# Deep learning

Dr. Aissa Boulmerka
a.boulmerka@centre-univ-mila.dz

2023-2024

# CHAPTER 8
# CONVOLUTIONAL NEURAL NETWORKS (CNNS)
# "FOUNDATIONS OF CNN"
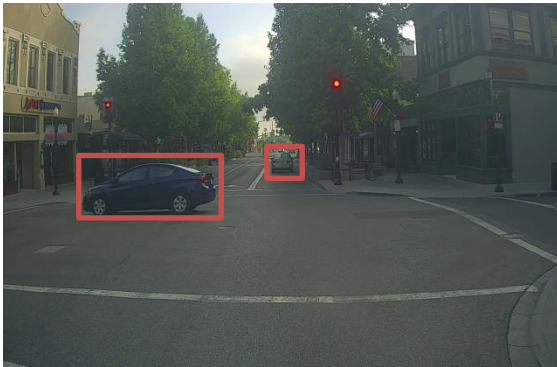
# Computer Vision Problems

## Image Classification



→ Cat? (0/1)

## Object detection



## Neural Style Transfer

# Deep Learning on large images

Image Classification



64x64x3



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

Cat? (0/1)



$x_1$
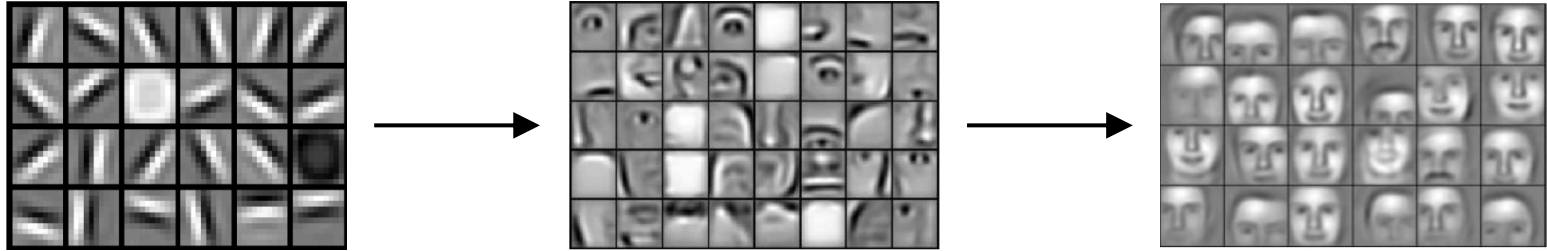
$x_2$

$\vdots$

$x_n$

$\hat{y}$

3M        1000        $\Rightarrow$    $3B$ parameters!

# Computer Vision Problem





vertical edges

horizontal edges

# Vertical edge detection



| 3[1] | 0[0] | 1[-0] | 2[-0] | 7[-0] | 4[-1] |
|---|---|---|---|---|---|
| 1[1] | 5[0] | 8[-0] | 9[-0] | 3[-0] | 1[-1] |
| 2[1] | 7[0] | 2[-0] | 5[-0] | 1[-0] | 3[-1] |
| 0[1] | 1[0] | 3[-0] | 1[-0] | 7[-0] | 8[-1] |
| 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**

6

# 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 |

$*$

| 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 |

\*

| 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 |

# 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 |

$*$

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

$=$

| | | | |
|---|---|---|---|
| | | | |
| | | | |
| | | | |

# Padding

$$6 \times 6$$
$$n \times n$$

$$\text{p} = \text{padding} = \ 1$$

$$*$$

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

$$=$$

$$6 \times 6$$

with padding:
$$(n + 2p - f + 1) \times (n + 2p - f + 1)$$
$$(6 + 2 - 3 + 1) \times (6 + 2 - 3 + 1)$$
$$= 6 \times 6$$

# Valid and Same convolutions

**"Valid" :**   No padding

$$n \times n \quad * \quad f \times f \quad \rightarrow \quad (n - f + 1) \times (n - f + 1)$$

