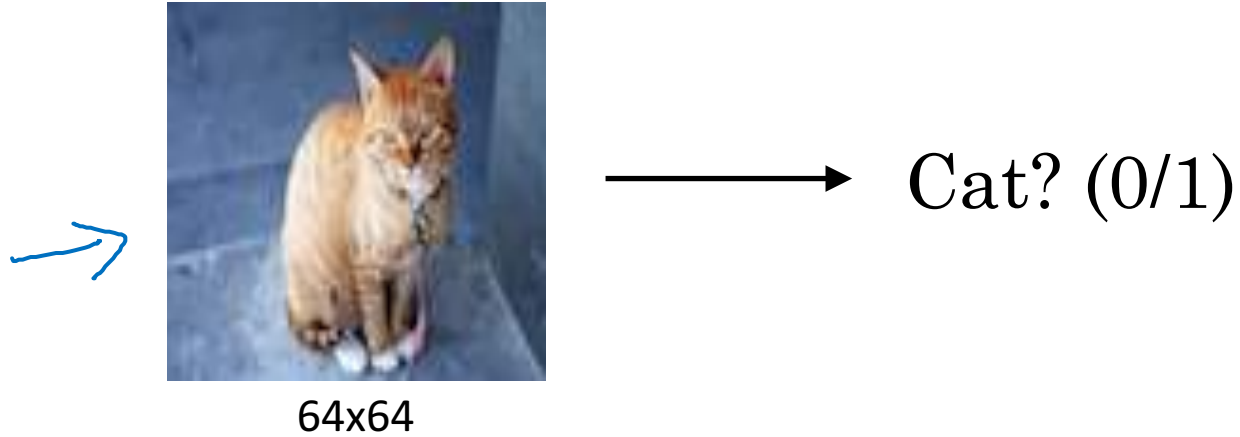# Convolutional Neural Networks

## Computer vision

# Computer Vision Problems

## Image Classification

Cat? (0/1)

