

Computer Vision problems

Image classification : Cat? 0/1

Object detection :- Identify cars in a given image

Neural Style transfer :- Some face image + some style (scenery) painting image



Face embedded in scenery

Convolution operation

In a given image - we would like to find edges
(horizontal as well as vertical edges)

Eg

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

Convolution operation



1	0	-1	$\begin{bmatrix} -5 & -4 & 0 \\ -10 & -2 & 2 \\ 0 & -2 & -4 \\ -3 & -2 & -3 \end{bmatrix}$
1	0	-1	4×4
1	0	-1	

3×3

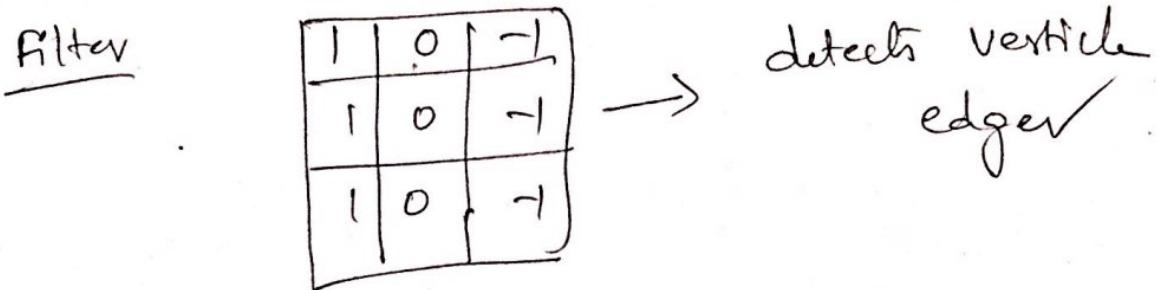
filter (or)

kernel

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

Summary

$$6 \times 6 \text{ image} * 3 \times 3 \text{ kernel or filter} = 4 \times 4 \text{ image}$$



Example

6×6 image

$$\begin{array}{|c|c|c|c|c|c|} \hline 10 & 10 & 10 & 0 & 0 & 0 \\ \hline 10 & 10 & 10 & 0 & 0 & 0 \\ \hline 10 & 10 & 10 & 0 & 0 & 0 \\ \hline 10 & 10 & 10 & 0 & 0 & 0 \\ \hline 10 & 0 & 0 & 0 & 0 & 0 \\ \hline 10 & 10 & 10 & 0 & 0 & 0 \\ \hline \end{array} * \begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline \end{array} = \begin{array}{|c|c|c|c|} \hline 0 & 30 & 30 & 0 \\ \hline 0 & 30 & 30 & 0 \\ \hline 0 & 30 & 80 & 0 \\ \hline 0 & 30 & 30 & 0 \\ \hline \end{array}$$

3×3

\downarrow

\downarrow

\downarrow

\downarrow

More on Verticle edge detection

$$\text{Eg} \quad \begin{pmatrix} 0 & 0 & 0 & 10 & 10 & 10 \\ 0 & 0 & 0 & 10 & 10 & 10 \\ 0 & 0 & 0 & 10 & 10 & 10 \\ 0 & 0 & 0 & 10 & 10 & 10 \\ 0 & 0 & 0 & 10 & 10 & 10 \\ 0 & 0 & 0 & 10 & 10 & 10 \end{pmatrix} * \begin{pmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{pmatrix} = \begin{pmatrix} 0 & -30 & -30 & 0 \\ 0 & -30 & -30 & 0 \\ 0 & -30 & -30 & 0 \\ 0 & -30 & -30 & 0 \end{pmatrix}$$

$$\begin{pmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{pmatrix}$$

detect Verticle edge

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

detect horizontal edge

$$\text{Eg} \quad \begin{pmatrix} 10 & 10 & 10 & 0 & 0 & 0 \\ 10 & 10 & 10 & 0 & 0 & 0 \\ 10 & 10 & 10 & 0 & 0 & 0 \\ 0 & 0 & 0 & 10 & 10 & 10 \\ 0 & 0 & 0 & 10 & 10 & 10 \\ 0 & 0 & 0 & 10 & 10 & 10 \end{pmatrix} * \begin{pmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{pmatrix} = \begin{pmatrix} 0 & 0 & 0 & 0 \\ 30 & 10 & -10 & -30 \\ 30 & 10 & -10 & -30 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

6x6

Note: for large images (other than 6x6) \rightarrow There filter will find horizontal edges more accurately

~~Learn~~

learning to detect edges

①

$$\begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix}$$

→ called Sobel filter

→ to detect

Vertical
edge in

Literature

②

$$\begin{bmatrix} 3 & 0 & -3 \\ 10 & 0 & -10 \\ 3 & 0 & -3 \end{bmatrix}$$

→ Scharr filter

Can we learn a new filter?

Aus

Yes

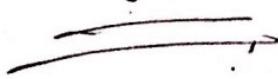
AI can learn

$$\begin{bmatrix} w_1 & w_2 & w_3 \\ w_4 & w_5 & w_6 \\ w_7 & w_8 & w_9 \end{bmatrix}$$

, this kind

of filter by using backpropagation

- which may find not only horizontal, & vertical
also can find tilted edges.



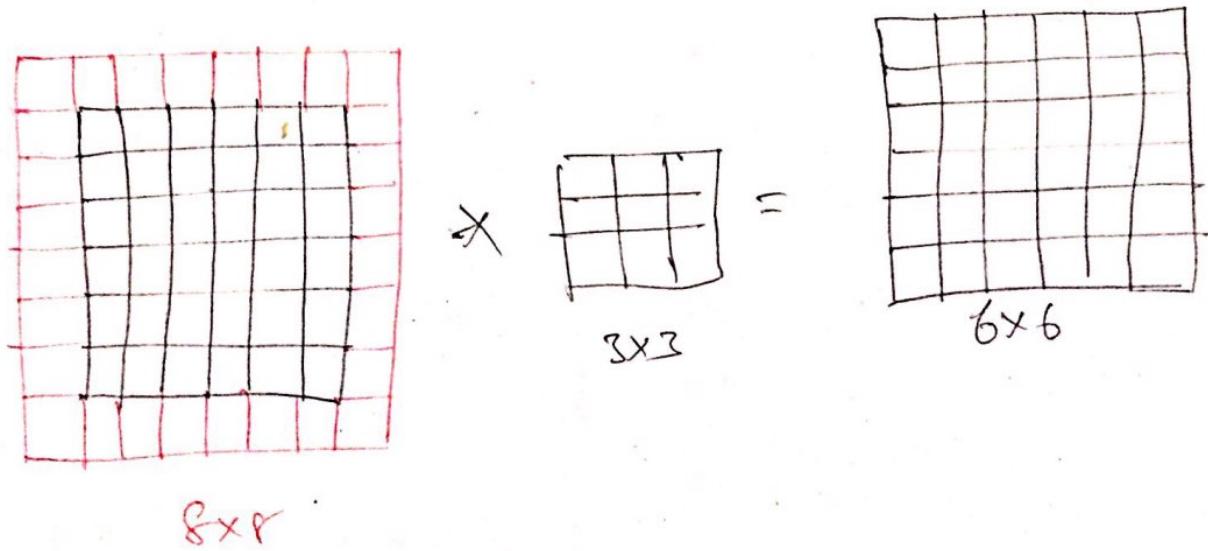
Padding ??

