

 intro

deeplearning.ai

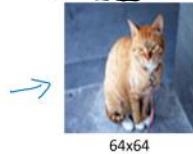
Convolutional Neural Networks

Computer vision

The literature is creative
can find inspiration from it
for other algorithms.

Computer Vision Problems

Image Classification



64x64

→ Cat? (0/1)

Neural Style Transfer ↗

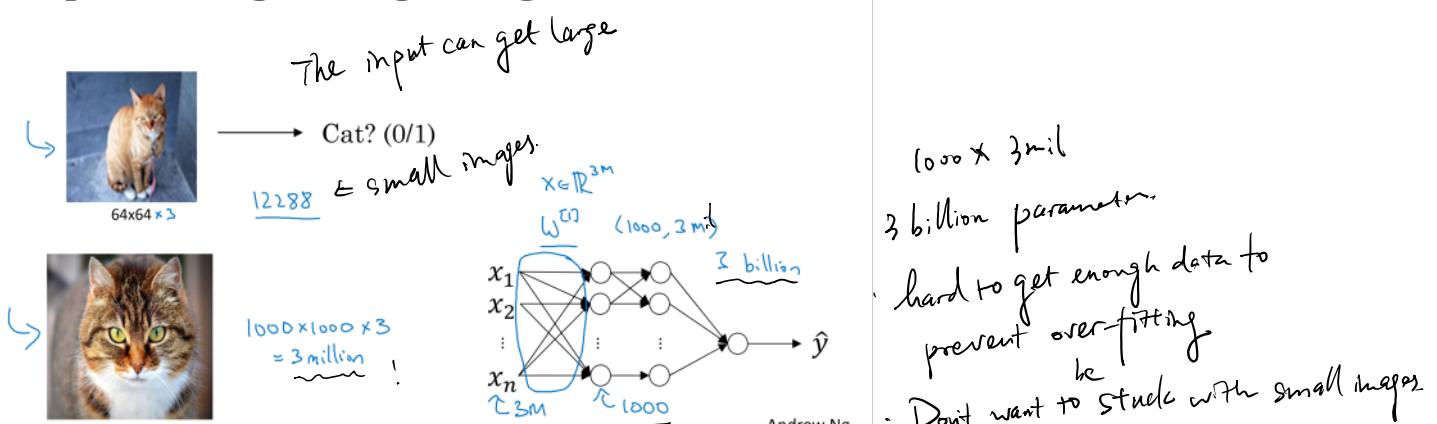


repaint the image
on the right with
the style on the right
Andrew Ng

Object detection



Deep Learning on large images



edge detect

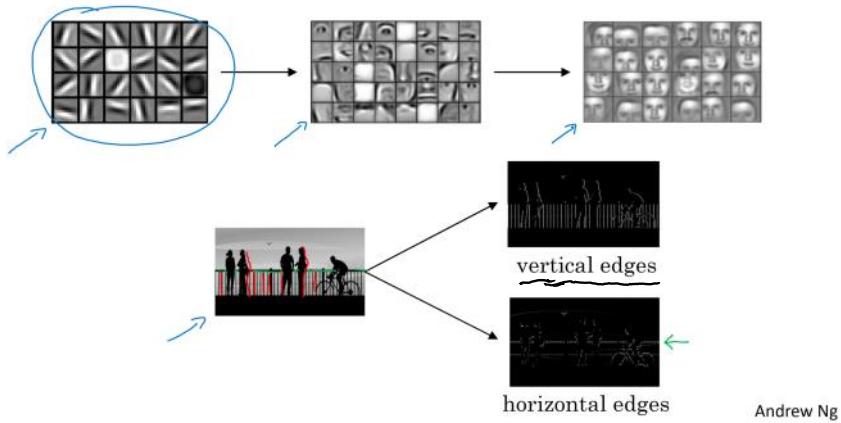


deeplearning.ai

Convolutional Neural Networks

Edge detection example

Computer Vision Problem



Vertical edge detection

$$\begin{array}{|c|c|c|c|c|c|} \hline
 3 & 0 & 1 & 2 & 7 & 4 \\ \hline
 1 & 5 & 8 & 3 & 0 & -1 \\ \hline
 2 & 7 & 2 & 1 & 0 & -1 \\ \hline
 0 & 1 & 3 & 1 & 7 & 8 \\ \hline
 4 & 2 & 1 & 6 & 2 & 8 \\ \hline
 2 & 4 & 5 & 2 & 3 & 9 \\ \hline
 \end{array}$$