64x64
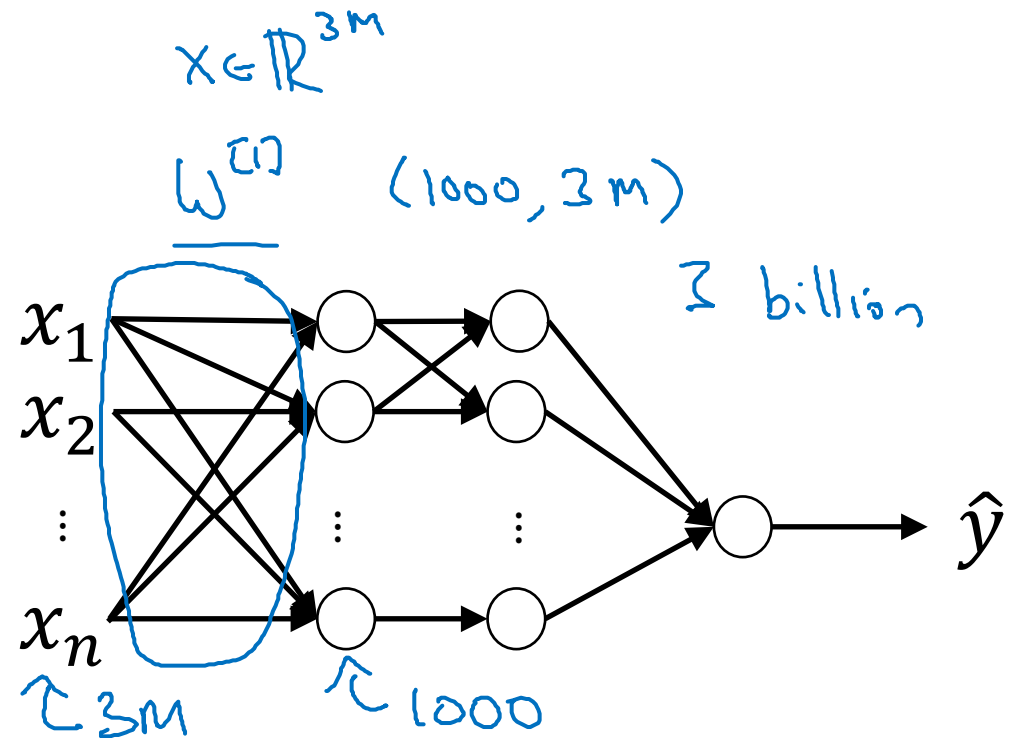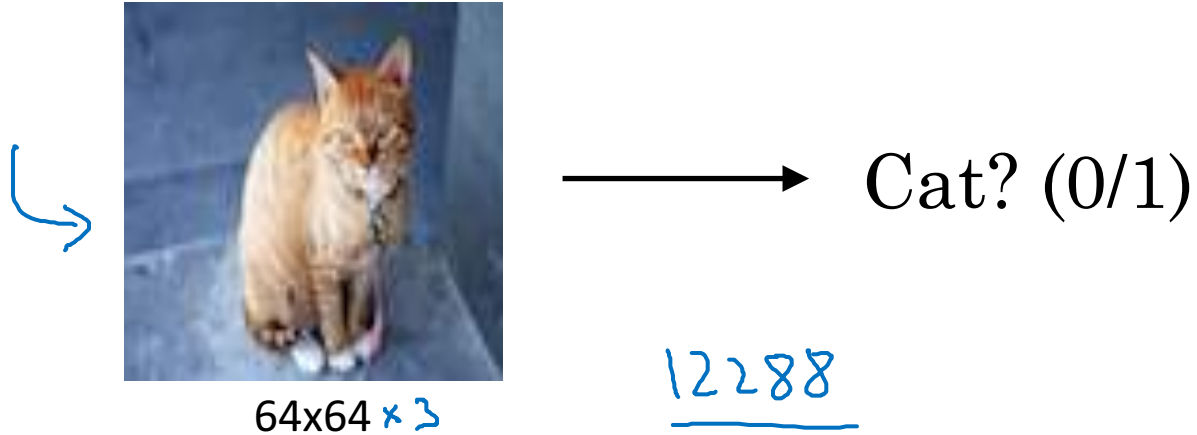
## Object detection

## Neural Style Transfer

# Deep Learning on large images



Cat? (0/1)