$$(6 \times 6) * (3 \times 3) = (4 \times 4)$$

in general

$$(n \times n) * (f \times f) = (n-f+1 \times n-f+1)$$

Note : In the prev filter moving over image
 (of size say 6×6) \rightarrow the corner pixels are
less contributing : (looking lot of info at corners)
 Image is also shrinking after convolution
 \hookrightarrow to address (each pixel to be equally considered
 in convolution)



Here $p = \text{padding} = 1$

After conv. image size will be

$$(n+2p-f+1) \times (n+2p-f+1)$$

$$6+2-3+1 \Rightarrow (6 \times 6)$$

→ Is it only padding $P=1$ always?

No

We can have padding $P=2$, or any integer



But, on one condition

Pad so that output size is the same
as the input size.

Eg:- (5×5) \times $(3 \times 3) = (3 \times 3)$

image filter image after
 padding



Image shrinking

How much padding we can apply?

$$n=5 \quad f=3 \quad P=?$$

$$\text{SA to get } n+2P-f+1 = n$$

$$8+2x-3+1 = 8$$

$$2x = 8 + 3 - 1 \\ = 10$$

$$\boxed{x=5}$$

If $n=5, f=5$

$$\text{Then } P = \frac{f-1}{2} = 2$$

Note:

Always f should be odd (usual convention)

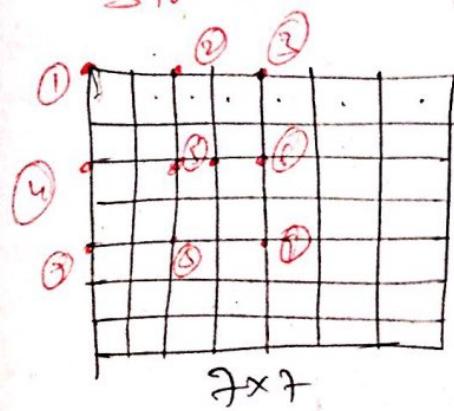
Remark

If $P=0$ (no padding)

is usually called "Valid"

"Same": pad so that output is the same
as the input size.

Strided Convolution



$$\begin{matrix} & X & \longrightarrow & Y & \\ \times & \begin{matrix} 3 & 4 & 4 \\ 1 & 0 & 2 \\ 1 & 0 & 3 \end{matrix} & = & \begin{matrix} & & \\ & & \\ & & \end{matrix} & 3 \times 3 \end{matrix}$$

Stride = 2 (measured by two pixels)

formula $(n \times n) * (f \times f)$ have $n=7$, $f=3$, $s=2$

padding: P :

Stride: S

result image will be

$$\text{floor} \left(\left[\frac{n+2P-f}{S} \right] + 1, \left[\frac{n+2P-f}{S} \right] + 1 \right)$$

$$\frac{7+0-3}{2} + 1 = 3$$

∴ The resultant image will be

(3x3)

$$\boxed{2.4} = 2 \text{ floor operation}$$

Summary of Convolution

(n × n) \times
padding P
Stride S

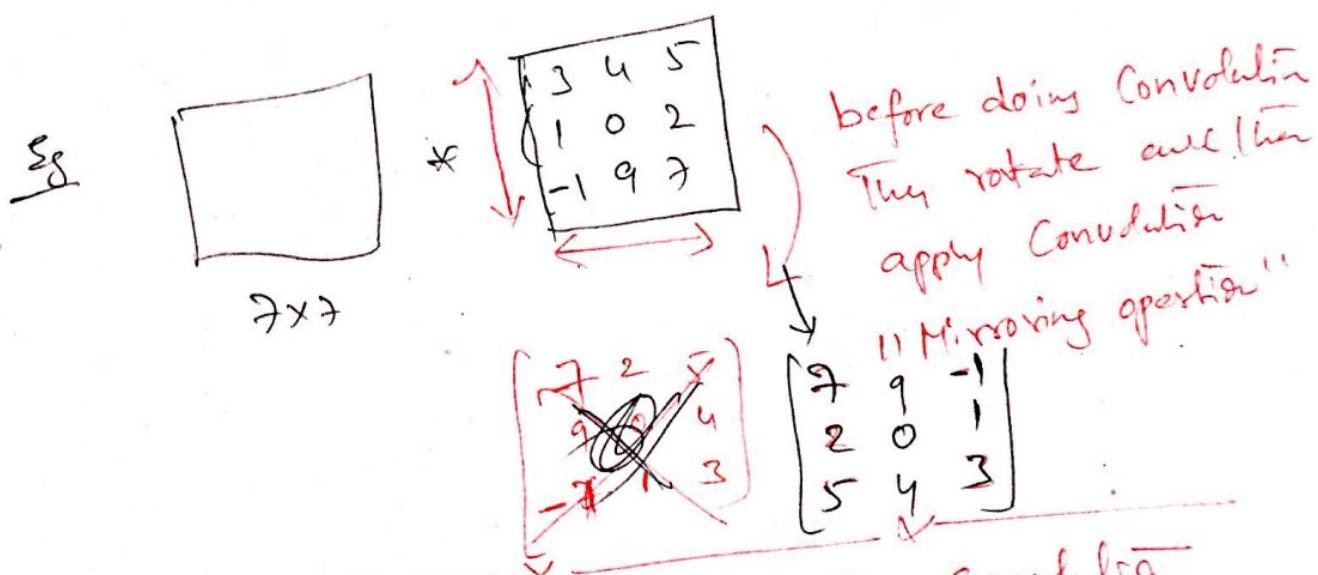
$$(f \times f) = \left(\left\lceil \frac{n+2P-f}{S} + 1 \right\rceil \times \left\lfloor \frac{n+2P-f}{S} + 1 \right\rfloor \right)$$

→ output Image size

Cross-Correlation

Vs. Convolution

In classical literature
who with mirroring operation
↑ what is this?



