

Convolutional Neural Networks

Computer Vision

Application

- Image classification
- Object detection
- Neural style transfer → re-painting a picture in a different style. Inputs 2 images converts one to style of the other

→ One of the challenges of CV problems is that the inputs can get really big.

Eg for a 1000×1000 image input will have $1000 \times 1000 \times 3$ input features \times as its in RGB. If 2nd layer has 1000 units then W_1 will be (1000, 3 million) dimensional matrix. $\therefore W_1$ will have 3 billion parameters

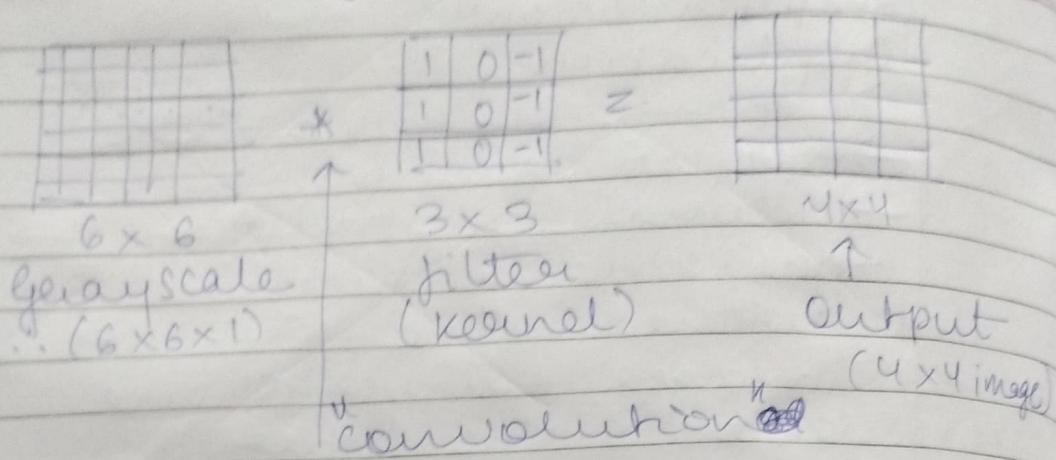
→ It's difficult to prevent overfitting also there are computational & memory requirements

→ Therefore we use convolutional networks

Edge detection

→ Convolution operation is a fundamental building block of convolutional neural networks

How edge detection works



- In math * is convolution
- Since in python * is multiplication it is overloaded

→ To calculate upper left of 4x4 what we do is take the 3x3 filter and paste it on top of original the 3x3 region of original input image

$$\begin{array}{|c|c|} \hline \text{Patch} & \text{to get this as O/P} \\ \hline \end{array} \quad * \quad \begin{array}{|c|c|} \hline \end{array} = \quad \begin{array}{|c|c|} \hline \end{array}$$

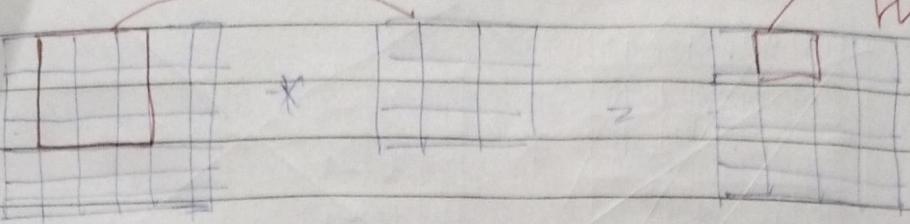
Take element wise product when we paste the filter