64x64 $\times 3$

12288

$1000 \times 1000 \times 3$
$= 3 \text{ million}$

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

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

3 billion

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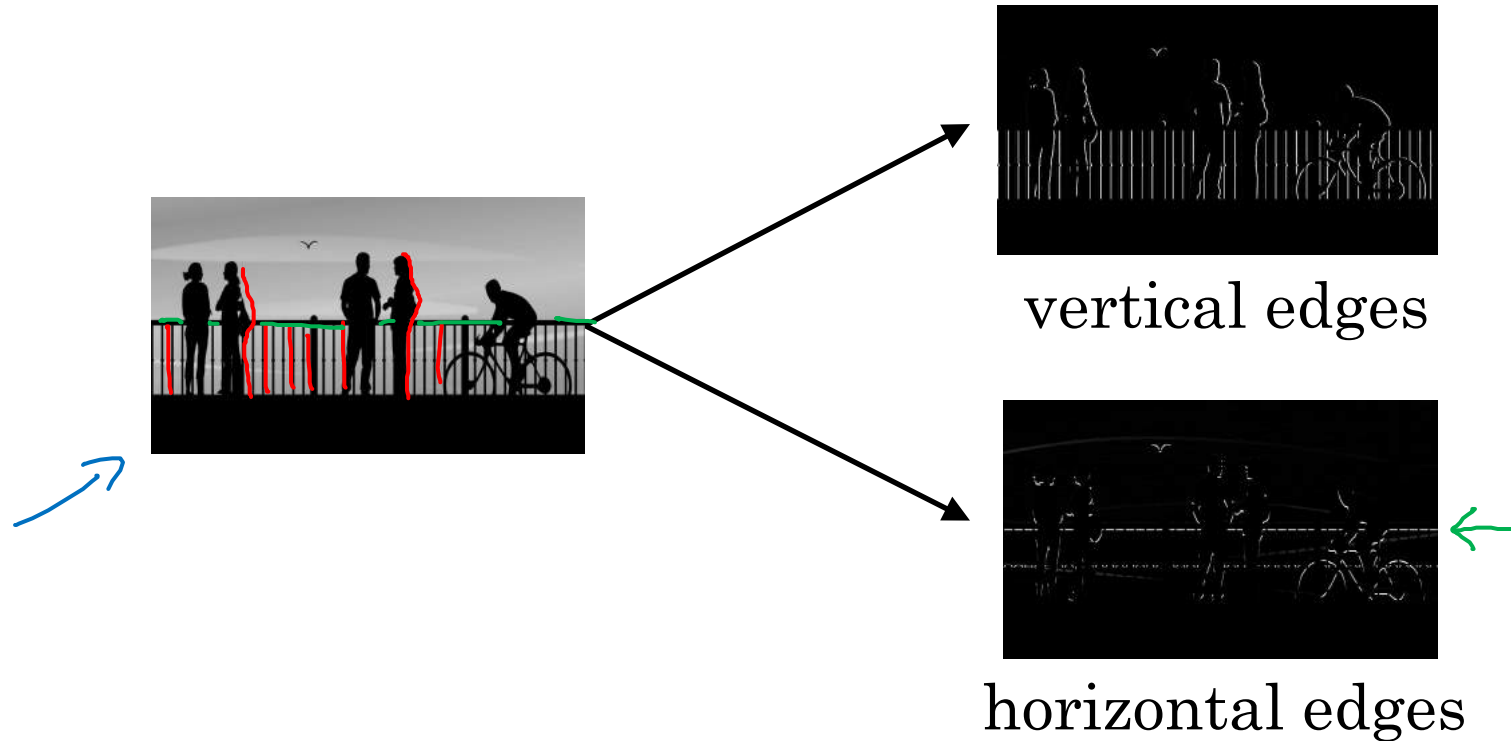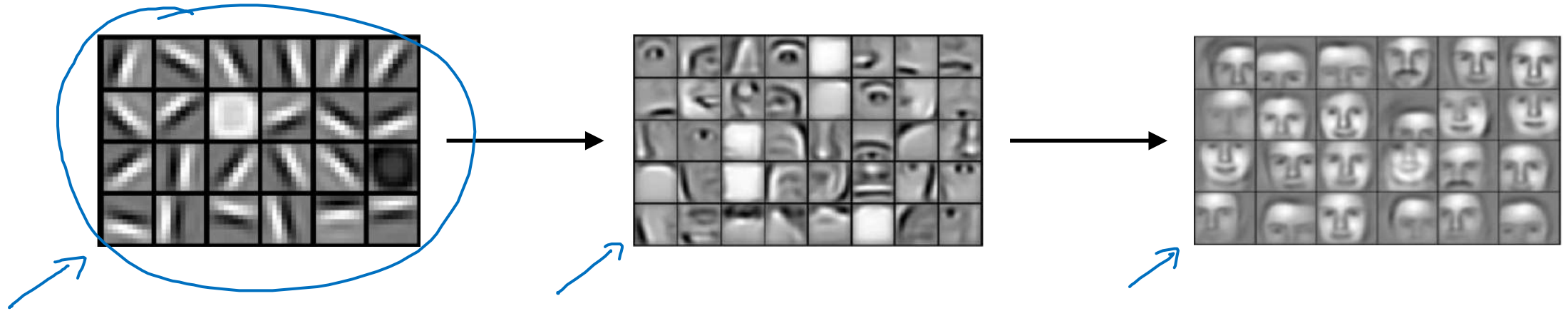$x_n$

