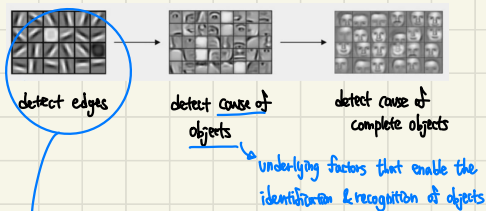


# <Course 4> - Convolutional Neural Network

## Edge Detection Example

### - Computer Vision Problem



### - Horizontal Edge Detection

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

 $\times$ 

bright	1	1	1
	0	0	0
horizontal dark filter	-1	-1	-1

 $=$ 

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

### - Vertical filters (make it horizontal by tilting 90°)

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

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

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

Sobel filter      Scharr filter

Learn these filters & use these 9#s as parameters

### - Vertical Edge Detection

$10 \times 1 + 10 \times 1 + 10 \times 0 + 10 \times 0 + 10 \times 0 + 10 \times (-1) + 10 \times (-1) + 10 \times (-1) = 0$

Convolution ( $\frac{1}{9} + \frac{1}{9}$ )

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

 $\times$ 

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

6x6      3x3 filter (kernel)      4x4

bright dark

$\times$   $=$

if pixels ↑  
gonna be more accurate

$w_1$	$w_2$	$w_3$
$w_4$	$w_5$	$w_6$
$w_7$	$w_8$	$w_9$

&  $\left. \begin{matrix} \text{tilting} \\ 45^\circ \\ 70^\circ \\ \vdots \end{matrix} \right\} \text{apply degrees}$

## Padding

- Downside of convolutional operator

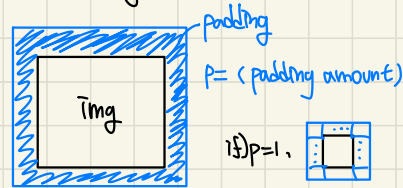
① apply convolutional op  $\rightarrow$  img shrinks (shrinking output)

$$(n \times n) * (f \times f) = \{(n-f+1) \times (n-f+1)\}$$

② corner pixels apply  $f$  only once  $\rightarrow$  throwing away corner info,

	10	10	0	0	0
10	10	10	0	0	0
10	10	10		0	0
0	0	0	10	10	10
0	0	0	10	10	10
0	0	0	10	10	10

- Solution: Padding



$$\therefore \text{Output pixel: } (n+2p-f+1) \times (n+2p-f+1)$$

- Valid & Same convolution

• "Valid", no padding

$$(n \times n) * (f \times f) = \{(n-f+1) \times (n-f+1)\}$$

• "Same"; pad so that output size is the same as input size

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

$\therefore p = \frac{f-1}{2}$  is usually odd

1x1  
3x3  
5x5  
7x7  
:

## Strided convolutions

Stride = 2

2	3	7	4	6	2	9
6	6	9	8	7	4	3
3	4	8	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

7x7  
padding p=0

Stride S = 2

3	4	4
1	0	2
-1	0	3

3x3

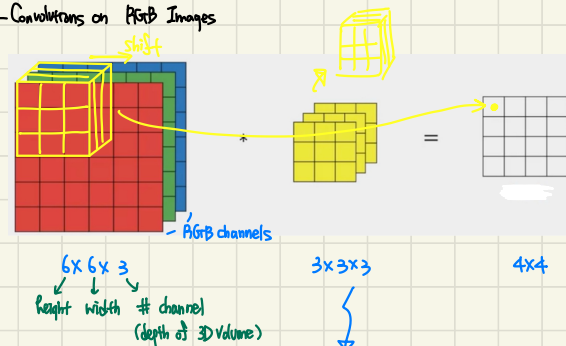
3x3

91 144 83  
69 91 127  
44 72 74

3x3

## Convolutions over volumes

- Convolutions on RGB Images



- Cross-correlation vs convolution

For common convolution, we should flip  $f$  horizontally & vertically

But in ML, we don't.

→ It's called "convolution" in ML convention,  
but it's more similar to "cross-correlation"!

