# Computer Vision Problems

### Image Classification

Cat? (0/1)

64x64

### Object detection

### Neural Style Transfer

---

# Deep Learning on large images

Cat? (0/1)

$64\times64\times3$ → 12288

$1000\times1000\times3$ = 3 million

$X \in \mathbb{R}^{3M}$

$W^{[1]}$  (1000, 3m)

3 billion

$x_1$
$x_2$
$\vdots$
$x_n$

↰ 3M    ↰ 1000

$\hat{y}$

---

# Computer Vision Problem

vertical edges

horizontal edges

# Vertical edge detection

$3\times1 + 1\times1 + 2\times1 + 0\times0 + 5\times0 + 7\times0 + 1\times-1 + 8\times-1 + 2\times-1 = -5$

| 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 |

6×6

"convolution"

$*$

| 1 | 0 | -1 |
|---|---|---|
| 1 | 0 | -1 |
| 1 | 0 | -1 |

3×3 filter    "Kernel"

$=$

| -5 | -4 | 0 | 8 |
|---|---|---|---|
| -10 | -2 | 2 | 3 |
| 0 | -2 | -4 | -7 |
| -3 | -2 | -3 | -16 |

4×4

python: conv-forward
tensorflow: tf.nn.conv2d
keras: Conv2D

Andrew Ng

# Vertical edge detection

edge  edge

| 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 |

6×6

$*$

| 1 | 0 | -1 |
|---|---|---|
| 1 | 0 | -1 |
| 1 | 0 | -1 |

3×3

$=$

| 0 | 30 | 30 | 0 |
|---|----|----|---|
| 0 | 30 | 30 | 0 |
| 0 | 30 | 30 | 0 |
| 0 | 30 | 30 | 0 |

4×4

$*$

Andrew Ng

# 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 |

Light ·········· Dark

$*$

| 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 |

| 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 |

$*$

| 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 |

Andrew Ng

# 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 |

6 x 6

*

| 1 | 1 | 1 |
|---|---|---|
| 0 | 0 | 0 |
| -1 | -1 | -1 |

=

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

Andrew Ng

# Learning to detect edges

| 1 | 0 | -1 |
|---|---|----|
| 1 | 0 | -1 |
| 1 | 0 | -1 |

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

"More Robust"

Sobel filter

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

Scharr filter

| 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 |

convolution

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

3×3

$$\begin{bmatrix} 45° \\ 70° \\ 73° \end{bmatrix}$$

Andrew Ng

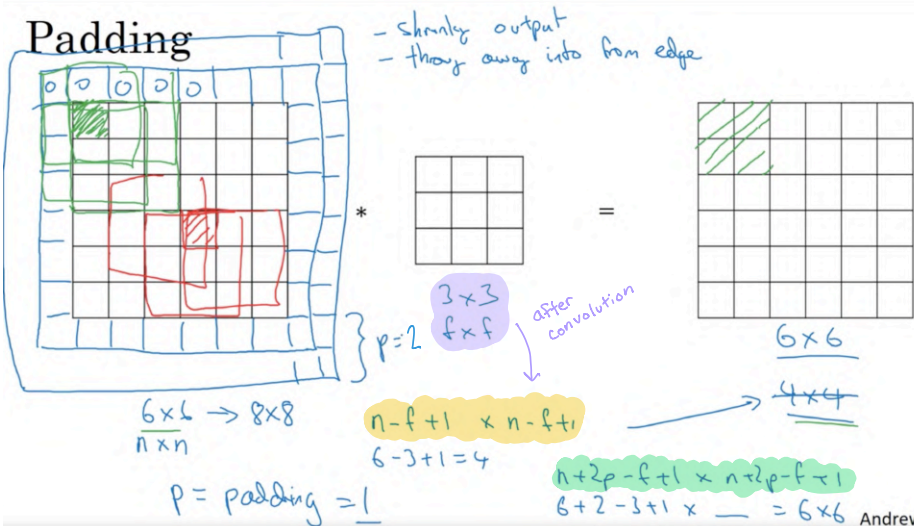# Padding

The two downsides to this; one is that, if every time you apply a convolutional operator, your image shrinks, so you come from six by six down to four by four then, you can only do this a few times before your image starts getting really small, maybe it shrinks down to one by one or something, so maybe, you don't want your image to shrink every time you detect edges or to set other features on it, so that's one downside is that, if you look the pixel at the corner or the edge, this little pixel is touched as used only in one of the outputs, because this touches that three by three region. Whereas, if you take a pixel in the middle, say this pixel, then there are a lot of three by three regions that overlap that pixel and so, is as if pixels on the corners or on the edges are use much less in the output. So you're throwing away a lot of the information near the edge of the image.