→ Then add up the sum of all elements

→ Move filter by 1 column to get next pixel and follow same steps

Paste

To get
this



→ and so on to fill the 4×4 matrix

→ This is an edge detector

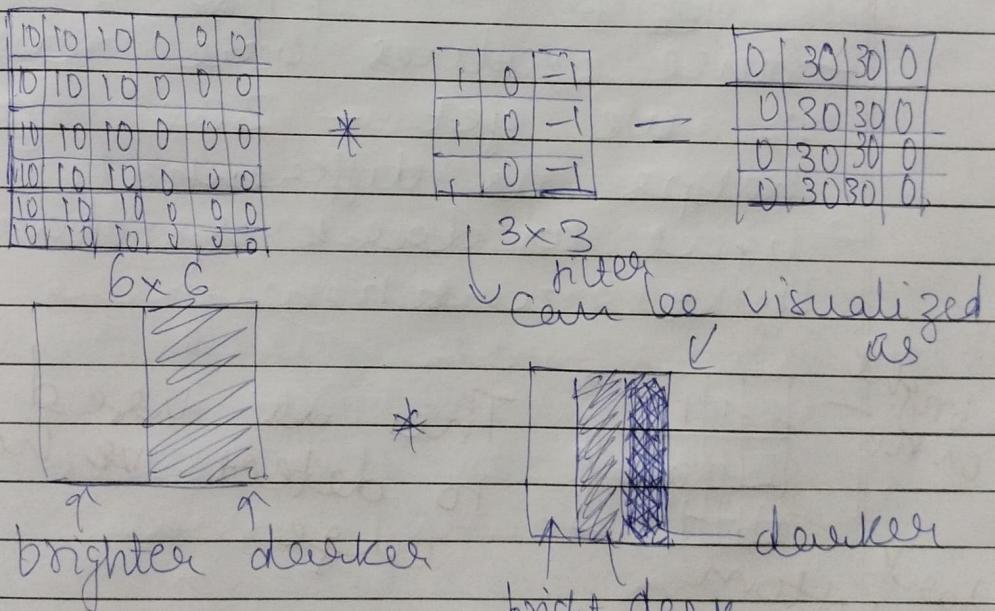
→ To denote the convolution various languages have various functions

Python: convolve2d

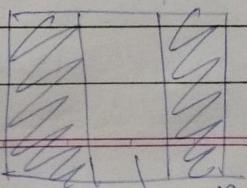
TensorFlow: tf.nn.conv2d

Keras: conv2D

Vertical edge detection



product is plotted (4×4 matrix):



Lighter region
in middle correspond
to having detected
an edge

Bright center down the middle
 6×6 image

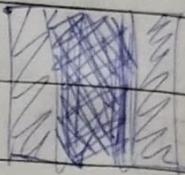
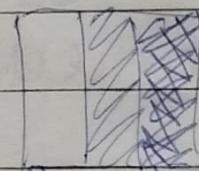
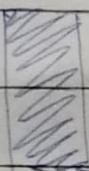
Eg 2)

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

*

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

0	-30	-300
0	-30	-300
0	-30	-300
0	-30	-300
0	-30	-300



— Since values are reversed
∴ we get -30 ~~for this~~
helps us differentiate b/w
light to dark & dark to
light transition

bright
on top
—
dark
on bottom

1	1	1
0	0	0
H	H	H

This is used
to detect ~~an~~ horizontal
edge

Ex. 3)

introduced
verticals
(copy right part
of image left
of image right)
y - ve on right

from lighter
to dark

10	10	10	0	0	0	0	0	0
10	10	10	0	0	0	1	1	1
10	10	10	0	0	0	0	0	= 30
0	0	0	10	10	10	-1	-1	-1
0	0	0	10	10	10			
0	0	0	10	10	10			

darker
to light

→ 3x3 filter is one possible choice

→ $\begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix}$ puts a little more weight to central pixel ∴ more robust Sobel filter

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

→ To detect edges in complicated image you don't need to handpick now just the filter we can treat them as parameters and learn using backpropagation to

w_1	w_2	w_3
w_4	w_5	w_6
w_7	w_8	w_9

get a good edge detector

→ This can be used not only for vertical or horizontal edges but also for inclined edges

II Padding → To build deep learning networks one modification to convolutional operation is padding

→ If we convolve an $n \times n$ image by $f \times f$ filter we get the O/P image by $(n-f+1) \times (n-f+1)$

→ We can only do convolution a few times as image shrinks & shifts every time and we don't want image to shrink every time we detect edges on it

→ Second downside is that pixels at corners or edges are used much less in output - lot of info near the edge of image is thrown away

This is solved by padding image by adding a border of 1 pixel around the image.

