

---

---

---

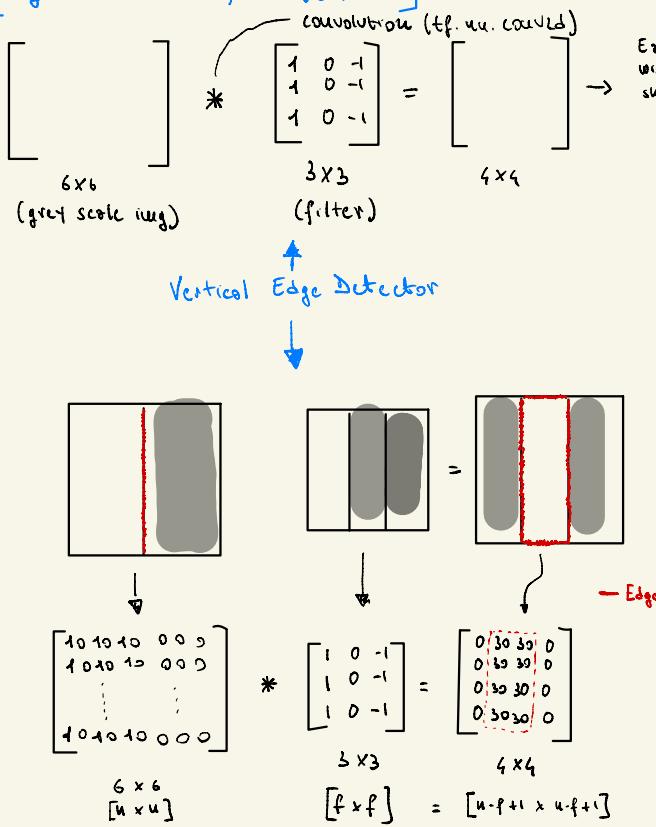
---

---



# CONV NETS - SETTING THE BASE

## [Edge Detection, Convolution]



- Light  $\rightarrow$  Dark, Dark  $\rightarrow$  Light edges
- Vertical vs Horizontal edge
- Different types of filters (ie Sobel filter)
- We can treat the 3 numbers in the filter as parameters and leave the values for these. This allows to learn different types of edges
- Valid vs Same convolution
- Strided convolution: skipping the filter by skipping steps with a stride of  $i \neq 1$
- Padding  $p = \frac{u+2p-f+1}{s} + 1 \times \frac{u+p-f+1}{s}$
- RGB images:  $\frac{\text{img}}{\text{filter}} = \frac{\text{img}}{\text{filter}}$   
 $6 \times 6 \times 3 \quad * \quad 3 \times 3 \times 3 \quad \# \text{channels} \rightarrow \text{must be the same}$
- We can even apply multiple filters to a img:

$$n \times u \times n_c * n \times u \times n_c = n-p+1 \times n-p+1 \times n_c \quad \text{~~~} \rightarrow \text{number of filters}$$

## [Padding]

Every time you apply convolution, your img shrinks  
 Edge pixels are used much less than other pixels } Cons of convolution

if padding with 1 pixel:  $p=1 \rightarrow [u+2p-f+1 \times u+2p-f+1]$

$\hookrightarrow$  We can pad the img before applying convolution: adding a border to the img of pixels with 0

## [Example of layer]

$(4 \times 4 \rightarrow u-p+1 \rightarrow 6-3+1 \rightarrow 4 \times 4)$

$$\begin{aligned} 6 \times 6 \times 3 &\ast \underbrace{3 \times 3 \times 3}_{w^{[1]}} \rightarrow \text{Relu}(\underbrace{4 \times 4 + b_1}_{w^{[1]}, b^{[1]}}) \rightarrow 4 \times 4 \\ &\ast 3 \times 3 \times 3 \rightarrow \text{Relu}(4 \times 4 + b_2) \rightarrow 4 \times 4 \end{aligned}$$

# of filters applied

$$\begin{aligned} Z^{[1]} &= W^{[1]} \cdot x^{[1]} + b^{[1]} \\ g^{[1]} &= g(Z^{[1]}) \end{aligned}$$

10 filters  $3 \times 3 \times 3$  in 1 layer, # of params:

$$3 \times 3 \times 3 + \text{bias} = 28 \text{ params} \quad | \quad 1 \text{ filter} \xrightarrow{\times 10} 280 \text{ params}$$

Notation  $\rightarrow$

## Notation

Convolution layer  $\ell$ :

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

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

$$s^{[\ell]} = \text{stride size}$$

$$u_c^{[\ell]} = \text{number of filters}$$

$$\rightarrow \text{Each filter is: } f^{[\ell]} \times f^{[\ell]} \times u_c^{[\ell-1]}$$

$$\rightarrow \text{Activations: } I^{[\ell]} \rightarrow u_h^{[\ell]} \times u_w^{[\ell]} \times u_c^{[\ell]}$$

$$\rightarrow \text{Weights: } f^{[\ell]} \times f^{[\ell]} \times u_c^{[\ell-1]} \times \underbrace{u_c^{[\ell]}}_{\text{number of filters in layer } \ell} \rightarrow \text{filters in layer } \ell$$