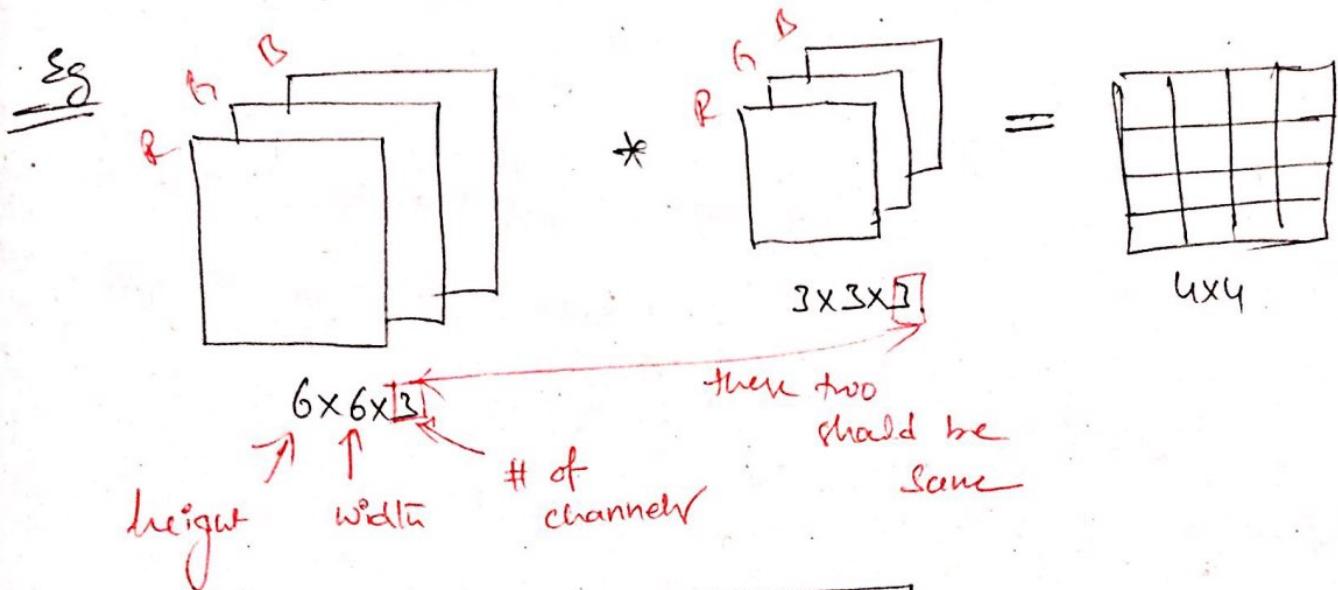
before doing Convolution
they rotate all then
apply Convolution

"Mirroring operation"

- Note:
- ① we use the terminology only convolution
 - ② In some texts "Conv" & "Cross-Correlation" are synonymous but we use "Convolution"
 - ③ Associative property $A * (B * C) = (A * B) * C$

Convolution over Volumes

we see not on grey scale images but on color images



Note: $\# \text{channels} = \# f$

Eg if our filter is

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

$\frac{1}{\sqrt{3}}$ R G B

It will find vehicle edge in Red channel

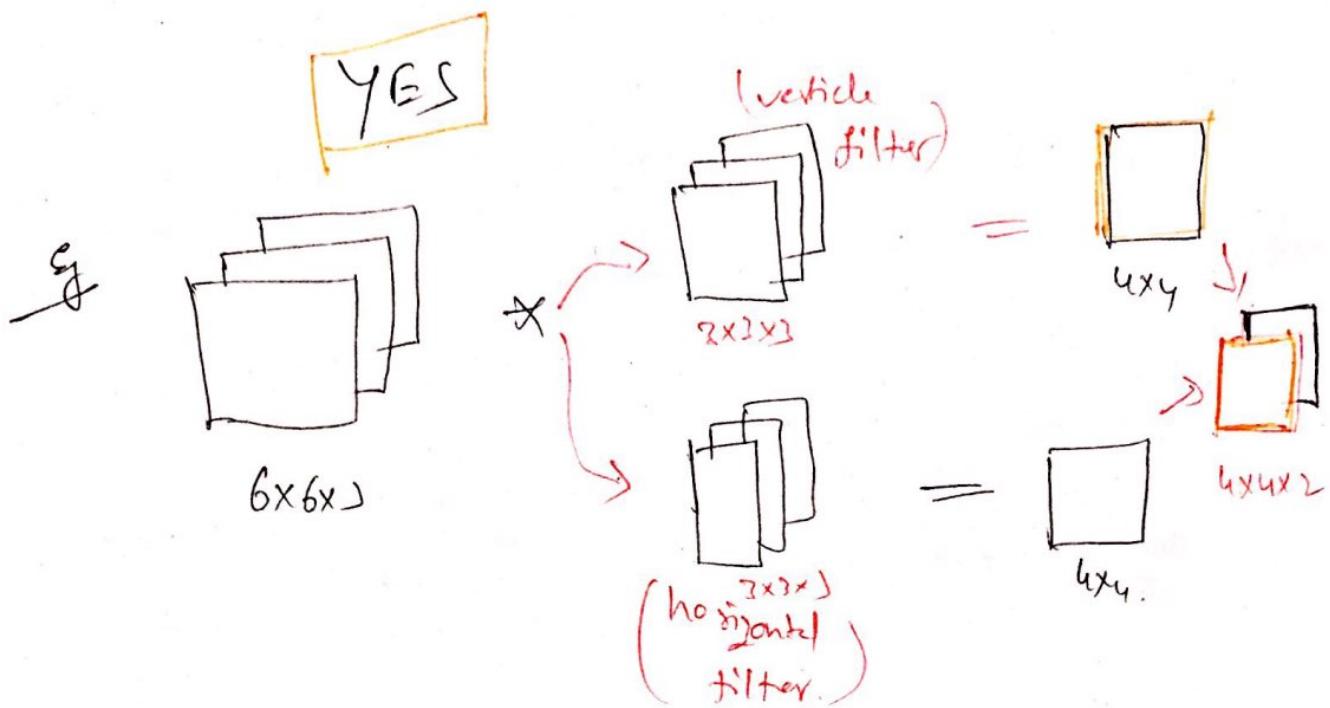
If our filter channels are

$$\begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix} \begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix} \begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix} \Rightarrow \text{detects vehicle edges in all 3 channels}$$

Can we apply multiple filter

filter 1 → to find horizontal edges

filter 2 → to find vertical filter
; 45° edges
; 75° edges (60° edges etc)



