

# CNNs

→ majorly used in computer vision problems

- ↳ image classification
- ↳ object detection
- ↳ Neural style transfer

A big challenge in CV is

[inputs can get really big]

for an image of 1000x600 pixel image.  
1 MB

The dimension of input features [x] becomes  $1000 \times 600 \times 3$   
 $= 3\text{million} \text{ } \text{ } \text{ } \text{ } \text{ }$

e.g:-

for a fully connected network, even if your 1<sup>st</sup> layer contains just 1000 neurons, the size of weight matrix is  $3\text{M} \times 1000 = 3\text{GB}$  ← enormous.

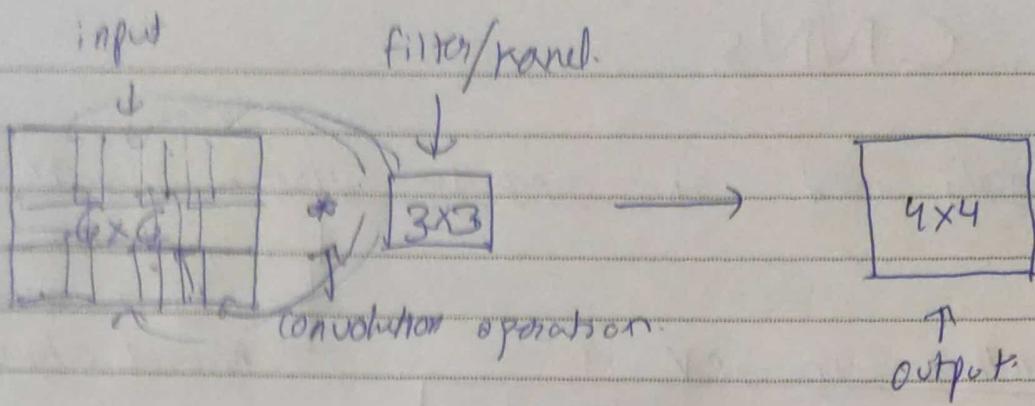
→ with that much size of data, chances of overfitting increase v much, + computer memory requirement increase

→ that doesn't mean we can't use big images with NN  
↳ CNNs allow us to do so

# CNNs example :- Edge Detection

Convolution operation  $\rightarrow \left\{ \begin{array}{l} \text{Big matrix} \times \text{small matrix} \\ \text{elt wise product over the entire matrix.} \end{array} \right.$

gives the output image.



\* Different types of filters/kernels detect different features

e.g.:

$$\begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix} \rightarrow \text{detects vertical edge.}$$

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \end{bmatrix} \star \begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix} = \begin{bmatrix} 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \\ 0 & 30 & 30 & 0 \end{bmatrix}$$

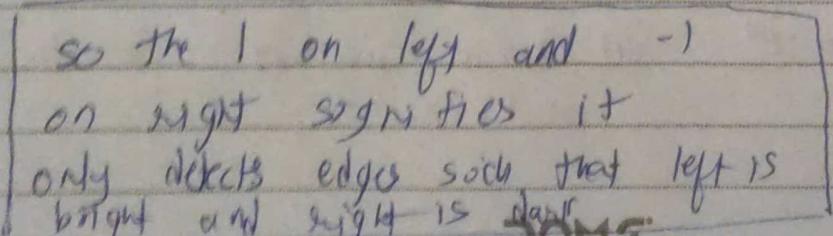
S.  verify 

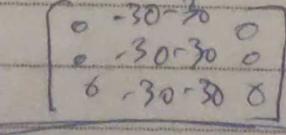
{ the middle part is white means there is an edge in middle }

The detected edge looks thick  
bcz 6x6 is a very small image

for larger size images its perf

→ what if input was  $\begin{bmatrix} 0 & 0 & 0 & 10 & 10 & 10 \\ 0 & 0 & 0 & 10 & 10 & 10 \\ 0 & 0 & 0 & 10 & 10 & 10 \end{bmatrix}$  → output becomes

\*  so the 1 on left and -1 on right signifies it only detects edges such that left is bright and right is dark



11

if left of image is bright then  $\begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{bmatrix}$  is the right kernel to use.

\* similarly  $\begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{bmatrix}$  detects horizontal edges such that ~~left~~ image is light on top and dark at bottom

\* To detect edges  $\begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix}$  is a common filter

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

Sobel filter

puts more weight to middle part

→ diff types of filters give us diff properties

\* → how long will you keep handpicking the right kernels?

→ you're not even sure what are "the right features" to detect to get the desired output

So you stop handpicking them and let the machine learn it.

\* In CNNs the kernels are your weights  
the machine "learns the right kernels" to apply which will give the required output

$$\text{Input image} \times \begin{bmatrix} w_1 & w_2 & w_3 \\ w_4 & w_5 & w_6 \\ w_7 & w_8 & w_9 \end{bmatrix} = \begin{bmatrix} \text{required output} \\ \dots \end{bmatrix}$$

DIMS weights of your model

## # PADDING.

$$\begin{array}{l} \boxed{6 \times 6 * 3 \times 3 \rightarrow 4 \times 4} \\ \boxed{n \times n * f \times f \rightarrow (n-f+1) \times (n-f+1)} \end{array}$$

{ → so your image shrinks after endng convolution operat<sup>n</sup>  
 { → edge pixels don't contribute much size  
 they are  $n \rightarrow$  you're throwing away relevant info from edges

→ you increase the size of input image!

if you add  $p$  layers on top and bottom  
 ↗ (left, right)

$$\text{non padding} \rightarrow (n+2p) * (n+2p) * f \times f \rightarrow (n+2p+f-1) * (n+2p+f-1)$$

what do you add in these  $p$  layers?  
 → 0

Valid convolution → no padding

Same Convolution →  $\text{input size} = \text{output size}$

$$n+2p+f-1 = n$$

$$p = \frac{f-1}{2}$$

If  $f$  is usually odd.

