

deci: computer vision

example:

① image classification

$$(x \rightarrow y(0,1) \text{ cat})$$

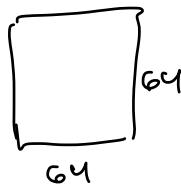
② object detection

$$(x \rightarrow \text{bounding box, labels})$$

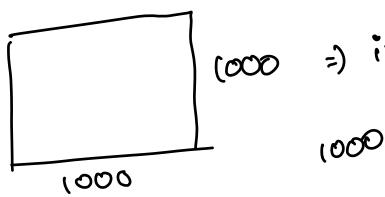
③ Neural style Transfer

$$(x + \text{style} \rightarrow x_{\text{style image}})$$

Deep learning on large images:



\Rightarrow input features ($64 \times 64 \times 3$) = 12288



\Rightarrow input features ($1000 \times 1000 \times 3$) = 3 million.

$$\begin{matrix} x_1 & \circ \\ \vdots & \vdots \\ x_n & \circ \\ & w^{(r)} \end{matrix} \quad X \in \mathbb{R}^{3^M} \quad \rightarrow$$

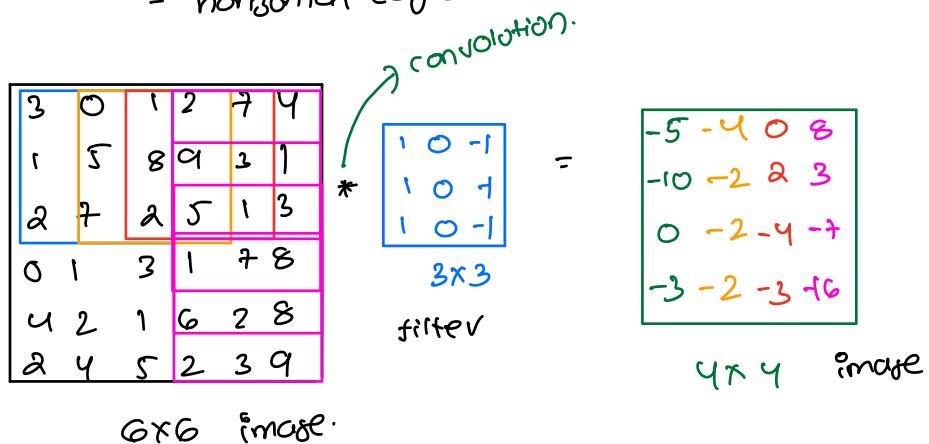
$$w^{(r)} = (1000, 3^M) \Rightarrow \text{Very huge matrix}$$

\Rightarrow Memory issue.

\Rightarrow need very large data.