In order to fix both of these problems, what you can do is the full apply of convolutional operation. You can pad the image.

- Shrink output
- throw away info from edge

*

=

$p = 2$

$3 \times 3$
$f \times f$

after convolution

6×6

4×4

$6 \times 6 \rightarrow 8 \times 8$
$n \times n$

$n - f + 1 \times n - f + 1$
$6 - 3 + 1 = 4$

$p = padding = 1$

$n + 2p - f + 1 \times n + 2p - f + 1$
$6 + 2 - 3 + 1 \times \_\_ = 6 \times 6$

Andrew Ng

# Valid and Same convolutions

→ no padding

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

$6 \times 6$ $*$ $3 \times 3$ → $4 \times 4$

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

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

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

$3 \times 3$ $p = \dfrac{3-1}{2} = 1$

$5 \times 5$ $f = 5$ $p = 2$

odd!

f is usually odd

$1 \times 1$
$3 \times 3$ ✓
$5 \times 5$ ✓
$7 \times 7$ ✓

Andrew Ng

---

# Strided convolution

Stride of 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 |

$7 \times 7$

$*$

| 3 | 4 | 4 |
|---|---|---|
| 1 | 0 | 2 |
| -1 | 0 | 3 |

$3 \times 3$

Stride = 2

$=$

| 91 | 100 | 83 |
|---|---|---|
| 69 | $a_1$ | 127 |
| 44 | 72 | 74 |

$3 \times 3$

$\lfloor z \rfloor = floor(z)$

$n \times n$ $*$ $f \times f$

padding $p$ stride $S$

$S = 2$

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

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

Andrew Ng

---

# 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 |

$=$

| 7 | 9 | -1 |
|---|---|---|
| 2 | 0 | 1 |
| 5 | 4 | 3 |

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

Andrew Ng

# Convolutions on RGB images



$6 \times 6 \times 3$ ······ Same ······ $3 \times 3 \times 3$

height width #channels

$4 \times 4$

# Convolutions on RGB image



height width channels

$6 \times 6 \times 3$

$3 \times 3 \times 3$

27 numbers

$4 \times 4$

R  G  B

$\rightarrow 3 \times 3 \times 3$

$\rightarrow 3 \times 3 \times 3$

Edge Detector

# Multiple filters



Vertical edge

$6 \times 6 \times 3$

channels

$3 \times 3 \times 3$

$= \quad 4 \times 4$

horizontal edge

$3 \times 3 \times 3$

$= \quad 4 \times 4$

$4 \times 4 \times 2$

$4 \times 4 \times 2$

$n_c$ of next layer

Summary: $n \times n \times n_c \quad * \quad f \times f \times n_c \quad \rightarrow \quad n-f+1 \times n-f+1 \times n_c'$

$6 \times 6 \times 3 \qquad 3 \times 3 \times 3 \qquad \quad 4 \quad \times \quad 4 \quad \times 2$ #filters

# Example of a layer



$\cancel{8}$ filters
10

"$W^{[1]} a^{[0]}$"

$R \rightarrow z^{[1]}$

$\rightarrow$ ReLU

$+ b_1$

$4 \times 4$

$\rightarrow$ ReLU

$+ b_2$

$4 \times 4$

$6 \times 6 \times 3$

$a^{[0]}$

$3 \times 3 \times 3$

$3 \times 3 \times 3$

"$W^{[1]}$"

