

Convolution Neural Network

Page No.

Date

CNN

- Image Recognition
- Object Detection & Localization
- Semantic Segmentation
- Medical Image Analysis

Application in Computer Vision & Image Analysis

for CNN we apply Image as input, different types of Image Gray Scale, RGB.

Input Image \rightarrow 2D matrix of pixels having particular Numerical value (8 bit \rightarrow 0 to 255) ($n_x \times n_y \rightarrow$ Input size)

$[n_x \times n_y \times 3]$ - Input size
A channels

CNN layers & operations

Q. What CNN consists of?

Like ANN, CNN contains stack of sequence layers followed by output layers - classifier layer

Each layer performs a Convolution operation or a Pooling operation

The conv. layer & pooling layers are alternatively stacked leading to a series of fully connected layers followed by an output layer.

Conv. Layer

(feature MAP / Activation MAP)

A G B

Output Volume

b/p have stack of feature maps

If IP don't contain more than 3 channels then also

Andrew Ng

Vertical Edge detection:

IP	filter	output	OIP
6x6	3x3	4x4	4x4
3 0 1 2 7 4 1 5 8 9 3 2 2 7 2 5 1 3 0 1 3 2 7 8 4 2 1 6 2 8 2 4 5 2 3 9	convolution	$\begin{array}{ c c c c } \hline & + & 0 & -1 \\ \hline + & & 0 & -1 \\ \hline & + & 0 & -1 \\ \hline \end{array}$	$\begin{array}{ c c c c } \hline -5 & -4 & 0 & 8 \\ \hline -10 & -2 & 2 & 3 \\ \hline 0 & -2 & -4 & -7 \\ \hline -3 & -2 & -3 & -16 \\ \hline \end{array}$
(3*1) + (0*0) + (1*-1) + (1*1) + (5*0) + (8*-1) + (2*1) + (7*0) + (2*-1) = 3 + 0 - 1 + 1 + 0 - 8 + 2 + 0 - 2 = -5 at (1,1) in OIP.			
(0*1) + (1*0) + (2*-1) + (5*1) + (8*0) + (9*-1) + (7*1) + (2*0) + (5*-1) = 0 + 0 - 2 + 5 + 0 - 9 + 7 + 0 - 5 = -4 at (1,2) in OIP			
(1*1) + (2*0) + (7*-1) + (8*1) + (9*0) + (3*-1) + (2*1) + (5*0) + (1*-1) = 1 + 0 - 7 + 8 + 0 - 3 + 2 + 0 - 1 = 0 at (1,3) in OIP			

similarly do all of this shift column by column IF Column end then shift ↓ row down side, do same process for this also
 OIP gives 4x4 Matrix

