

CONV NETS

COURSE III

EDGE DETECTION EXAMPLE

* Vertical edge detection:

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

$n \times n$

"convolution"

$*$

FILTER
(KERNEL)
 $f \times f$
usually odd

$=$

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

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

$$3 \cdot 1 + 1 \cdot 1 + 2 \cdot 1 + 0 \cdot 0 + 0 \cdot 5 + 0 \cdot 7 + (-1) \cdot 1 + (-1) \cdot 3 + (-1) \cdot 2 = -5$$

Curved arrow from the first row of the input to the first row of the output.

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

$*$

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

$=$

-5		

* you apply the filter and perform element-wise multiplication:

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

$*$

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

$=$

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

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



BRIGHT TO DARK

$*$

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

$=$

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



* horizontal edge detection:

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

* what filter is the best filter?

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

SOBEL FILTER: it gives a higher weight to the middle row and makes it more robust.

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

SCHARR FILTER

w_1	w_2	w_3
w_4	w_5	w_6
w_7	w_8	w_9

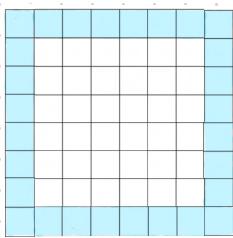
* you can treat each element of the filter as a parameter and use backprop to calculate them.

PADDING

* it helps with shrinking output when applying filters
 ↘ throwing away info from the edge

* how much to pad?
 ↘ VAUD: no padding
 ↘ SAME: pad so the output size = input size

PADDING



$$6 \times 6 \rightarrow 8 \times 8$$

$$(n+2p) \times (n+2p)$$

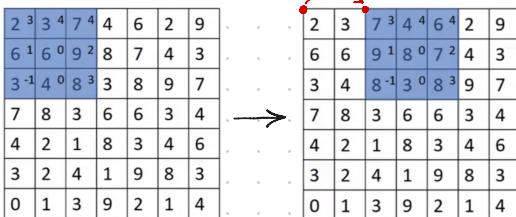
$$\begin{array}{ccc} * & \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} & = \end{array} \begin{array}{|c|c|c|c|} \hline & & & \\ \hline \end{array}$$

3×3
 $f \times f$

$4 \times 4 \rightarrow 6 \times 6$
 $(n+2p-f+1) \times (n+2p-f+1)$

STRIDED CONVOLUTIONS

- * how many steps the filter moves.



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

$n \times n$

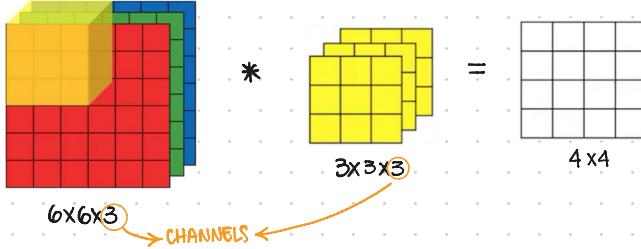
$$* \quad \begin{matrix} 3 & 4 & 4 \\ 1 & 0 & 2 \\ -1 & 0 & 3 \end{matrix} = \begin{matrix} & & \\ & & \\ & & \end{matrix}$$

$$f \times f \quad S=2, P=1 \quad \left(\frac{n+2p-f}{s} + 1 \right) \times \left(\frac{n+2p-f}{s} + 1 \right)$$

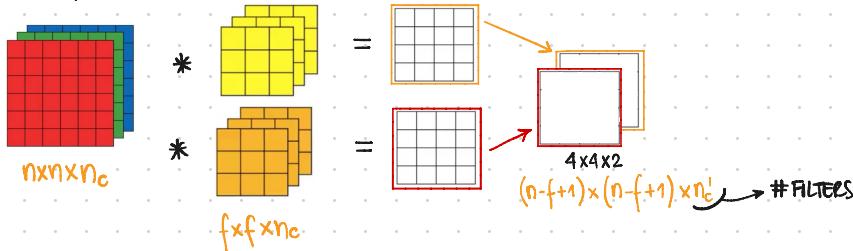
if it's not an integer,
we round down.