↪ If even  $f$  doesn't give same convolution

↪ If even  $f$  doesn't have a "centre"

in CV it's good to have a centre to your filter

## # Strided Convolutions.

convolution operation means taking a kernel and sliding it over the input image.

at every step you take element-wise product.

→ but how much do you "slide" your kernel by → stride normally its 1  
but it can be whatever you want

11 stride is same in both horizontal and vertical direction.

$$n \times n \text{ padding } p \quad * \quad f \times f \quad \rightarrow? \quad \left( \frac{n+2p-f}{s} + 1 \right) \times$$

if not integer, take floor.  
means, you never take your filter out of your image.

11 in some math textbooks the convolution operation means rotating/flipping the kernel horizontally than vertically than taking the element-wise product over the entire image, but in all its just the last part 11

$$\text{also } (A * B) * C = A * (B * C) \quad \checkmark \text{ associativity}$$

## # CONVOLUTIONS OVER VOLUME.

what if you've an RGB image instead of grayscale

height  $\sqrt{6 \times 6 \times 3}$  with no. of channels

DOMS

3 must equal 3

→ now your filter changes from  $3 \times 3$  to  $3 \times 3 \times 3$  !!  
with  $\uparrow \uparrow \uparrow$   
rest with no. of channels

no. of channels in image must equal no. of channels  
in kernel

$$6 \times 6 \times 3 * 3 \times 3 \times 3 \rightarrow \underline{4 \times 4}$$

like you had 9 elts in one elt wise product  
 $\rightarrow$  not  $6 \times 6 \times 3$

$\rightarrow$  you summed up those 9 elt to get output  
elt in output

Now you'll have 27 elts in each elt wise product  
you sum up those 27 elts to get corresponding elt for  
output.

Now multiply the 1<sup>st</sup> layer of next with 1<sup>st</sup> layer of input  
2<sup>nd</sup> layer with 2<sup>nd</sup> layer of input  
etc gives 9 elt gives 9 elt etc with 3<sup>rd</sup>  
gives 9 elt.

## # Multiple filters

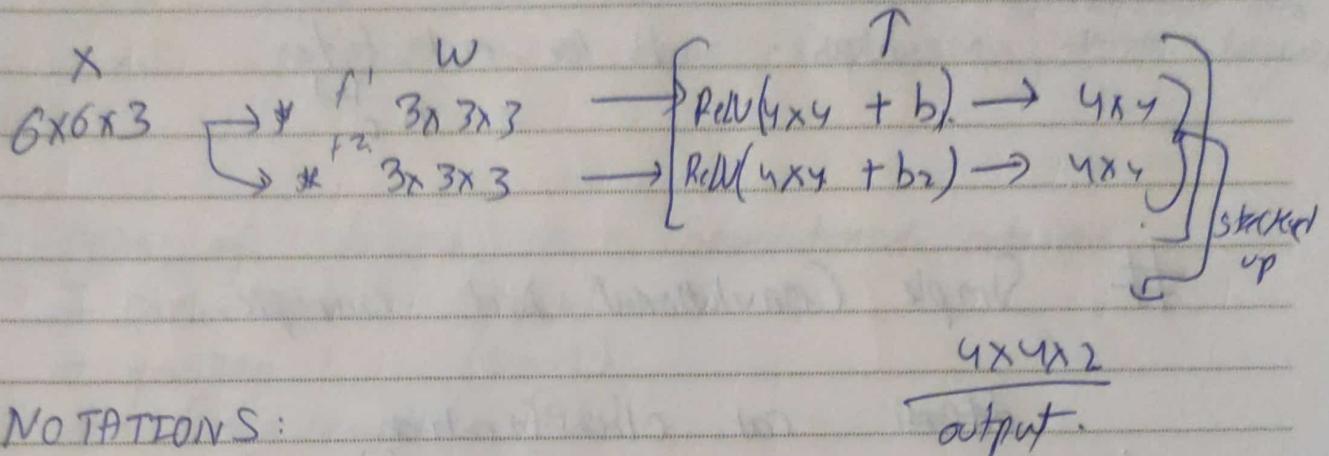
obv you want to detect multiple features,  
so you convolve the input with multiple filters

how do you get the output?

To detect Vedges or anything else

$$6 \times 6 \times 3 * \overbrace{3 \times 3 \times 3}^{\text{to detect vertically}} \rightarrow 4 \times 4 \quad \left. \begin{array}{l} \text{stacked up} \\ \text{output} \end{array} \right\} \begin{array}{l} 4 \times M \times 2 \\ \text{can be generalized to } n \end{array}$$

# # ONE LAYER OF CNN      $\text{ReLU}(Wx + b)$



NOTATIONS :

If layer  $l$  is a convolution layer.

for layer  $l$

$f^{[l]}$	= 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 of prev layer.