*6x6 x 1
gray scale
w/ separate RGB
channels.*

"convolution" in mathematics - this denote convolution. But in Python `*` means multiplication -

Convolute
with a 3×3
filter -

$\begin{array}{|c|c|c|} \hline
 1 & 0 & -1 \\ \hline
 1 & 0 & -1 \\ \hline
 1 & 0 & -1 \\ \hline
 \end{array}$ $\quad = \quad \begin{array}{|c|c|c|c|} \hline
 -5 & -4 & 0 & 8 \\ \hline
 -10 & -2 & 2 & 3 \\ \hline
 0 & -2 & -4 & -7 \\ \hline
 -3 & -2 & -3 & -16 \\ \hline
 \end{array}$
 3×3
 \rightarrow filter kernel

python: `conv-forward`
tensorflow: `tf.nn.conv2d`

Andrew Ng

Vertical edge detection

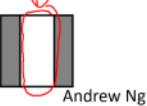
$$\begin{array}{|c|c|c|c|c|c|c|}\hline 10 & 10 & 10 & 0 & 0 & 0 & 0 \\ \hline 10 & 10 & 10 & 0 & 0 & 0 & 0 \\ \hline 10 & 10 & 10 & 0 & 0 & 0 & 0 \\ \hline 10 & 10 & 10 & 0 & 0 & 0 & 0 \\ \hline 10 & 10 & 10 & 0 & 0 & 0 & 0 \\ \hline 10 & 10 & 10 & 0 & 0 & 0 & 0 \\ \hline \end{array} \quad * \quad \begin{array}{|c|c|c|}\hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline 1 & 0 & -1 \\ \hline \end{array}_{3 \times 3} = \begin{array}{|c|c|c|c|c|c|}\hline 0 & 30 & 30 & 0 & & \\ \hline 0 & 30 & 30 & 0 & & \\ \hline 0 & 30 & 30 & 0 & & \\ \hline 0 & 30 & 30 & 0 & & \\ \hline \end{array}_{4 \times 4}$$

\downarrow

$\uparrow 4 \times 4$

*





\downarrow

$\uparrow 4 \times 4$

Andrew Ng

6×6

$\uparrow \uparrow \uparrow$

 more edge



deeplearning.ai

Convolutional Neural Networks

More edge
detection

light to dark

dark to light.

Ex:

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

*

-J	-J	-J
0	0	0
J	J	J

=

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

Sign weights

Vertical vs horizontal edge

Vertical

J	0	-J
J	0	-J
J	0	-J

Horizontal

-J	-J	-J
0	0	0
J	J	J

← negative sign

Vertical and Horizontal Edge Detection

~~~~~

Sobel filter: put more weights to  
the central weight

$\Rightarrow$  more robust

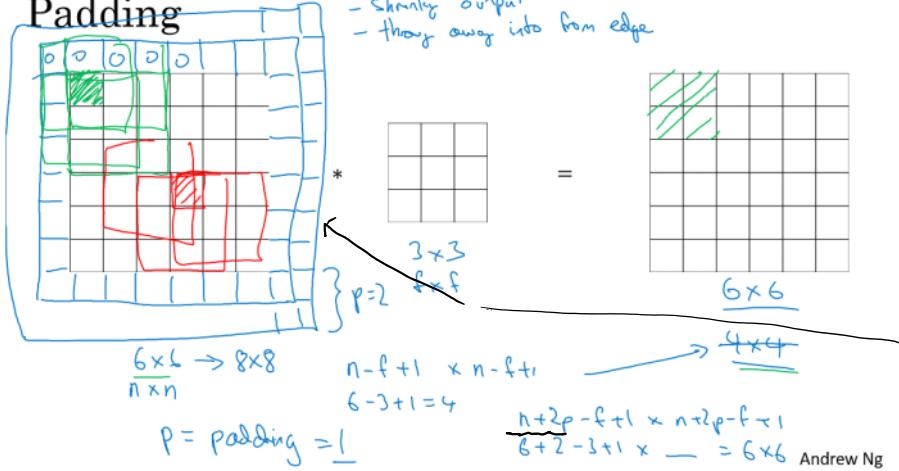
Scharr filter: another filter

with ML, filters need not be hand-picked  
it can be learned!

Treating the  $\gamma$  numbers  $w_i, i \in \{1, 9\}$   
as parameters,

 padding

## Padding



$n \times l$  picture

Dimension  $f \times f$  filter size  
 $(n-f+1) \times (n-f+1)$

Problem ① Image shrink every time  
you detect edges, etc.

② pixels at the corners are used less  
⇒ throwing away information

Solution pad the image.

$$p = \text{padding}$$

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

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

## Valid and Same convolutions

→ no padding

"Valid":  $n \times n * f \times f \rightarrow n-f+1 \times n-f+1$

$6 \times 6 * 3 \times 3 \rightarrow 4 \times 4$

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

$$n+2p-f+1 \times n+2p-f+1$$

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

$$3 \times 3 \quad p = \frac{3-1}{2} = 1 \quad \begin{array}{c} f \text{ is usually odd} \\ \boxed{3 \times 3} \\ 5 \times 5 \\ 7 \times 7 \end{array}$$

Andrew Ng

Valid conv: no padding.

Same conv: pad  $\Rightarrow$  same size.

odd num filters convention  
(x) is possible

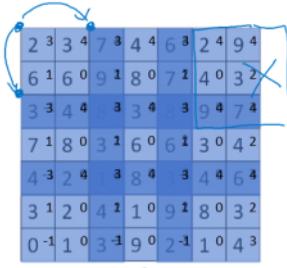


deeplearning.ai