for Implementing CNN convolution layer

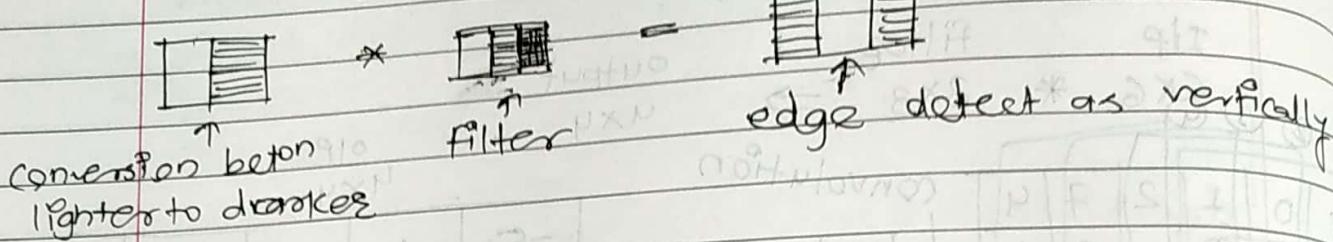
In python: conv-forward

tensorflow: tf.nn.conv2d

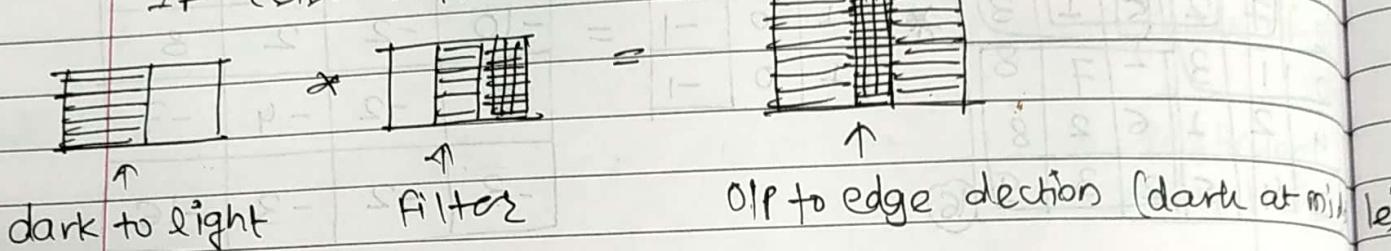
keras: conv2D

$> I - \text{white}$ \rightarrow darker $I - \text{lighter}$ \rightarrow lighter $I - \text{more dark}$ (black)

Vertical Edge detection.



IF color flipped



(1) So different filters allow to find out vertical as well as horizontal edges

1	0	-1
2	0	-2
1	0	-1

Advantage

3	0	-3
10	0	-10
3	0	-3

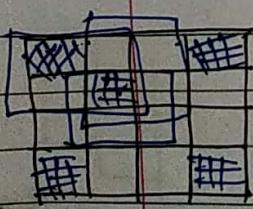
Sobel filter \rightarrow more robust

\leftarrow Sobel filter

Vertical & Horizontal edge detection

Padding

Input dimension $(n \times n)$ Img \star Filter dimension $(F \times F)$ = Output Img dimension $(n - F + 1) \times (n - F + 1)$



4 or 8 corner pixels Images get evaluated once means it gives much less olp if middle pixels give evaluated much more time so to evaluate this problems

you have to know both the shrinking outputs &

Throwing away information from edges (corners) of img

To fix both of these problems you can before applying convolution operation you can pad the image, with additional one border of one pixel around the edges, so you can do that then (6x6 img) $\xrightarrow{\text{padded to}} (8 \times 8 \text{ img})$

In convolution

(8x8) img * (3x3) filter $\xrightarrow{\text{get}} (6 \times 6) \text{ img o/p}$

so you padded 0's

$\therefore p$ is padding amount

i.e. $p=1$ in this case

you add extra 1 pixel around img

0	0	0	0	0	0	0
0	-	-	-	-	-	-
0	-	-	-	-	-	-
0	-	-	-	-	-	-
0	-	-	-	-	-	-
0	-	-	-	-	-	-
0	0	0	0	0	0	0

In o/p become $\rightarrow (n+2P-f+1) \times (n+2P-f+1)$

$$\text{Ex} - (6+2(1)-3+1) \times (6+2(1)-3+1)$$

$$= (6+2-2) \times (6+2-2)$$

$$= (6) \times (6) \text{ Imp o/p}$$

Valid & same convolutions -

Valid \rightarrow No padding

$$(n \times n) \text{ img} * (f \times f) \text{ filter} \xrightarrow{\text{o/p}} (n-f+1) \times (n-f+1)$$

same \rightarrow Pad so that output size is the same as the input size then,

$$(n+2P-f+1) * (n+2P-f+1) \Rightarrow n \times n \text{ then,}$$

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

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

$$F = 3$$

$$\text{Ex} - p = \frac{3-1}{2} = \frac{2}{2} = 1$$

\downarrow Pad around image

usually, F is always odd.

same padding gives you a natural padding as same image at ($I/P = O/P$).