Output :  $n_h^{[l]} \times n_w^{[l]} \times n_c^{[l]}$

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

serve to  $n_w^{[l]}$

$n_c^{[l]}$  = no. of filters you're using

Each filter is :  $f^{[l]} \times f^{[l]} \times n_c^{[l-1]}$  → must match no. of channels in input.

Activations .  $a^{[l]} \rightarrow n_h^{[l]} \times n_w^{[l]} \times n_c^{[l]}$   
 for  $m$  examples.

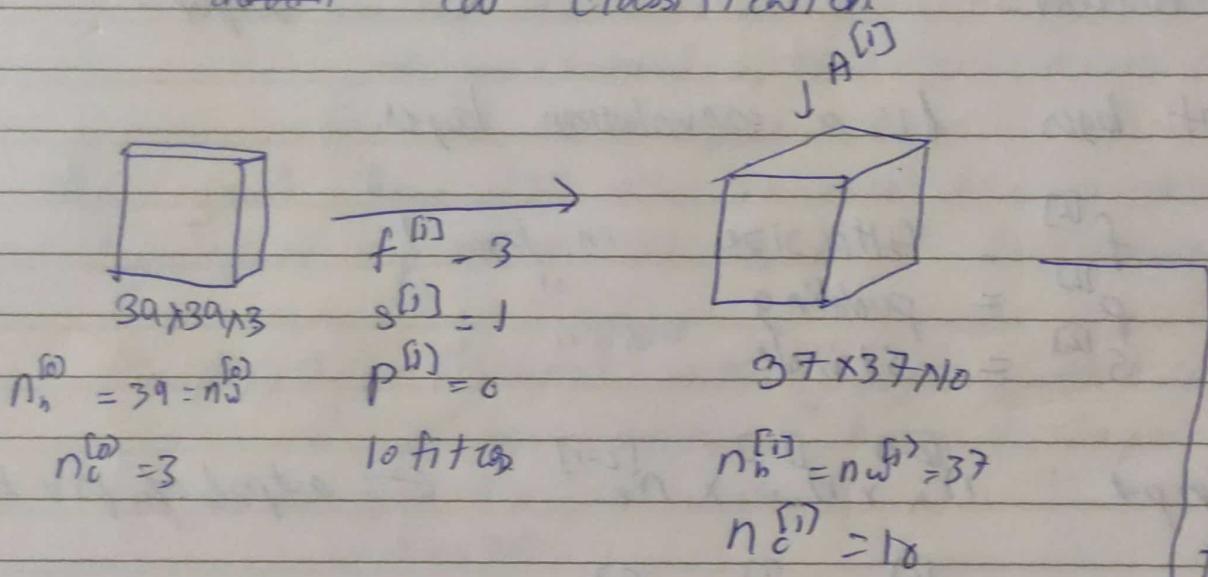
Wagns

$$\underbrace{f^{(1)} \times f^{(1)}}_{\text{size of each filter}} \times n_c^{(1)} \times \text{no. of filters}$$

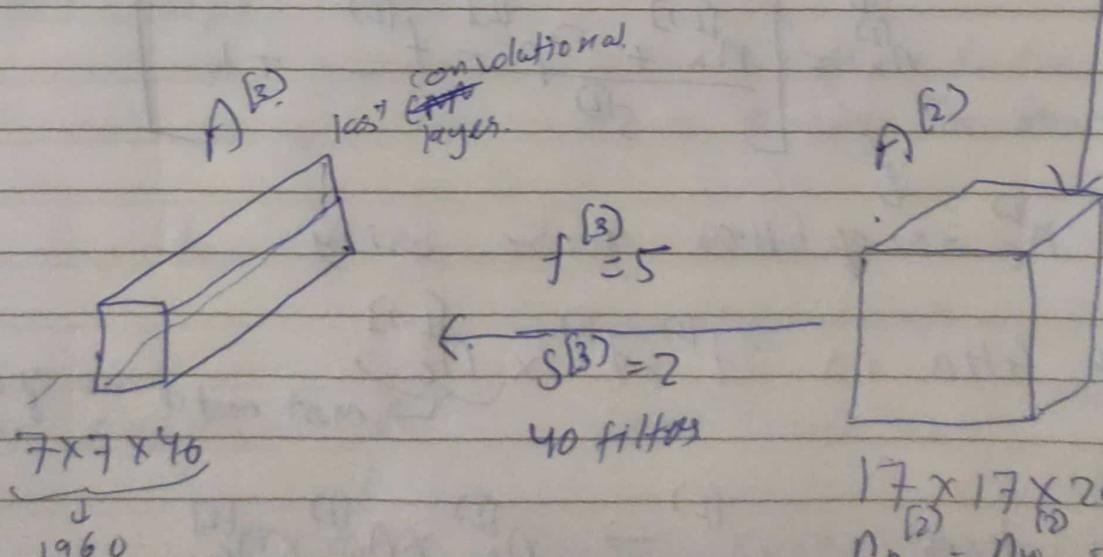
bias : 1 for each filter  
 $= n_c^{(1)}$

## # Simple Convolutional Net example.

~~Input~~ cat classification



$f^{(2)} = 5$ ,  $s^{(2)} = 2$ ,  $P^{(2)} = 0$ , 20 filters.



$$n_h^{(2)} = n_w^{(2)} = 17, n_c^{(2)} = 20$$

convert into vector  $\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_{1960} \end{bmatrix}$  feed  $\rightarrow$   $y \rightarrow$  final output.  
 Softmax  
 Logistic regression unit  
 50MS

Usually in convolutional networks

your input size keeps decreasing } through each  
and channel size keeps increasing } as we layer

# types of layers, in convolutional network.

- convolution (conv)
  - pooling (pool.)
  - Fully connected (FC)

# POOLING LAYERS.

Max pooling → you break your input into certain sections, then take each section's max

4x4 input grid

2x2 max pooling filter

certain sections = filter size

for  $4 \times 4 \rightarrow 2 \times 2$  noting to  
~~hyperparameters~~ { filters size = 2 // there are no <sup>non.</sup>  
 stride = 2 parameters in  
 Max pooling layers //

But what is max pooling doing?