# Convolutional Neural Networks

## Strided convolutions

### Strided convolution

 \* 

$\underbrace{n \times n}_{\text{input}} \quad * \quad \underbrace{f \times f}_{\text{filter}}$   
 $\underbrace{p}_{\text{padding}} = 1 \quad \underbrace{s}_{\text{stride}} = 2$

$\begin{bmatrix} 3 & 4 & 4 \\ 1 & 0 & 2 \\ -1 & 0 & 3 \end{bmatrix} = \begin{bmatrix} 9_1 & 10_0 & 8_3 \\ 6_9 & 9_1 & 12_7 \\ 4_4 & 7_2 & 7_4 \end{bmatrix}$

$\boxed{\text{stride } = 2} \quad \boxed{\lfloor \frac{2}{2} \rfloor = \text{floor}(\frac{2}{2})}$

$\left[ \frac{n+2p-f}{s} + 1 \right]^{f \times f} \times \left[ \frac{n+2p-f}{s} + 1 \right]$

$\frac{7+0-3}{2} + 1 = \frac{4}{2} + 1 = 3$

Andrew Ng

$\left[ \frac{n+2p-f}{s} + 1 \right]$  round down.

## Summary of convolutions

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

padding  $p$       stride  $s$

$$\boxed{\text{Output size:} \left\lceil \frac{n+2p-f}{s} + 1 \right\rceil \quad \times \quad \left\lceil \frac{n+2p-f}{s} + 1 \right\rceil}$$

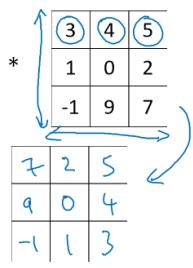
Andrew Ng

→ stride > 1 output  $\frac{1}{s}$  sample...

Technical note on cross-correlation vs. convolution

Convolution in math textbook: flip before operation

|   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|
| 2 | 3 | 2 | 7 | 5 | 4 | 6 | 2 |
| 6 | 6 | 9 | 4 | 8 | 7 | 4 |   |
| 3 | 4 | 8 | 3 | 3 | 8 | 9 |   |
| 7 | 8 | 3 | 6 | 6 | 6 | 3 |   |
| 4 | 2 | 1 | 8 | 3 | 4 |   |   |
| 3 | 2 | 4 | 1 | 9 | 8 |   |   |



$$= \begin{array}{|c|c|c|} \hline 1 & 0 & 2 \\ \hline 9 & 0 & 4 \\ \hline -1 & 1 & 3 \\ \hline \end{array}$$

$$(A * B) * C = A * (B * C)$$

Andrew Ng

## Cross-correlation

this is Cross-correlation  
in DL called "convolution"  
in convention