$\hat{y}$

$\curvearrowright 3M$    $\curvearrowright 1000$

Andrew Ng

deeplearning.ai

Convolutional Neural Networks

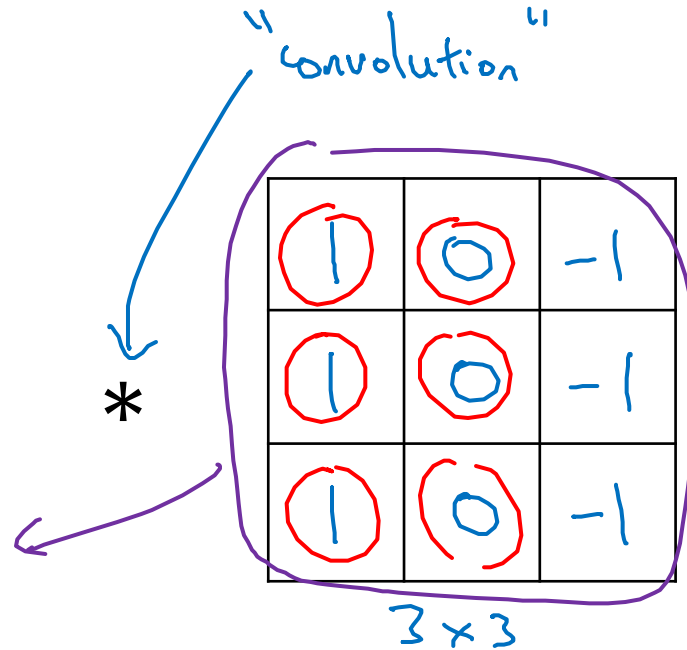Edge detection example

# Computer Vision Problem



vertical edges

horizontal edges

Andrew Ng

# Vertical edge detection

$3\times1 + 1\times1 + 2\times1 + 0\times0 + 5\times0 + 7\times0 + 1\times7 + 8\times-1 + 2\times-1 = -5$



"convolution"

$*$

3×3

→ filter
kernel

$=$

6×6

4×4

# 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 |

6 x 6

\*

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

3 x 3

=

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

4 x 4

\*

Andrew Ng

deeplearning.ai

Convolutional Neural Networks

More edge detection

# 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 |

\*

| 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 |

# 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 |

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

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

$3 \times 3$

$*$

$=$

$45°$
$70°$
$73°$

Convolutional
Neural Networks

Padding

deeplearning.ai

# Padding



$*$  3×3  $f \times f$

$=$

$p=2$

— Shrink output
— throw away info from edge

6×6

$\dfrac{6 \times 6}{n \times n} \rightarrow 8 \times 8$

$\dfrac{4 \times 4}{}$

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

$6 - 3 + 1 = 4$

$P = \text{padding} = \underline{1}$

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

$6 + 2 - 3 + 1 \times \underline{\phantom{xx}} = 6 \times 6$

Andrew Ng

# Valid and Same convolutions

$\rightarrow$ no padding

"Valid": $n \times n$ $*$ $f \times f$ $\longrightarrow$ $\underline{n - f + 1} \times n - f + 1$

$6 \times 6$ $*$ $3 \times 3$ $\longrightarrow$ $4 \times 4$

"Same": Pad so that output size is the <u>same</u> as the input size.