Filters

$\begin{matrix} R \\ G \\ B \end{matrix}$

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

→ 3x3x3 filter detect vertical edges in Red channel

$\begin{matrix} R \\ G \\ B \end{matrix}$

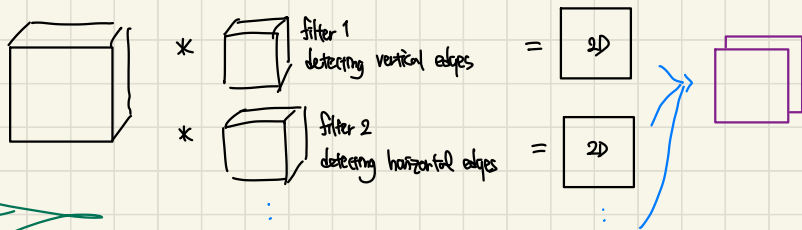
$\begin{matrix} R & \begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix} \\ G & \begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix} \\ B & \begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix} \end{matrix}$

→ 3x3x3 filter detect vertical edges in any color.

∴ tons of variation

- Multiple filters

applying multiple filters simultaneously



\* Summary

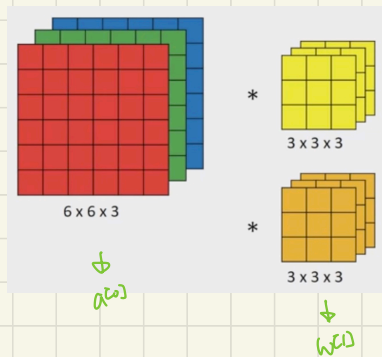
$$(n \times n \times n_c) * (f \times f \times n_c) * \rightarrow \{(n-f+1) \times (n-f+1) \times n'_c\}$$

↪ # channel (same)

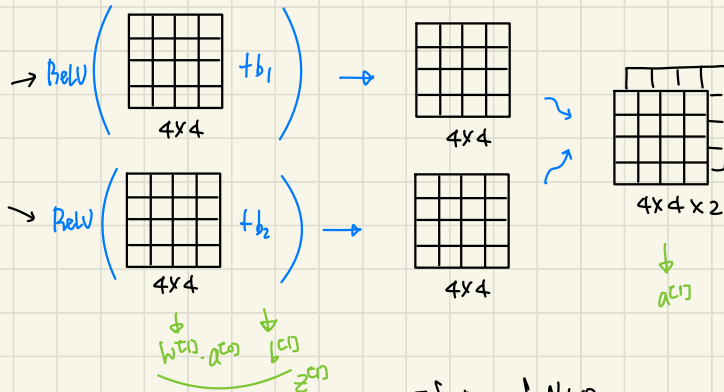
↪ # filters

# One layer of convolutional network

- Example of a layer



- ① apply bias ( $b \in \mathbb{R}$ ) w/ Python broadcasting, apply bias to all params
- ② apply non-linearity ex ReLU

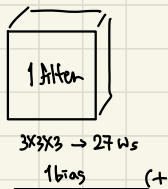


transform  $6 \times 6 \times 3 \rightarrow 4 \times 4 \times 2$   
 ⇒ one layer of  
convolutional neural network

$$z^{(1)} = w^{(1)} a^{(0)} + b^{(1)}$$

$$a^{(1)} = g(z^{(1)})$$

Q) 3x3x3 matrix of 10 filters  
 in one layer → how many parameters?



⊕ Even though input size ↑,  
 # param is not affected!