$\therefore 6 \times 6$ becomes 8×8

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

\rightarrow Padding is done with 0's conventionally

\rightarrow New O/P image size:

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

Valid & same convolutions

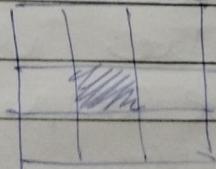
\rightarrow These are 2 common choices in terms of how much to pad

~~no padding~~

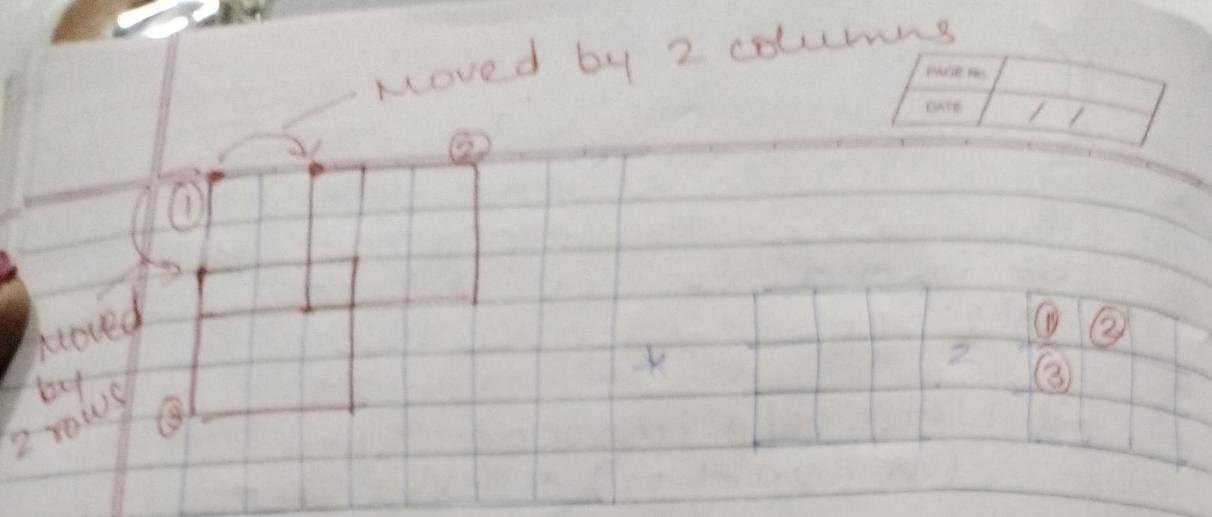
$$\text{"Valid": } n \times n * f \times f \rightarrow (n - f + 1) \times (n - f + 1)$$

\rightarrow "Same": Pad so O/P size is same as I/P size

filters are usually of odd order $f \times f$ & f is odd
 (helps give asymmetric patterns without asymmetric padding)
 also has a central position



= Strided convolutions



7×7
if stride is 2 instead of
moving by 1 column of row
we move by 2

order of O/P matrix:-

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

If this is not a whole no we round off to nearest integer

$\lfloor z \rfloor = \text{Floor}(z)$ - means taking z & rounding to nearest integer

→ It means that if the filter does not lie entirely in the image we skip that computation

→ In main to before passing the filter and doing all the steps there is one step to be done before