in signal process: flipping helps in getting associativity property  
the But this doesn't matter for DL  
This is just a caution about definition

↳ notation -  
May matter when you read <sup>some</sup> papers

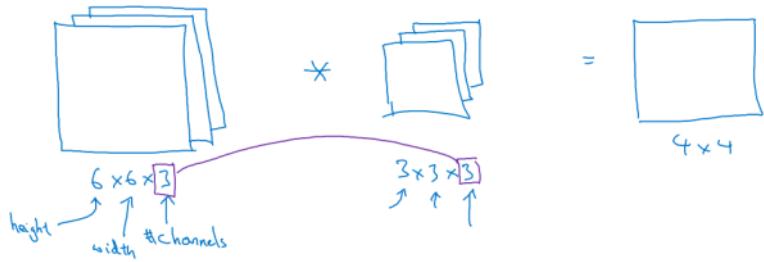


deeplearning.ai

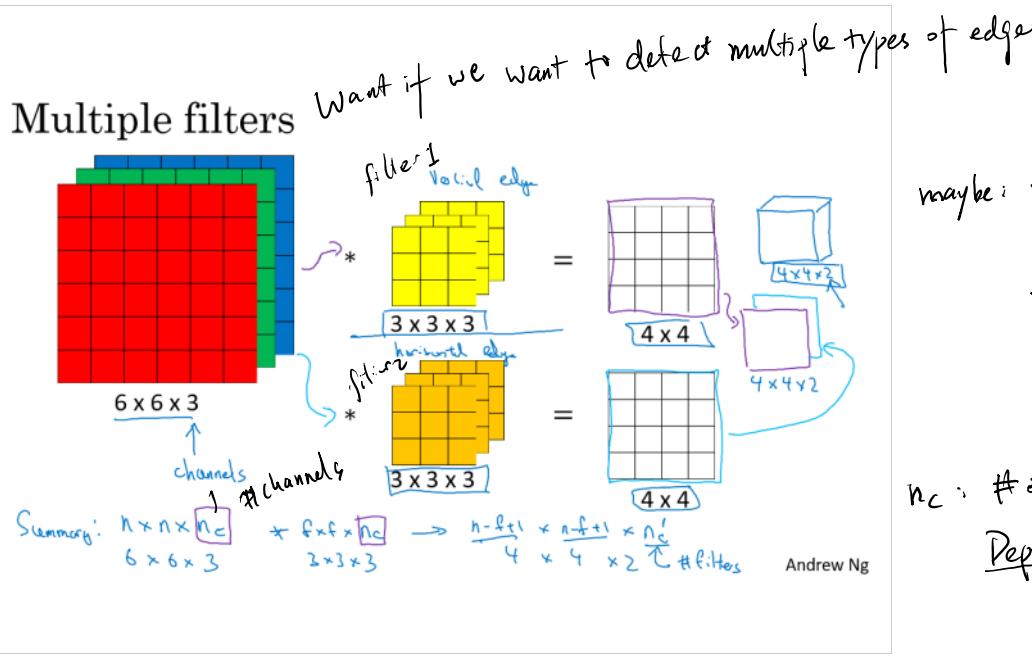
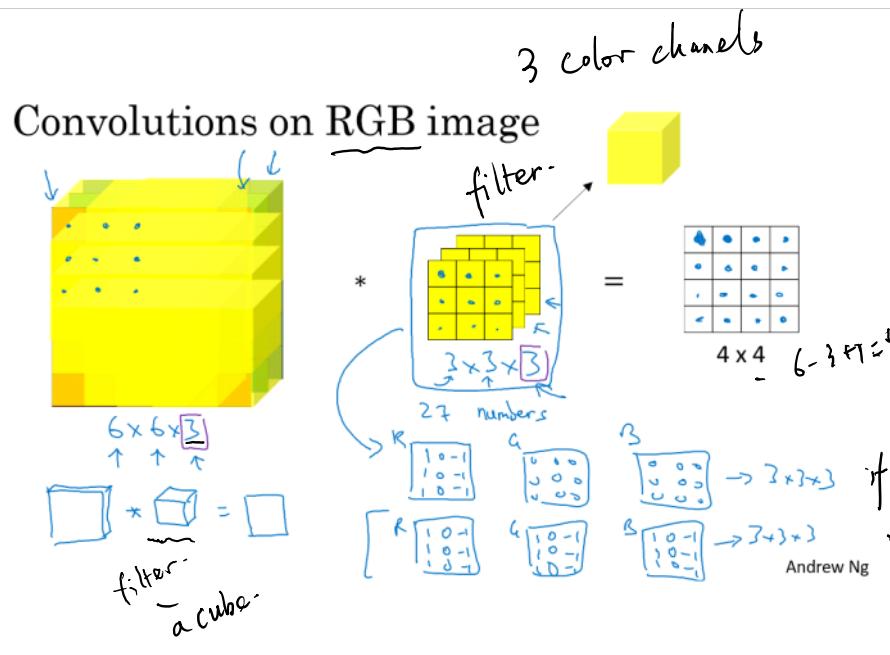
# Convolutional Neural Networks

## Convolutions over volumes

### Convolutions on RGB images

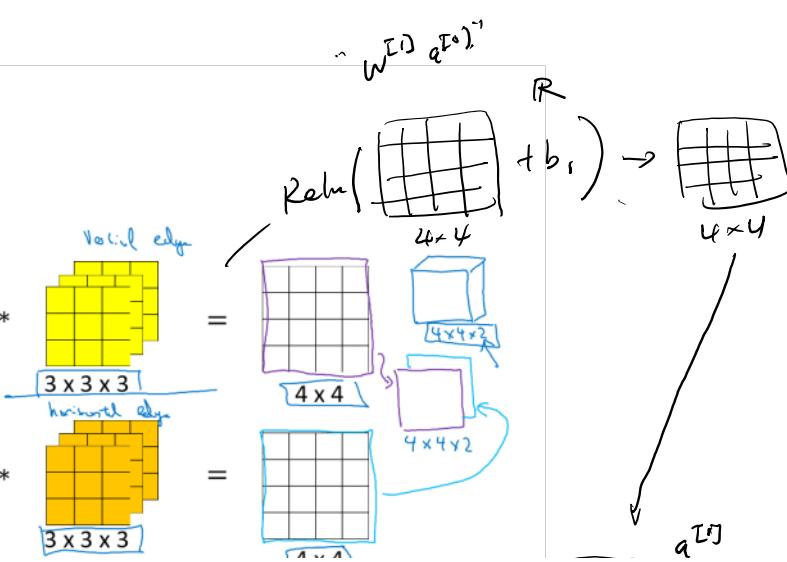
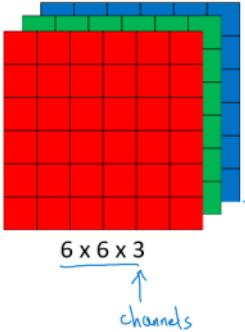


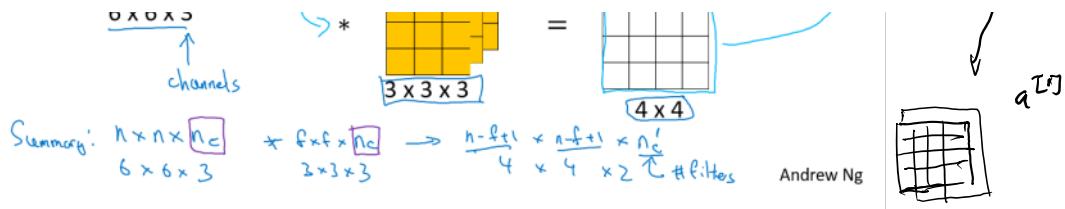
Andrew Ng



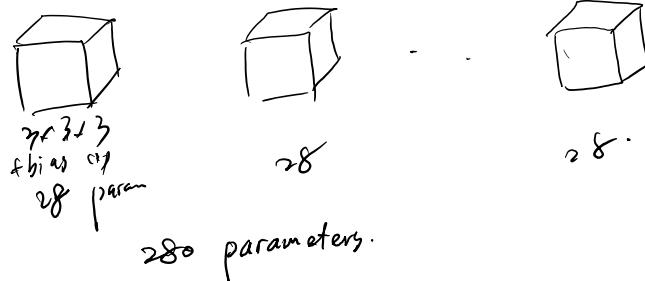
one layer of CNN

Multiple filters





If we have 10 filters that are  $3 \times 3 \times 3$  in one layer.  
How many param?



10 filters that detect different edges  
can work for even large images  
 $\Rightarrow$  significantly less parameters  
than not using CNN (described in the notes)  
 $\Rightarrow$  avoid overfitting problem.

If layer  $l$  is a convolution layer

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

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

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

$$n_c^{[l]} = \text{number of filters.}$$

$$\text{Each filter is: } f^{[l]} \times f^{[l]} \times n_c^{[l-1]}$$

$$\text{Activation: } a^{[l]} \rightarrow n_H^{[l]} \times n_w^{[l]} \times n_c^{[l]}$$

$$\text{Weights: } f^{[l]} \times f^{[l]} \times n_c^{[l-1]} \times n_c^{[l]}$$

$\curvearrowright$  # of filters in layer  $l$

$$\text{bias: } - (1, 1, 1, n_c^{[l]}) \curvearrowleft \text{by channel } \rightarrow \text{bias.}$$

Caution: in open source implementation, no convention about the order of dimensions  
Some frameworks put channel first.