28 parameters per filter

→ 10 filter x 28 params each = 280 params in total

## Summary of Notation

if layer  $l$  is a convolutional layer:

$f^{(l)}$ : filter size  
 $p^{(l)}$ : padding  
 $s^{(l)}$ : stride  
 $n_c^{(l)}$ : # filter

Input:  $n_H^{(l-1)} \times n_W^{(l-1)} \times n_C^{(l-1)}$   
 Output:  $n_H^{(l)} \times n_W^{(l)} \times n_C^{(l)}$   

$$n_{H \times W}^{(l)} = \left\lfloor \frac{n_H^{(l-1)} + 2p^{(l)} - f^{(l)}}{s^{(l)}} + 1 \right\rfloor$$

Each filter is:  $f^{(l)} \times f^{(l)} \times n_C^{(l-1)}$

activation:  $a^{(l)} \rightarrow n_H^{(l)} \times n_W^{(l)} \times n_C^{(l)} \rightarrow A^{(l)} \rightarrow m \times n_H^{(l)} \times n_W^{(l)} \times n_C^{(l)}$

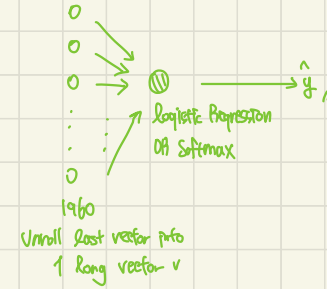
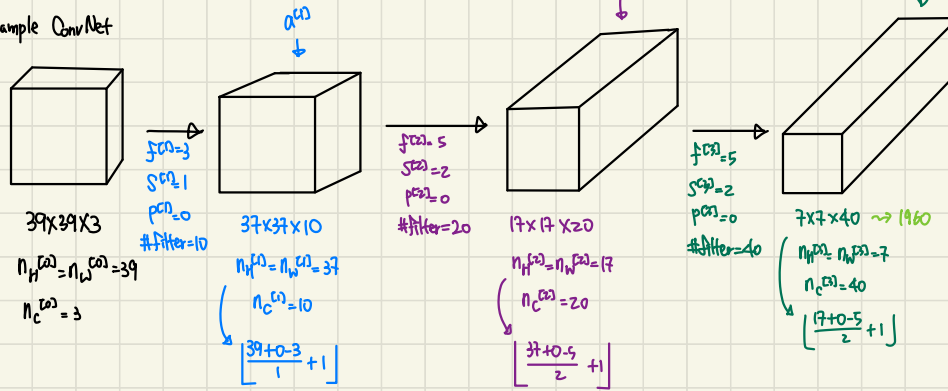
weights:  $f^{(l)} \times f^{(l)} \times n_C^{(l-1)} \times n_C^{(l)}$

each filter      # filter

bias:  $n_C^{(l)}$

# Simple Convolutional Network Example

## - Example ConvNet



## - Types of layer in a convolutional network

- ① Convolution (Conv)
- ② Pooling (POOL)
- ③ Fully Connected (FC)

## Pooling Layer

- Pooling Layer: Max pooling

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

4x4

$f=2$   
 $S=2$

hyperparameters of max pooling

→ each output will be Max from the corresponding region

9	2
6	3

2x2

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

5x5x $n_c$

$f=3$   
 $S=1$

9	9	5
9	9	5
8	6	9

3x3x $n_c$

- Pooling Layer: average pooling → not used much,,

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

4x4

$f=2$   
 $S=2$

3.75	1.25
4	2

2x2

- Summary of pooling

Hyperparams:

$f$ : filter size

$S$ : stride

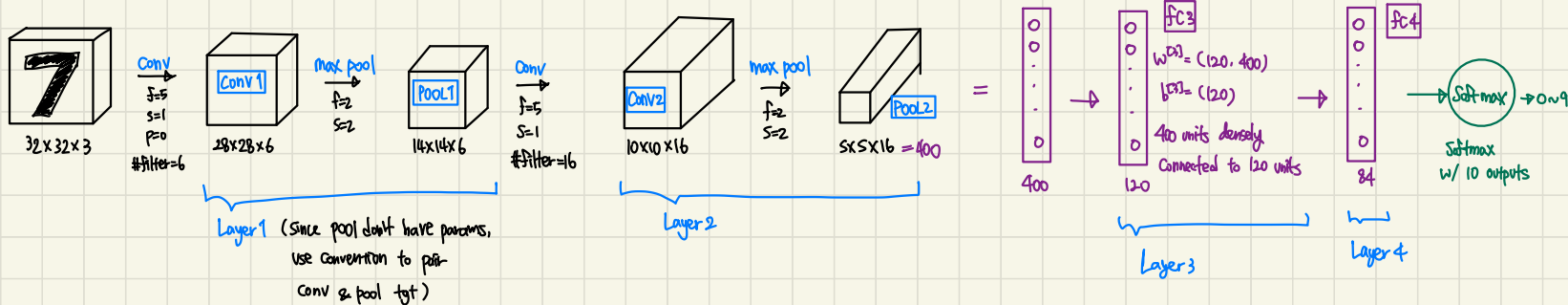
usually no padding ( $p=0$ )