Strided Convolutions

It is another basic building block of convolutions as used in CNN.

In this when we slide the window of convolution, instead of sliding one by one slide 2 steps so, we can skip over two steps like so on for rows as well as columns.

So, we get (3×3) O/P Imp. {less dimension}

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

Padding P

$f = P = 0, s = 2$

$$(7 \times 7) * (3 \times 3) \Rightarrow \left[\frac{7+2(0)-3}{2} + 1 \right] \times \left[\frac{7+2(0)-3}{2} + 1 \right]$$

Take floor

$$= (2+1) \times (2+1)$$

$$= 3 \times 3$$

Technical Notes on Cross-correlation Vs Convolution

- Most convolution technicians before applying convolution flip the matrix of filter first then apply convolution operation.

But to simplify code

3	4	5
1	0	2
7	9	7

7	2	5
9	0	4
-1	1	3

flip

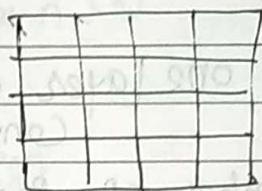
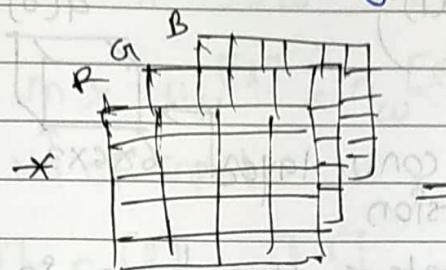
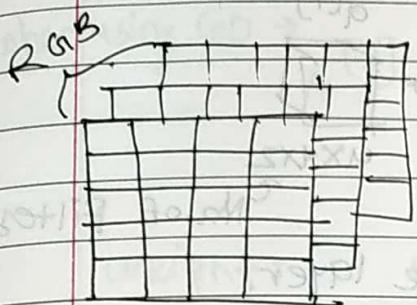
Cross-correlation

In CNN we omit this operation. & done convolution

(should not effect anything)