$z^{[1]} = W^{[1]} a^{[0]} + b^{[1]}$

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

stack

$6 \times 6 \times 3$   $a^{[0]}$

$a^{[1]}$  $4 \times 4 \times 2$

$4 \times 4 \times 2$

$4 \times 4 \times 10$  $\leftarrow a^{[1]}$

One layer

# Number of parameters in one layer

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


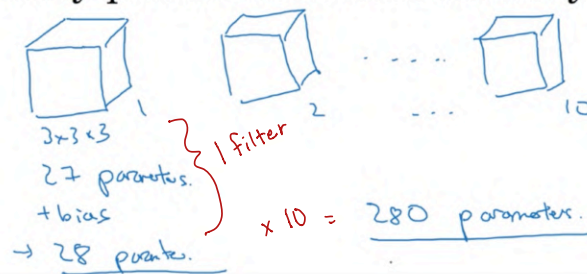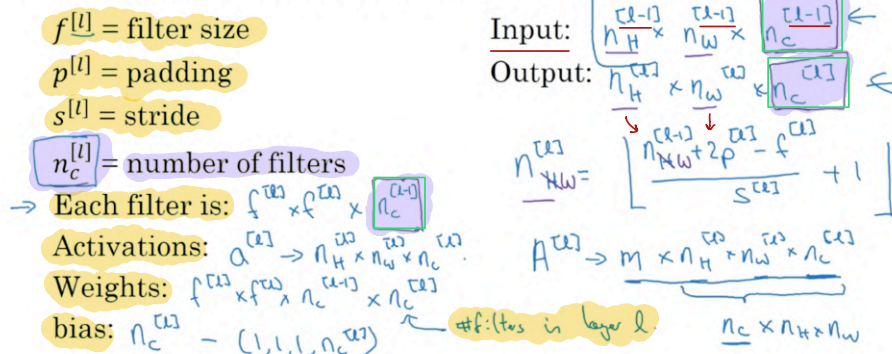
2

10

$3 \times 3 \times 3$
27 parameters.
+ bias
$\rightarrow$ 28 parameters.

1 filter

$\times 10 = $ <u>280</u> parameters.

# Summary 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
$\rightarrow$ Each filter is: $f^{[l]} \times f^{[l]} \times n_c^{[l-1]}$
Activations: $a^{[l]} \rightarrow n_H^{[l]} \times n_W^{[l]} \times n_c^{[l]}$
Weights: $f^{[l]} \times f^{[l]} \times n_c^{[l-1]} \times n_c^{[l]}$
bias: $n_c^{[l]} - (1,1,1,n_c^{[l]})$

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

$n_{H,W}^{[l]} = \left[ \dfrac{n_{H,W}^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1 \right]$

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

$\leftarrow$ #filters in layer $l$.

$n_c \times n_H \times n_W$

# Example ConvNet



$a^{[0]}$     $a^{[2]}$

$x$   $39 \times 39 \times 3$

$n_H^{[0]} = n_W^{[0]} = 39$

$n_c^{[0]} = 3$

$\rightarrow f^{[1]} = 3$
$\rightarrow s^{[1]} = 1$
$\rightarrow p^{[1]} = 0$
$\rightarrow$ 10 filters