→ V flip the filter on horizontal axis as well as vertical axis

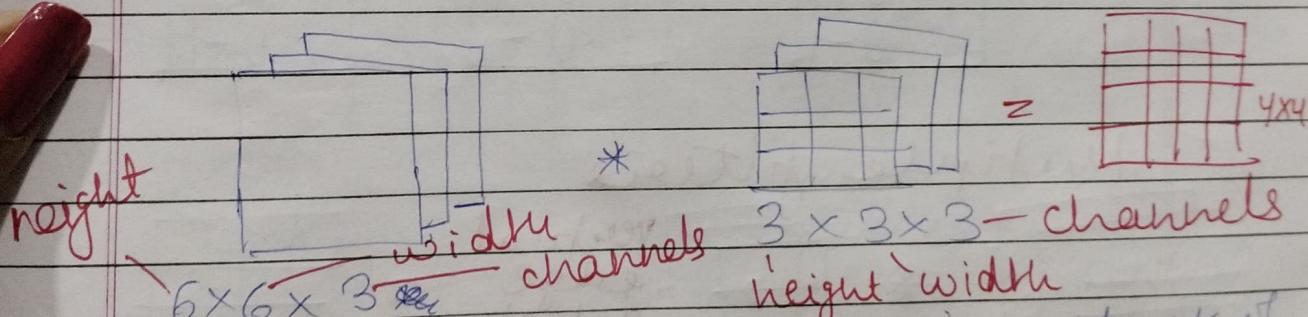
Eg

3	4	5		7	9	-1
1	0	2	→	2	0	1
5	1	9		5	4	3

- If we skip the above operation it is technically cross-correlation instead of convolution
- This operation allows associativity $(A * B) * C = A * (B * C)$
- This is not needed for convolutional neural networks so we skip it

Convolutions over volume

Convolutions on RGB images



→ An RGB image is a stack of 3 6x6 images so the filter is a stack of 3 3x3 matrices

→ channels of image of filter must be equal

- You can consider the $3 \times 3 \times 3$ filter as a 3D cube
- The $3 \times 3 \times 3$ filter has 27 parameters
- We multiply each layer of filter with corresponding R/G/B channel then add all 27 multiplication (9 for each channel)
- Each filter can have diff parameters

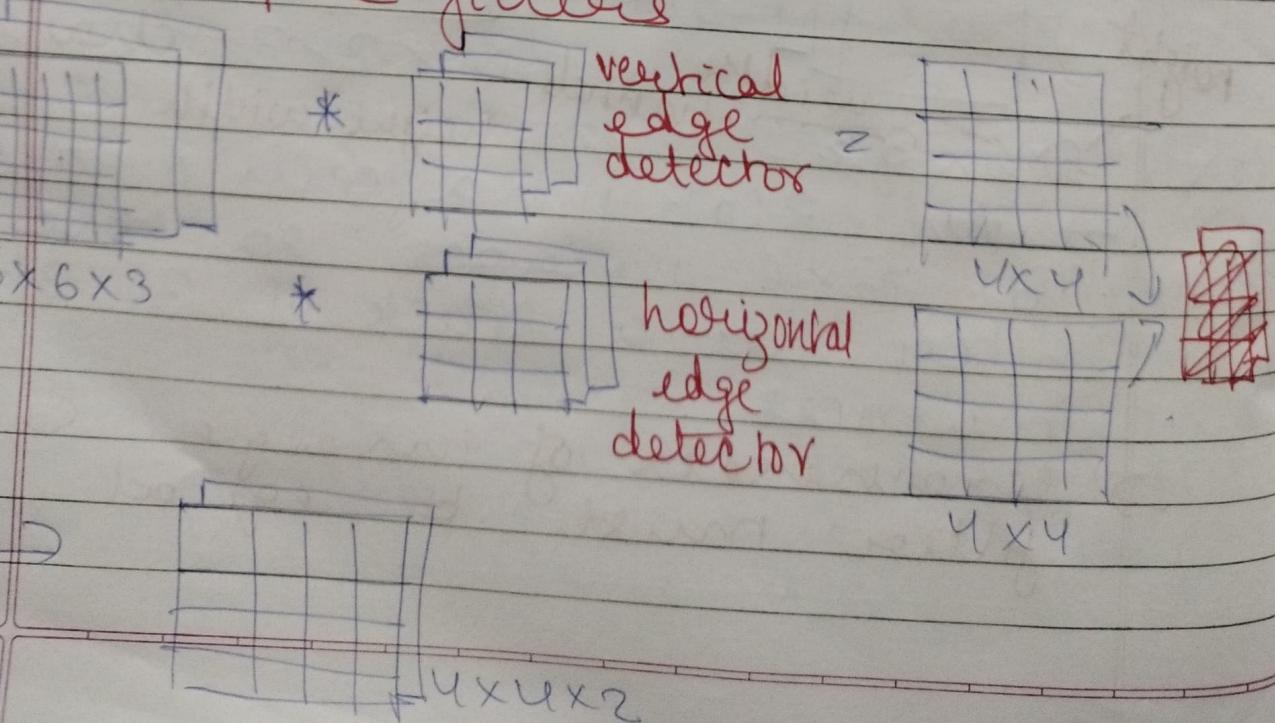
$$\text{Eq} \quad \begin{matrix} R \\ G \\ B \end{matrix}$$

