

Course 4: CNN

WEEK 1

Computer Vision Problems:

- Image classification
- Object detection
- Neural style transfer

One of the problem is that the inputs will be a lot. If we have a 1000×1000 px image then it'll be $1000 \times 1000 \times 3$ (RGB) input values = 3M
So ~~the~~ if the 1st hidden layer has 1000 hidden values, then $w^{[1]}$ will be $[1000, 3M]$ which is 3 billion values!

CNN: Edge detection example,

We will use convolution to get the vertical edge of an image.

3	0	1	2	7	4
1	5	8	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

6x6
image

convolution

↓

*

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

3x3
filter/
kernel

=

-5	-4	0	8
-10	-2	2	3
0	-2	-4	-7
-3	-2	-3	-16

4x4
output image

For -5,

3	0	1
1	5	8
2	7	2

$$3 \times 1 + 1 \times 1 + 2 \times 1 + 0 \times 0 + 5 \times 0 + 7 \times 0 + 1 \times -1 + 8 \times -1 + 2 \times -1 = -5$$

For -4,

0	1	2
5	8	9
7	2	5

$$0 \times 1 + 5 \times 1 + 7 \times 1 + 1 \times 0 + 8 \times 0 + 2 \times 0 + 2 \times -1 + 9 \times -1 + 5 \times -1 = -4$$

For -10,

1	5	8
2	7	2
0	1	3

$$1 \times 1 + 1 \times 2 + 0 \times 1 + 0 \times 5 + 7 \times 0 + 1 \times 0 + 8 \times -1 + 2 \times -1 + 3 \times -1 = -10$$

This is how its done.