Summary

$n \times n \times n_c$ input image

$n_c = \# \text{ channels}$

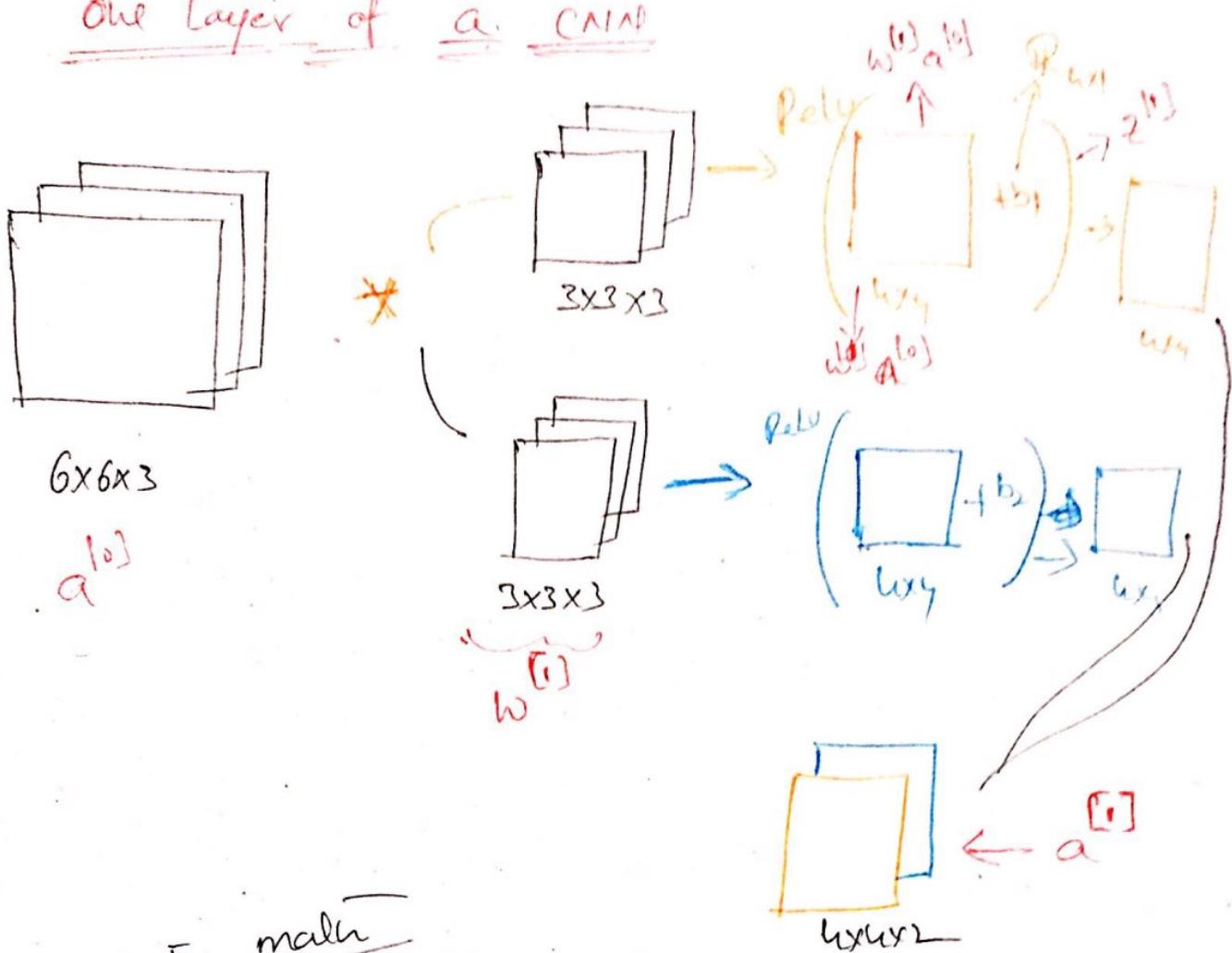
$f \times f \times n_c$ filter

$$(n \times n \times n_c) \times (f \times f \times n_c) \rightarrow ((n-f+1) \times (n-f+1) \times n'_c)$$

$$\boxed{n'_c = \# \text{ filters}}$$

Note that n_c , n'_c are different
+ $\Rightarrow \# \text{ filters}$ Some time n_c is also called or DEPTH

One Layer of a CNN

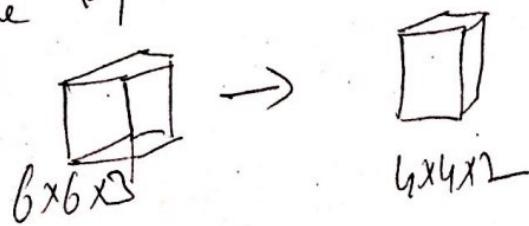


Recall the math

$$z^{(1)} = w^{(1)} a^{(0)} + b^{(1)}$$

$$a^{(1)} = g(z^{(1)})$$

In one layer of a CNN



6x4x2

(each filter 3x3x3)

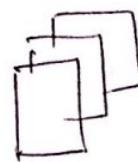
If there are 10 layers filter

Then 6x6x3 \rightarrow 6x6x10 ✓

Number of parameters in one layer?

If 10 filters each $3 \times 3 \times 3$

for each filter



$$3 \times 3 \times 3 = 27$$

+ bias

28 parameters

∴ for 10 filters $28 \times 10 = 280$ parameters

Summary Notation

l : layer number

$f^{[l]}$ = filter size in layer l

$p^{[l]}$ = padding

$s^{[l]}$ = stride

Input

$n_H^{[l-1]} \times n_W^{[l-1]} \times n_C^{[l-1]}$

Output:

$n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$

$n_C^{[l]} = \# \text{ filters}$

$$n_H^{[l]} = \left\lfloor \frac{n_H^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1 \right\rfloor$$

④ Each filter size

$$f^{(l)} \times f^{(l)} \times n_c^{(l)}$$

④ activations: $a^{(l)} \Rightarrow n_H^{(l)} \times n_W^{(l)} \times n_G^{(l)}$

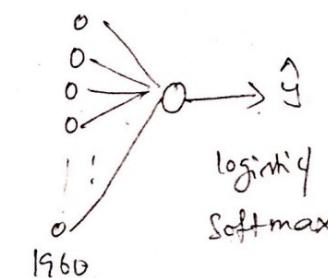
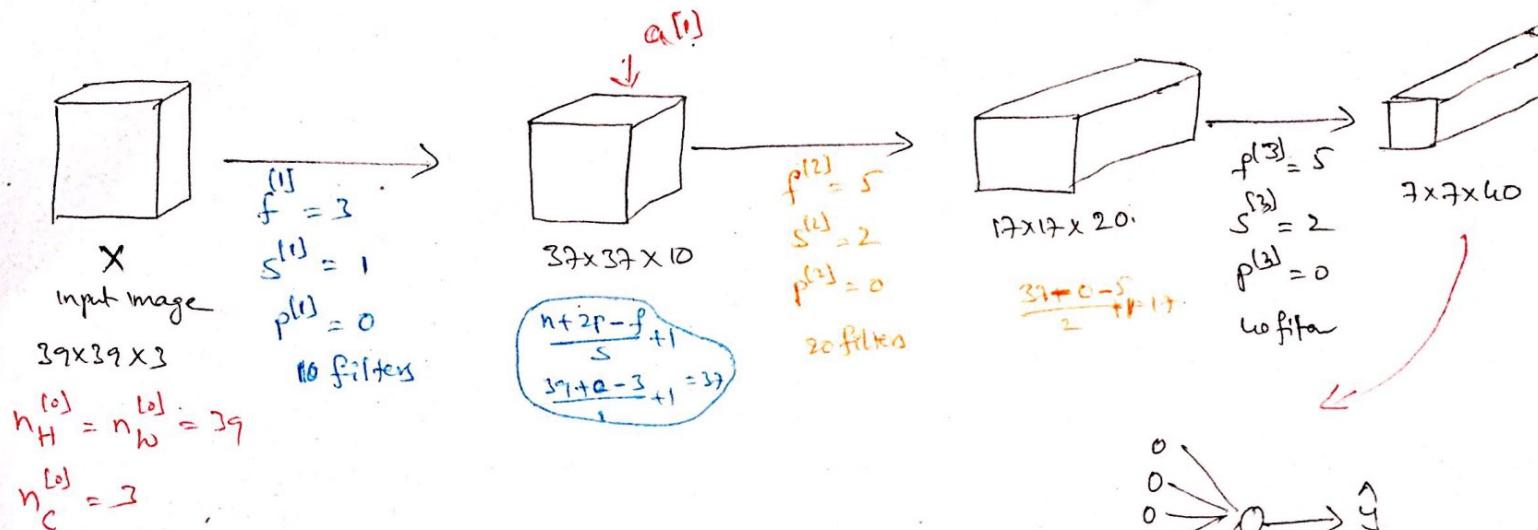
Vectorized implement $A^{(l)} = m \times n_H^{(l)} \times n_W^{(l)} \times n_G^{(l)}$