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

$-3 \times 3 \times 3$

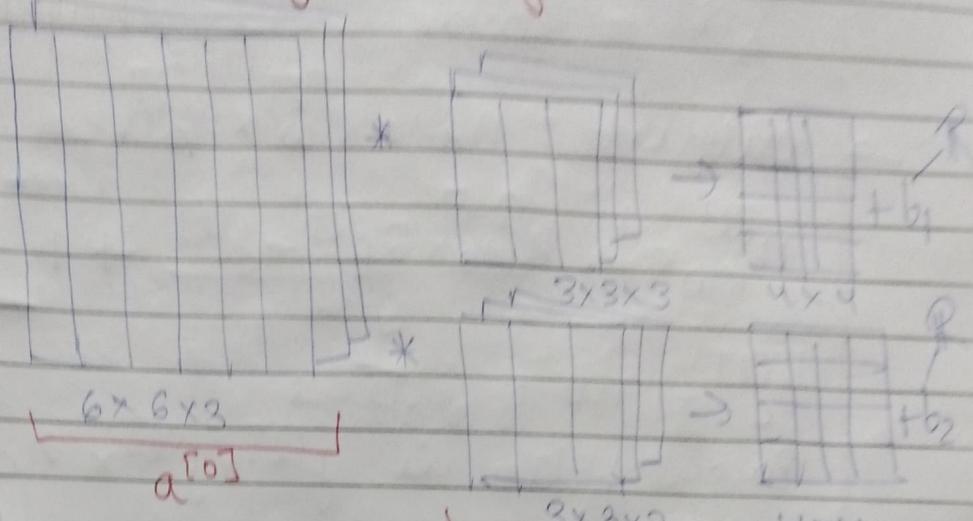
$$\begin{matrix} \begin{vmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{vmatrix} & \begin{vmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{vmatrix} & \begin{vmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{vmatrix} \end{matrix} - 3 \times 3$$

Multiple filters



→ This is used to find horizontal and vertical edges simultaneously

III Example of a layer



$$\rightarrow \text{ReLU} (\begin{matrix} \boxed{\quad} & \boxed{\quad} & \boxed{\quad} \\ \boxed{\quad} & \boxed{\quad} & \boxed{\quad} \\ \boxed{\quad} & \boxed{\quad} & \boxed{\quad} \end{matrix} + b_1) \rightarrow \begin{matrix} \boxed{\quad} & \boxed{\quad} & \boxed{\quad} \\ \boxed{\quad} & \boxed{\quad} & \boxed{\quad} \\ \boxed{\quad} & \boxed{\quad} & \boxed{\quad} \end{matrix}$$

$$\rightarrow \text{ReLU} (\begin{matrix} \boxed{\quad} & \boxed{\quad} & \boxed{\quad} \\ \boxed{\quad} & \boxed{\quad} & \boxed{\quad} \\ \boxed{\quad} & \boxed{\quad} & \boxed{\quad} \end{matrix} + b_2) \rightarrow \begin{matrix} \boxed{\quad} & \boxed{\quad} & \boxed{\quad} \\ \boxed{\quad} & \boxed{\quad} & \boxed{\quad} \\ \boxed{\quad} & \boxed{\quad} & \boxed{\quad} \end{matrix}$$

This is 1 layer of convolutional network

→ The parameters depend on size of filter & no of filters. No matter how big an image the no of parameters are unaffected. This prevents overfitting.

