(1) HW 3 Due Next <u>Friday</u> <span style="color:red">Start Early</span>

(2) Next Thursday : all sections $\longrightarrow$ office hours

# Convolutional Neural Network

**W**

# Multi-layer Neural Network

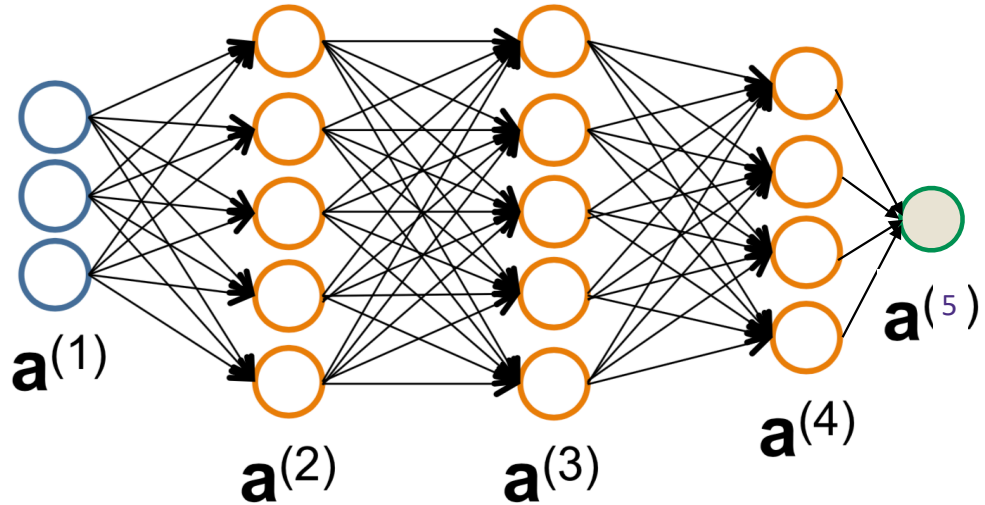$$a^{(1)} = x$$

$$z^{(2)} = \Theta^{(1)}a^{(1)}$$

$$a^{(2)} = g\left(z^{(2)}\right)$$

$$\vdots$$

$$z^{(l+1)} = \Theta^{(l)}a^{(l)}$$

$$a^{(l+1)} = g\left(z^{(l+1)}\right)$$

$$\vdots$$

$$\widehat{y} = a^{(L+1)}$$



$\mathbf{a}^{(1)}$

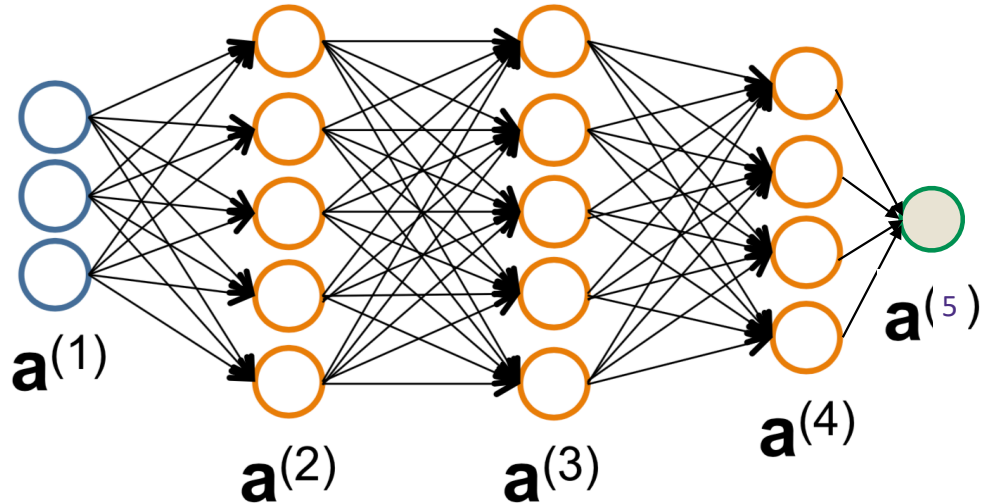$\mathbf{a}^{(2)}$    $\mathbf{a}^{(3)}$    $\mathbf{a}^{(4)}$    $\mathbf{a}^{(5)}$