dec2: Edge Detection Example:-

- vertical edges
- horizontal edges

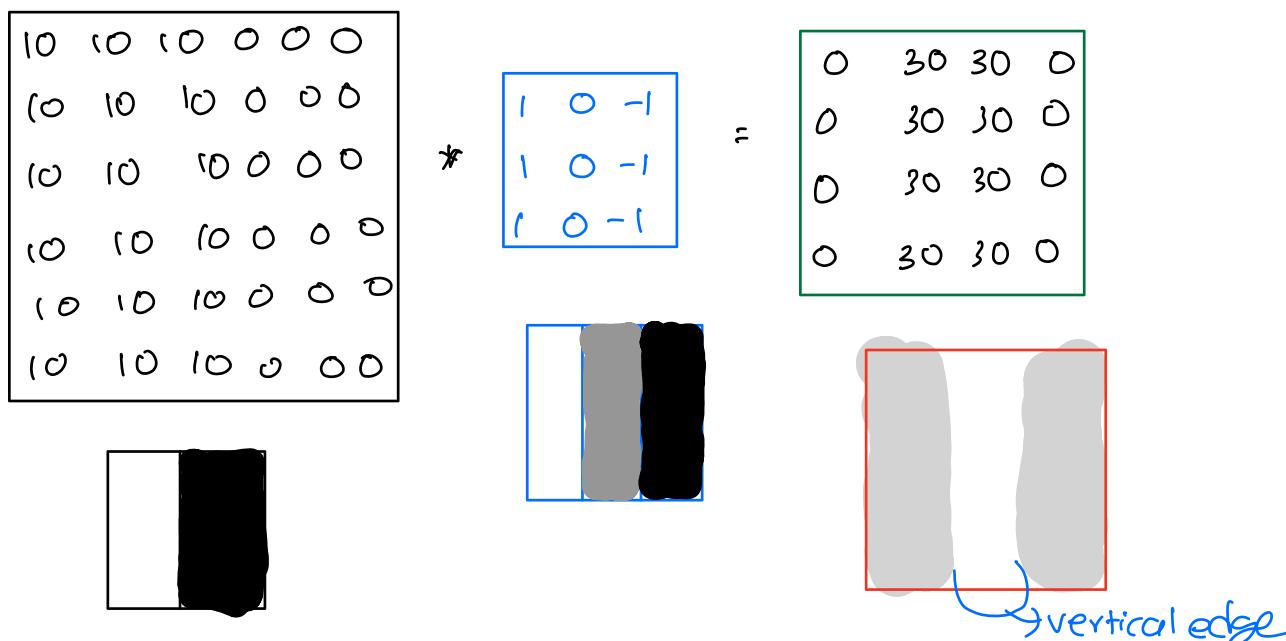


$$(0,0) = 3 \times 1 + 1 \times 1 + 2 \times 1 + 0 \times 0 + 5 \times 0 + 7 \times 0 + (-1 + 8 \times 1 + 2 \times 1) = -5$$

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

Python: conv-forward

Tensorflow: tf.nn.conv2d



dec 3: More edge detection:-

$$\begin{bmatrix}
 1 & 0 & 1 & 0 & 0 & 0 & 0 \\
 0 & 1 & 0 & 1 & 0 & 0 & 0 \\
 1 & 0 & 1 & 0 & 0 & 0 & 0 \\
 0 & 1 & 0 & 1 & 0 & 0 & 0 \\
 1 & 0 & 1 & 0 & 0 & 0 & 0 \\
 0 & 1 & 0 & 1 & 0 & 0 & 0 \\
 0 & 1 & 0 & 1 & 0 & 0 & 0
 \end{bmatrix} * \begin{bmatrix}
 1 & 0 & -1 \\
 1 & 0 & -1 \\
 1 & 0 & -1
 \end{bmatrix} = \begin{bmatrix}
 0 & 30 & 30 & 0 \\
 0 & 30 & 30 & 0 \\
 0 & 30 & 30 & 0 \\
 0 & 30 & 30 & 0
 \end{bmatrix} \quad (\text{dark to light})$$

$$\begin{bmatrix}
 0 & 0 & 0 & 10 & 10 & 10 \\
 0 & 0 & 0 & 10 & 10 & 10 \\
 0 & 0 & 0 & 10 & 10 & 10 \\
 0 & 0 & 0 & 10 & 10 & 10 \\
 0 & 0 & 0 & 10 & 10 & 10 \\
 0 & 0 & 0 & 10 & 10 & 10
 \end{bmatrix} * \begin{bmatrix}
 1 & 0 & -1 \\
 1 & 0 & -1 \\
 1 & 0 & -1
 \end{bmatrix} = \begin{bmatrix}
 0 & -30 & -30 & 0 \\
 0 & -30 & -30 & 0 \\
 0 & -30 & -30 & 0 \\
 0 & -30 & -30 & 0
 \end{bmatrix} \quad (\text{right to dark})$$