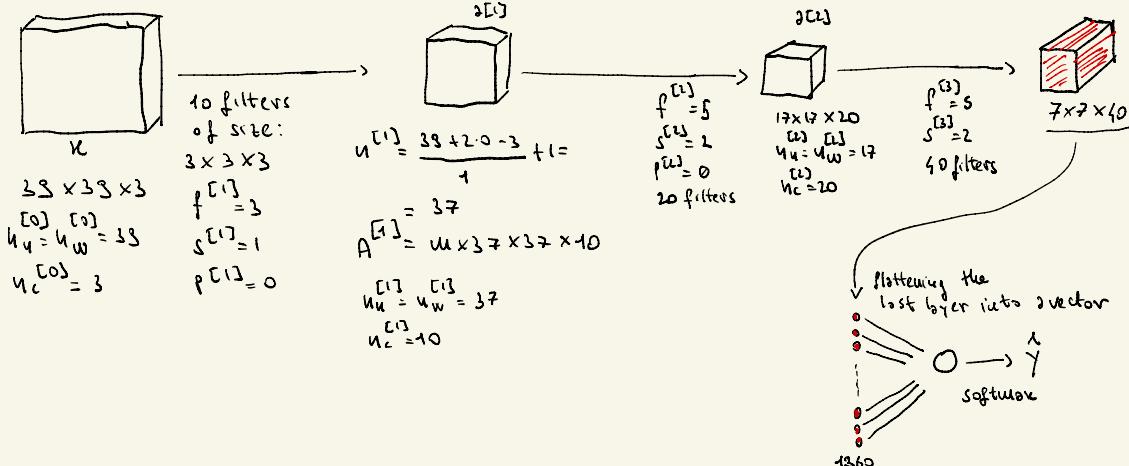
$$\rightarrow \text{Bias: } u_b^{[\ell]} = (1, 1, 1, u_c^{[\ell]})$$

$$\text{Input: } u_h^{[\ell-1]} \times u_w^{[\ell-1]} \times u_c^{[\ell-1]}$$

$$\text{Output: } u_h^{[\ell]} \times u_w^{[\ell]} \times u_c^{[\ell]}$$

$$u_h^{[\ell]} = \frac{u_h^{[\ell-1]} + 2p^{[\ell]} - f^{[\ell]}}{s^{[\ell]}} + 1$$

$$I^{[\ell]} \rightarrow \underbrace{u_h^{[\ell]} \times u_w^{[\ell]} \times u_c^{[\ell]}}_{\# \text{ of examples}}$$



## [Type of layers in a conv net]

- 1) Convolution 2) Pooling 3) Fully Connected
- (conv) (pool) (FC)

## 2) Pooling

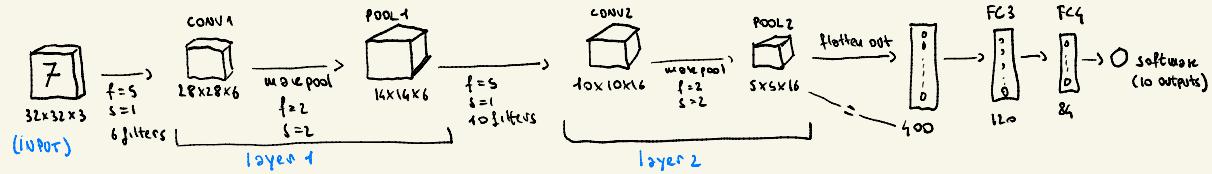
• Max pooling:  $\begin{bmatrix} 4 \times 4 \end{bmatrix} \xrightarrow{\text{Report only the max value}} \begin{bmatrix} 2 \times 2 \end{bmatrix}$

• Average pooling: Report only the avg value.  
It's less used

- it helps detecting a feature in the different quadrants where it's applied
- it has no parameters to learn, just hyperparameters ( $f = \text{size}, s = \text{stride}$ )

## [CNN Example]

→ Recognition of hand digits



Network Size

	Activation shape	Activation Size	# parameters
Input:	$(32, 32, 3)$	<u>3,072</u>	0
CONV1 ( $f=5, s=1$ )	$(28, 28, 8)$	<u>6,272</u>	208
POOL1	$(14, 14, 8)$	<u>1,568</u>	0
CONV2 ( $f=5, s=1$ )	$(10, 10, 16)$	<u>1,600</u>	416
POOL2	$(5, 5, 16)$	<u>400</u>	0
FC3	$(120, 1)$	<u>120</u>	48,001
FC4	$(84, 1)$	<u>84</u>	10,081
Softmax	$(10, 1)$	<u>10</u>	841

## [Why Convolutions]

⇒ Less parameters to learn if compared to FC layers

A) Parameter sharing: a feature detector that's useful in one part of the img is probably useful in another part of the img

B) Sparsity of connections: in each layer, each output value depends only on a small number of inputs

$A + B \rightarrow$  Traversing with smaller traversal sets

- Less prone to overfitting
- Good at capture translation in variance