Convolutions over Volume

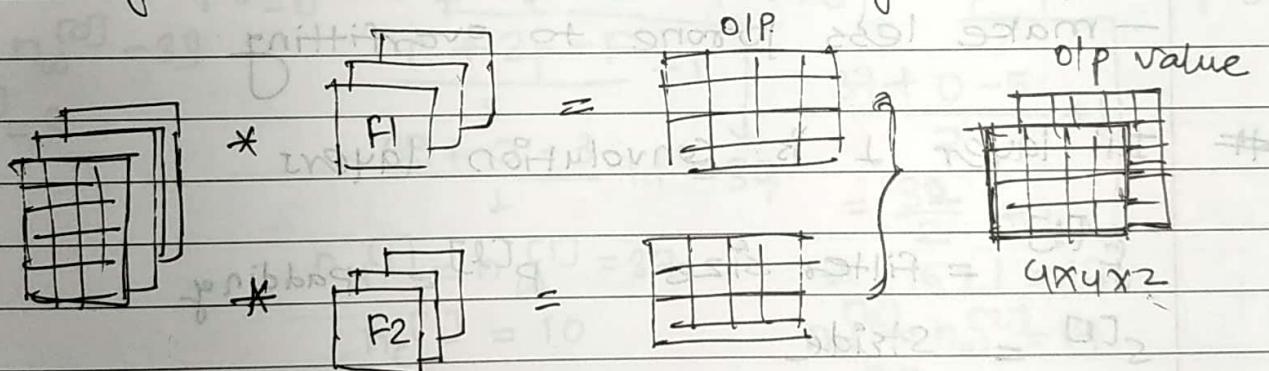
Convolution on RGB Image -

for 1st side small 27 No. = get single cell value in o/p

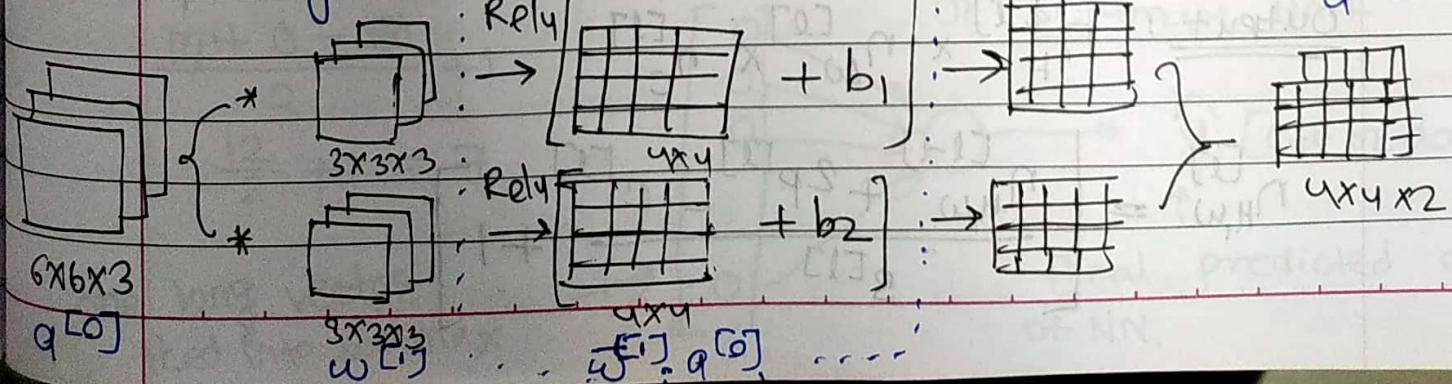
Same for all cells.

Multiple filters -

Apply more than one inputs filters to the inputs.

we get that much of uxy o/p Img.

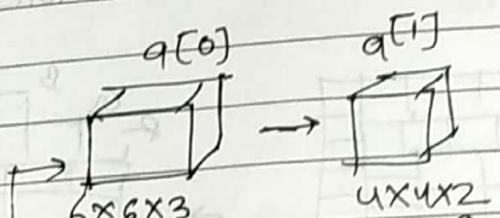
$$n \times n \times [n_c] * f \times f * [n_c] \rightarrow (n-f+1) \times (n-f+1) \times n_c$$

one layer of CNN $Z^{[0]}$ 

$$z^{[l]} = w^{[l]} \cdot a^{[l-1]} + b^{[l]}$$

$$a^{[l]} = g(z^{[l]})$$

one layer of conv. layer
Conversion



u x u x 2
No. of Filters

Number of parameters you have in one layer.

Q. If you have 10 filters that are 3x3x3 in one layer of NN, how many parameters does that layer have?

\rightarrow filters $(3 \times 3 \times 3) = 27 + \text{bias} = 28 \times 10 \text{ filters} = 280 \text{ parameters}$

- No matter how big your IP Img is \rightarrow OLP Parameters is 280 only
This detects features of IP Img

property of CNN
- make less prone to overfitting

If layer l is convolution layers

$f^{[l]}$ = filter size $p^{[l]}$ = padding

$s^{[l]}$ = stride

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

H - Height W - width C - channels

IP to l layers has OLP from l-1 Layer

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

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

$n_c^{[l]}$ - No. of Filters Filter $\rightarrow F^{[l]} \times F^{[l]} \times n_c^{[l-1]}$