$$L(y, \widehat{y}) = y\log(\widehat{y}) + (1-y)\log(1-\widehat{y})$$

$$g(z) = \frac{1}{1 + e^{-z}}$$
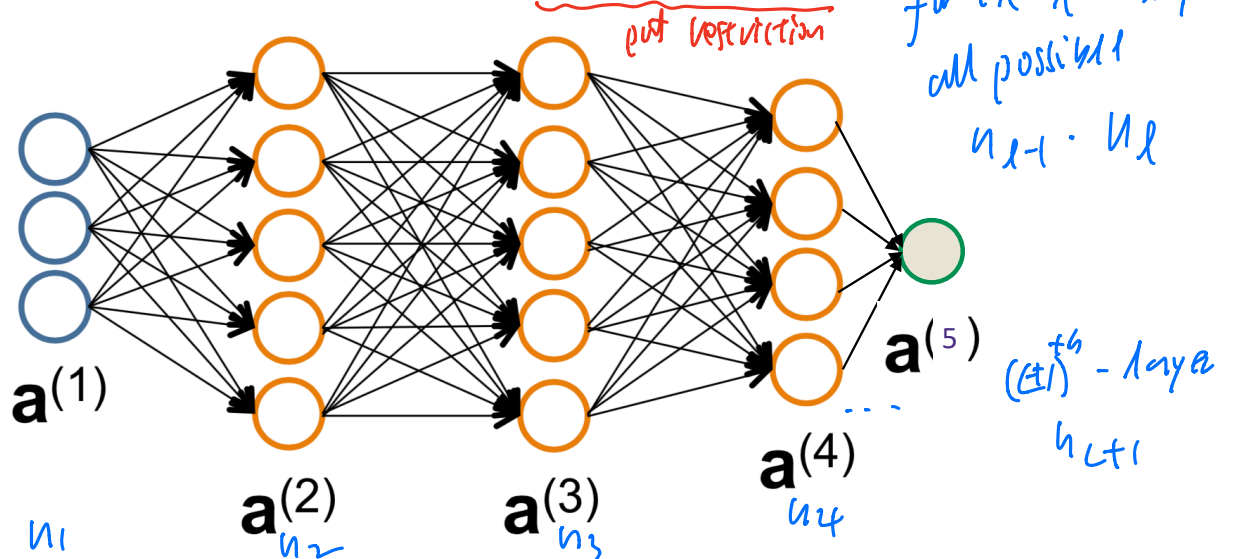
Binary Logistic Regression

# Neural Network Architecture

*depth*

The neural network architecture is defined by the number of layers, and the number of nodes in each layer, but also by **allowable edges**.

*width*

$\mathbf{a}^{(1)}$

$\mathbf{a}^{(2)}$

$\mathbf{a}^{(3)}$

$\mathbf{a}^{(4)}$

$\mathbf{a}^{(5)}$

# Neural Network Architecture

The neural network architecture is defined by the number of layers, and the number of nodes in each layer, but also by **allowable edges**.
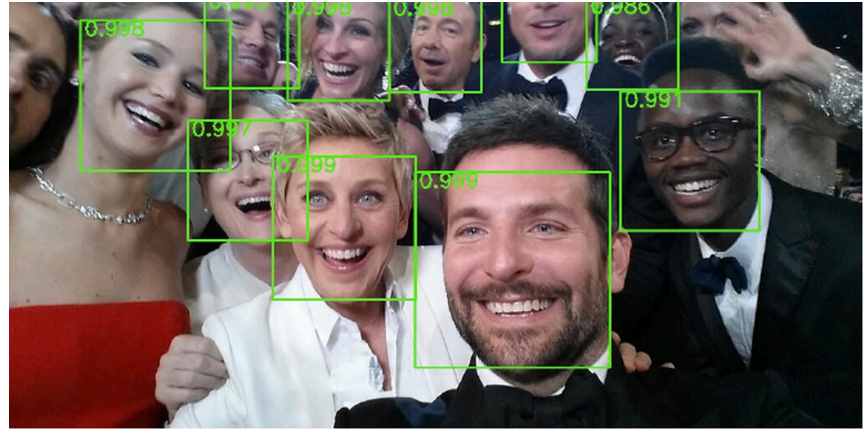
*put restriction*

