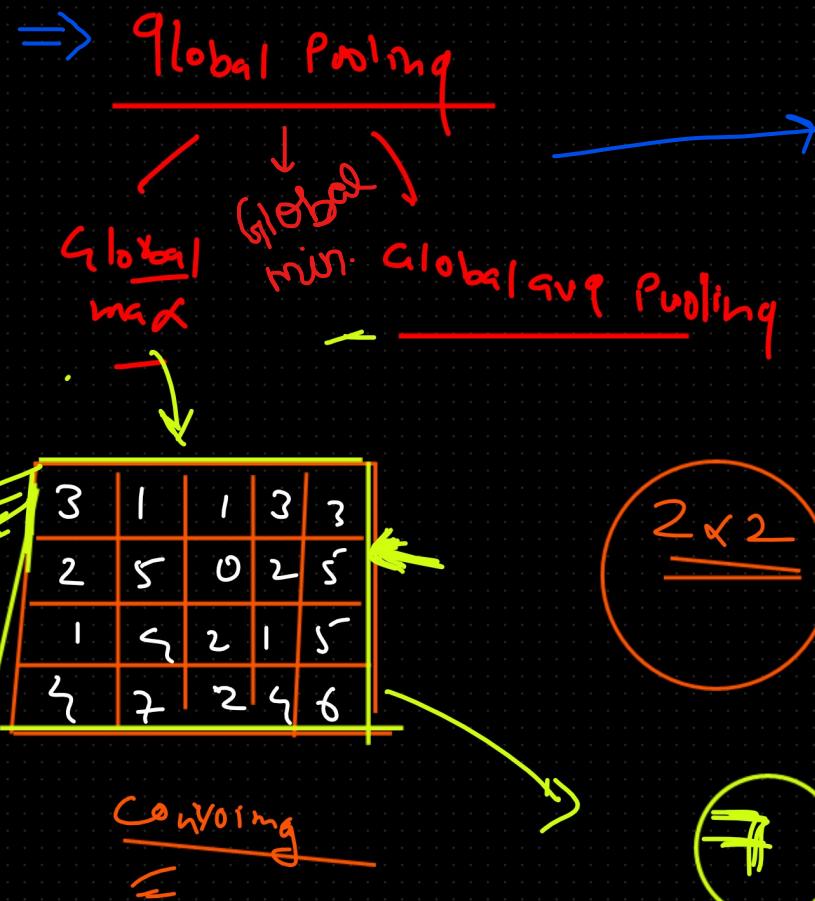
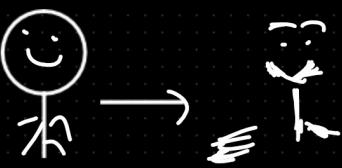
Section where Dog is present (irrespective of left or right position), will be having High intensity value whereas, the background of the Dog will be having Low intensity value. So, for such a use case we can easily pool out important features by using the Max Pooling approach.



Ques: Like how machine knows which are the important features. How machine makes that distinction?

Ans: That's why we have different variants of Pooling which will be implemented as per use case, but mostly Max pooling is used frequently for the selection of the important features.





2×2

7

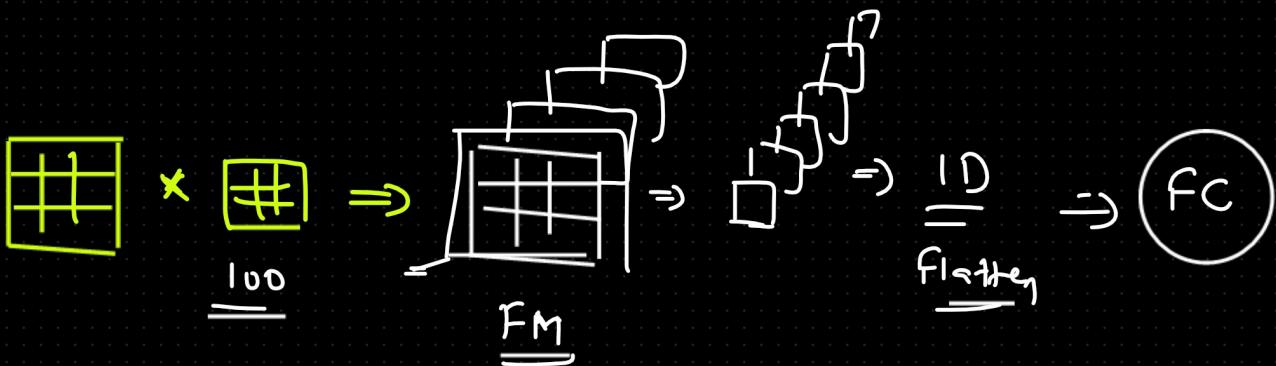
Scalar value

Global max we are converting the whole into single scalar which is maximum of the entire feature map

Whereas, in Global avg we are taking the average of entire feature map and form a single scalar

Similarly, Global min .. so on

\max, \min, avg



Global max pooling

Intense
max

low
intensity
↓
binary