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

$f$ is usually odd

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

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

$3 \times 3$ $\quad p = \dfrac{3-1}{2} = 1$ $\quad\Big|\quad$ $5 \times 5 \quad p = 2$

$f = 5$

Convolutional
Neural Networks

Strided
convolutions

deeplearning.ai

# Strided convolution

| 2³ | 3⁴ | 7³ | 4⁴ | 6³ | 2⁴ | 9⁴ |
|----|----|----|----|----|----|----|
| 6¹ | 6⁰ | 9¹ | 8⁰ | 7¹ | 4⁰ | 3² |
| 3³ | 4⁴ | 8³ | 3⁴ | 8³ | 9⁴ | 7⁴ |
| 7¹ | 8⁰ | 3¹ | 6⁰ | 6¹ | 3⁰ | 4² |
| 4³ | 2⁴ | 1³ | 8⁴ | 3³ | 4⁴ | 6⁴ |
| 3¹ | 2⁰ | 4¹ | 1⁰ | 9¹ | 8⁰ | 3² |
| 0⁻¹ | 1⁰ | 3⁻¹ | 9⁰ | 2⁻¹ | 1⁰ | 4³ |

7×7

*

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

3×3

Stride = 2

=

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

3×3

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

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

padding $p$    stride $s$

$s = 2$

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

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

Andrew Ng

# Summary of convolutions

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

padding $p$ $\qquad$ stride $s$

Output Size:

$$\left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor$$
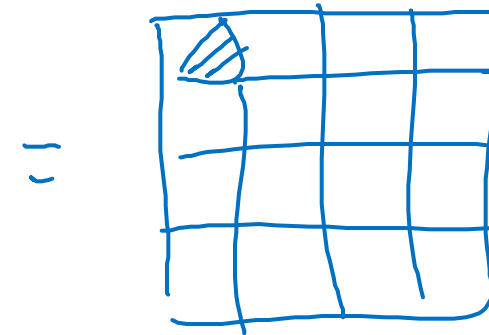
# Technical note on cross-correlation vs. convolution

## Convolution in math textbook:



* In mathematics, the convolution operation requires preliminarily flipping the filter. In ML this step is omitted and what is called "convolution" is actually a "cross-correlation" between 2 matrices.

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

← the filter's mirroring, which we don't use in ML is used to have this property, useful in signal processing but not necessarily in ML.
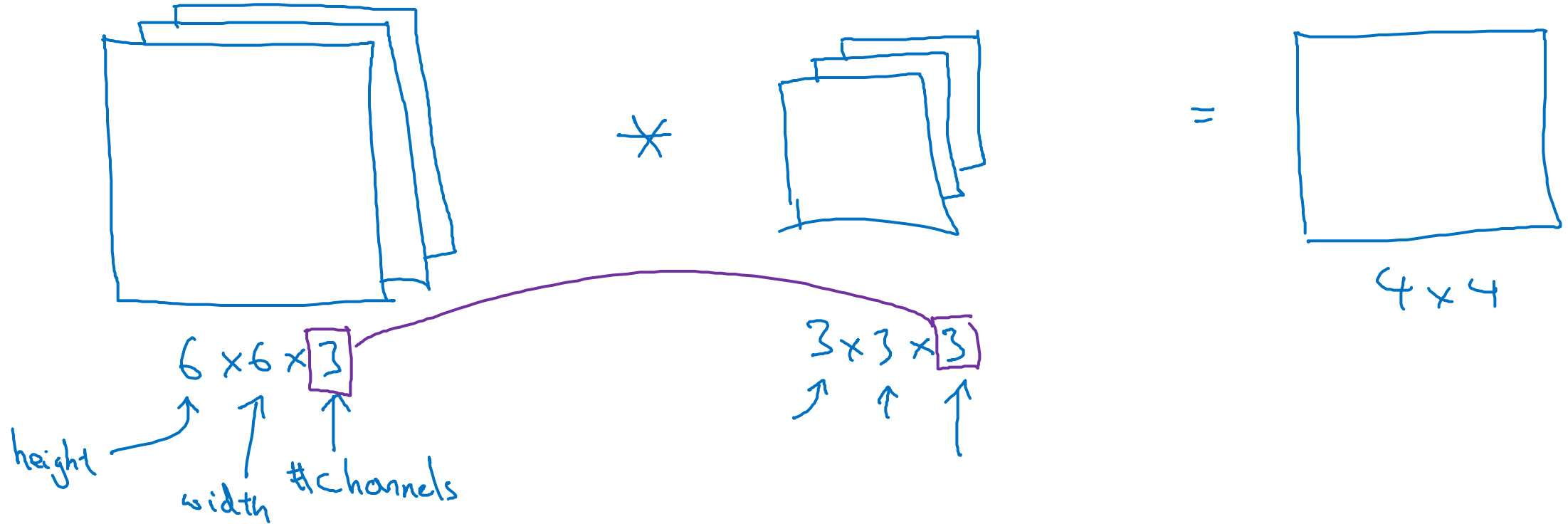
Andrew Ng