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

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

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

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

$f$ is usually odd

$$3 \times 3 \Longrightarrow p = 1$$
$$5 \times 5 \Longrightarrow p = 2$$
$$7 \times 7 \Longrightarrow p = 3$$
$$\vdots$$

# Strided convolution

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

$7 \times 7$

$*$

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

$3 \times 3$

$=$

$$\begin{array}{|c|c|c|}
\hline
91 & 100 & 83 \\
\hline
69 & 91 & 127 \\
\hline
44 & 72 & 74 \\
\hline
\end{array}$$

$3 \times 3$

$$\begin{array}{c} n \times n \\ \text{Padding } p \end{array} \quad * \quad \begin{array}{c} \text{Stride } s \\ s = 2 \end{array} \qquad \left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor$$

$$\lfloor z \rfloor = floor(z) \qquad\qquad \frac{7+0-3}{2} + 1 = \frac{4}{2} + 1 = 3$$

# Summary of convolutions

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

padding $p$          stride $s$

Output size:

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

# 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 |

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

$=$

|   |   |   |
|---|---|---|
|   |   |   |
|   |   |   |

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

# Convolutions on RGB images

$$6 \times 6 \times 3 \qquad * \qquad 3 \times 3 \times 3 \qquad = \qquad 4 \times 4$$

**Height** $\times$ **Witdh** $\times$ **#channels**

# Convolutions on RGB image

$6 \times 6 \times 3$ * $3 \times 3 \times 3$ = $4 \times 4$

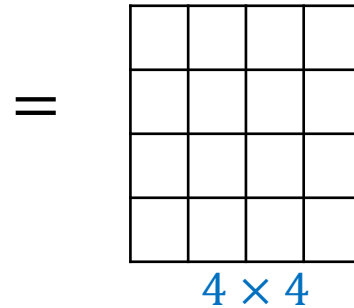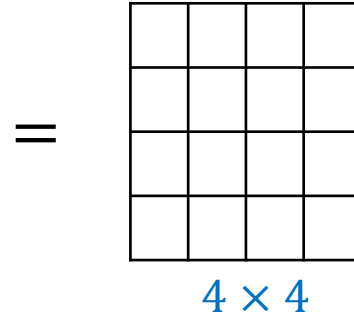# Multiple filters



Vertical edges
$3 \times 3 \times 3$

Horizontal edges
$3 \times 3 \times 3$

$6 \times 6 \times 3$

$*$

$=$

$4 \times 4$

$4 \times 4$

$4 \times 4 \times 2$

**Summary:** $n \times n \times nc \quad * \quad f \times f \times nc \quad \rightarrow \quad (n - f + 1) \times (n - f + 1) \times nc'$

$6 \times 6 \times 3 \quad * \quad 3 \times 3 \times 3 \quad \rightarrow \quad 4 \times 4 \times 2$

**# filters**

# Example of a layer

$w^{[1]}$

$w^{[1]}a^{[0]}$

$*$

$3 \times 3 \times 3$

$\longrightarrow$ RELU(     $+b_1$)

$4 \times 4$

$6 \times 6 \times 3$

$a^{[0]}$

$*$

$3 \times 3 \times 3$

$\longrightarrow$ RELU(     $+b_2$)

$4 \times 4$

$4 \times 4 \times 2$

$a^{[1]}$

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

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

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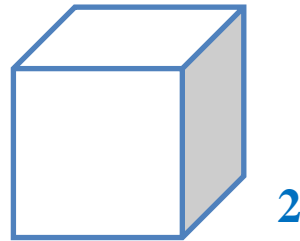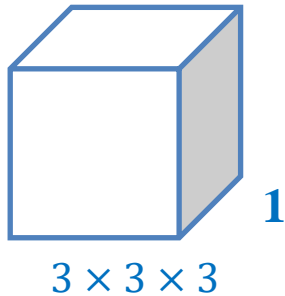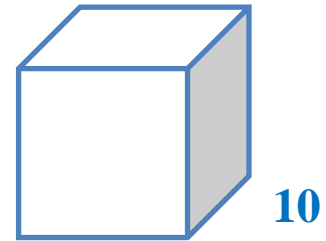
$6 \times 6 \times 3$      $4 \times 4 \times 2$

# Number of parameters in one layer

- If you have **10 filters** that are $\mathbf{3 \times 3 \times 3}$ in one layer of a neural network, how many **parameters** does that layer have?



