

# Convolutional Neural Networks

Most of this material is from Prof. Andrew Ng and Chang's slides.

## Computer Vision Problems

### Image Classification



64x64

→ Cat? (0/1)

### Neural Style Transfer



### Object detection

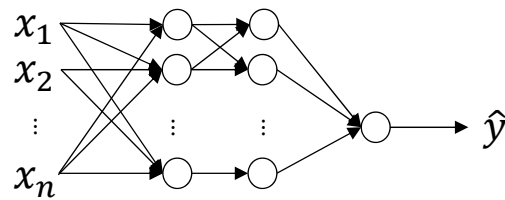


## Deep Learning on large images

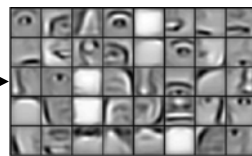
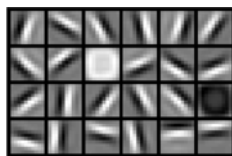


64x64

→ Cat? (0/1)



## Computer Vision Problem



vertical edges



horizontal edges

## Vertical edge detection

3 <sup>1</sup>	0 <sup>0</sup>	1 <sup>0</sup>	2 <sup>0</sup>	7 <sup>0</sup>	4 <sup>-1</sup>
1 <sup>1</sup>	5 <sup>0</sup>	8 <sup>-1</sup>	9 <sup>0</sup>	3 <sup>-0</sup>	1 <sup>-1</sup>
2 <sup>1</sup>	7 <sup>0</sup>	2 <sup>-1</sup>	5 <sup>0</sup>	1 <sup>-0</sup>	3 <sup>-1</sup>
0 <sup>1</sup>	1 <sup>0</sup>	3 <sup>-1</sup>	1 <sup>0</sup>	7 <sup>-0</sup>	8 <sup>-1</sup>
4	2	1	6	2	8
2	4	5	2	3	9

 $*$ 


 $=$ 

0	-2	-4	-7
-3	-2	-3	-16

## 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

 $*$ 

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

$*$  $=$

## Vertical edge detection examples

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

\*

= 

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

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

\*

= 

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

## Vertical and Horizontal Edge Detection

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

Vertical

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

Horizontal

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

\*

= 

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

## Learning to detect edges

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



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

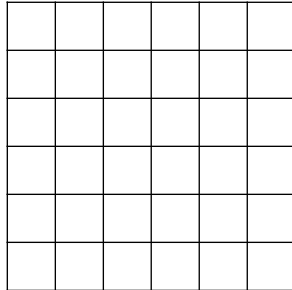
$w_1$	$w_2$	$w_3$
$w_4$	$w_5$	$w_6$
$w_7$	$w_8$	$w_9$


## Padding

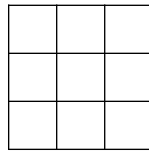

\*


=

## Padding



\*



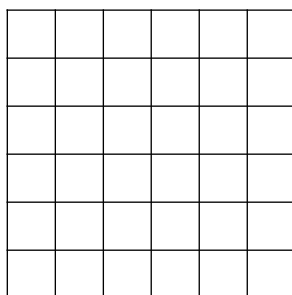
=

input:  $n \times n$

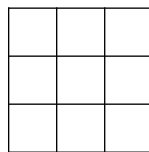
filter:  $f \times f$

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

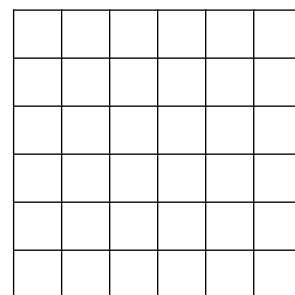
## Padding



\*



=



input:  $n \times n$

filter:  $f \times f$

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

If we use padding  $p$



output:  
 $(n + 2p - f + 1)$   
 $\times (n + 2p - f + 1)$

## Valid and Same convolutions

“Valid”: No padding

“Same”: Pad so that output size is the same as the input size.