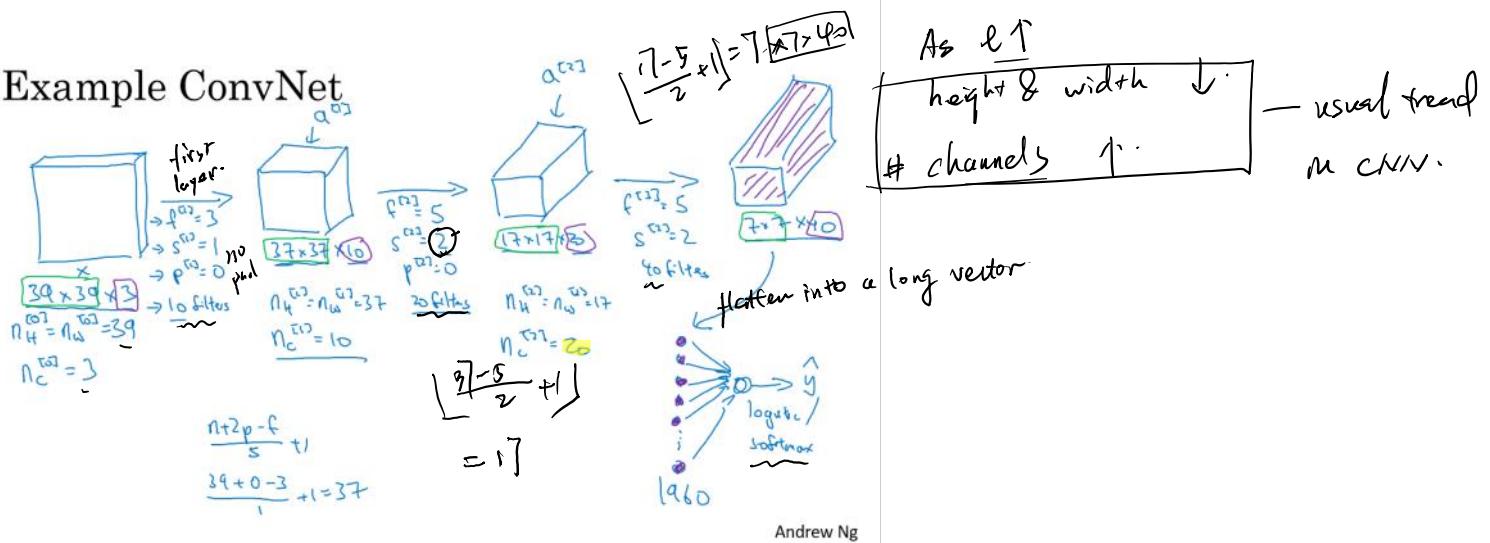


deeplearning.ai

# Convolutional Neural Networks

## A simple convolution network example

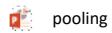
### Example ConvNet



Types of layer in a convolutional network:

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

Andrew Ng



pooling



deeplearning.ai

Convolutional  
Neural Networks

Pooling layers

## Pooling layer: Max pooling

|   |   |   |   |
|---|---|---|---|
| 1 | 3 | 2 | 1 |
| 2 | 9 | 1 | 1 |
| 1 | 3 | 2 | 3 |
| 5 | 6 | 1 | 2 |

|   |   |
|---|---|
| 9 | 2 |
| 6 | 3 |

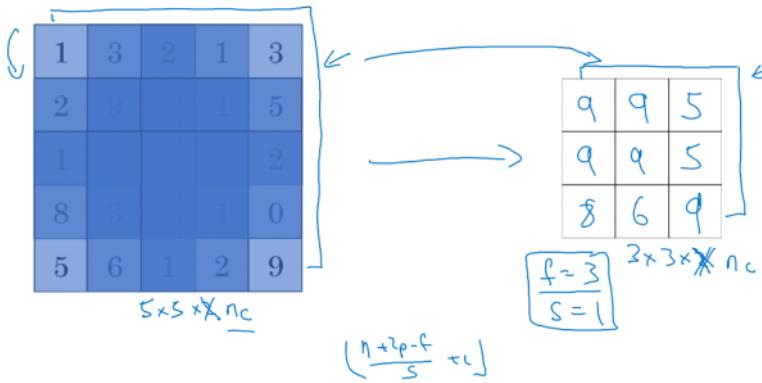
(large value can be viewed as detecting some feature.)

Max Pooling: if this feature is detected in this region → keep the large number  
 Andrew Ng (Intuition)  
 But why it is used? Proved to work.

No proof  
 that this is the real underlying mechanism

## Pooling layer: Max pooling

这里我们用 3x3 filter, filter size f  
 stride s



Andrew Ng

## Pooling layer: Average pooling

|   |   |   |   |
|---|---|---|---|
| 1 | 3 | 2 | 1 |
| 2 | 9 | 1 | 1 |
| 1 | 4 | 2 | 3 |
| 5 | 6 | 1 | 2 |



|      |      |
|------|------|
| 3.75 | 1.25 |
| 4    | 2    |

$$f=2$$

$$s=2$$

$$\underline{7 \times 7 \times 1000} \rightarrow 1 \times 1 \times 1000$$

Andrew Ng

Max pooling is used more often than avg pooling.  
In very deep NN, may need to use avg pooling

→ collapses represen  $7 \times 7 \times 1000$

$$\Rightarrow 1 \times 1 \times 1000$$

## Summary of pooling

Hyperparameters:

$f$ : filter size       $f=2, s=2$   
 $s$ : stride               $f=3, s=2$   
Max or average pooling  
 $\rightarrow p$ : padding

No parameters to learn!

$$\begin{aligned}