1

$3 \times 3 \times 3$

2

…

10

27 parameters + 1 bias
=> **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** $: a^{[l]} \rightarrow n_H^{[l]} \times n_W^{[l]} \times n_c^{[l]}$

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

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

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

# Example ConvNet



$39 \times 39 \times 3$
$n_H^{[0]} = n_W^{[0]} = 39$
$n_C^{[0]} = 3$

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

$37 \times 37 \times 10$
$n_H^{[1]} = n_W^{[1]} = 37$
$n_C^{[1]} = 10$

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

$17 \times 17 \times 20$
$n_H^{[2]} = n_W^{[2]} = 17$
$n_C^{[2]} = 20$

$f^{[3]} = 5$
$s^{[3]} = 2$
$p^{[3]} = 0$
$40 \, filters$

$7 \times 7 \times 40$

Flatten

1960

$\hat{y}$

**Logistic Softmax**

# Types of layer in a convolutional network

- Convolution (CONV)

- Pooling (POOL)

- Fully connected (FC)

# Pooling layer: Max pooling



$4 \times 4$

$2 \times 2$

Hyperparameters:
$$f = 2$$
$$s = 2$$

# 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 |

$$5 \times 5 \times 2$$

$\longrightarrow$

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

$$3 \times 3 \times 2$$

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

Hyperparameters:
$$f = 3$$
$$s = 1$$

25

# Pooling layer: Average pooling



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

# Summary of pooling

Hyperparameters:

   f : filter size

   s : stride

Max or average pooling

$$n_H \times n_W \times n_C$$

$$\downarrow$$

$$\left\lfloor \frac{n_H - f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n_W - f}{s} + 1 \right\rfloor \times nc$$
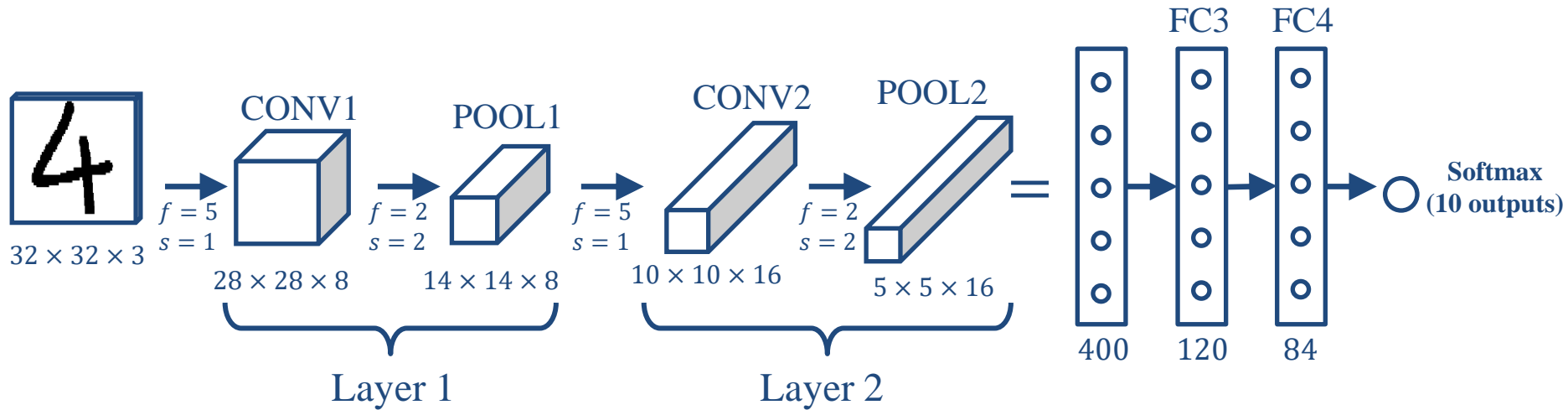
No parameters to learn!