$37 \times 37 \times 10$

$n_H^{[1]} = n_W^{[1]} = 37$
$n_c^{[1]} = 10$

$f^{[2]} = 5$
$s^{[2]} = 2$
$p^{[2]} = 0$
20 filters

$17 \times 17 \times 20$

$n_H^{[2]} = n_W^{[2]} = 17$
$n_c^{[2]} = 20$

$f^{[3]} = 5$
$s^{[3]} = 2$
40 filters

$7 \times 7 \times 40$

flatten

$\rightarrow \hat{y}$   logistic / softmax

1960

$\dfrac{n + 2p - f}{s} + 1$

$\dfrac{39 + 0 - 3}{1} + 1 = 37$

**Trend**

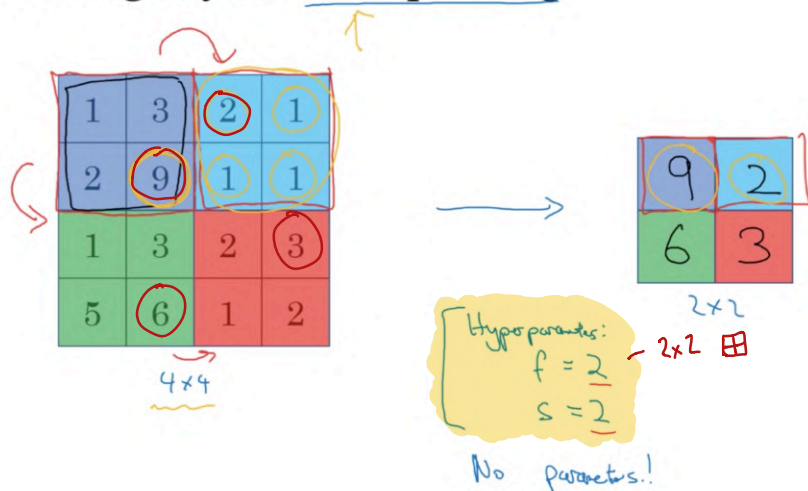$n_H, n_W$   $n_c$

Trend down   Trend Up

Andrew Ng

---

# Types of layer in a convolutional network:

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

Andrew Ng
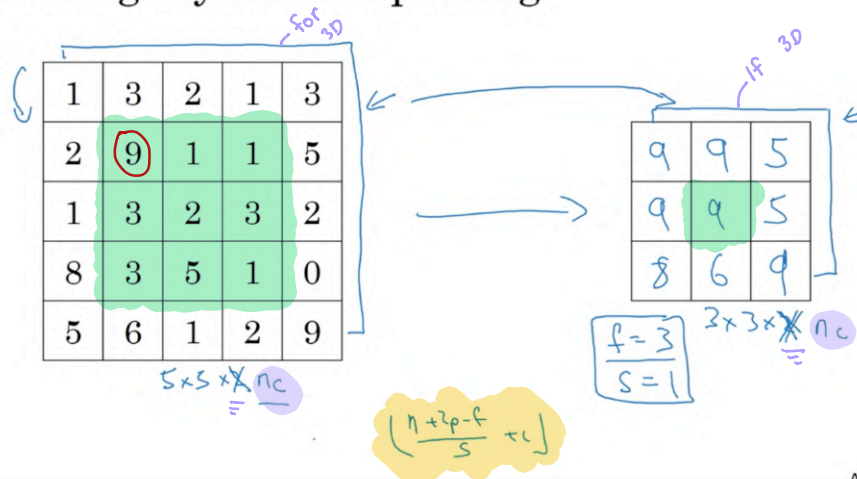
---

# Pooling layer: Max pooling



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

$4 \times 4$

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

$2 \times 2$

Hyperparameters:
$f = 2$   ← $2 \times 2$
$s = 2$

No parameters!

So here's the intuition behind what max pooling is doing. If you think of this four by four region as some set of features, the activations in some layer of the neural network, then a large number, it means that it's maybe detected a particular feature. So, the upper left-hand quadrant has this particular feature. It maybe a vertical edge or maybe a eye or whisker if you trying to detect a cat. Clearly, that feature exists in the upper left-hand quadrant. Whereas this feature, maybe it isn't cat eye detector. Whereas this feature, it doesn't really exist in the upper right-hand quadrant. So what the max operation does is a lots of features detected anywhere, and one of these quadrants , it then remains preserved in the output of max pooling. So, what the max operates to does is really to say, if these features detected anywhere in this filter, then keep a high number. But if this feature is not detected, so maybe this feature doesn't exist in the upper right-hand quadrant. Then the max of all those numbers is still itself quite small. So maybe that's the intuition behind max pooling.

Andrew Ng

## Pooling layer: Max pooling

*for 3D*

| 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 |

$5 \times 5 \times n_c$

*If 3D*

| 9 | 9 | 5 |
|---|---|---|
| 9 | 9 | 5 |
| 8 | 6 | 9 |

$3 \times 3 \times n_c$

$f = 3$
$s = 1$

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

## Pooling layer: Average pooling

Not Used Often

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

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

**You might use average pooling to collapse your representation from say, 7 by 7 by 1,000. An average over all the [inaudible] , you get 1 by 1 by 1,000.**