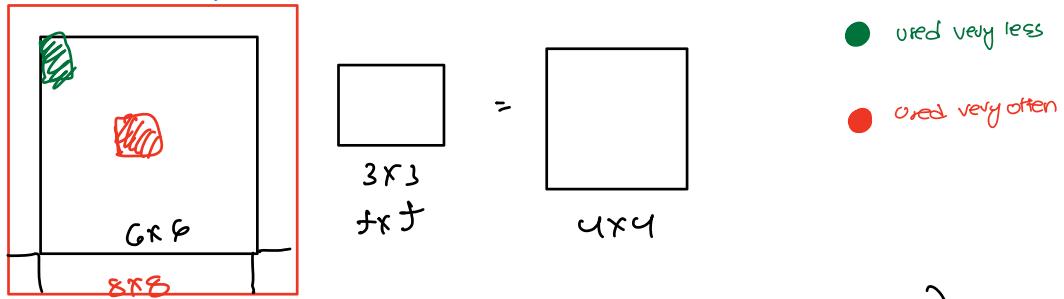
Vertical edge filter $\begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix}$ Horizontal edge filter $\begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{bmatrix}$

Sobel filter $\begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix}$ Scharr filter $\begin{bmatrix} 3 & 0 & -3 \\ 10 & 0 & -10 \\ 3 & 0 & -3 \end{bmatrix}$

Goal is learn this matrix, by treating its parameters.

$$\begin{bmatrix} w_1 & w_2 & w_3 \\ w_4 & w_5 & w_6 \\ w_7 & w_8 & w_9 \end{bmatrix}$$

dec4: padding:



$n \times n$

$$\begin{aligned} \text{output dim} &= (n-f+1) \times (n-f+1) \\ \text{without pad} &= (6-3+1) \times (6-3+1) \\ &= (4 \times 4) \end{aligned}$$

$$\begin{aligned} \text{Output dim} &= p=1 \text{ padding} \\ \text{with pad} &= (n+2p-f+1) \times (n+2p-f+1) \\ &= (6+2 \times 1-3+1) \times (6+2 \times 1-3+1) \\ &= (6 \times 6) \text{ (original dim is preserved)} \end{aligned}$$

● used bit more with padding

Valid and Same Convolutions:-

valid: No padding $(n, n) * (f, f) = (n-f+1, n-f+1)$

Same: padding Input size = Output size

(n, n) with p padding $* (f, f) = (n, n)$

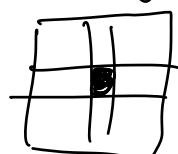
$$(n+2p-f+1) = n$$

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

f is usually odd.

if $f=(1, 3)$ then $p=1$

if $f=(5, 5)$ then $p=2$



• Position of filter.

dec 5: Strided Convolutions:-

Diagram illustrating a stride-2 strided convolution:

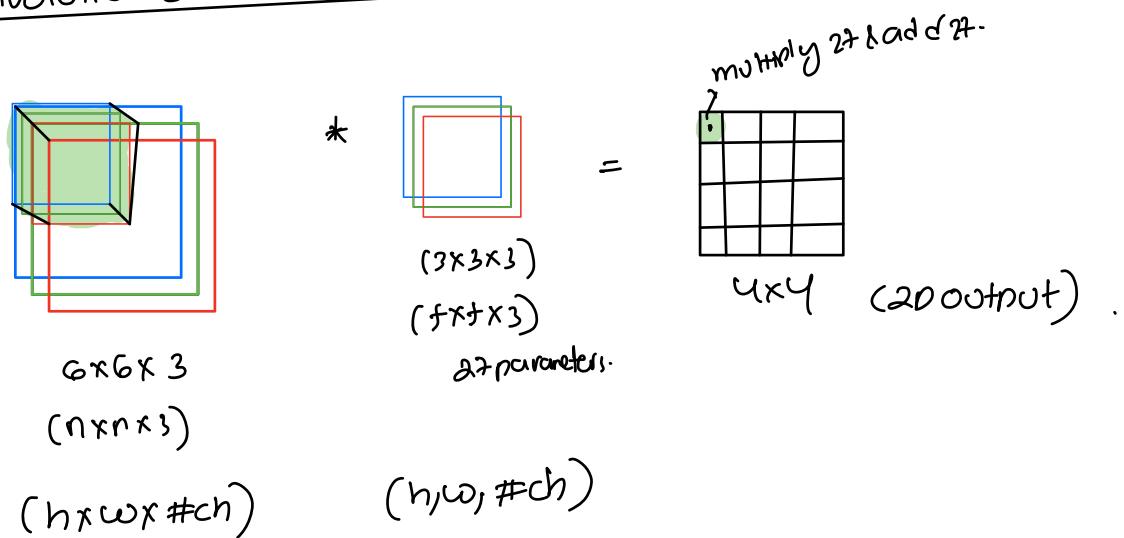
- Input Matrix:** 7×7 (labeled $n \times n$)
- Filter:** 3×3
- Stride:** $s=2$
- Output Matrix:** (3×3)
- Calculation:** $(n \times n) * (f \times f)$ with pad p & stride s