## Strided convolution

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

\*

3	4	4
1	0	2
-1	0	3

=


## Strided convolution

2 <sup>3</sup>	3 <sup>4</sup>	7 <sup>3</sup>	4 <sup>4</sup>	6 <sup>3</sup>	2 <sup>4</sup>	9 <sup>4</sup>			
6 <sup>1</sup>	6 <sup>0</sup>	9 <sup>1</sup>	8 <sup>0</sup>	7 <sup>1</sup>	4 <sup>0</sup>	3 <sup>2</sup>			
3 <sup>3</sup>	4 <sup>4</sup>	3 <sup>3</sup>	3 <sup>4</sup>	3 <sup>3</sup>	9 <sup>4</sup>	7 <sup>4</sup>			
7 <sup>1</sup>	8 <sup>0</sup>	3 <sup>1</sup>	6 <sup>0</sup>	6 <sup>1</sup>	3 <sup>0</sup>	4 <sup>2</sup>			
4 <sup>3</sup>	2 <sup>4</sup>	1 <sup>3</sup>	8 <sup>4</sup>	3 <sup>3</sup>	4 <sup>4</sup>	6 <sup>4</sup>			
3 <sup>1</sup>	2 <sup>0</sup>	4 <sup>1</sup>	1 <sup>0</sup>	9 <sup>1</sup>	8 <sup>0</sup>	3 <sup>2</sup>			
0 <sup>-1</sup>	1 <sup>0</sup>	3 <sup>-1</sup>	9 <sup>0</sup>	2 <sup>-1</sup>	1 <sup>0</sup>	4 <sup>3</sup>			

\*

3	4	4
1	0	2
-1	0	3

=


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

## Summary of convolutions

$n \times n$  image       $f \times f$  filter

padding  $p$       stride  $s$

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



## Technical note on cross-correlation vs. convolution

Convolution in math textbook:

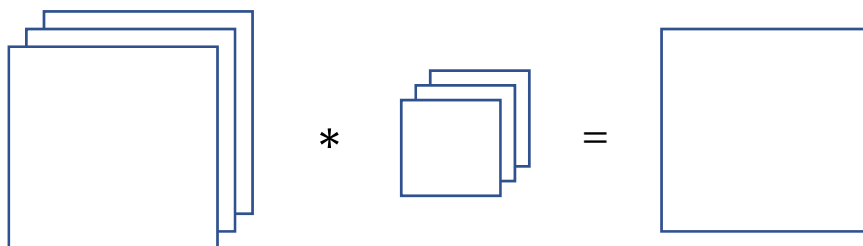
2	3	7	4	6	2
6	6	9	8	7	4
3	4	8	3	8	9
7	8	3	6	6	3
4	2	1	8	3	4
3	2	4	1	9	8