- if a feature is detected in any filter  
keep it preserved
- also it "just works well"

11

Average pooling → you take avg, instead of max  
 ↳ maybe not used much  
 ↳ used to collapse your network & in very deep networks

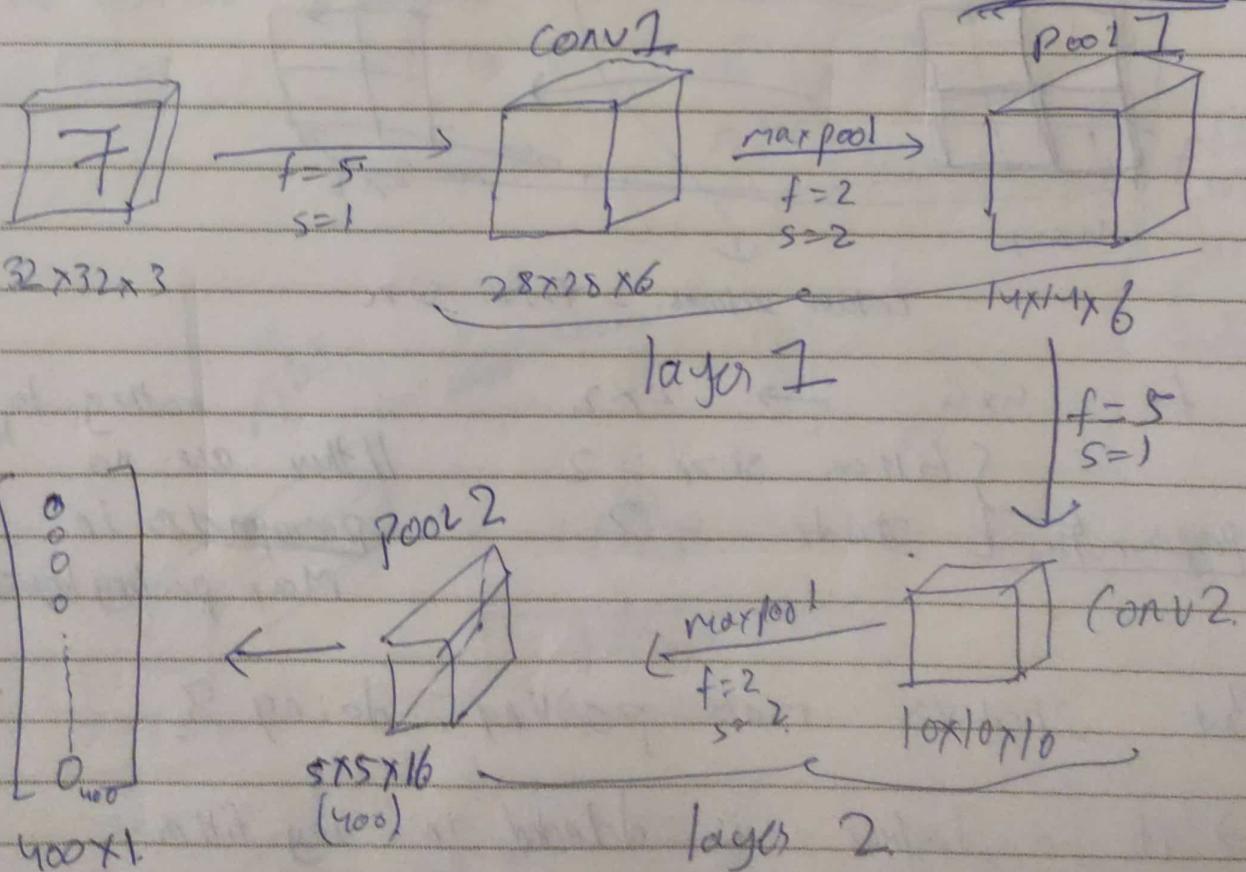
$$\frac{N_h \times N_w \times N_c}{\downarrow} \xrightarrow{\text{padding}} \left[ \frac{N_h - f + 1}{s} \right] \times \left[ \frac{N_w - f + 1}{s} \right] \times N_c$$