filter must lie completely inside the padded image.

$$\begin{aligned}
 &= \left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n+2p-1}{s} + 1 \right\rfloor \text{ floor } [z] \\
 &= \left\lfloor \frac{7+0-3}{2} + 1 \right\rfloor \times \left\lfloor \frac{7-3}{2} + 1 \right\rfloor \\
 &= (3 \times 3)
 \end{aligned}$$

dec6: convolution Over Volume :-

convolution over RGB image :-



#ch-image = #ch-filters

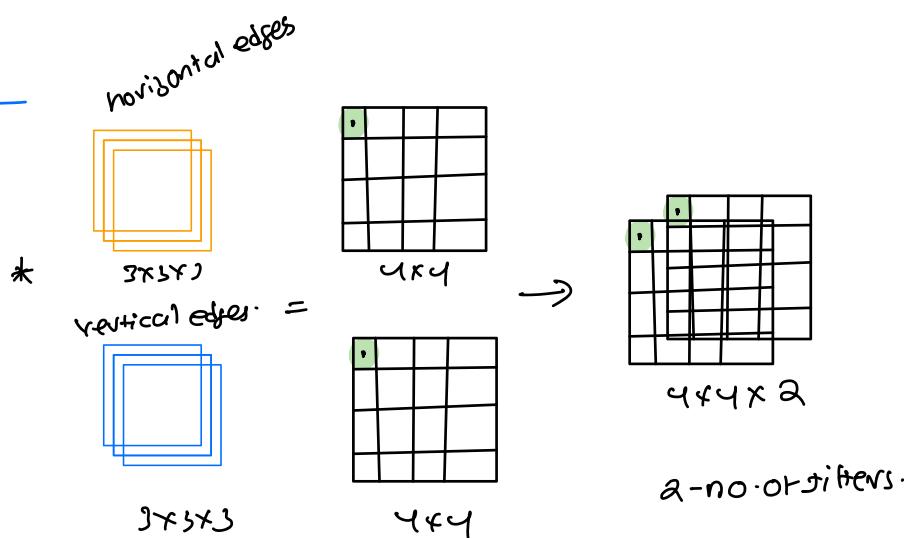
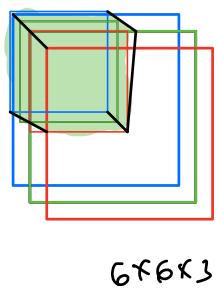
$$\begin{array}{ccc}
 R & \begin{matrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{matrix} & G \begin{matrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{matrix} \\
 & & B \begin{matrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{matrix}
 \end{array}$$

\rightarrow edges in red channel.

$$\begin{array}{ccc}
 R & \begin{matrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{matrix} & G \begin{matrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{matrix} \\
 & & B \begin{matrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{matrix}
 \end{array}$$

\rightarrow edges in any colors.