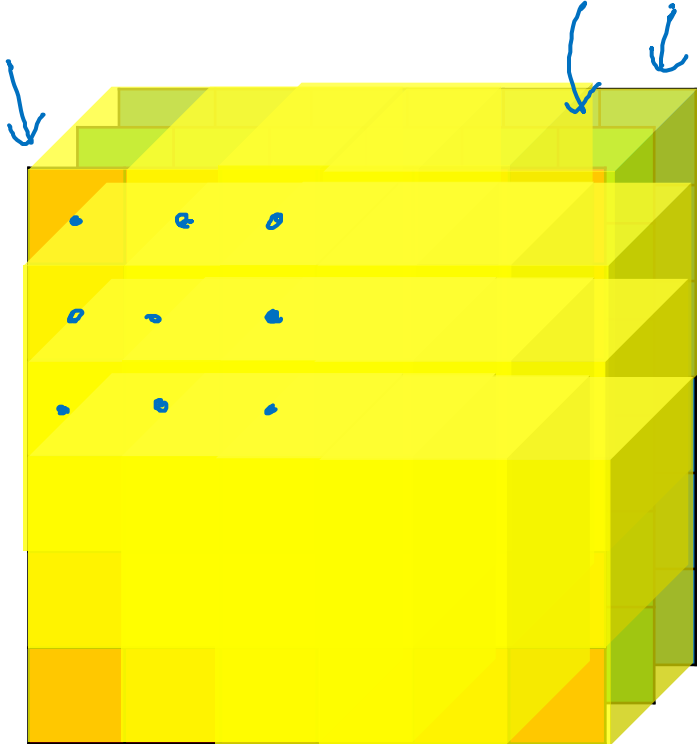deeplearning.ai

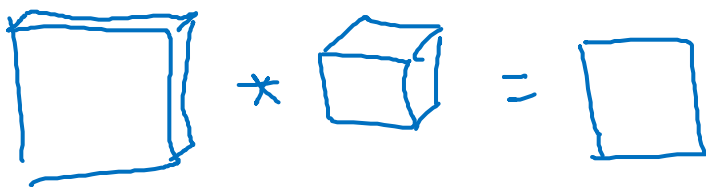Convolutional Neural Networks

Convolutions over volumes

# Convolutions on RGB images



$6 \times 6 \times \boxed{3}$

height

width

#channels

$3 \times 3 \times \boxed{3}$

$*$

$=$

$4 \times 4$
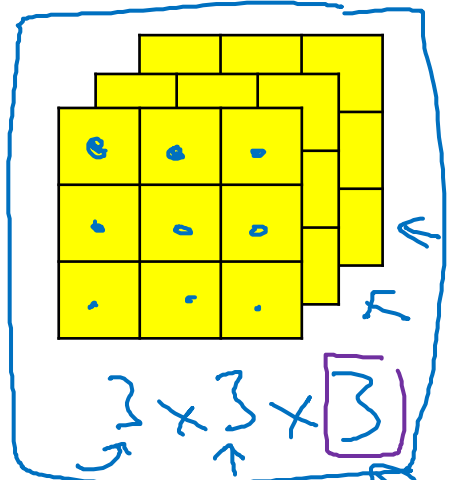
# Convolutions on RGB image