 $*$ 

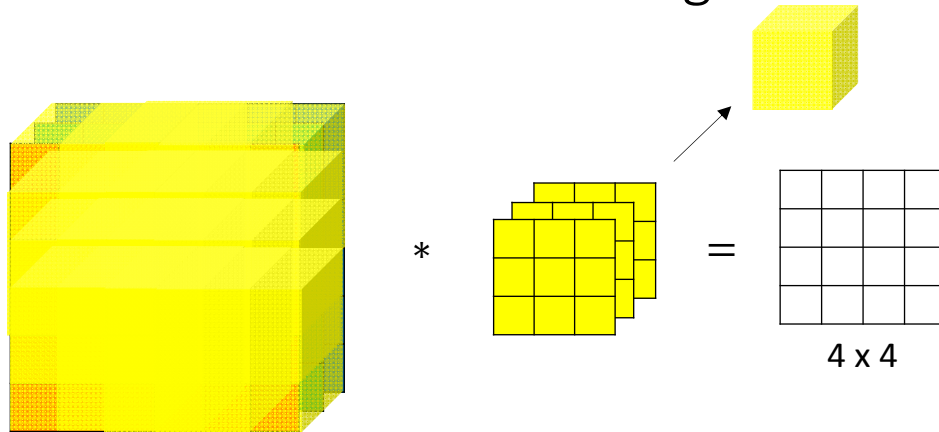
3	4	5
1	0	2
-1	9	7

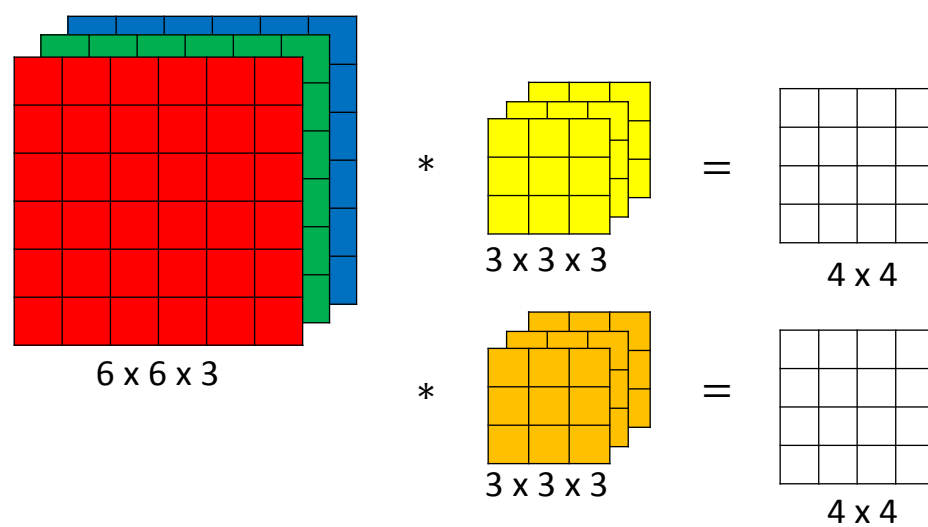

## Convolutions on RGB images



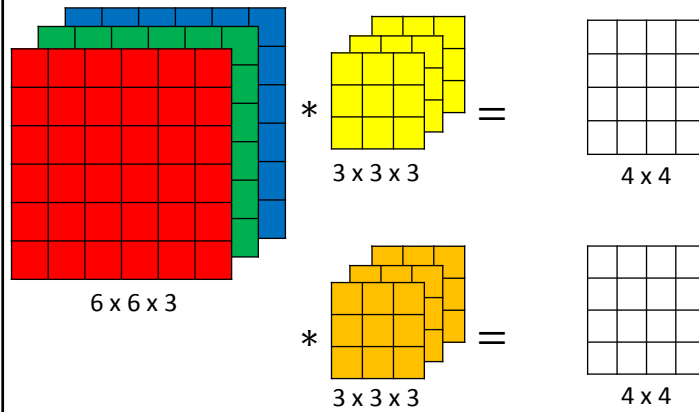
## Convolutions on RGB images



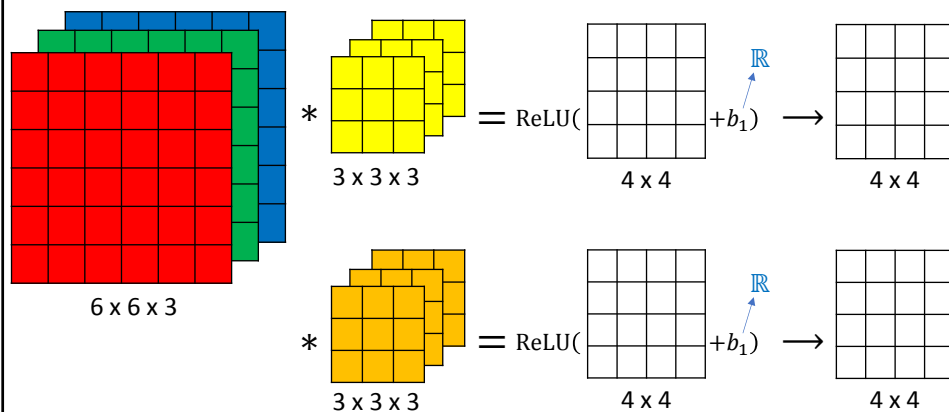
## Multiple filters



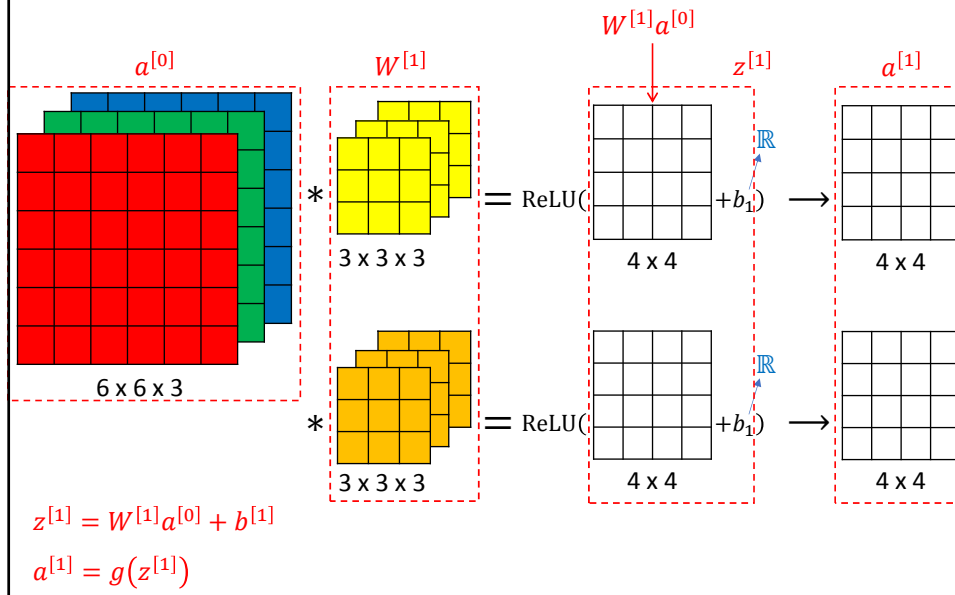
## Example of a layer



## Example of a layer



## Example of a layer



## Number of parameters in one layer

- If you have 10 filters that are 3 x 3 x 3 in one layer of a neural network, how many parameters does that layer have?

## Summay of notation

If layer  $l$  is a convolution layer:

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

$p^{[l]}$  = padding

$s^{[l]}$  = stride

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

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

## Summay of notation

If layer  $l$  is a convolution layer:

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

$p^{[l]}$  = padding

$s^{[l]}$  = stride

$n_C^{[l]}$  = number of filters

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

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

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

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

## Summay of notation

If layer  $l$  is a convolution layer:

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

$p^{[l]}$  = padding

$s^{[l]}$  = stride

$n_C^{[l]}$  = number of filters

Each filter is  $f^{[l]} \times f^{[l]} \times n_C^{[l-1]}$

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

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