④ weights: $f^{(l)} \times f^{(l)} \times n_c^{(l-1)} \times \underline{n_c^{(l)}} \Rightarrow \# \text{ filters in layer } l.$

④ bias: $n_c^{(l)} \rightarrow (1, 1, 1, n_c^{(l)})$

A Single Convolution at A1

ConvNet



Hyperparameters
 f, s, p , #filter ($n_C^{(l)}$)
 each layer

Types of layer in a ConvNet

- ① Convolution layers (CONV)
- ② Pooling layers (POOL)
- ③ Fully Connected layer (FC)

MAX POOLING

1	3	2	1
2	9	1	1
1	3	2	3
5	6	1	2

4x4

9	2
6	3

Ex: 1

Hyperparameters

$$f = 2$$

$$s = 2$$

(No parameters to learn)

2x2

1	3	2	1	3
2	9	1	1	5
1	3	2	3	2
8	3	2	3	2
5	6	1	2	9

5x5x2

$$n_H = n_D = 5$$

$$\begin{aligned} f &= 3 \\ s &= 1 \\ p &= 0 \end{aligned}$$

9	9	5
9	9	5
8	6	9

3x3x2

$$\frac{n + 2p + f}{s} + 1 = \frac{5+0+3}{1} + 1 = ③$$

Input image size

$$n_H \times n_W \times n_C$$

Output image size

$$\left(\frac{n_H + 2p + f}{s} + 1 \times \frac{n_W + 2p + f}{s} + 1 \times n_C \right)$$

Hyperparameters
f
s

Average pooling

1	3	21	
2	9	1	1
1	4	2	3
5	6	1	2

4x4.

filter 2x2

3.75	1.25
4	2

2x2

↑
output

$$\frac{4+0-2}{2} + 1 = 2$$

Input

SUMMARY

Hyper parameters

Output image size

f : filter size

$$\left(\left[\frac{n_H+2p-f}{s} + 1 \right] \times \left[\frac{n_W+f-p}{s} + 1 \right] \times n_C \right)$$

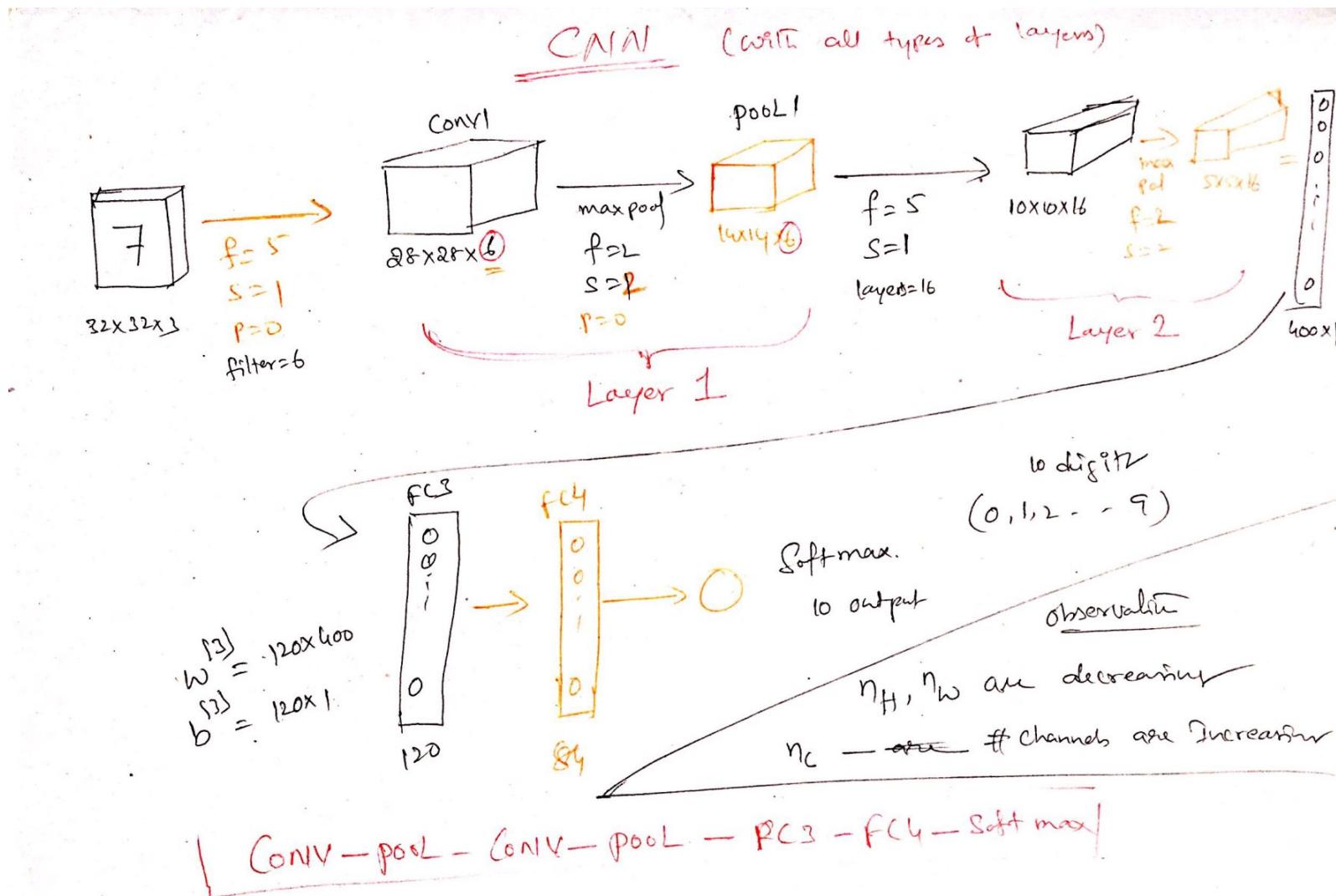
s : stride

max or Avg pool

(Usually padding doesn't make sense)