 & n_H \times n_W \times n_C \\
 & \downarrow \\
 & \left\lfloor \frac{n_H-f+1}{s} \right\rfloor \times \left\lfloor \frac{n_W-f}{s} + 1 \right\rfloor \\
 & \times n_C
 \end{aligned}$$

usually no padding.

Andrew Ng

cnn example



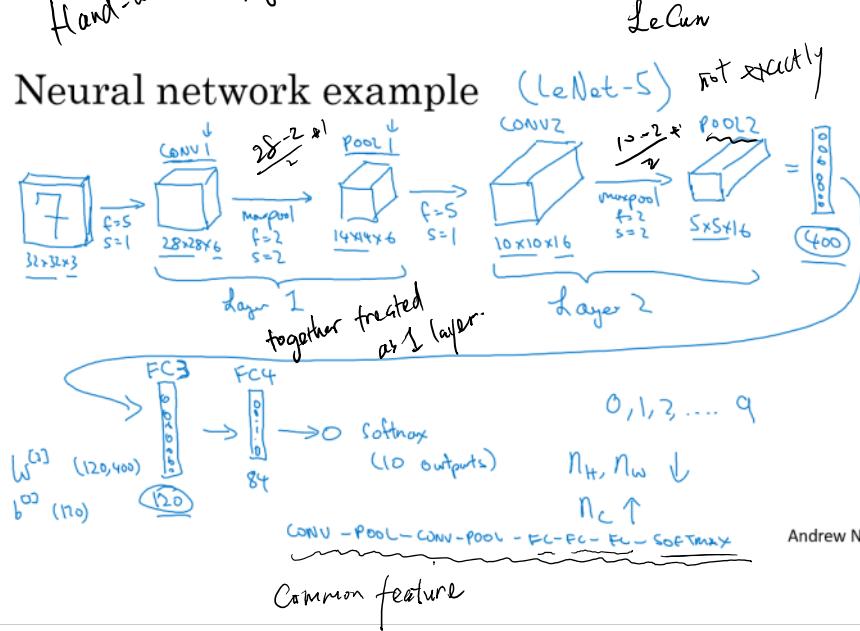
deeplearning.ai

# Convolutional Neural Networks

## Convolutional neural network example

Hand-written digit recognition

### Neural network example



Notation thing: counting layers that have weights

search for hyper-parameters in the literature

Don't blindly search for own hyperparam  
usually: channel ( $n_c$ )  $\uparrow \Rightarrow$  full connected

## Neural network example

|        | Activation shape | Activation Size | # parameters |
|--------|------------------|-----------------|--------------|
| Input: | (32,32,3)        | 3,072 $a^{(1)}$ | 0            |
|        |                  | —               | ←            |
|        |                  | —               | ←            |
|        |                  | —               | ←            |
|        |                  | —               | ←            |
|        |                  | —               | —            |
|        |                  | —               | —            |
|        |                  | —               | —            |
|        |                  | —               | —            |
|        |                  | —               | —            |

Andrew Ng

• Max Pooling layer have no params.

Fully connected layers have lots of parameters.

lots of conv net share similar structure

putting layers together requires insight

→ It'll be good to build on prior work.

Why convolution

 why conv