Below is another example where you can see vertical edge detection

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	0	0	0

6x6

*

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

3x3
Vertical

=

0	30	30	0
0	30	30	0
0	30	30	0
0	30	30	0

4x4

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

6x6

*

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

3x3
Vertical

=

0	-30	-30	0
0	-30	-30	0
0	-30	-30	0
0	-30	-30	0

4x4

Similarly

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

is horizontal filter

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

\times

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

$=$

0	0	0	0
30	10	-10	-30
30	10	-10	-30
0	0	0	0

\therefore Different filters help us to find horizontal and vertical edges

There are other filters,

1	0	-1
2	0	2
1	0	-1

Sobel filter

3	0	3
10	0	-10
3	0	-3

Scharr filter

We can even learn it as parameters using back propagation

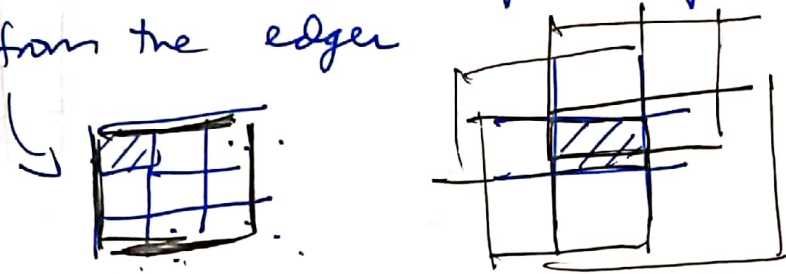
w_1	w_2	w_3
w_4	w_5	w_6
w_7	w_8	w_9

And it'll learn the best edge filter for that image data

Padding

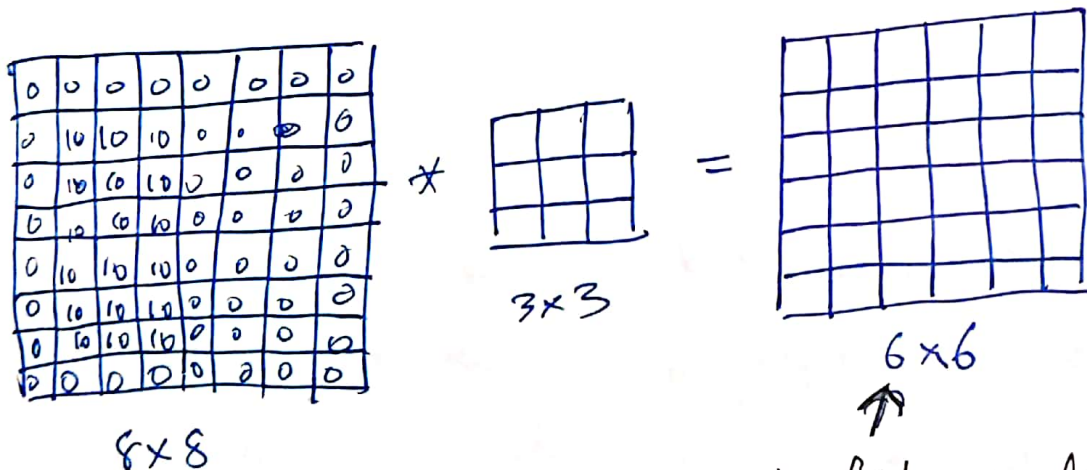
There are ~~two~~ 2 problems with the previous method:

- 1) The output shrinks from 6×6 to 4×4
- 2) We are throwing away information from the edges



If it's at edge we only consider it in 1 box, but if it's at centre, we consider it in more boxes

The solution is padding (padding 0s on all sides so it becomes 8×8)



Here padding, $p = 1$

How to choose p :

- Valid convolution - no padding

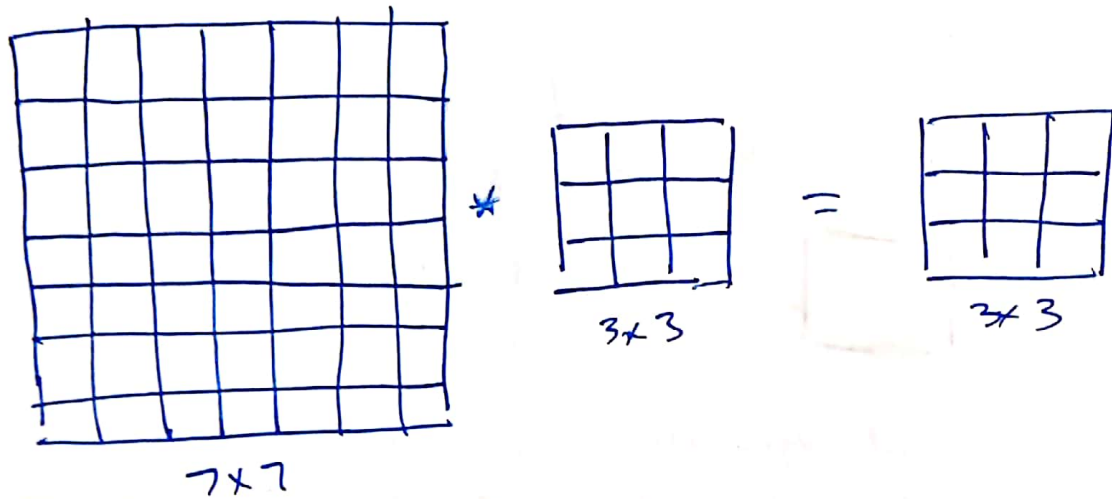
- Same convolution - $p = \frac{f-1}{2}$ where f is filter size (usually f is odd)

$$\begin{matrix} \text{input} & \text{filter} \\ \text{size} & \text{size} \end{matrix} \quad n-f+1 \times n-f+1$$

$$n \times n * f \times f \rightarrow n-f+1 \times n-f+1$$

Strided Convolutions

Instead of ~~shift~~ shifting (striding) the square by 1, you can select the stride s .