Activations = $a^{[l]} \rightarrow n_H^{[l]} \times n_W^{[l]} \times n_c^{[l]}$

when using GD \rightarrow

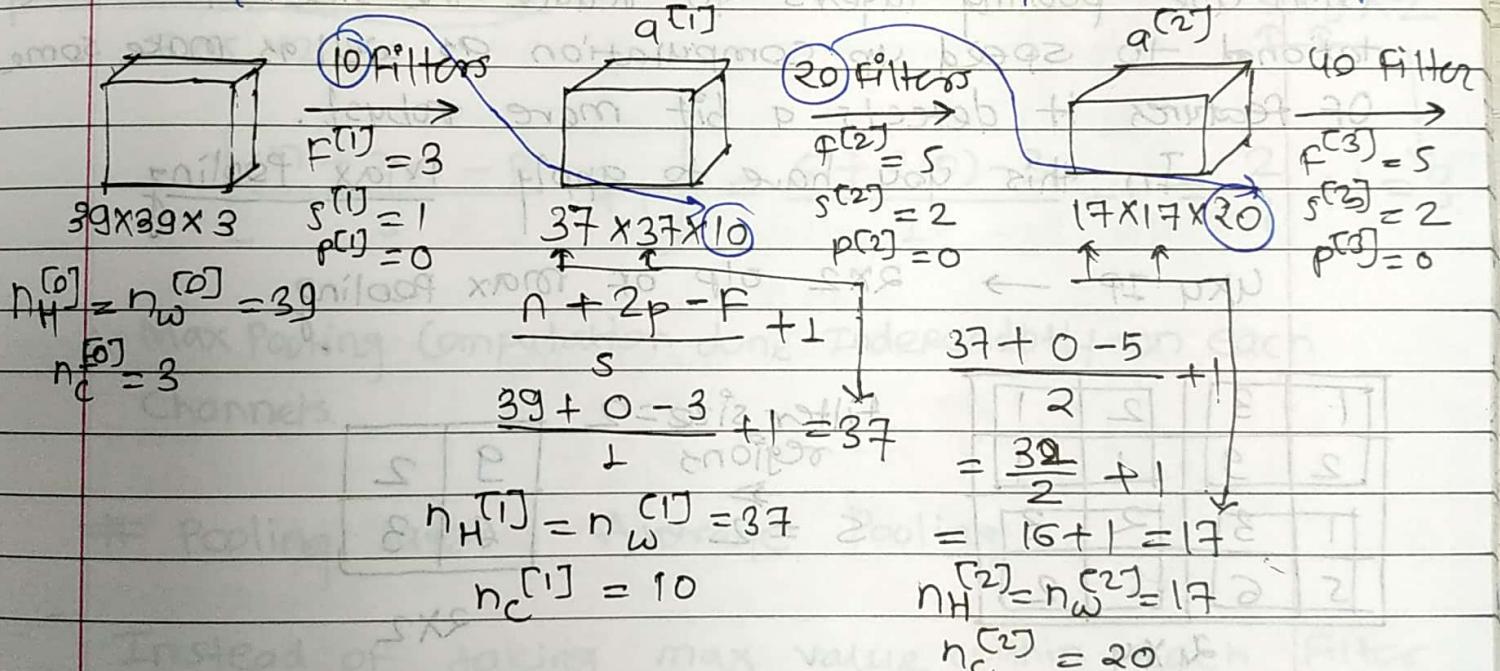
$$A_H^{[l]} = m \times n_H^{[l]} \times n_W^{[l]} \times n_c^{[l]}$$

\uparrow
No. of Input

$$\text{weights} = f^{[l]} \times f^{[l]} \times n_c^{[l-1]} \times n_c^{[l]}$$

$$\text{bias} = n_c^{[l]} - [1, 1, 1, n_c^{[l]}] \quad \# \text{No. of filters in layer}$$

simple Convolution Network Example. [cat or dog]