Summary of notation

→ if l is a convolutional layer,
 $f^{[l]}$ = filter size Input = $n_H^{[l-1]} \times n_W^{[l-1]} \times n_C^{[l-1]}$
 $p^{[l]}$ = padding height $h^{[l]}$, width $w^{[l]}$
 $s^{[l]}$ = stride

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

$$n_{H/W}^{[l]} = \left\lfloor \frac{n_{H/W}^{[l-1]} \times 2^{p^{[l]}} - f^{[l]}}{s^{[l]}} + 1 \right\rfloor$$

$n_C^{[l]}$ = no of filters
 Each filter is $f^{[l]} \times f^{[l]} \times n_C^{[l-1]}$

→ Activations : $a^{[l]} \rightarrow n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$

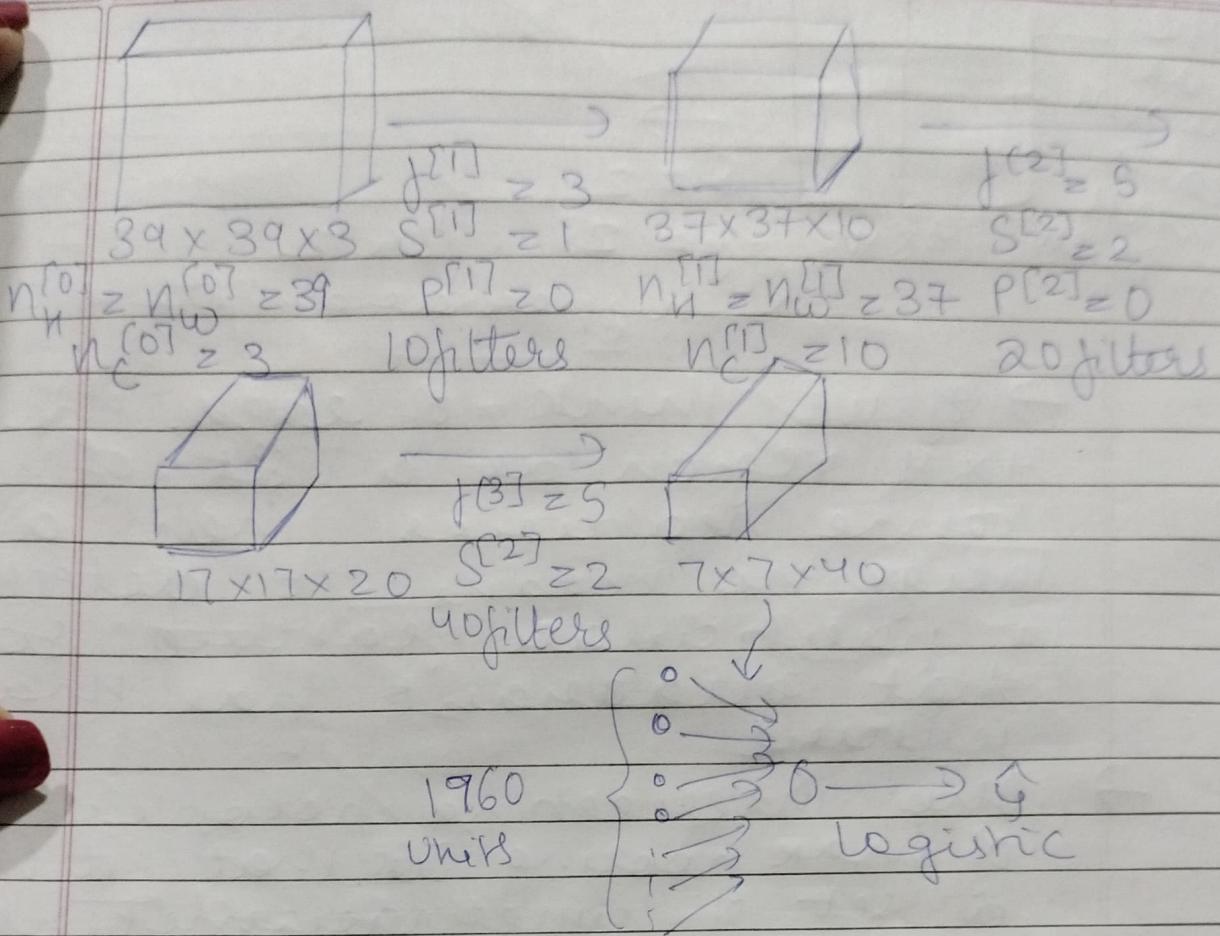
↑
training

examples

→ Weights : $f^{[l]} \times f^{[l]} \times n_C^{[l-1]} \times n_C^{[l]}$

→ bias : 1 no for
 each filter no of
 $= n_C^{[l]} = (1, 1, 1, \dots, n_C^{[l]})$ filters in
 layer l

Simple convolutional Network



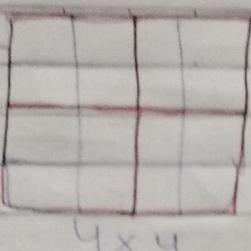
→ We usually see that as no of channels ↑ the height & width ↓

- ★ Types of layers in convolutional network
- Convolution (CONV)
- Pooling (POOL)