*for the $l^{th}$ layer*
*all possible*
*$n_{l-1} \cdot n_l$*



$\mathbf{a}^{(1)}$
$\mathbf{a}^{(2)}$
$\mathbf{a}^{(3)}$
$\mathbf{a}^{(4)}$
$\mathbf{a}^{(5)}$

*$(L+1)^{th}$ - layer*
*$n_{L+1}$*

$n_1$
$n_2$
$n_3$
$n_4$

We say a layer is **Fully Connected (FC)** if all linear mappings from the current layer to the next layer are permissible.

$$\mathbf{a}^{(k+1)} = g(\Theta \mathbf{a}^{(k)}) \quad \text{for any } \Theta \in \mathbb{R}^{n_{k+1} \times n_k}$$
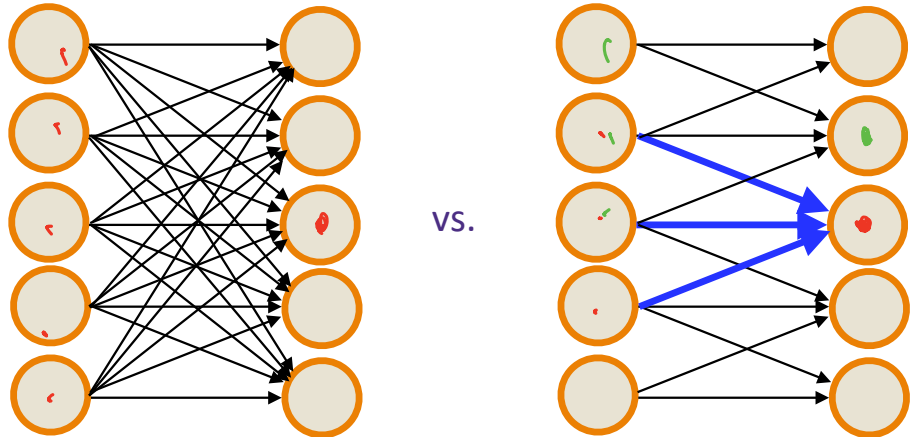
A lot of parameters!!    $n_1 n_2 + n_2 n_3 + \cdots + n_L n_{L+1}$

# Neural Network Architecture

Objects are often **localized in space** so to find the faces in an image, not every pixel is important for classification—makes sense to drag a window across an image.

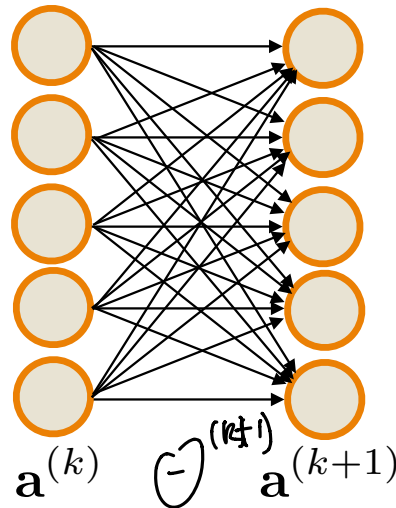

recognition

# Neural Network Architecture

Objects are often **localized in space** so to find the faces in an image, not every pixel is important for classification—makes sense to drag a window across an image.
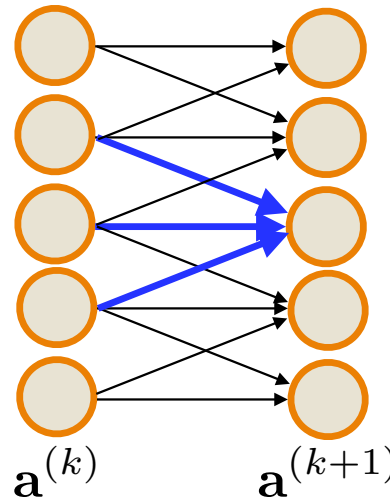


Similarly, to identify edges or other local structure, it makes sense to only look at **local information**

vs.