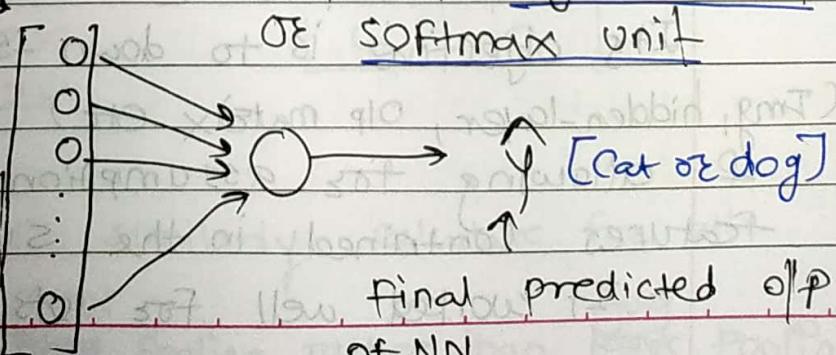


Now take this volume (1960) size & flatten it (unroll it)

$7 \times 7 \times 40 = 1960 \rightarrow$ into & fit into logistic unit.

$$\frac{17 + 0 - 5}{2} + 1 = 7$$

$$\frac{12}{2} + 1 = 7$$



one long vector
feed into softmax

as you go deeper in NN
 Img $(39 \times 39) \rightarrow (37 \times 37) \rightarrow (17 \times 17) \rightarrow (7 \times 7)$

& No of channels increasing as from
 $3 \rightarrow 10 \rightarrow 20 \rightarrow 40$

Types of layers in convolutional NN -

- Convolution Layer ✓
- Pooling Layer
- Fully connected layer

Pooling Layer - Max Pooling

other than convolutional layers CNN often also use pooling layers to reduce the size of representation to speed up computation as well as make some of features it detects a bit more robust.

In this you have to apply max pooling

4×4 IP \rightarrow 2×2 OLP of max pooling

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

2×2

filter size = 2
regions
$S = 2$

9	2
6	3

2×2

Take IP break into diff. regions each of this output will be just max from correspondingly shaded region.

maxpooling hyperparameters

$$f = 2, S = 2$$

Its objective is to down-sample an input representation (Img, hidden-layer, OLP matrix etc) reducing its dimensionality & allowing for assumptions to be made about features contained in the sub-regions.

- It worked well for lots of examples

property of max pooling -

- It has no parameters to learn right there's actually nothing for gradient descent to learn.
- Once you fix 'f' & 's' then GD doesn't change anything

Ex -

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

$$F = 3$$

$$S = 1$$

9	9	5
9	9	5
8	6	9

$3 \times 3 \times 2$

$5 \times 5 \times 2$

$$\left\lfloor \frac{n+2p-F}{s} + 1 \right\rfloor = \frac{5+2(0)-3}{1} + 1 = \frac{2}{1} + 1 = 3.$$

Max Pooling Computation done Independently on each channels.