padding is applied  
to each channel individually

II padding is v rarely used in pooling

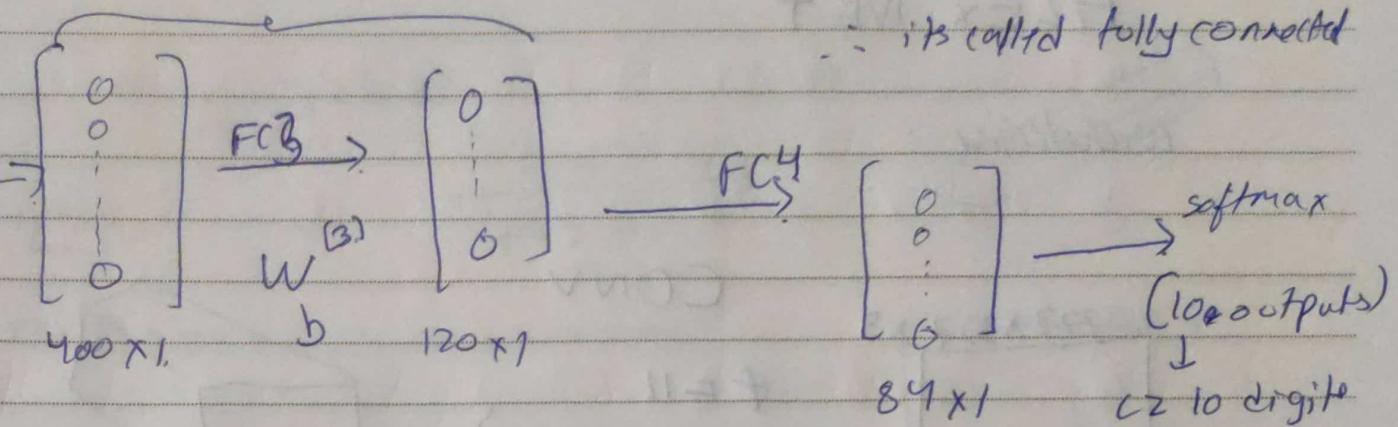
usually pooling layers are not considered layers  
bcz they don't have any parameters

## # NEURAL NETWORK EXAMPLE — LeNet-5



normal layers like we  
saw in NNs

$$W^{(B)} = \underbrace{120 \times 400}_{\text{bcz each unit TS connected}} \\ \therefore \text{it's called fully connected}$$



|| a lot of hyperparameters !!

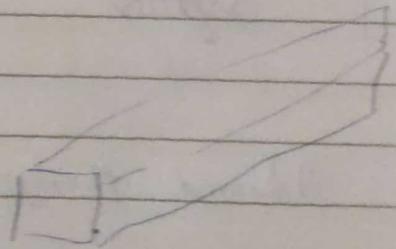
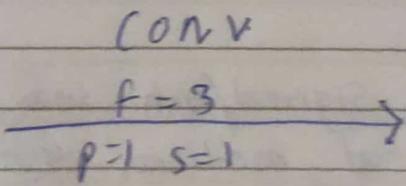
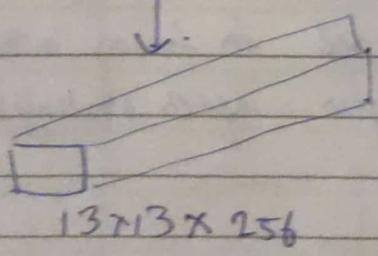
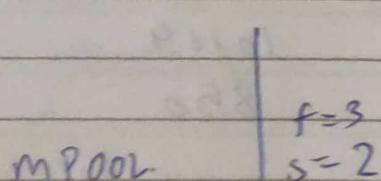
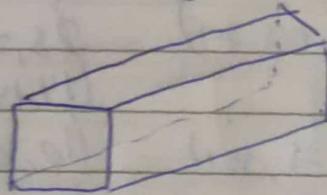
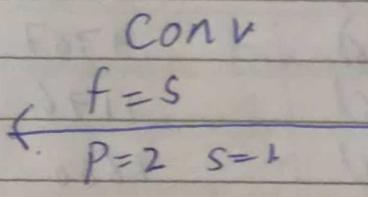
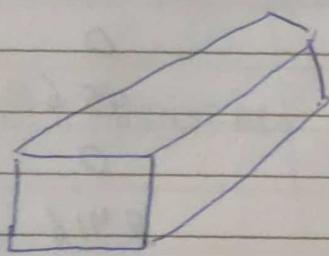
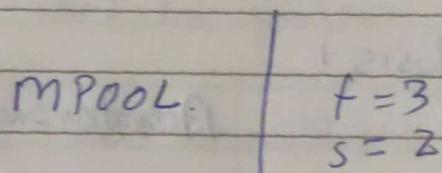
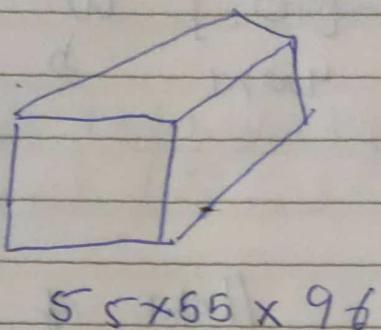
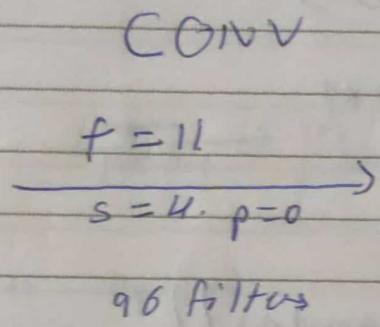
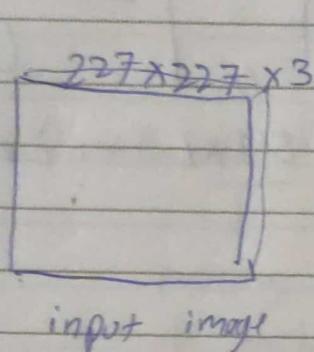
### \*Analysis:

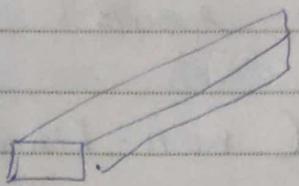
	Activation shape	Activation Siz [A]	# parameters
Input	(32, 32, 3)	3072	0
CONV1 ( $f=5, s=1$ )	(8, 20, 6)	4707	456
POOL 1	(14, 14, 6)	1176	0
CONV2 ( $f=5, s=1$ )	(10, 10, 16)	1600	2416
POOL 2	(5, 5, 16)	400	0
FC3	(120, 1)	120	48210
FC4	(84, 1)	84	10164
Softmax	(10, 1)	16	850

Advanced things:-  $\rightarrow$  sigmoid/tanh was used as activation.  
 $\rightarrow$  had non-linear activation after poolings!

# ALEX NET

architecture :

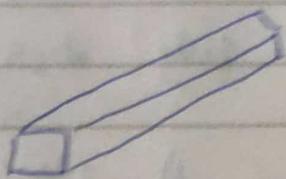




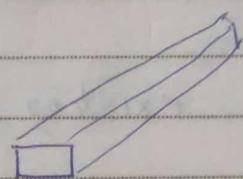
$13 \times 13 \times 384$

$\xrightarrow{\text{CONV}}$   
 $f=3 \quad p=1$   
 $s=1$

384 filters

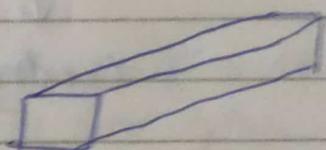


$13 \times 13 \times 128$



$6 \times 6 \times 256 = 4096$

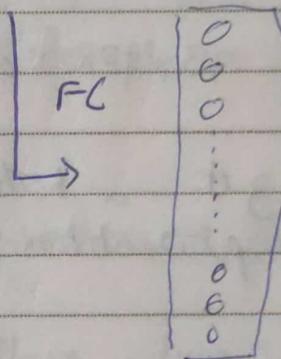
$\xleftarrow{\text{MPOOL}}$   
 $f=3$   
 $s=2 \quad p=1$



$13 \times 13 \times 128$

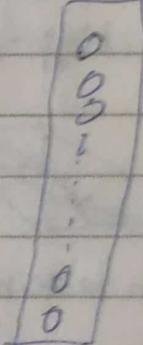
$256 \text{ filters}$   
 $f=3$   
 $p=1$   
 $s=1$

$\downarrow \text{CONV}$



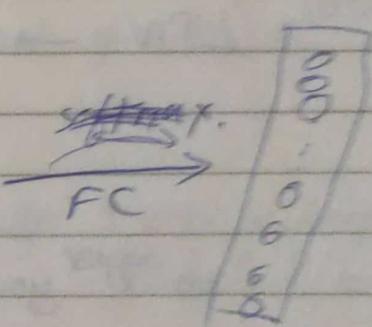
4096

$\xrightarrow{\text{FC}}$



~~1000~~

4096

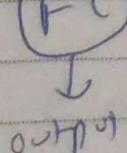
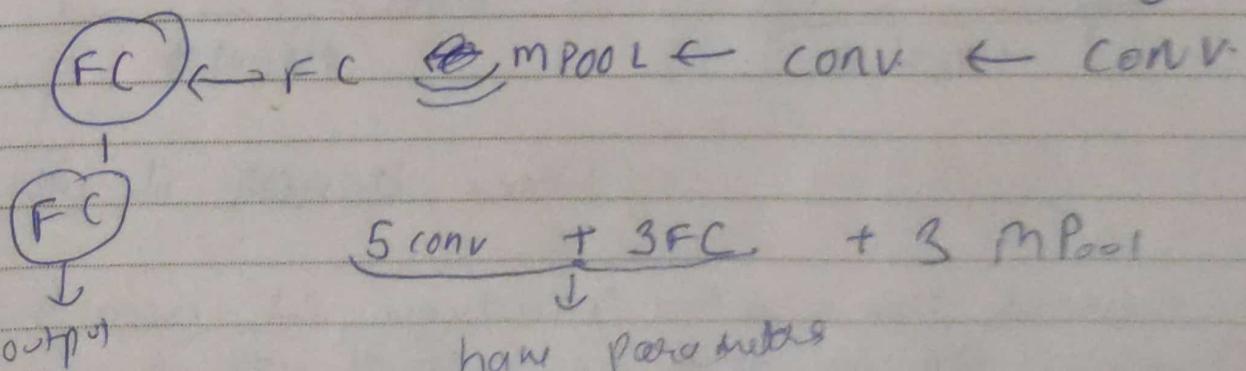


1000

(output)

$\text{CONV} \rightarrow \text{MPOOL} \rightarrow \text{CONV} \rightarrow \text{MPOOL} \rightarrow \text{CONV}$

$\downarrow$



$\underbrace{5 \text{ conv}}_{\text{how many}} + \underbrace{3 \text{ FC}}_{\text{parameters}} + 3 \text{ MPool}$

parameters

DOMS

Why does Alexnet achieve better results?