Multiple Filters :-



summary :-

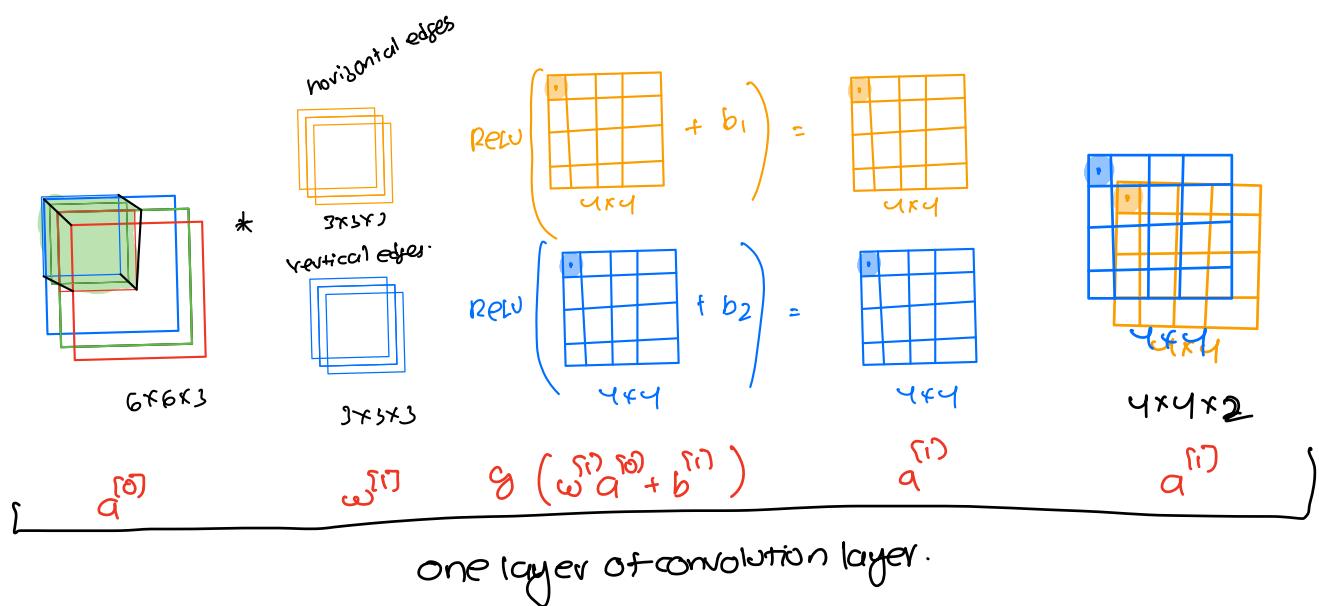
$$(n \times n \times [n_c]) * (f \times f \times [n_c]) \rightarrow (n-f+1) \times (n-f+1) \times n_c^1$$

$n_c^1 = \text{no. of filters}$

$S=1$
 $P=0$

$6 \times 6 \times 3$ $(3 \times 3 \times 3)$

dec 7: One layer of Convolution Network:



non convol

$$z^{(i)} = w^{(i)} a^{(0)} + b^{(i)}$$

$$a^{(i)} = g(z^{(i)})$$

$$a^{(0)} \rightarrow a^{(i)}$$

6x6x3 4x4x2

2 - no. of filters.

if we have 10 filters, then

$$a^{(f)} = 4x4x10$$

Question) If you have 10 filters that are $3 \times 3 \times 3$ in one layer of NN,
how many parameters does the layer have?

each filter has

$$\begin{array}{c} \text{filter} \\ \text{size} 3 \times 3 \times 3 \end{array} + \text{bias} = (3 \times 3 \times 3) + 1 = 28$$

$$\therefore \text{no. of parameter of NN} = 280$$

which is independent to no. of input image, which can
be very large.