# Neural Network Architecture



vs.

$\mathbf{a}^{(k)}$     $\Theta^{(k+1)}$     $\mathbf{a}^{(k+1)}$                    $\mathbf{a}^{(k)}$           $\mathbf{a}^{(k+1)}$

each node has

only 3 connections

$$\begin{bmatrix} \Theta_{0,0} & \Theta_{0,1} & \Theta_{0,2} & \Theta_{0,3} & \Theta_{0,4} \\ \Theta_{1,0} & \Theta_{1,1} & \Theta_{1,2} & \Theta_{1,3} & \Theta_{1,4} \\ \Theta_{2,0} & \Theta_{2,1} & \Theta_{2,2} & \Theta_{2,3} & \Theta_{2,4} \\ \Theta_{3,0} & \Theta_{3,1} & \Theta_{3,2} & \Theta_{3,3} & \Theta_{3,4} \\ \Theta_{4,0} & \Theta_{4,1} & \T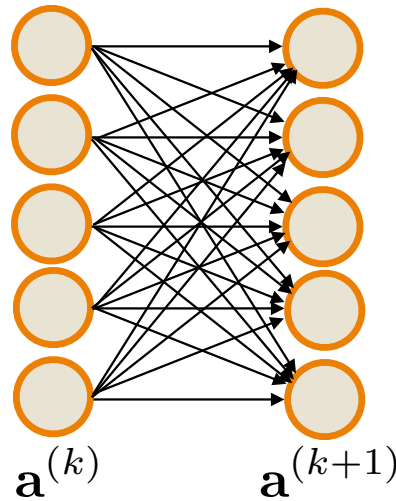heta_{4,2} & \Theta_{4,3} & \Theta_{4,4} \end{bmatrix} \qquad \begin{bmatrix} \Theta_{0,0} & \Theta_{0,1} & 0 & 0 & 0 \\ \Theta_{1,0} & \Theta_{1,1} & \Theta_{1,2} & 0 & 0 \\ 0 & \Theta_{2,1} & \Theta_{2,2} & \Theta_{2,3} & 0 \\ 0 & 0 & \Theta_{3,2} & \Theta_{3,3} & \Theta_{3,4} \\ 0 & 0 & 0 & \Theta_{4,3} & \Theta_{4,4} \end{bmatrix}$$
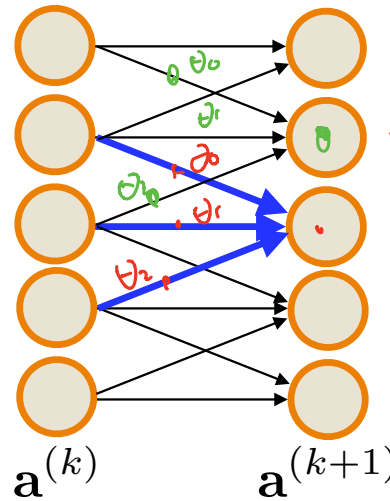
Parameters:          $n^2$                                    $3n - 2$

$$\mathbf{a}_i^{(k+1)} = g\left( \sum_{j=0}^{n-1} \Theta_{i,j} \mathbf{a}_j^{(k)} \right)$$

$$a_j^{(k+1)} = g\left( \Theta_{i,j-1}^{(k+1)} a_{j-1}^{(k)} + \Theta_{j,j}^{(k+1)} a_j^{(k)} + \Theta_{j,j+1}^{(k+1)} a_{j+1}^{(k)} \right)$$

# Neural Network Architecture



vs.

$\mathbf{a}^{(k)}$ $\mathbf{a}^{(k+1)}$ $\mathbf{a}^{(k)}$ $\mathbf{a}^{(k+1)}$

m: # of connecting

m=3

**Mirror/share local weights everywhere (e.g., structure equally likely to be anywhere in image)**

$$\begin{bmatrix} \Theta_{0,0} & \Theta_{0,1} & \Theta_{0,2} & \Theta_{0,3} & \Theta_{0,4} \\ \Theta_{1,0} & \Theta_{1,1} & \Theta_{1,2} & \Theta_{1,3} & \Theta_{1,4} \\ \Theta_{2,0} & \Theta_{2,1} & \Theta_{2,2} & \Theta_{2,3} & \Theta_{2,4} \\ \Theta_{3,0} & \Theta_{3,1} & \Theta_{3,2} & \Theta_{3,3} & \Theta_{3,4} \\ \Theta_{4,0} & \Theta_{4,1} & \Theta_{4,2} & \Theta_{4,3} & \Theta_{4,4} \end{bmatrix}$$