$6 \times 6 \times \boxed{3}$

$*$

$3 \times 3 \times \boxed{3}$

27 numbers

$=$

4 x 4

example 1:
"filter that detects vertical edges in the R channel"

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

G

B
$\rightarrow 3 \times 3 \times 3$

example 2:
"filter that detects vertical edges in any channel"

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

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

B
| 1 | 0 | -1 |
| 1 | 0 | -1 |
| 1 | 0 | -1 |
$\rightarrow 3 \times 3 \times 3$

Andrew Ng

# Multiple filters



Vertical edge

horizontal edge

$6 \times 6 \times 3$

channels

$3 \times 3 \times 3$

$3 \times 3 \times 3$

$4 \times 4$

$4 \times 4$

$4 \times 4 \times 2$

$4 \times 4 \times 2$

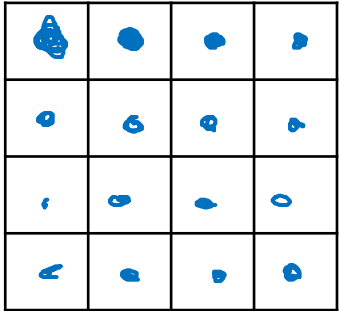Summary: $n \times n \times n_c$ $*$ $f \times f \times n_c$ $\longrightarrow$ $\dfrac{n-f+1}{} \times \dfrac{n-f+1}{} \times n_c'$

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

Andrew Ng

deeplearning.ai

Convolutional
Neural Networks

One layer of a
convolutional
network

# Example of a layer



$6 \times 6 \times 3$

$a^{[0]}$

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

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

\# filters
10

$3 \times 3 \times 3$

$3 \times 3 \times 3$

"$W^{[1]}$"

$\rightarrow$ ReLU

$\rightarrow$ ReLU

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

$+ b_1$   $4 \times 4$

$+ b_2$   $4 \times 4$

$z^{[1]}$

$4 \times 4$

$4 \times 4$

$6 \times 6 \times 3$   $a^{[0]}$ $\rightarrow$ $a^{[1]}$   $4 \times 4 \times 2$

$4 \times 4 \times 2$   $a^{[1]}$
$4 \times 4 \times 10$

Andrew Ng

# 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?



$3 \times 3 \times 3$

27 parameters.

+ bias

$\rightarrow$ 28 parameters.

280 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

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

Activations: $a^{[l]} \to n_H^{[l]} \times n_W^{[l]} \times n_c^{[l]}$

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

bias: $n_c^{[l]} - (1,1,1,n_c^{[l]})$

← #filters in layer l.

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

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

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

#examples

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

$n_c \times n_H \times n_W$

deeplearning.ai

Convolutional
Neural Networks

A simple convolution
network example

# Example ConvNet

$a^{[1]}$

$a^{[2]}$

$x$

$39 \times 39 \times 3$

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

$n_c^{[0]} = 3$

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

$7 \times 7 \times 40$

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

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

General trend of the
activations through a CNN:
* size   shrinks
* no. channels increases

$1960$

$\hat{y}$

logistic/
softmax

Andrew Ng

# Types of layer in a convolutional network:

- Convolution $(CONV)$ ←

- Pooling $(POOL)$ ←

- Fully connected $(FC)$ ←

# Convolutional Neural Networks

# Pooling layers

# Pooling layer: Max pooling

Shrinks activation maps by only retaining important info.

INTUITION: the activation of a filter indicates how much the detected feature appears in the input image.

By maxpooling the activation we discard some information while retaining where in the image the feature is more prominent.

- performed on each channel independently
- no parameters to learn!