- Fully connected (FC)

Pooling layer :

Max Pooling

The max no from each region is taken



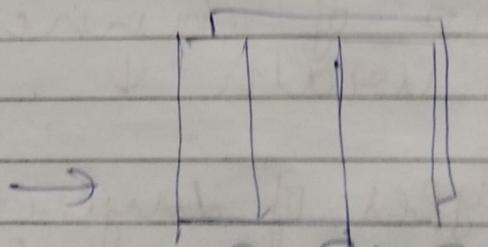
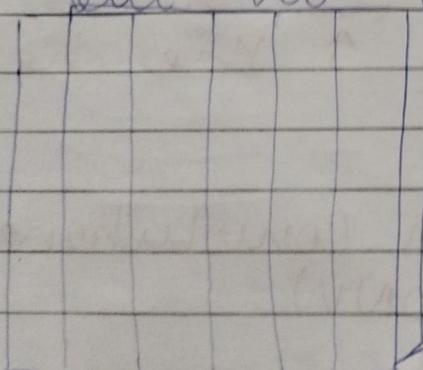
$$\begin{matrix} \rightarrow \\ S=2 \\ f=2 \end{matrix}$$

(hyperparameters)

→ if we think of each the 4x4 region as a set of some features then a large no means it detected a particular feature so that quadrant has that feature

→ One interesting property of max pooling is that it has a set of hyperparameters but no parameters to learn

Ex



$$\begin{matrix} 6 \times 8 \times 2 \\ \rightarrow \\ 3 \times 3 \times 2 \end{matrix}$$

$$f=3 \quad S=1$$

- Max pooling is done independently on each channel
- It is not used very often

Average pooling - we take avg instead of max

- sometimes in deep network we use avg pooling to collapse representation from say $4 \times 7 \times 1000$ to $1 \times 1 \times 1000$
- Max pooling is used more than avg

→ Hyperparameters for Pooling:

f : filter size ($f=2, s=2$
 $s = \text{stride}$ are common)