Summary:

$n \times n$ image, $f \times f$ filter, padding p , stride s

output size:

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

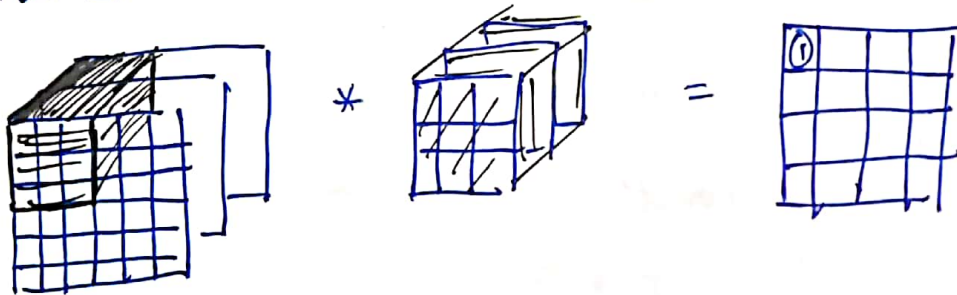
— floor value

In actual convolution, we flip the filter before multiplying, so we are actually doing cross-convolution (no flipping) but we just call it convolution.

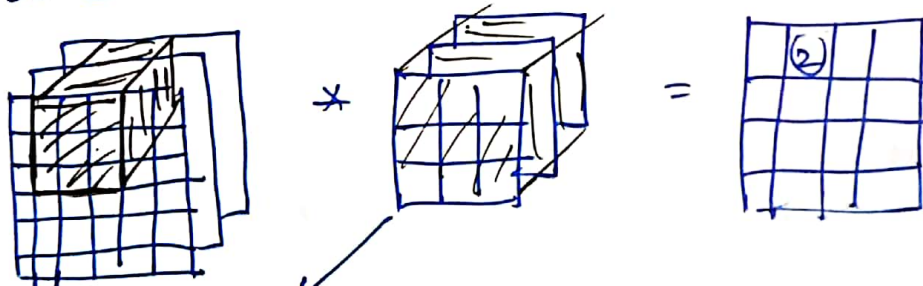
convolutions over RGB image

Diagram illustrating the multiplication of three 6×6 matrices, resulting in a 4×4 matrix. The first matrix is labeled $6 \times 6 \times 3$. The second matrix is labeled $3 \times 3 \times 3$. The result is a 4×4 matrix, with the first two columns circled and numbered 1 and 2 respectively.

For ①

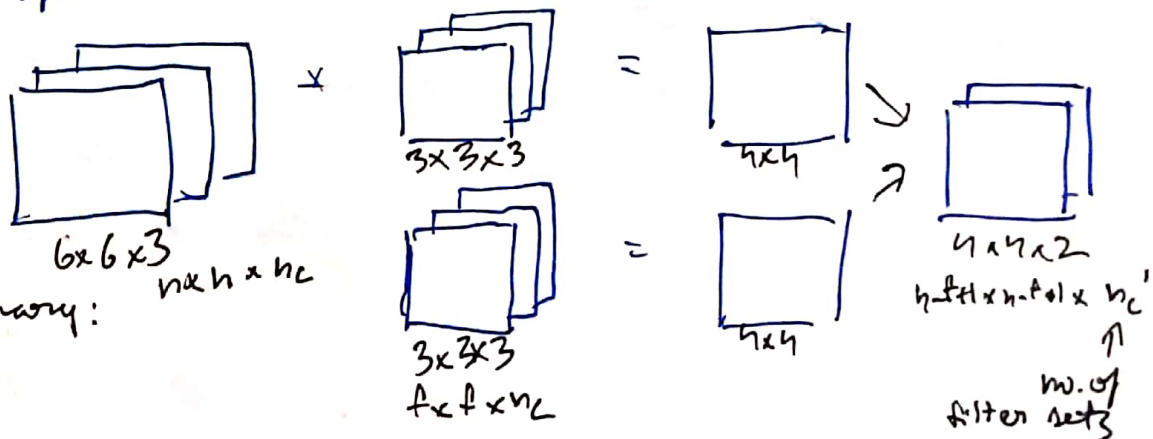


For (2)