# Pooling layer: Max pooling

| 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 |

$5 \times 5 \times n_c$

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

$3 \times 3 \times n_c$

$\dfrac{f = 3}{s = 1}$

$$\left( \frac{n + 2p - f}{s} + 1 \right)$$

# Pooling layer: Average pooling



$f = 2$

$s = 2$

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

# Summary of pooling

Hyperparameters:

$f$ : filter size

$s$ : stride

Max or average pooling

$f=2, s=2$

$f=3, s=2$

$\Rightarrow$ $p$ : padding.

No parameters to learn!

$$n_H \times n_W \times n_C$$

$$\downarrow$$

$$\left\lfloor \frac{n_H - f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n_W - f}{s} + 1 \right\rfloor$$

$$\times n_C$$

Convolutional
Neural Networks

Convolutional neural
network example

deeplearning.ai

# Neural network example (LeNet-5)



CONV1
↓
POOL1
↓
CONV2
POOL2

32×32×3 → f=5, s=1 → 28×28×6 → maxpool f=2, s=2 → 14×14×6 → f=5, s=1 → 10×10×16 → maxpool f=2, s=2 → 5×5×16 = 400

Layer 1

Layer 2

By convention, pooling layer is associated with the previous layer because it does not have learnable params

FC3 → 120
FC4 → 84
→ 0 Softmax (10 outputs)

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

0, 1, 2, .... 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)  8 filters of size 5x5x3 | (28,28,8) | 6,272 | 608 ← |
| POOL1 | (14,14,8) | 1,568 | 0 ← |
| CONV2 (f=5, s=1) | (10,10,16) | 1,600 | 3216 ← |
| POOL2 | (5,5,16) | 400 | 0 ← |
| FC3 | (120,1) | 120 | 48120 |
| FC4 | (84,1) | 84 | 10164 |
| Softmax | (10,1) | 10 | 850 |

Convolutional
Neural Networks

Why convolutions?

deeplearning.ai

# Why convolutions



$5 \times 5 \quad — \quad 25$

$26$

$6 \times 26 = 156$ parameters

$f = 5$

6 filters

$32 \times 32 \times 3$

$28 \times 28 \times 6$

3,072

4,704

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

3072

4704

Andrew Ng

# Why convolutions

translation ~~invariance~~ Invariance

| 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 |

$3 \times 3$

=

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

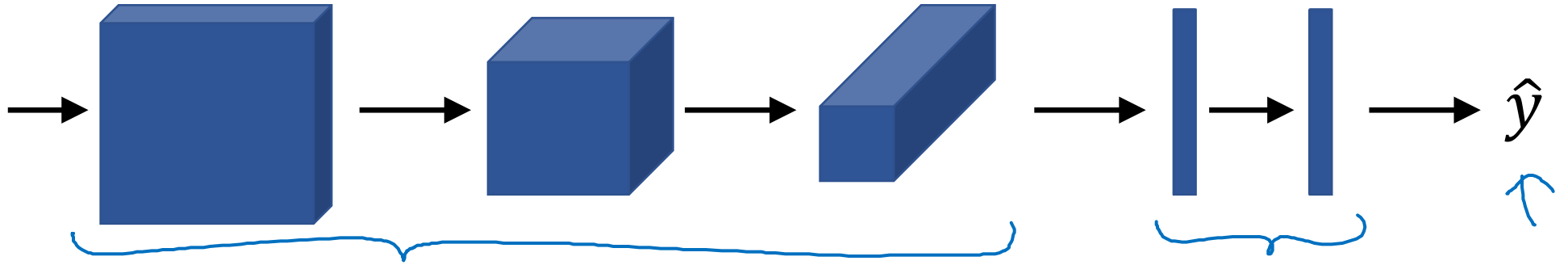**Parameter sharing:** A feature detector (such as a vertical edge detector) that's useful in one part of the image is probably useful in another part of the image. (As opposed to a FC layer connected to an image, which would probably relearn the same feature multiple times in different units). → => also favors translation invariance. & reduces overfitting.

SPATIAL LOCALITY

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

FC

# Putting it together

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



$w, b$

$\hat{y}$

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 $J$