note
=

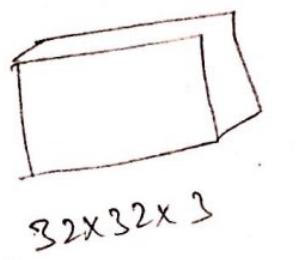
no parameters to learn



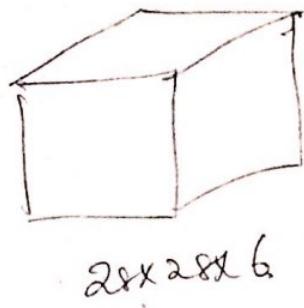
Summary

	<u>Activatn shape</u>	<u>activtn size</u>	<u># parameters</u>
Input	(32x32x3)	3072 a^3	0
Conv1 ($f=5, s=1$)	(28, 28, 8)	6272 $8 \times (5 \times 5 \times 3)$	608 (608)
pool1	(14, 14, 8)	1568	0 ←
Conv2 ($f=5, s=1$)	(10, 10, 16)	1600 $(5 \times 5 \times 8)$ * 16	416 3216
pool2	(5, 5, 16)	600	0 ←
FC3	(120, 1)	120 $600 \times 120 + 120$	48,000 / 48(?)
FC4	(84, 1)	84 $\cancel{600 \times 120 + 120}$ $120 \times 84 + 84$ $12,084$	12,084 - 6000 → 10,484
Softmax	(10, 1)	10 $84 \times 10 + 10$	850

Why Convolutions?



f_{25}
6 filters



parameters =
each filter

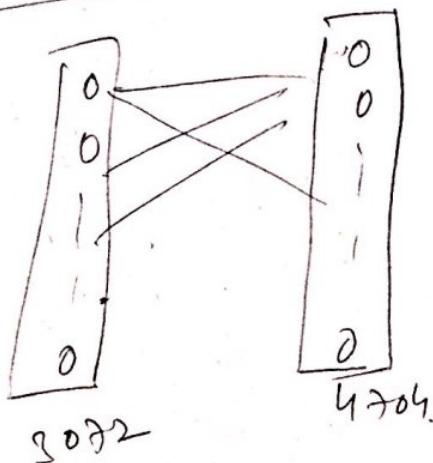
$$5 \times 5 = 25 + 1 \text{ (bias)} = 26$$

$$6 \times 26 = 156 \text{ parameters}$$

In general fully connected DNN

$$32 \times 32 \times 3 = 3072$$

$$28 \times 28 \times 6 = 4704$$



weight matrix

$$3072 \times 4704$$

≈ 14 Million

Summary No parameters in Conv layers are small

When compare to full fully connected layer in NN

Parameter Sharing:- A feature detector (such as vertical edge detector) that's useful in one part of the image is probably useful in another part of the image.

Sparcity of Connections

In each layer, each output value depends only on a small number of inputs.

(Recall filter for edge detection)