Summary of notation:-

If layer ℓ is a convolution layer:

f^{ℓ} = filter size ($f \times f$)

p^{ℓ} = padding (valid or same)

s^{ℓ} = stride

n_c^{ℓ} = number of filters

Each filter is: $f \times f \times n_c^{\ell}$

Activation: $a \rightarrow (n_h \times n_w \times n_c^{\ell})$

$A \rightarrow (m \times n_h \times n_w \times n_c^{\ell})$

Weights: $f \times f \times n_c^{\ell} \times n_c^{\ell}$

bias: $n_c^{\ell} - (1, 1, 1, n_c^{\ell})$

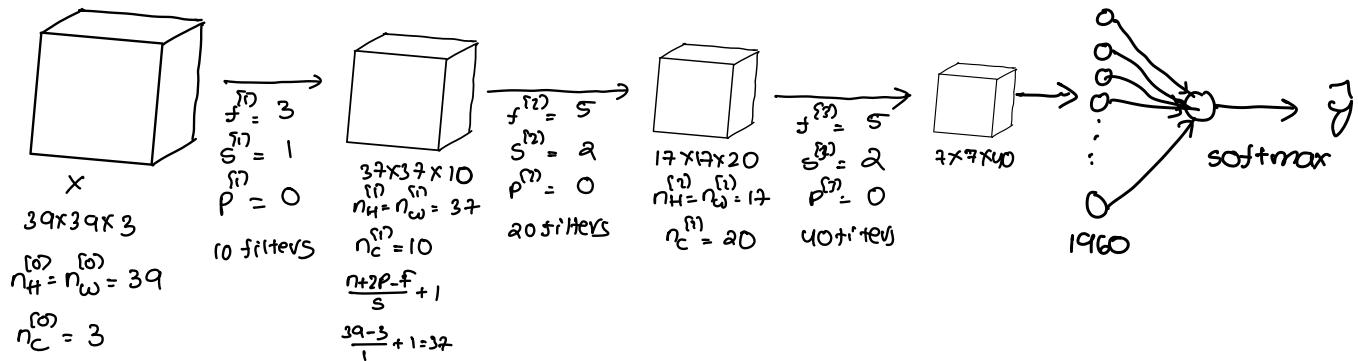
Input: $n_h \times n_w \times n_c^{\ell-1}$
 Output: $n_h \times n_w \times n_c^{\ell}$

$$n_h^{\ell} = \left\lfloor \frac{n_h^{\ell-1} + 2p - f}{s^{\ell}} + 1 \right\rfloor$$

$$n_w^{\ell} = \left\lfloor \frac{n_w^{\ell-1} + 2p - f}{s^{\ell}} + 1 \right\rfloor$$

XEC 8: A simple convolution network:

ConNet:



Types of layers in convnet:

- convolution (conv)
- pooling (pool)
- Fully connected (FC)

dec9: Pooling layer: Max Pooling:

$$\begin{array}{c|cc} 1 & 3 & 2 \\ \hline 2 & 9 & 1 \\ 1 & 3 & 2 \\ 5 & 6 & 1 \end{array} \xrightarrow[f=2]{s=2} \begin{matrix} 9 & 2 \\ 6 & 3 \end{matrix}$$

No parameters to learn.

Hyperparameters:

f: filter size

s: stride

p: padding

Max, avg pooling

$$\begin{array}{ccccc} 1 & 3 & 2 & 1 & 3 \\ 2 & 9 & 1 & 1 & 5 \\ 1 & 3 & 2 & 3 & 2 \\ 8 & 3 & 5 & 1 & 0 \\ 5 & 6 & 1 & 2 & 9 \end{array} \xrightarrow[f=3]{s=1} \begin{matrix} 9 & 9 & 5 \\ 9 & 9 & 5 \\ 8 & 6 & 9 \end{matrix}_{3 \times 3}$$