$$\begin{bmatrix} \Theta_{0,0} & \Theta_{0,1} & 0 & 0 & 0 \\ \Theta_{1,0} & \Theta_{1,1} & \Theta_{1,2} & 0 & 0 \\ 0 & \Theta_{2,1} & \Theta_{2,2} & \Theta_{2,3} & 0 \\ 0 & 0 & \Theta_{3,2} & \Theta_{3,3} & \Theta_{3,4} \\ 0 & 0 & 0 & \Theta_{4,3} & \Theta_{4,4} \end{bmatrix}$$

$$\begin{bmatrix} \theta_1 & \theta_2 & 0 & 0 & 0 \\ \theta_0 & \theta_1 & \theta_2 & 0 & 0 \\ 0 & \theta_0 & \theta_1 & \theta_2 & 0 \\ 0 & 0 & \theta_0 & \theta_1 & \theta_2 \\ 0 & 0 & 0 & \theta_0 & \theta_1 \end{bmatrix}$$

Parameters:   $n^2$         $3n - 2$         $3$

$$\mathbf{a}_i^{(k+1)} = g\left( \sum_{j=0}^{n-1} \Theta_{i,j} \mathbf{a}_j^{(k)} \right)$$

$$\mathbf{a}_i^{(k+1)} = g\left( \sum_{j=0}^{m-1} \theta_j \mathbf{a}_{i+j}^{(k)} \right)$$

# Neural Network Architecture

**Fully Connected (FC) Layer**

$$\begin{bmatrix} \Theta_{0,0} & \Theta_{0,1} & \Theta_{0,2} & \Theta_{0,3} & \Theta_{0,4} \\ \Theta_{1,0} & \Theta_{1,1} & \Theta_{1,2} & \Theta_{1,3} & \Theta_{1,4} \\ \Theta_{2,0} & \Theta_{2,1} & \Theta_{2,2} & \Theta_{2,3} & \Theta_{2,4} \\ \Theta_{3,0} & \Theta_{3,1} & \Theta_{3,2} & \Theta_{3,3} & \Theta_{3,4} \\ \Theta_{4,0} & \Theta_{4,1} & \Theta_{4,2} & \Theta_{4,3} & \Theta_{4,4} \end{bmatrix}$$

**Convolutional (CONV) Layer (1 filter)**

$$\begin{bmatrix} \theta_1 & \theta_2 & 0 & 0 & 0 \\ \theta_0 & \theta_1 & \theta_2 & 0 & 0 \\ 0 & \theta_0 & \theta_1 & \theta_2 & 0 \\ 0 & 0 & \theta_0 & \theta_1 & \theta_2 \\ 0 & 0 & 0 & \theta_0 & \theta_1 \end{bmatrix} \quad \text{m=3}$$

$$\mathbf{a}_i^{(k+1)} = g\left( \sum_{j=0}^{n-1} \Theta_{i,j} \mathbf{a}_j^{(k)} \right)$$

$$\mathbf{a}_i^{(k+1)} = g\left( \sum_{j=0}^{m-1} \theta_j \mathbf{a}_{i+j}^{(k)} \right) = g([\theta * \mathbf{a}^{(k)}]_i)$$

Convolution*

$$\theta = (\theta_0, \ldots, \theta_{m-1}) \in \mathbb{R}^m \text{ is referred to as a "filter"}$$