Translation Invariance

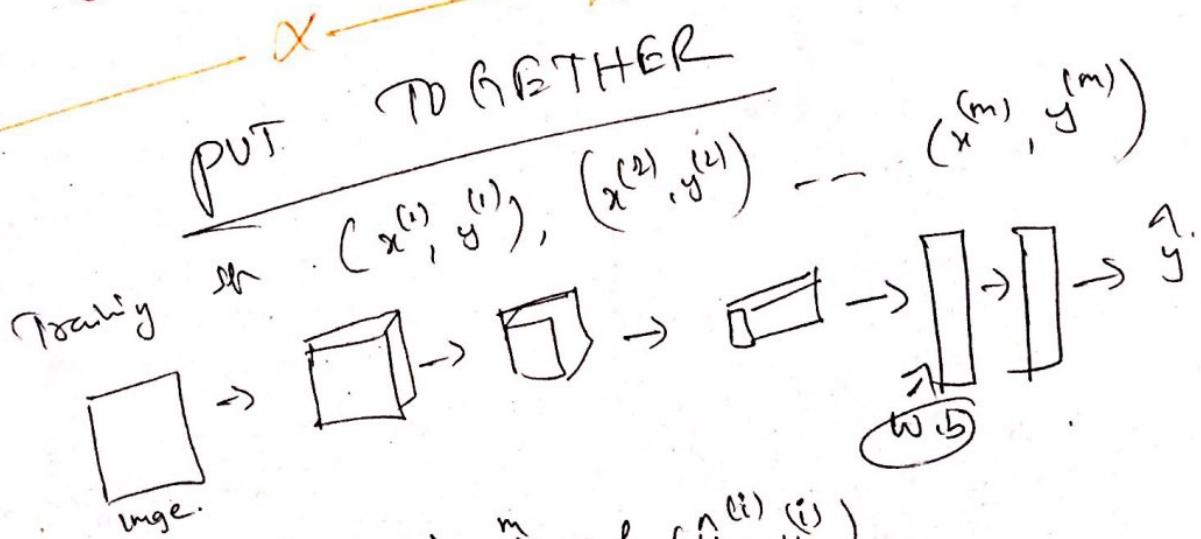


Why Convolution

① parameter sharing

② sparsity of connections

③ Translation Invariance



$$\text{Cost } J = \frac{1}{m} \sum_{i=1}^m L(g^{(i)}, y^{(i)})$$

Use Grad. Descent to optimize parameters to reduce J

Care Studies - CNNs

Deep Convolutional Model

Classification Networks

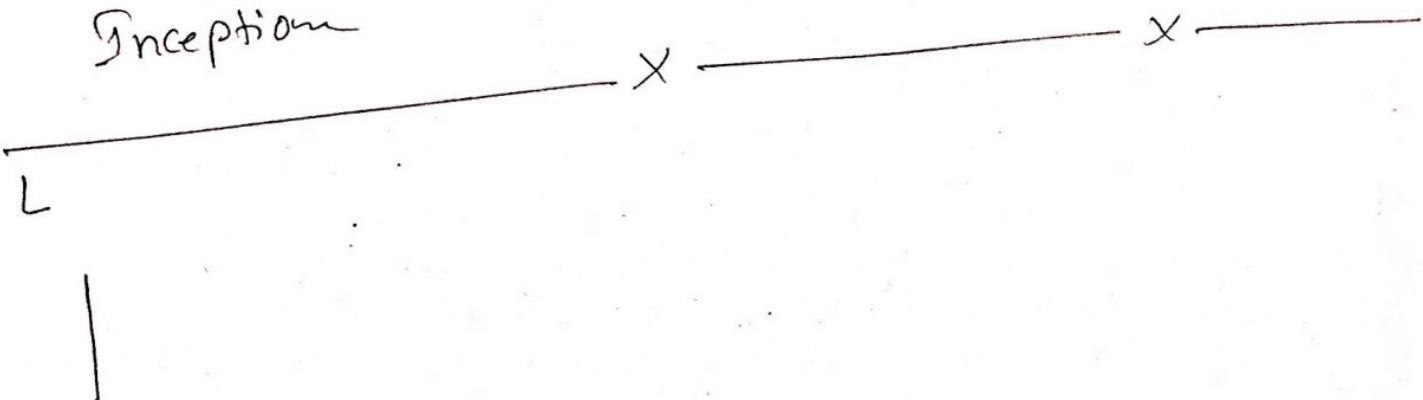
* LeNet - 5

* AlexNet ←

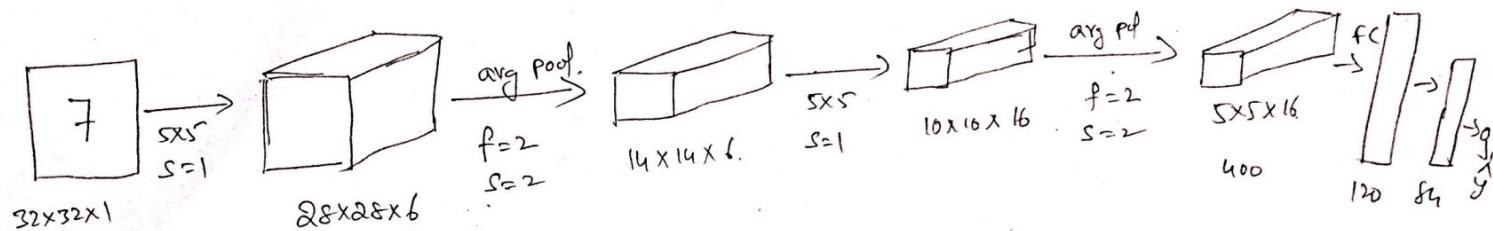
* VGG

ResNet (Residual Network)

Inception

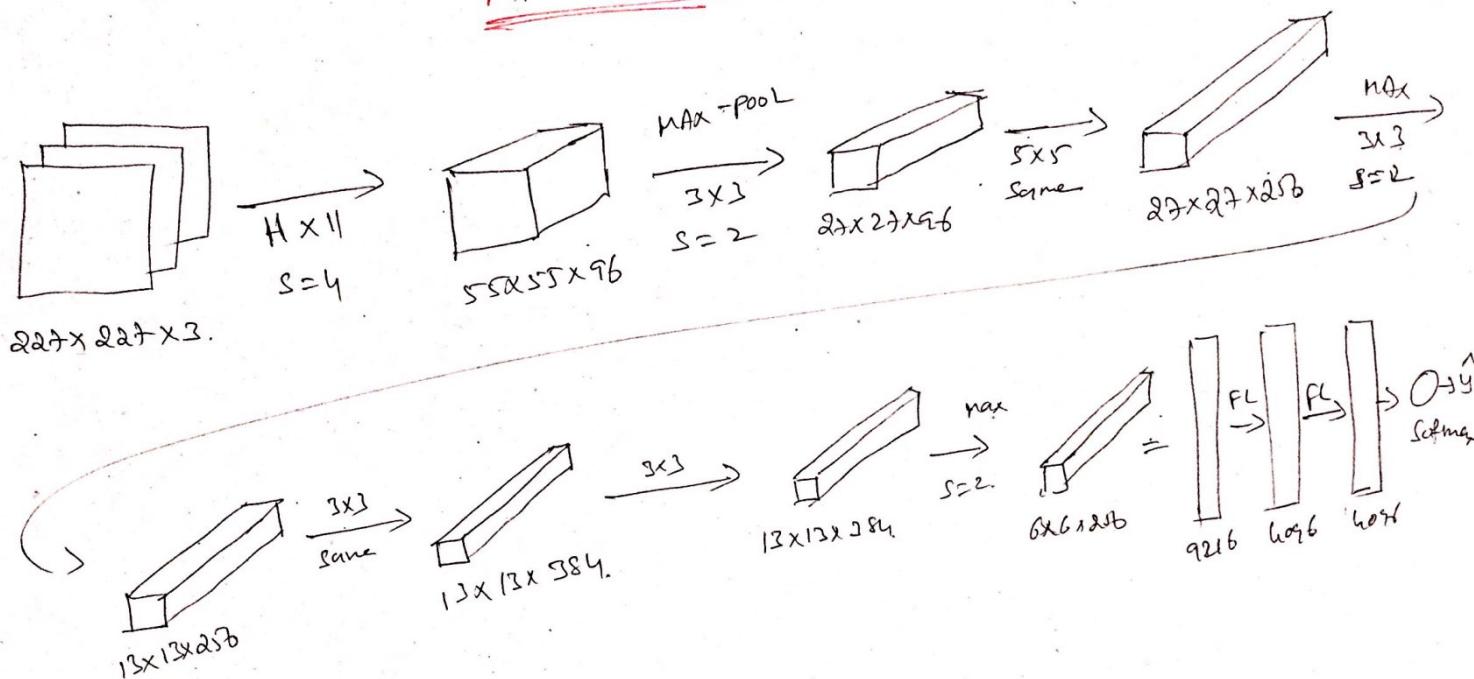


LeNet - 5



here $n_H, n_W \downarrow$ and $n_C \uparrow$.