No parameters to learn!

$$\text{Input } n_H \times n_W \times n_C \rightarrow \text{Output } \left\lfloor \frac{n_H - f}{S} + 1 \right\rfloor \times \left\lfloor \frac{n_W - f}{S} + 1 \right\rfloor \times n_C$$

# CNN example

-Neural Network Example : handwritten digit recognition



	Activation shape	Activation Size	# parameters
Input:	$(32, 32, 3)$	3,072	0
CONV1 (f=5, s=1)	$(28, 28, 6)$	4,704	456
POOL1	$(14, 14, 6)$	1,176	0
CONV2 (f=5, s=1)	$(10, 10, 16)$	1,600	2,416
POOL2	$(5, 5, 16)$	400	0
FC3	$(120, 1)$	120	48,120
FC4	$(84, 1)$	84	10,164
Softmax	$(10, 1)$	10	850

- ① Max pooling don't have params
- ② Conv layer tend to have few params
- ③ lots of params for fc
- ④ activation size gradually ↓
- ⑤  $\begin{cases} n_H \& n_W \text{ gradually } \uparrow \\ n_C \text{ gradually } \downarrow \end{cases}$

⑥ Common pattern;

CONV → POOL → CONV → POOL → FC → FC → Softmax

\* fc

In fully connected layer, each neuron is connected to every neuron in the prev layer

It takes high level features extracted by conv & pool, combines them to predict the final output.

→ feature combination

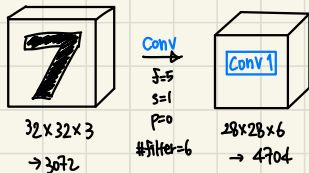
→ classification

Why convolutions?

① parameter sharing

② sparsity of connection

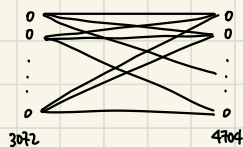
- Why convolutional over fully connected layers?



no convolutional

each filter:  $5 \times 5 \times 3 + 1 = 76 \times 6 \times 10^3 + 1$

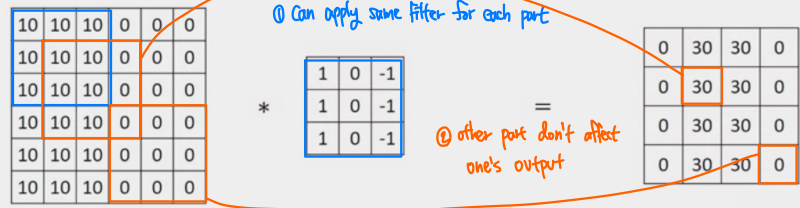
6 filters  $\rightarrow (5 \times 5 \times 3 + 1) \times 6 \times 10^3 = 456$



fully connected network

$3072 \times 4704 \approx 13M$

lots of params to train..!

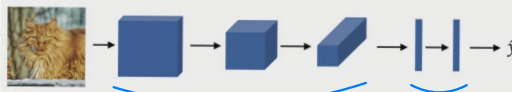


① Parameter Sharing: A feature detector (such as a vertical edge detector) that's useful in one part of the image is probably useful in another part of the image.

② Sparsity of Connection: In each layer, each output value depends only on a small number of inputs.

- Putting it together

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



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

→ use gradient descent to optimize params to reduce J.