Pooling layer : Average Pooling

Instead of taking max value within each filter you take the average

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

Hyperparameters

$$F = 2$$

$$S = 2$$

3.75	1.25
4	2

$$\frac{1+3+2+9}{4} = \frac{15}{4} = 3.75 \text{ similarly calculate for all}$$

Now a days we use Avg Pooling rather than ~~Max Pooling~~

Except, sometimes very deep in NN you might use Avg. pooling

Max pooling use much more in NN than Avg. pooling

Hyperparameters of Pooling

f: filter size s: stride

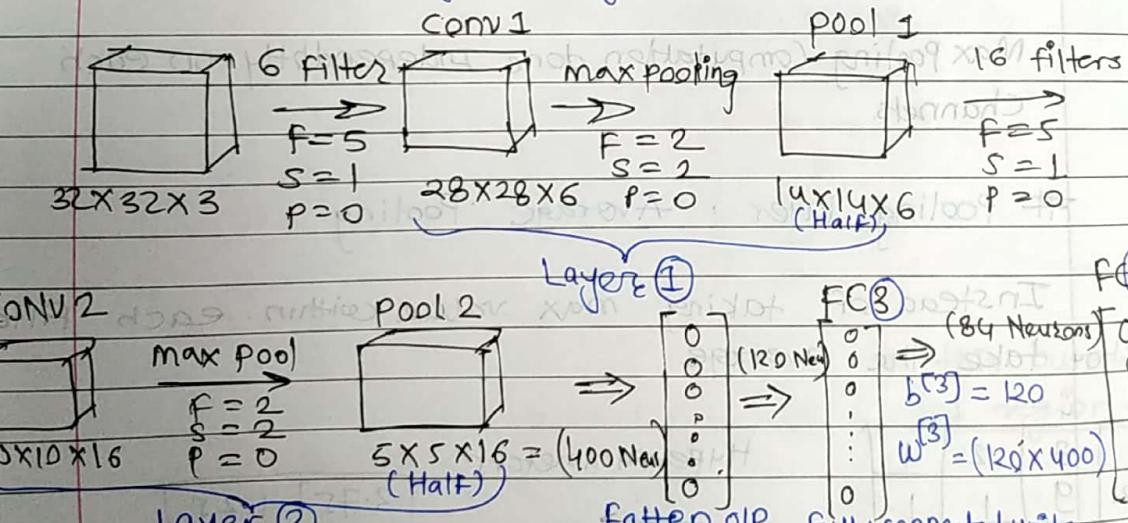
Max or Avg. pooling

p: padding (mostly p=0)

$$\text{Input} - n_H \times n_W \times n_C \rightarrow \text{Olp} \rightarrow \left[\frac{n_H - F}{s} + 1 \right] \times \left[\frac{n_W - F}{s} + 1 \right] \times n_C$$

No parameters to learn \rightarrow when backpropagate No parameters to backprop. adapt through max pooling
i.e. Nothing to learn \rightarrow fix parameters

NN. fx - Recognize digit (0 to 9)



Generally the layer which has OLP as weight & parameters then calculate (said) layer but Pooling don't have any parameters & weight so (conv1 + pooling) \rightarrow layer 1 called.
Layers = layers that have weight.

Fully connected
Fully connected each of the Neurons just like that you saw in you have weight also have

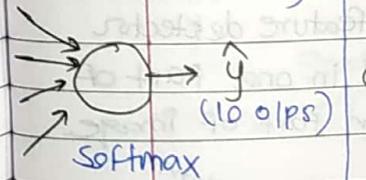
As go [32 \rightarrow where [3 \rightarrow]

Pattern [conv]

[see in]

NN

Finally recognize OLP [0 to 9]



Fully Connected Layer

Fully Connected :-

each of the 600 Neurons connected to each of 120

Neurons just like the Single Neural Network layer
that you saw in cases 1 & 2 Standard NN, where
you have weight matrix

also have Bias parameter.

As go deep you $n_h \times n_w$ decreases[32 \rightarrow 28 \rightarrow 14 \rightarrow 10 \rightarrow 5]

where No. of channels will increase

[3 \rightarrow 6 \rightarrow 16 \rightarrow Fully connected layer]

Pattern of NN (common pattern)

[conv + pool + conv + pool ... + FC + FC ... + softmax]

to see in NN

NN Ex -

	Input	Activation shape	Activation size	Params
	Input \rightarrow	(32, 32, 3)	3,072	0
finally recognize obj [0 to 9]	conv 1 (F=5, S=1)	(28, 28, 6)	6,272	208
	pool 1	(14, 14, 6)	1,568	0
	conv 2 (F=5, S=1)	(10, 10, 16)	1,600	416
	pool 2	(5, 5, 16)	400	0
	FC 3	(120, 1)	120	48,001
	FC 4	(84, 1)	84	10,081
	softmax	(10, 1)	10	841

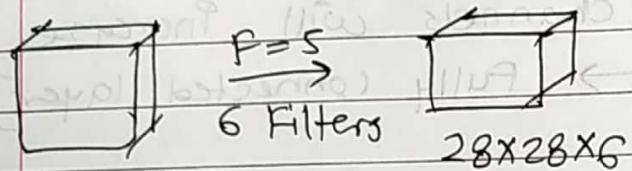
Activation size decrease gradually as you go in deeper \rightarrow finally softmax o/p

Why Convolutions are so useful?