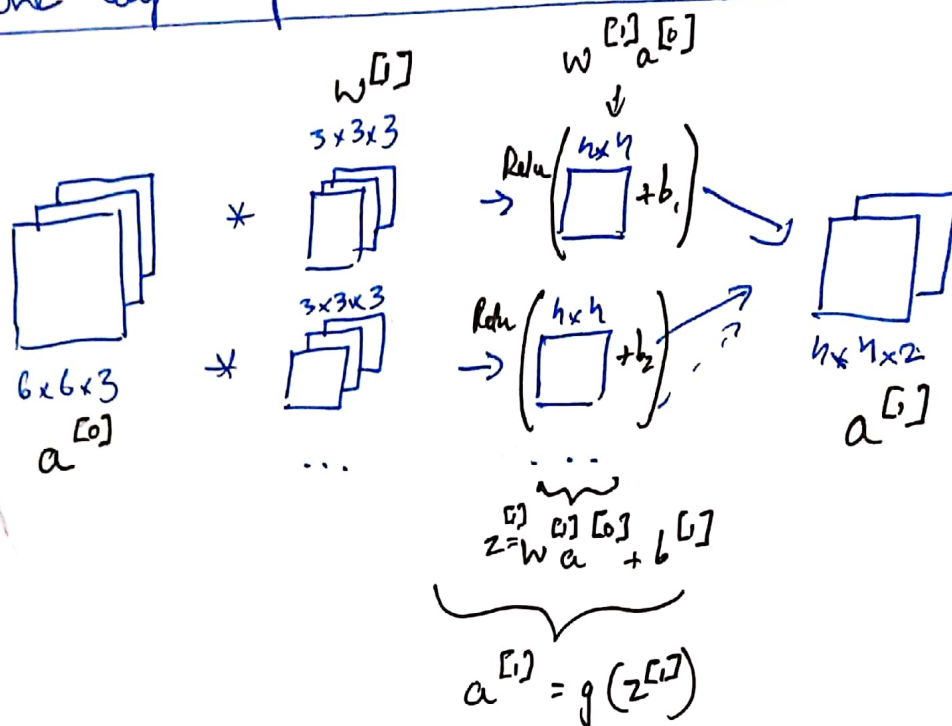


only check edges in red channel
multiple filters

Multiple Filters



One layer of a Convolutional Network



Number of parameters in one layer:

The no. of parameters remain the same for a NN. For example, for 10 filters that are $3 \times 3 \times 3$,

$$1 \text{ filter} = 3 \times 3 \times 3 = 27$$

$$+ 1 \text{ bias} = 27 + 1 = 28$$

$$10 \text{ filters} = 28 \times 10 = \underline{280 \text{ parameters}}$$

Notation summary:

If layer l is a convolution layer:

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

$n_c^{[l]}$ = no. of filters, input: $n_H^{[l-1]} \times n_W^{[l-1]} \times n_C^{[l-1]}$