Max or Avg

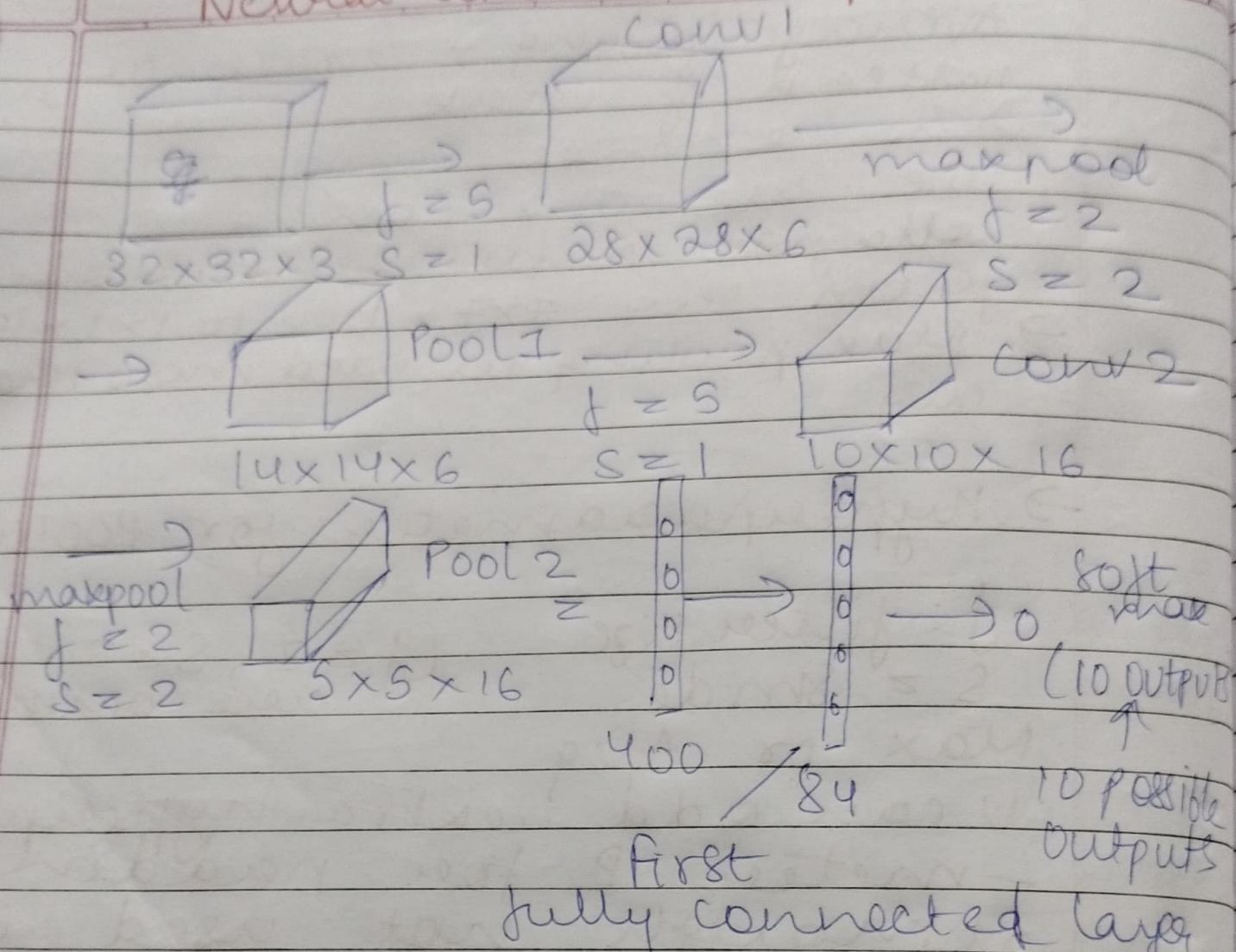
we can add extra hyperparameter - metter P for padding although not used in pooling mostly.

$$n_h \times n_w \times n_c \rightarrow \left[\frac{n_h-f+1}{s} \right] \times \left[\frac{n_w-f+1}{s} \right] n_c$$

Neural Network Example

PAGE NO.

DATE



Conv1 to Pool1 is layer 1

Conv2 to Pool2 is layer 2

Conv → Pool → Conv → Pool → FC → FC

→ FC → SOFTMAX

- Max pooling layers don't have any parameters
- Size (ie no of activations in a layer) yes as network progresses

Why convolutions

- No of parameters to train is less as they only depend on filter size (and 1 bias) and no of filters
- Reasons for less Parameters
 - 1) Parameter sharing -
A feature detector (such as vertical edge detector) that's useful in one part is also useful in another part of the image (e.g. the same filter is pasted on different sections of the image)
 - 2) Sparsity of connections -
In each layer each O/P value depends only on small no of inputs (value of pixel to decide in O/P layer only depends on those of S/I P layer used to derive it)