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

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

## Summay of notation

If layer  $l$  is a convolution layer:

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

$p^{[l]}$  = padding

$s^{[l]}$  = stride

$n_C^{[l]}$  = number of filters

Each filter is  $f^{[l]} \times f^{[l]} \times n_C^{[l-1]}$

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

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

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

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

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

## Summay of notation

If layer  $l$  is a convolution layer:

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

$p^{[l]}$  = padding

$s^{[l]}$  = stride

$n_C^{[l]}$  = number of filters

Each filter is  $f^{[l]} \times f^{[l]} \times n_C^{[l-1]}$

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

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

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

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

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

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

## Summay of notation

If layer  $l$  is a convolution layer:

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

$p^{[l]}$  = padding

$s^{[l]}$  = stride

$n_C^{[l]}$  = number of filters

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:  $f^{[l]} \times f^{[l]} \times n_C^{[l-1]} \times n_C^{[l]}$

Bias:  $n_C^{[l]}$

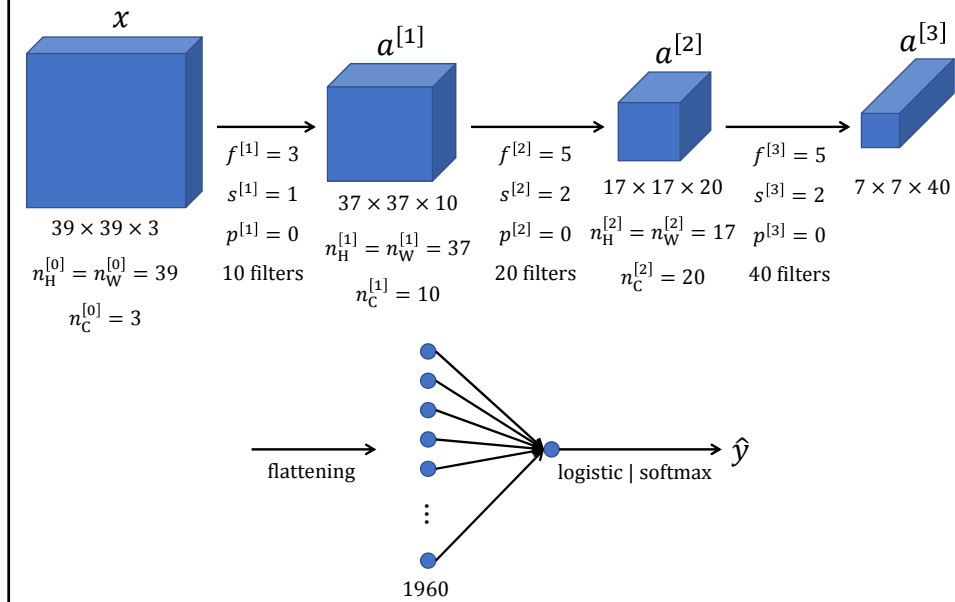
Input:  $n_H^{[l-1]} \times n_W^{[l-1]} \times n_C^{[l-1]}$