$f = 2$
$s = 2$

$7 \times 7 \times 1000 \longrightarrow 1 \times 1 \times 1000$

## Summary of pooling

Hyperparameters:

$$\begin{cases} f : \text{filter size} \\ s : \text{stride} \\ \text{Max or average pooling} \end{cases}$$
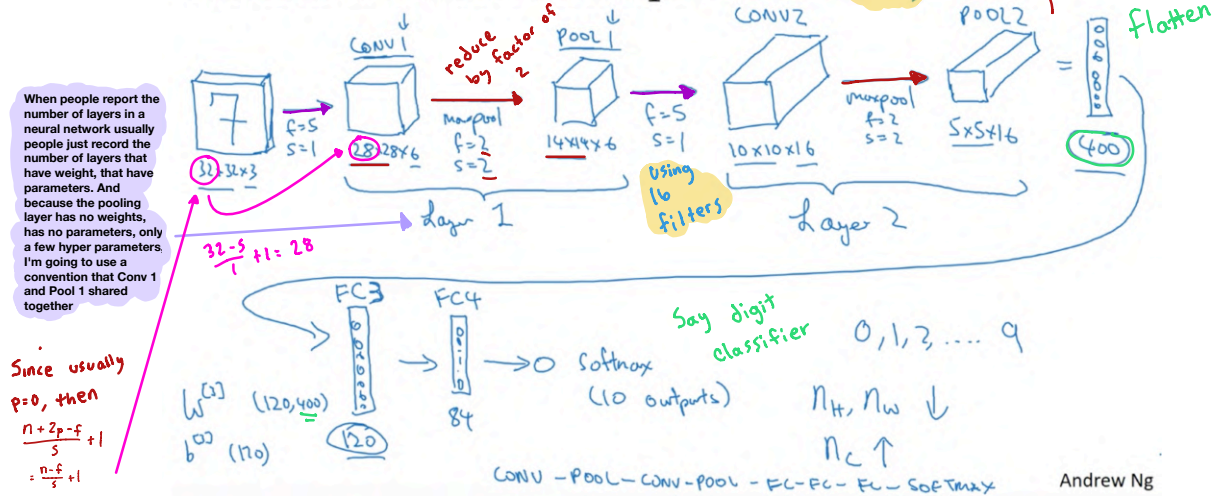
$f = 2, s = 2$
$f = 3, s = 2$

$\longrightarrow$ p: padding  *rarely used; usually p = 0*

No parameters to learn!

$n_H \times n_w \times n_c$

$\downarrow$

$$\left[ \frac{n_H - f}{s} + 1 \right] \times \left[ \frac{n_w - f}{s} + 1 \right] \times n_c$$

# Neural network example

(LeNet-5)

inspired by

flatten

CONV1 ↓  reduce by factor of 2  POOL1 ↓  CONV2  POOL2

$f=5$
$s=1$

maxpool
$f=2$
$s=2$

$14 \times 14 \times 6$

$f=5$
$s=1$

$10 \times 10 \times 16$

maxpool
$f=2$
$s=2$

$5 \times 5 \times 16$

$32 \times 32 \times 3$
$28 \times 28 \times 6$

$\frac{32-5}{1} + 1 = 28$

Layer 1

Using 16 filters

Layer 2

400

Since usually $p=0$, then

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

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

FC3   FC4

$W^{[3]}$  $(120, 400)$
$b^{[3]}$  $(120)$

120

84

→O Softmax
(10 outputs)

Say digit classifier  $0, 1, 3, \ldots 9$

$n_H, n_W \downarrow$

$n_C \uparrow$

CONV — POOL — CONV — POOL — FC — FC — FC — SOFTMAX

Andrew Ng

## Neural network example

|            | Activation shape | Activation Size | # parameters |
|------------|------------------|-----------------|--------------|
| Input:     | (32,32,3)        | — 3,072   $a^{[0]}$ | 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             | 48,001       |
| FC4        | (84,1)           | 84              | 10,081       |
| Softmax    | (10,1)           | 10              | 841          |

Andrew Ng

## Neural network example

Here are the 5 typos:

1. 208 should be (5*5*3 + 1) * 8 = 608
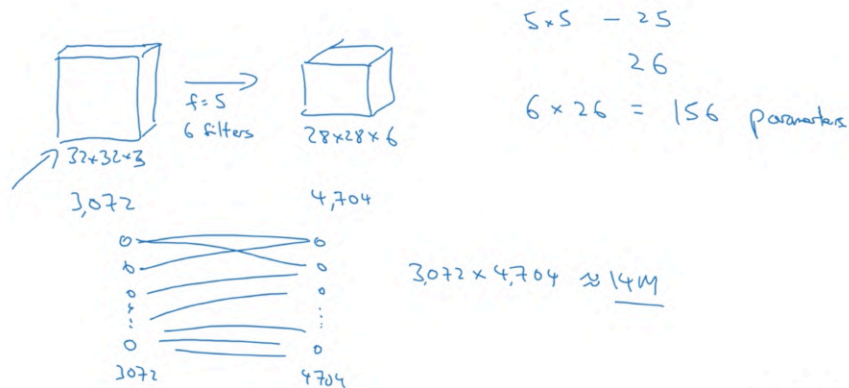
2. 416 should be (5*5*8 + 1) * 16 = 3216