CONVOLUTIONS OVER VOLUME

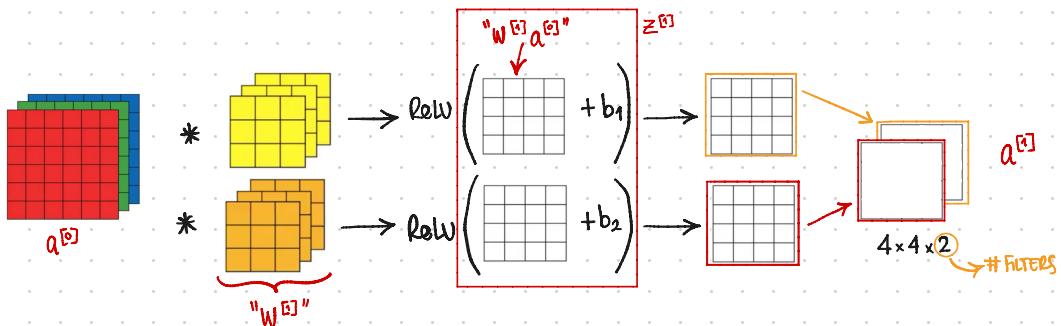
- * convolutions on RGB images:



- * can you use multiple filters at the same time?



ONE LAYER OF A CNN NET



* if you have 10 filters of size $3 \times 3 \times 3$, how many parameters you have?

$$3 \times 3 \times 3 = 27 \text{ (size filter)} + \text{bias} = 28 \times 10 \text{ (#filters)}$$

NOTATION:

$$f^{[l]} = \text{filter size}$$

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

$$p^{[l]} = \text{padding size}$$

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

$$S^{[l]} = \text{stride}$$

$$n_H^{[l]} = \left\lceil \frac{n_H^{[l-1]} + 2p^{[l]} - f^{[l]} + 1}{S^{[l]}} \right\rceil$$

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

$$f^{[l]} \times f^{[l]} \times n_C^{[l-1]} = \text{filter size}$$

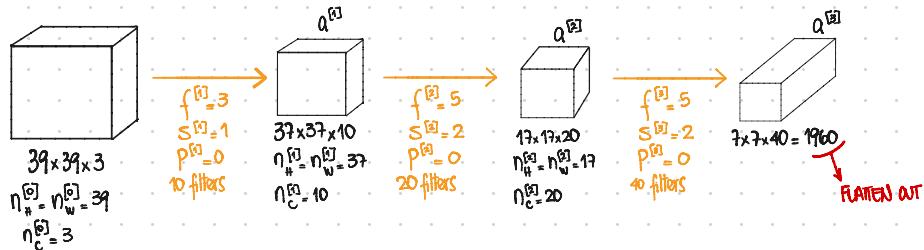
$$a^{[l]} \rightarrow n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$$

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

$$\text{weights} \rightarrow f^{[l]} \times f^{[l]} \times n_C^{[l-1]} \times n_C^{[l]}$$

$$\text{bias} \rightarrow n_C^{[l]}$$

EXAMPLE

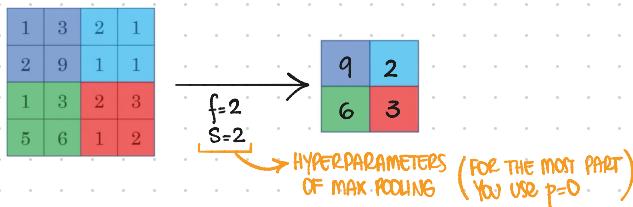


- * types of layers in a CNN
 - Conv (convolution)
 - Pool (pooling)
 - FC (fully connected) **simpler to define**

POOLING LAYERS

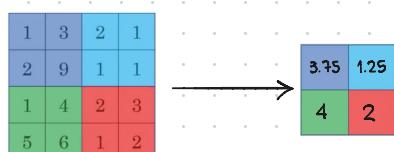
- * useful for
 - reduce size of the representation
 - speed up the computation
 - make some of the features that detect more robust

MAX POOLING:



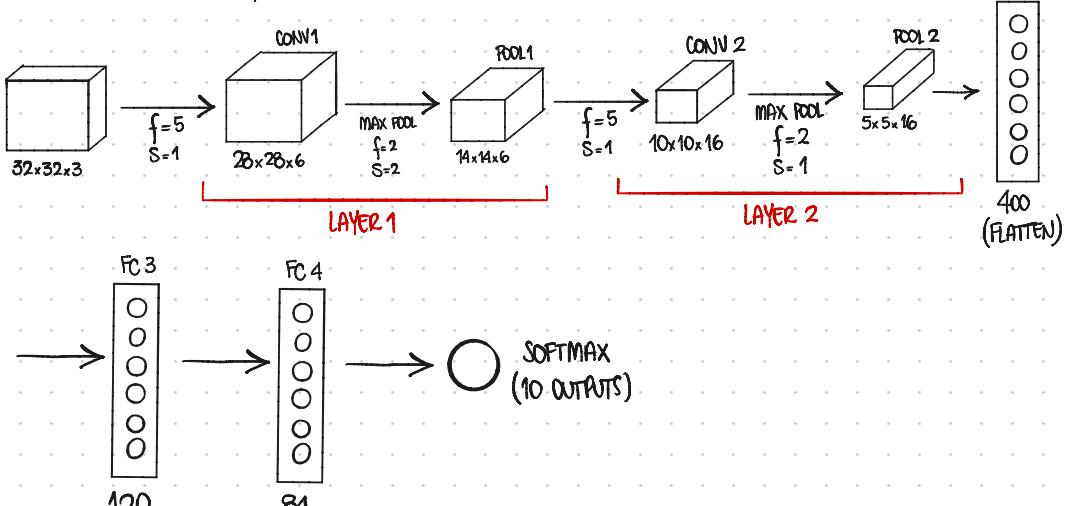
- * you take the highest value in the region
- * the important features detected in the quadrants (regions) remains preserved in the output.
- * it has hyperparameters but no parameters to learn, unlike convolution.

AVERAGE POOLING:



* not used as often as MAX pool. Sometimes used in very deep (NN) to collapse your representation, e.g. from $7 \times 7 \times 1000$ to $1 \times 1 \times 1000$.

EXAMPLE (LeNet-5 inspired)



* AS you move deep into the network... $n_h, n_w \downarrow$
 $n_c \uparrow$

	Activation shape	Activation Size	# parameters
Input:	(32,32,3)	3,072	0
CONV1 ($f=5, s=1$)	(28,28,8)	6,272	608
POOL1	(14,14,8)	1,568	0
CONV2 ($f=5, s=1$)	(10,10,16)	1,600	3216
POOL2	(5,5,16)	400	0
FC3	(120,1)	120	48120
FC4	(84,1)	84	10164
Softmax	(10,1)	10	850

$$(5 \times 5 \times 3 + 1) \times 8$$

$$(5 \times 5 \times 8 + 1) \times 16$$

IT DROPS GRADUALLY
AS THE (NN) GETS DEEPER

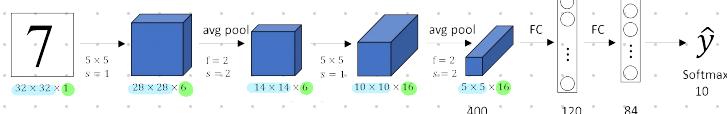
MOST PARAMETERS TEND TO BE
IN THE FC LAYERS
 $84 \times 10 + 10$

WHY CONN NETS?

- * advantages
 - parameter sharing: a feature detector that's useful in one part of the image is probably useful in another part of the image.
 - sparsity of connections: in each layer, each output value depends only on a small number of inputs.

CLASSIC NETWORKS

* LeNet 5:



* it has 60k parameters → relatively small compared to other networks.

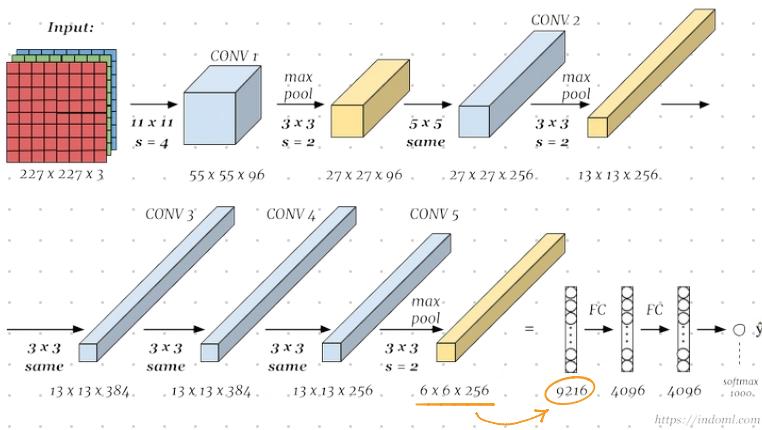
* as you move deeper into the network, n_{in} and n_{out} tend to go ↓

fc tends to go ↑

* general structure: conv pool + conv pool + fc + fc + output

* section II and III of the paper.

* AlexNet:



* Similar to LeNet, but much bigger] → ~60M parameters
60k parameters

* it has more hidden layers
been trained on much more data (ImageNet)

* it uses ReLU activations and Local Response Normalization

normalize inputs based on one location and across all input channels. } You don't want too many active neurons with a very high activation.