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

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

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

## Example ConvNet



## Types of layer in a convolutional network

- Convolution (CONV)
- Pooling (POOL)
- Fully connected (FC)



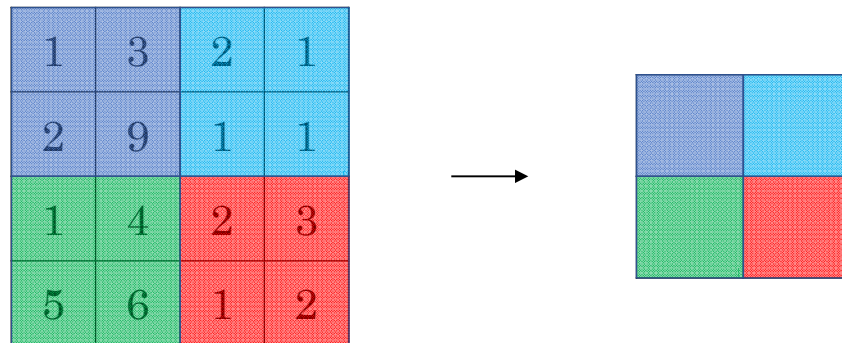
## Pooling layer: Max pooling

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


## Pooling layer: Max pooling

1	3	2	1	3
2	9	1	1	5
1	3	1	1	2
8	3	1	1	0
5	6	1	2	9


## Pooling layer: Average pooling



## Summary of pooling

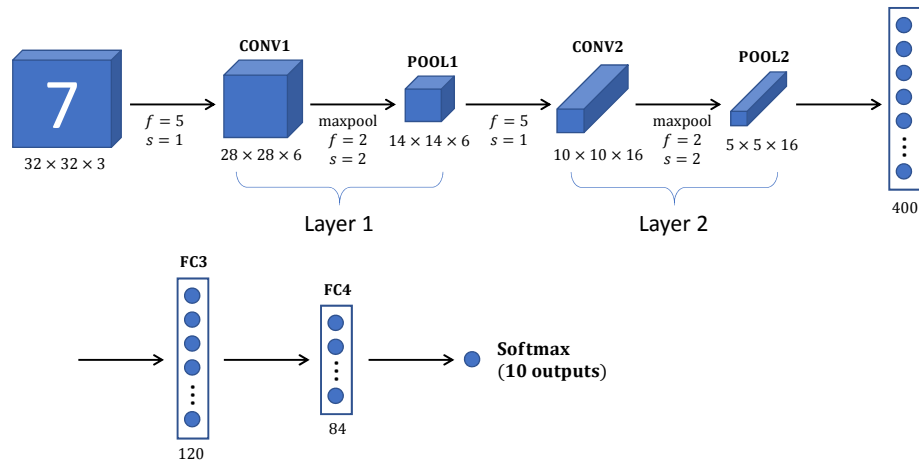
Hyperparameters:

$f$  : filter size

$s$  : stride

Max or average pooling

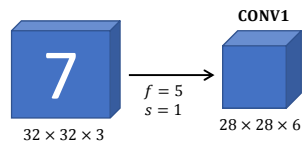
## Neural network example (LeNet-5)



## Neural network example (LeNet-5)

	Activation shape	Activation Size	# parameters
Input:	(32,32,3)	3,072	0
CONV1 ( $f=5, s=1$ )	(28,28,8)	6,272	208
POOL1	(14,14,8)	1,568	0
CONV2 ( $f=5, s=1$ )	(10,10,16)	1,600	416
POOL2	(5,5,16)	400	0
FC3	(120,1)	120	48,120
FC4	(84,1)	84	10,164
Softmax	(10,1)	10	850

## Why convolutions?



## Why convolutions?

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

\*

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

=

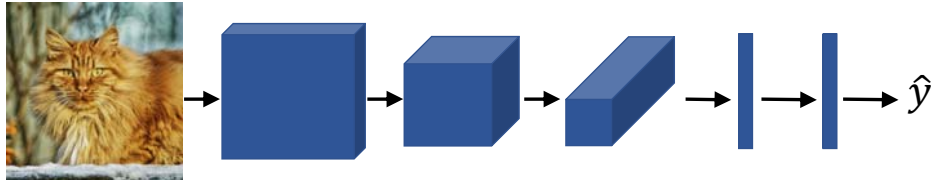
0	30	30	0
0	30	30	0
0	30	30	0
0	30	30	0

**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.

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

## Putting it together

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



$$\text{Cost } J = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

Use gradient descent to optimize parameters  
to reduce  $J$