# Neural network example

## (LeNet-5)



$$g_k(\tilde{x}) = \frac{\exp(-\widetilde{\mathbf{w}}_k^T \tilde{x})}{\sum_j \exp(-\widetilde{\mathbf{w}}_j^T \tilde{x})}$$

$$g_k(\tilde{x}) \in [0,1]$$
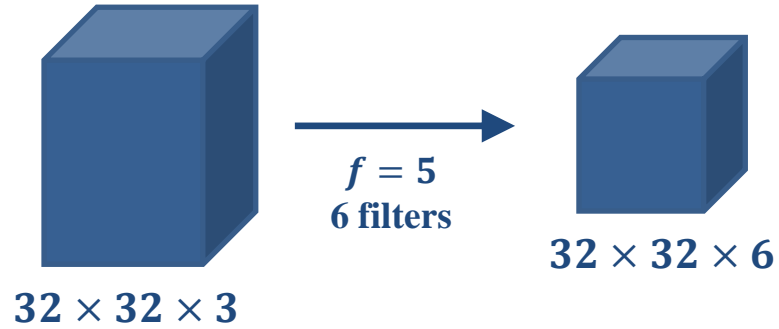
CONV-POLL-CONV-POOL-FC-FC-SOFTMAX

# Neural network example

| | Activation shape | Activation Size | # parameters |
|---|---|---|---|
| Input: | (32,32,3) | 3,072 | 0 |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |

# Neural network example

| | Activation shape | Activation Size | # parameters |
|---|---|---|---|
| Input: | (32,32,3) | 3,072 | 0 |
| CONV1 (f=5, s=1) | (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 | 48,120 |
| FC4 | (84,1) | 84 | 10,164 |
| Softmax | (10,1) | 10 | 850 |

# Why convolutions



$$5 \times 5 = 25 + 1$$
$$= 26 \; parameters \; per \; filter$$
$$6 \times 26 = 156 \; parameters$$

$32 \times 32 \times 3$

$f = 5$
6 filters

$32 \times 32 \times 6$

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

3072          4704

# 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 |

=

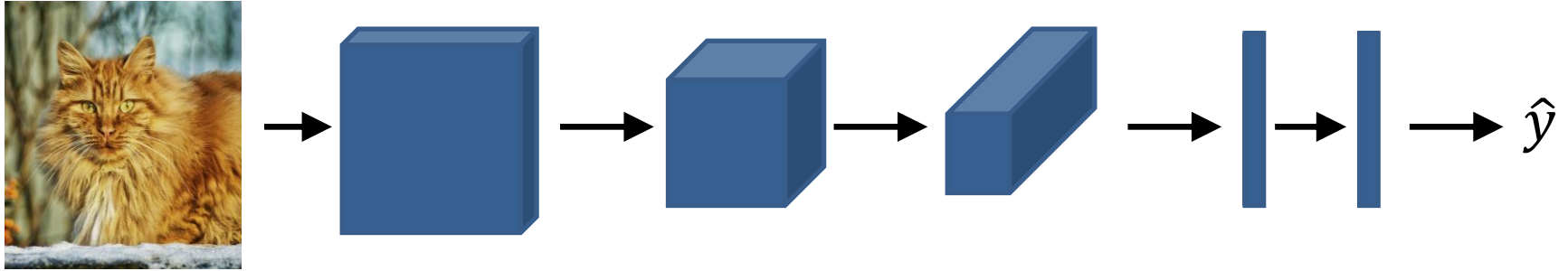| 0 | 30 | 30 | 0 |
|---|----|----|---|
| 0 | 30 | 30 | 0 |
| 0 | 30 | 30 | 0 |
| 0 | 30 | 30 | 0 |

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



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 $J$

# References

- Andrew Ng. Deep learning. Coursera.

- Geoffrey Hinton. Neural Networks for Machine Learning.

- Kevin P. Murphy. Probabilistic Machine Learning An Introduction. MIT Press, 2022.

- MIT Deep Learning 6.S191 (http://introtodeeplearning.com/)