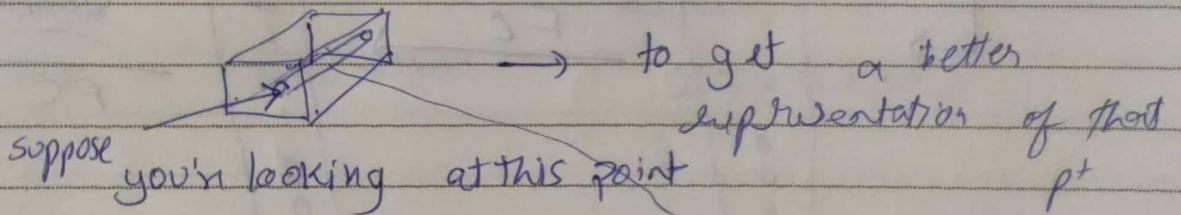
- 1) Applied ReLU activation (non-linearity)
- 2) Standardization (local response normalization)
- 3) Dropout regularization

<sup>some advanced  
CNN stuff</sup> (→ y) Data Augmentation (Enhanced Data)  
 $\downarrow$

Creating more data from existing data

- 5) Local Response Normalization.

LRN → suppose you're in a  $3 \times 13 \times 256$  volume



We normalize that point's value over this entire volume

Alex has lots of hyperparameters

## RESNETS

↳ build out of residual blocks

$$a^{(L)} \xrightarrow{\text{linear}} \text{ReLU} \xrightarrow{\text{linear}} \text{ReLU} \xrightarrow{\dots} a^{(L+1)}$$

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

$$a^{(L+2)} = g(z^{(L+1)})$$

by residual blocks what we do is

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

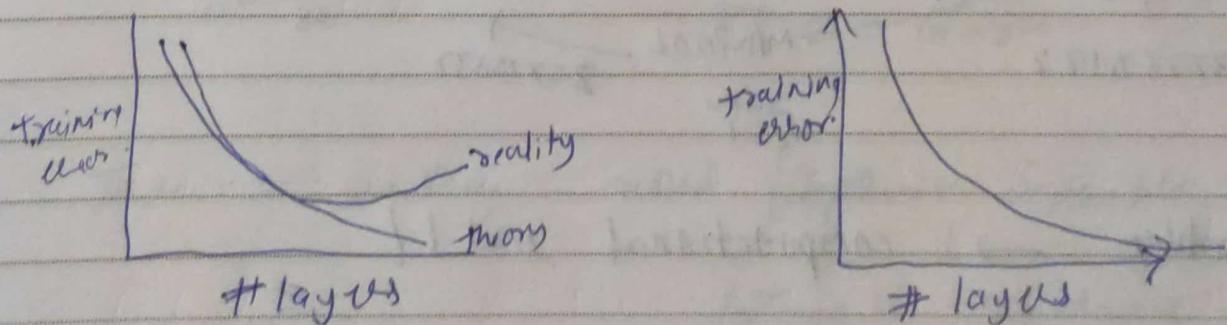
↳ we pick a previous layer and add it to some future layer

this is also called "shortcut/skip connection"  
 $a^{(L)}$  skips a few middle layers and goes to  $a^{(L+2)}$

→ Observation → Residual Networks allow to train very deep models (100+)

Plain networks

## RESNETS



\* Why do RESNETS work?

If using L2 regularization  $\rightarrow w, b$  will decrease  $\propto \epsilon$  for some layers so res blocks act like identity f'

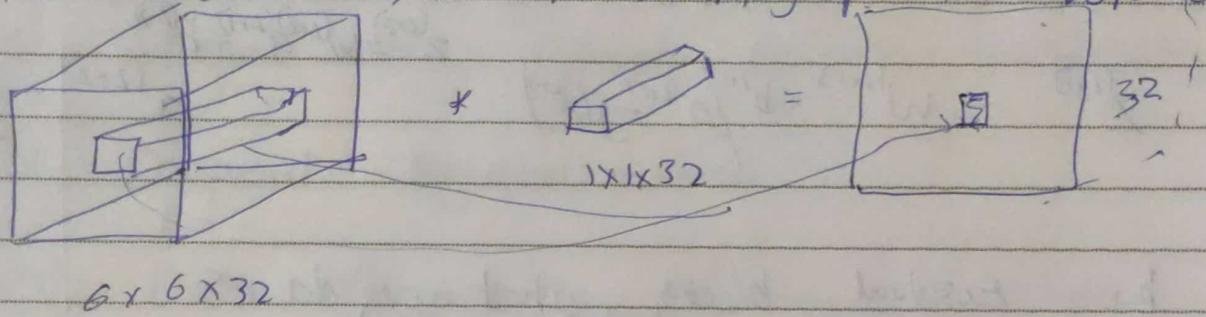
DOHS Keeping  $a^{(L+1)} \approx a^{(L)}$