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

where

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

- floor value

Each filter is : $f^{[L]} \times f^{[L]} \times n_c^{[L-1]}$

Activations : $a^{[L]} \rightarrow n_H^{[L]} \times n_W^{[L]} \times n_c^{[L]}$

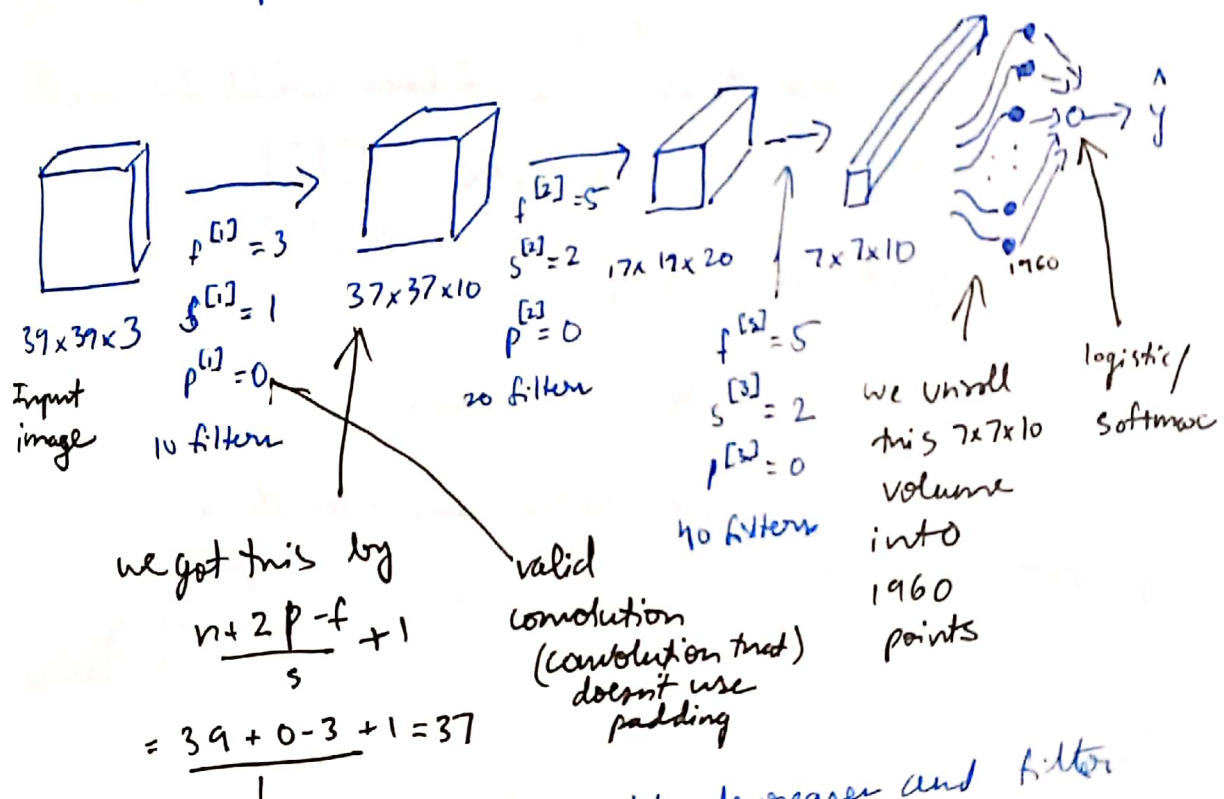
$A^{[L]} \rightarrow m \times n_H^{[L]} \times n_W^{[L]} \times n_c^{[L]}$

Weights : $f^{[L]} \times f^{[L]} \times n_c^{[L-1]} \times n_c^{[L]}$

↑
#filter in layer L

bias : $n_c^{[L]}$ but in programming we use
 $(1, 1, 1, n_c^{[L]})$

Example of a convolutional network



→ Notice how the height & width decreases and filter size increases, this is a usual trend

Types of layer in a convolutional network:

- Convolution (conv) - we have discussed this
- Pooling (pool)
- Fully Connected (FC) } we will discuss this now

Pooling

It is a fixed function layer. No parameters are learnt here.

→ Max pooling - pick the largest no. in the window
Here largest is 9


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

→
 $f = 2$

$s = 2$

↑

9	2
6	3

Since we use this, the filter will be 2×2
and we will move like 

∴ output is 4×4

Another example (for all channels):

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

$5 \times 5 \times n_c$

↑

no. of channels

we do the same for other layers

9 is max

→

$f = 3$

$s = 1$

9	9	5
9	9	5
8	6	9

$3 \times 3 \times n_c$

There is also average pooling - where the average is taken instead of max. But this isn't used usually.