$\frac{n+2p-f}{s} + 1$
 $s-3+1=3$

5×5

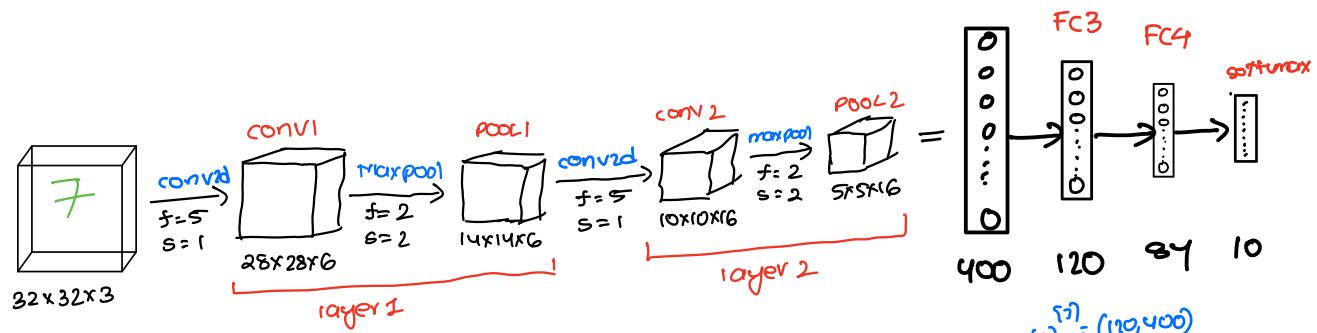
$5 \times 5 \times n_C \xrightarrow[f=3]{s=1} 3 \times 3 \times n_C$

Average pooling:-

$$\begin{array}{cccc} 1 & 3 & 2 & 1 \\ 2 & 9 & 1 & 1 \\ 1 & 4 & 2 & 3 \\ 5 & 6 & 1 & 2 \end{array} \xrightarrow[f=2]{s=2} \begin{array}{cc} 2.75 & 1.25 \\ 4 & 2 \end{array}$$

$7 \times 7 \times 1000 \xrightarrow{\text{avg pool}} 1 \times 1 \times 1000$

dec10: Neural Network Example: (deNet-5)



$$\text{CONV1 filter} = (5 \times 5 \times 3)$$

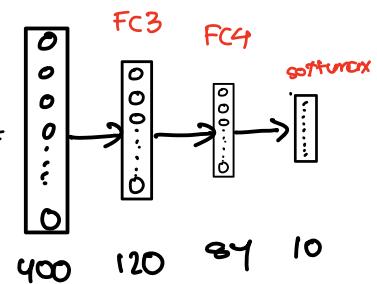
$$\text{para} = (5 \times 5 \times 1) + 1 = 26$$

$$6 \text{ filters} = 456$$

$$\text{CONV2 filter} = (5 \times 5 \times 6)$$

$$\text{para} = (5 \times 5 \times 6 + 1) = 151$$

$$16 \text{ filters} = 2416$$



$$\omega^{(1)} = (110, 400)$$

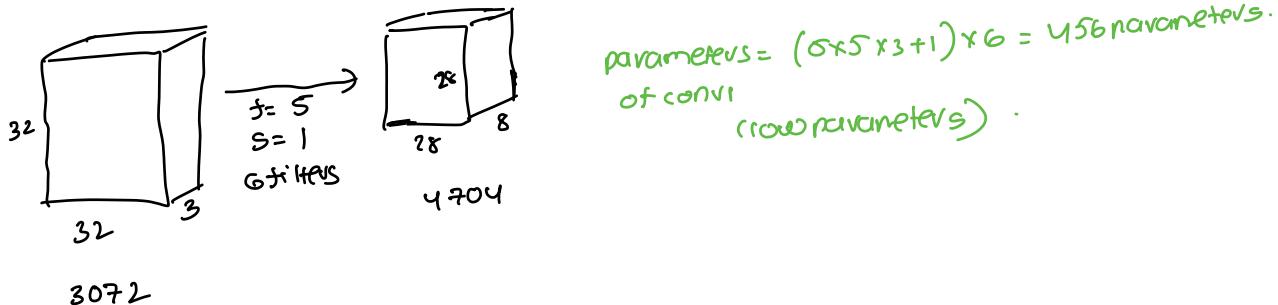
$$\omega^{(2)} = (84, 120)$$

	Activation shape	Activation Size	# parameters
Input:	(32, 32, 3)	3,072	0
CONV1 ($f=5, s=1$)	(28, 28, 6)	4704	456
POOL1 ($f=2, s=2$)	(14, 14, 6)	1176	0
CONV2 ($f=5, s=1$)	(10, 10, 16)	1600	2416
POOL2 ($f=2, s=2$)	(5, 5, 16)	400	0
FC3	(120, 1)	120	$(400 \times 120 + 120) = 48,120$
FC4	(84, 1)	84	$(120 \times 84 + 84) = 10164$
Softmax	(10, 1)	10	$(84 \times 10 + 10) = 850$

lec11: Why convolution layers:

advantages:

- parameters sharing
- sparsity of connections



3072

if you use FC layer

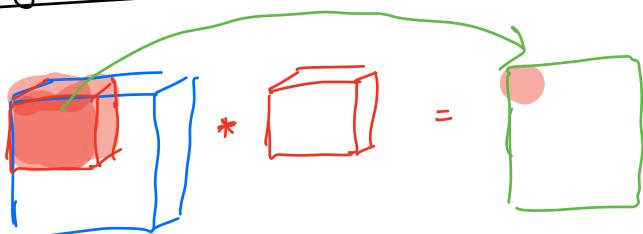
$$W = (4704, 3072)$$

$$\text{parameters} = (4704 \times 3072) \approx 10^{10}$$

parameter sharing:-

A feature detector (such as a vertical edge detector) that's useful in one part of image is probably useful in another part of image.

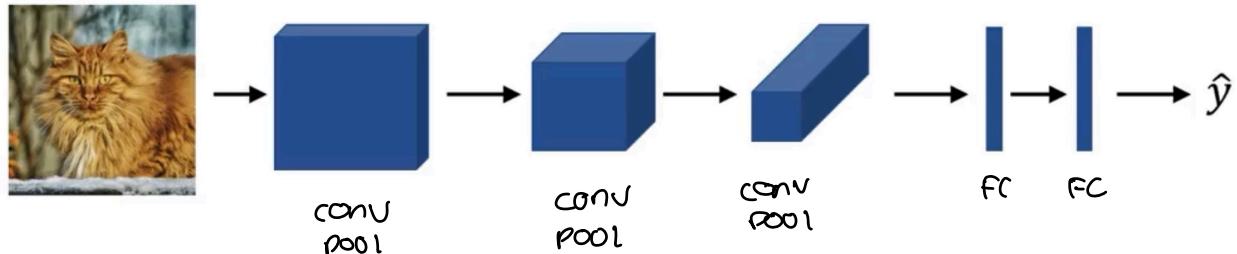
sparsity of connection:



In each layer, each output value depends only on a small number of inputs. (less overfitting)

Training set $(x^{(1)}, y^{(1)}) \dots (x^{(m)}, y^{(m)})$.

ω, b



$x^{(i)}$ - image

$y^{(i)}$ - binary output or k classes

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

use gradient descent to optimize parameters to reduce J

$$\textcircled{2} \quad w = [256, 16384] = 4194304$$

$$b = [256]$$

+

4194560

$$\textcircled{7} \quad f = (7 \times 7 \times 3)$$

$$\text{para} = (7 \times 7 \times 1) \times 128 = 18944$$

$$\textcircled{9} \quad 63, 67, 16 \xrightarrow{\begin{matrix} s=7 \\ s=2 \\ p=0 \\ 32 \text{ filters} \end{matrix}} \frac{n+2p-s}{s} + 1 = \frac{63-7}{2} + 1 = 29 \Rightarrow 29 \times 29 \times 3 \leftarrow$$

$$\textcircled{5} \quad 61 \times 61 \times 32 \quad p = 3$$

$$61+6 = (67 \times 67 \times 32)$$

$$\textcircled{6} \quad (121, 121, 32) \xrightarrow{\begin{matrix} f=5 \\ p=? \\ \#filters=32 \end{matrix}} (121, 121, 32)$$

$$\frac{n+2p-f}{s} + 1 = 121 \quad \frac{121+2p-5}{1} + 1 = 120$$

$$121+2p = 125 \quad p = 2$$

$$2p = 4$$