#  $1 \times 1$  convolutions

→ looks like just multiplying by a number  
 $\downarrow$

only in 2D

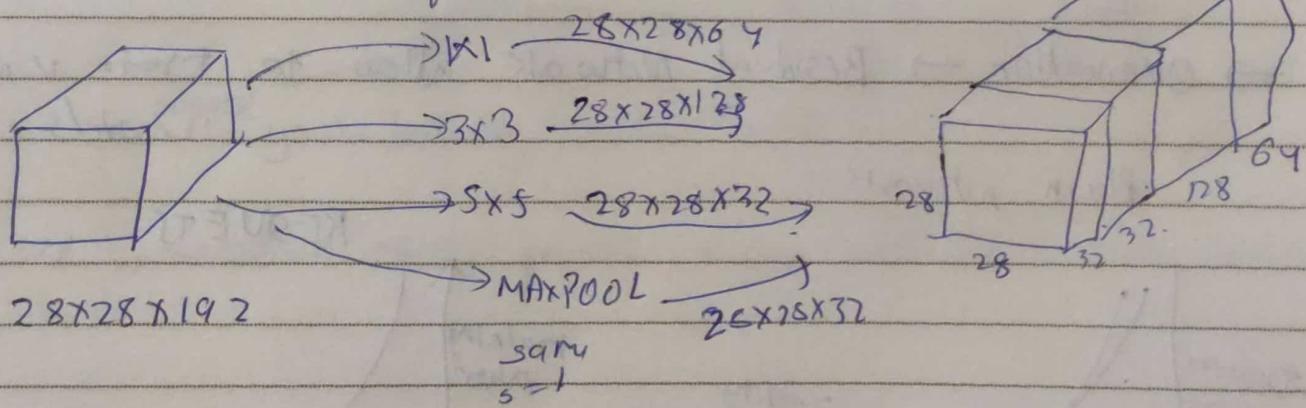
over a volume → it's like summing of entire volumes



Also called network in network.

## # Inception Network

"instead of choosing a specific filter size ( $1 \times 1$  or  $3 \times 3$ ...) just try all of them"

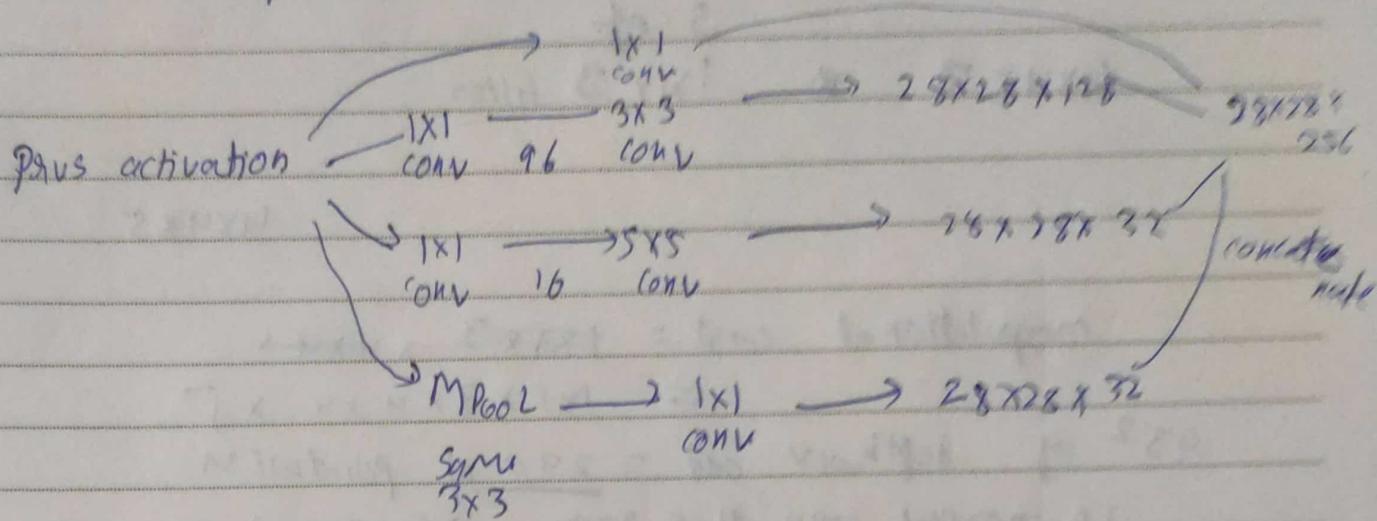


Problem → computational cost !!

$1 \times 1$  convolution can be used to reduce volume. "Bottleneck layer"

$$28 \times 28 \times 192 \xrightarrow{1 \times 1 \text{ conv}} 28 \times 28 \times 16$$

So for inception network



## # MOBILENET

→ run NN architectures on less powerful computational devices.

In a normal convolution

$$\text{computational cost} = \# \text{filter planes} \times \# \text{filter positions} \times \# \text{of filters}$$

$$= \underset{\substack{\uparrow \\ \text{input}}}{{\text{fxf}} \times \underset{\substack{\uparrow \\ \text{input}}}{{n_c}}} \times \underset{\substack{\downarrow \\ \text{input}}}{{n_h \times n_w}} \times \underset{\substack{\downarrow \\ \text{input}}}{{n_c}}$$

There is something called [Depthwise Separable Convolution]

depthwise + pointwise

depth wise.



$$* \quad 3 \times 3 \times 3 \quad = \quad 4 \times 4 \times 3$$

$$6 \times 6 \times 3$$

and all same !!

this is the intuition convolution  
DOMS

Pointwise.

$4 \times 4 \times 3$  \*  $3 \times 1 \times 3$  filters

3 of

$4 \times 4 \times 5$

$$\text{computational cost} = 4 \times 4 \times 3 \times 4 \times 4 \times 5$$

$$= 1 \times 1 \times 3 \times 4 \times 4 \times 5$$

$43^2$  of depthwise and = 290 of pointwise

in normal conv this same  $6 \times 6 \times 3 \rightarrow 4 \times 4 \times 5$

2160 computations were done

672 vs 2160

The ratio of depthwise to normal is

$$\frac{1}{n_c} + \frac{1}{f^3}$$