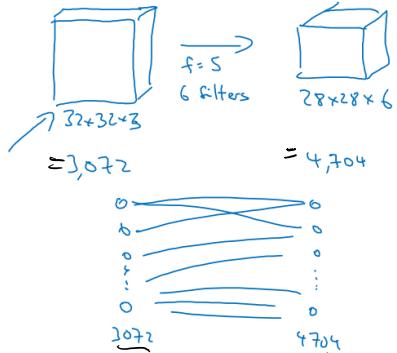


deeplearning.ai

## Convolutional Neural Networks

### Why convolutions?

## Why convolutions



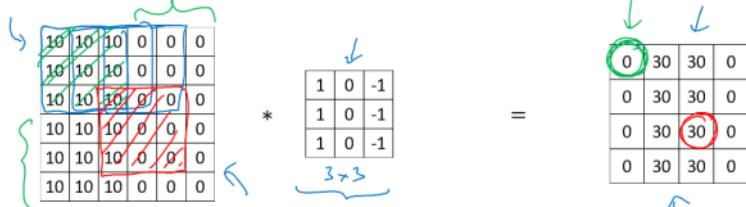
$$\begin{aligned} 5 \times 5 &= 25 \\ 26 & \\ 6 \times 26 &= 156 \text{ parameters} \end{aligned}$$

$$3072 \times 4704 \approx 14M$$

Andrew Ng

~~If fully connected to start~~  
 $\Rightarrow$  too many parameters!

## Why convolutions



**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.

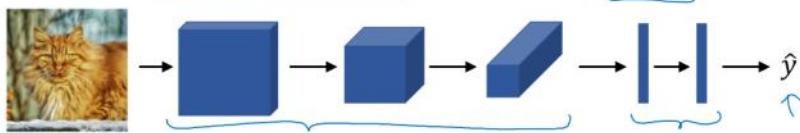
← How this works  
local features.  
output depends on only a small # of inputs

Andrew Ng

## Putting it together

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

$w, b$



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

Use gradient descent to optimize parameters to reduce  $J$

Andrew Ng

$$\begin{aligned} w^{(2)} \times n \times n \times n \times n &= v^{(2)} \\ p = \frac{1}{n} \times \frac{1}{n} \times \frac{1}{n} \times \frac{1}{n} &= \frac{1}{n^4} \end{aligned}$$

✓

Andrew Ng

✓



3. Suppose your input is a 300 by 300 color (RGB) image, and you use a convolutional layer with 100 filters that are each 5x5. How many parameters does this hidden layer have (including the bias parameters)?

7600 -- 所以自动乘以三?

1 / 1  
points



2. Suppose your input is a 300 by 300 color (RGB) image, and you are not using a convolutional network. If the first hidden layer has 100 neurons, each one fully connected to the input, how many parameters does this hidden layer have (including the bias parameters)?

27000100

1 / 1  
points



8. Because pooling layers do not have parameters, they do not affect the backpropagation (derivatives) calculation.

False!

1 / 1  
points