3. In the FC3, 48001 should be 400*120 + 120 = 48120, since the bias should have 120 parameters, not 1

4. Similarly, in the FC4, 10081 should be 120*84 + 84 (not 1) = 10164

(Here, the bias is for the fully connected layer. In fully connected layers, there will be one bias for each neuron, so the bias become In FC3 there were 120 neurons so 120 biases.)
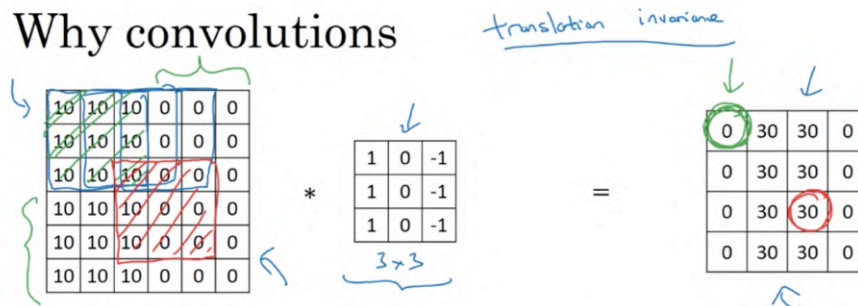
5. Finally, in the softmax, 841 should be 84*10 + 10 = 850

# Why convolutions



$5 \times 5 \quad - \quad 25$

$26$

$6 \times 26 = 156 \text{ parameters}$

$f = 5$
6 filters

$32 \times 32 \times 3$ → $28 \times 28 \times 6$

$3,072$ → $4,704$

$3,072 \times 4,704 \approx 14M$

Andrew Ng

# Why convolutions

translation invariance



$$
\begin{array}{cccccc}
10 & 10 & 10 & 0 & 0 & 0 \\
10 & 10 & 10 & 0 & 0 & 0 \\
10 & 10 & 10 & 0 & 0 & 0 \\
10 & 10 & 10 & 0 & 0 & 0 \\
10 & 10 & 10 & 0 & 0 & 0 \\
10 & 10 & 10 & 0 & 0 & 0 \\
\end{array}
\quad * \quad
\begin{array}{ccc}
1 & 0 & -1 \\
1 & 0 & -1 \\
1 & 0 & -1 \\
\end{array}
\quad = \quad
\begin{array}{cccc}
0 & 30 & 30 & 0 \\
0 & 30 & 30 & 0 \\
0 & 30 & 30 & 0 \\
0 & 30 & 30 & 0 \\
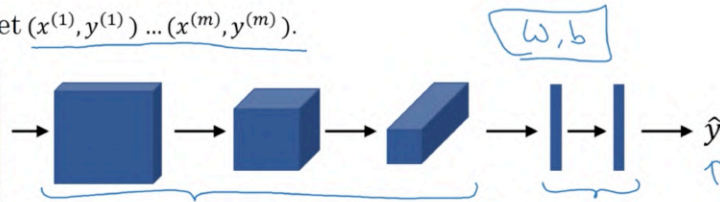\end{array}
$$

$3 \times 3$

**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.

Andrew Ng

# Putting it together

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



$W, b$

$(5*5*3+1)*6 = 456$

This is based on the equation:

$(f^{[l]} \times f^{[l]} \times n_c^{[l-1]} + 1) \times n_c^{[l]}$.

$f^{[l]}$ is the filter height (and width).

$n_c^{[l-1]}$ is the number of channels in the previous layer.

$n_c^{[l]}$ is the number of channels in the current layer.

The "1" is the bias term.

(It was (5*5+1)*6=156 in the video.)

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

Use gradient descent to optimize parameters to reduce $J$

Andrew Ng