Advantages of CNN -

① Parameters sharing

② Sparsity of connections



$32 \times 32 \times 3$

$P=2$

6 Filters

$28 \times 28 \times 6$

3072

4704

multi
NN

3072

$$W = 3072 \times 4704$$

$$W \approx 14M$$

Not so that is lots of parameters to train

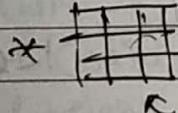
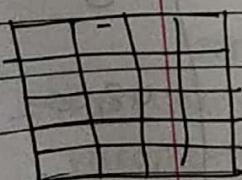
$$\text{Conv} \Rightarrow (5 \times 5) \text{ filters} = 25 + 1 \text{ bias} = 26$$

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

Reason for Conv NN have small parameters

① Parameter sharing

- motivated by observation that a feature detector such as vertical edge that's useful in one part of the image is probably useful in another part of image



Some filter (parameters) are used for whole image for feature extraction. There is probability that the upper left corner has same filters (parameters). Used & same as lower bottom filters (parameters), that's why having less parameters

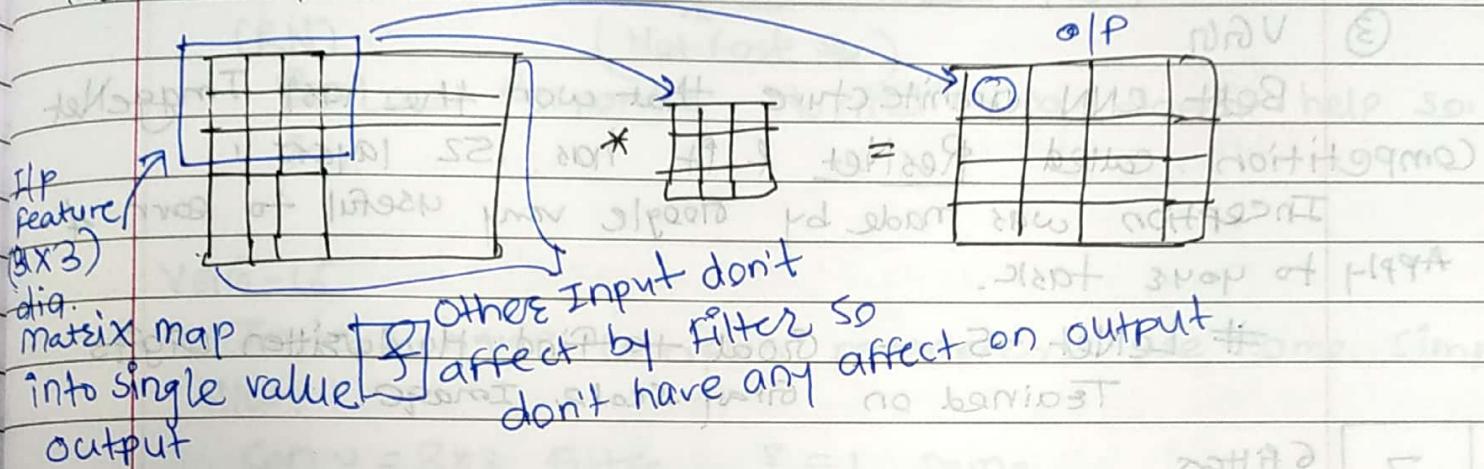
In each layer, each output value depends only on a small number of inputs which makes it translation invariance

Page No.

Date

② Sparsity connection:

In each layer, each o/p value depends only on a small no of inputs.



that why called sparse of connections.

(loosely connected). so

through this two mechanism a CNN has a few parameters which allows it to be trained to smaller training sets & less prone to be overfitting.

After putting together.

fully connected NN at the end has weight & bias also, & works same as other NN well.

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

use gradient descent to optimize parameters to reduce $J(\text{cost function})$