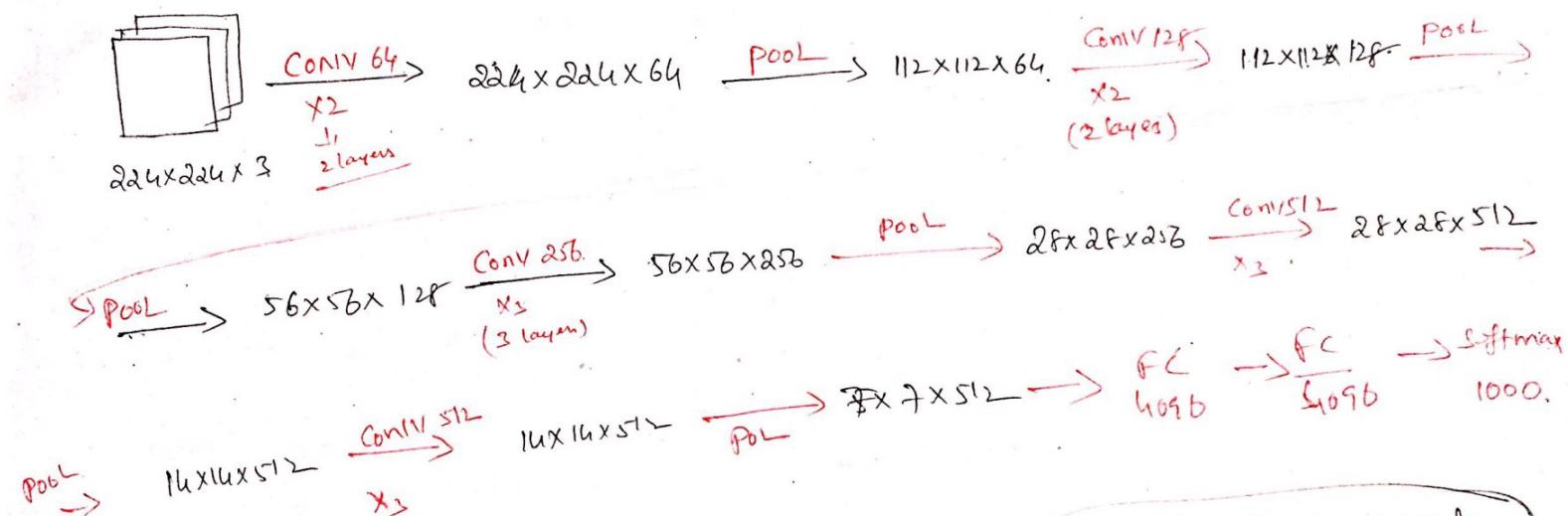
Alex Net



VGG-16

Conv = 3x3 filter, $s=1$, Same

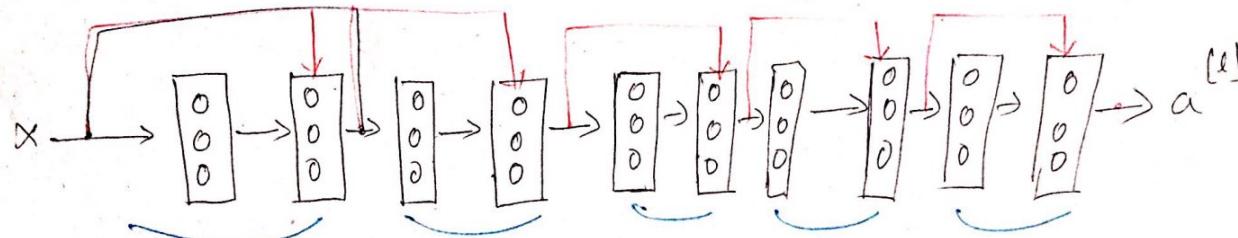
MAX-POOL = 2×2 , $s=2$



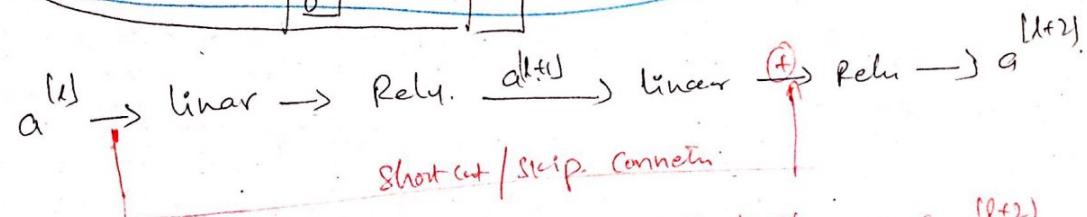
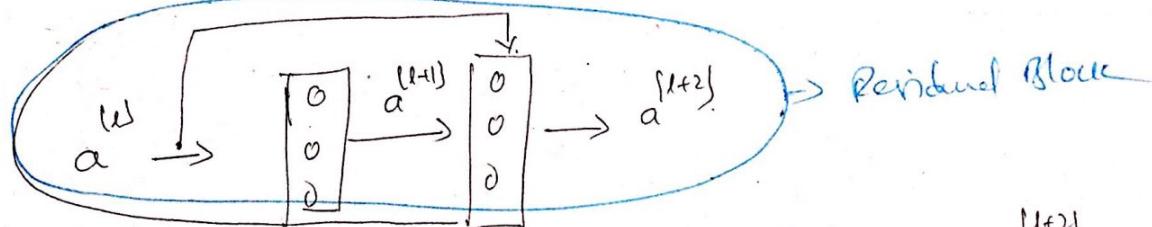
≈ 138 million parameters

Andri improvement
VGG-19

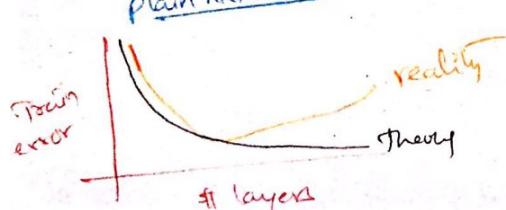
Residual Network (ResNet)



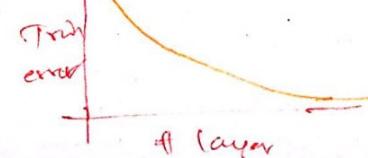
plain network
is without
short cut
connection/
residual



plain network

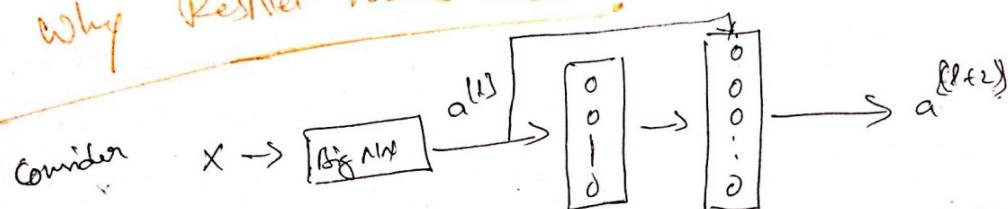


ResNet



$$a^{(l+2)} = g(z^{(l+2)} + a^{(l)})$$

Why ResNet works well?



as ReLU results values $\geq 0 \therefore a \geq 0$.

$$a^{(l+2)} = g\left(z^{(l+2)} + a^{(l)}\right) = g\left(w^{(l+2)} \frac{a^{(l+1)}}{a} + b^{(l+2)} + a^{(l)}\right)$$

because if $w^{(l+2)} = 0$, $g^{(l+2)} = 0$.

$$a^{(l+3)} = g(a^{(l)}) = a^{(l)}$$

Summary → Identity function is easy for residual block to learn.

→ easy if $z^{(l+2)}$ (or $a^{(l+2)}$) and $a^{(l)}$ are of same dimension

→ If not we multiply $a^{(l)}$ by another weight matrix W

$$\text{eg } a^{(l+2)} = 256 \times 1 \quad a^{(l)}_{128 \times 1} \quad W_{256 \times 128} \xrightarrow{\text{New weight}} \underline{\text{matrix}}$$