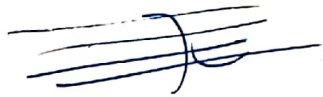
Hyperparameters for pooling:

→ f : filter size

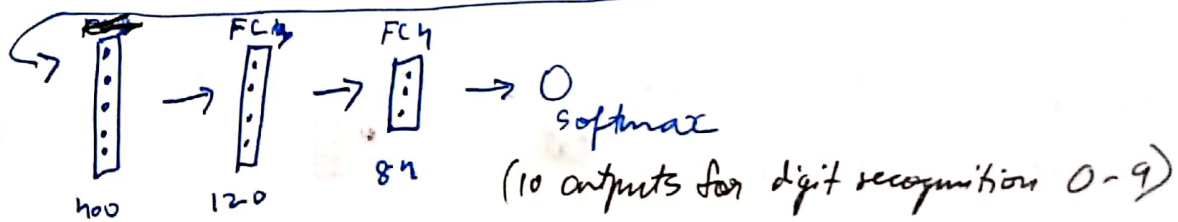
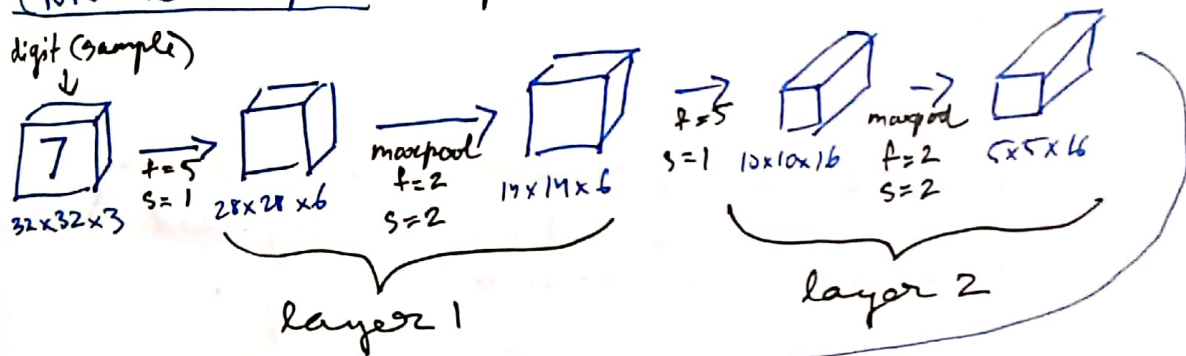
→ s : stride

→ Max or average pooling?

$$n_H \times n_W \times n_C \rightarrow \left\lfloor \frac{n_H - f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n_W - f}{s} + 1 \right\rfloor \times n_C$$



CNN example - digit recognition



CONV - POOL - CONV - POOL - FC - FC - FC - SOFTMAX



Layer 1 Layer 2

Table of above layer,

	Activation shape	Activation size	# parameters
Input:	(32, 32, 3)	3072	0
CONV1 (5, 5, 1)	(28, 28, 8)	6272	608 $(5 \times 5 \times 3 + 1) \times 8$
POOL1	(14, 14, 8)	1568	0
CONV2 (5, 5, 1)	(10, 10, 16)	1600	3216 $(5 \times 5 \times 8 + 1) \times 16$
POOL2	(5, 5, 16)	400	0
FCS	(120, 1)	120	48120 $160 \times 120 + 120$
FC4	(84, 1)	84	48120 10164 $120 \times 84 + 84$
Softmax	(10, 1)	10	84 850 $84 \times 10 + 10$

Why convolution?

- Parameter sharing: A feature detector (such as a vertical edge detector) that's useful in one part of the image is probably useful in another part of the image
 - Therefore less parameters are required (If we used a normal NN, we may need millions of parameters)
- Sparsity of connections: In each layer, each output value depends only on a small no. of inputs


→


 - allows translation invariance - if we shift the position of a cat in an image, it'll still recognise the feature