# Example (1d convolution)

$$(\theta * x)_i = \sum_{j=0}^{m-1} \theta_j x_{i+j}$$

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

Input $x \in \mathbb{R}^n$

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

Filter $\theta \in \mathbb{R}^m$

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

Output $\theta * x$
$\in \mathbb{R}^{n-2}$

# Example (1d convolution)



Input $x \in \mathbb{R}^n$

$$(\theta * x)_i = \sum_{j=0}^{m-1} \theta_j x_{i+j}$$

Filter $\theta \in \mathbb{R}^m$

Output $\theta * x$

# Example (1d convolution)

$$(\theta * x)_i = \sum_{j=0}^{m-1} \theta_j x_{i+j}$$

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

Input $x \in \mathbb{R}^n$

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

Filter $\theta \in \mathbb{R}^m$

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

/  0  |

| 2 | 1 | / |
|---|---|---|

Output $\theta * x$

# Example (1d convolution)

$$(\theta * x)_i = \sum_{j=0}^{m-1} \theta_j x_{i+j}$$

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

Input $x \in \mathbb{R}^n$

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

Filter $\theta \in \mathbb{R}^m$

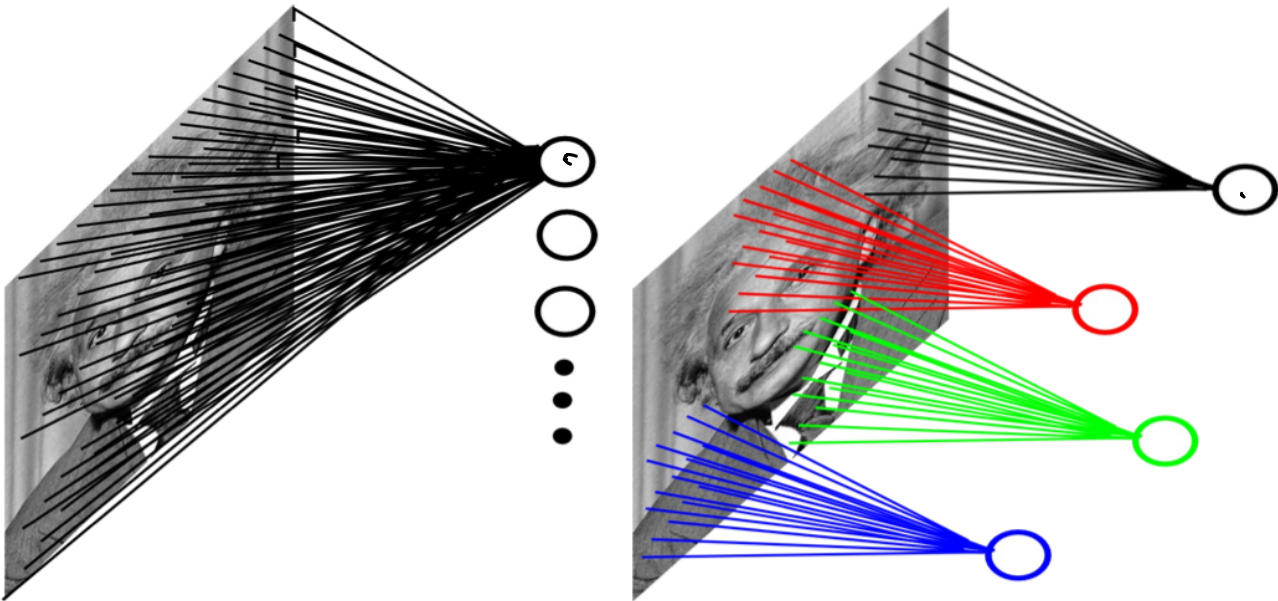| 1 | 1 | 1 | 0 | 0 |
|---|---|---|---|---|
|   |   | ×1 | ×0 | ×1 |

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

Output $\theta * x$

# 2d Convolution Layer



**Example: 200x200 image**
- Fully-connected, 400,000 hidden units = 16 billion parameters
- Locally-connected, 400,000 hidden units 10x10 fields = 40 million params
- Local connections capture local dependencies

# Convolution of images (2d convolution)

two directions

$$\left( i, j \right) = \left( 0, 0 \right)$$

$$(I * K)(i,j) = \sum_m \sum_n \underbrace{I(i+m, j+n)}\, K(m,n)$$

$$m, n = 0, 0$$
$$\rightarrow (2, 2)$$

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

Image $I$

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

Filter $K$

$3 \times 3$

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

Image

| 4 | | |
|---|---|---|
| | 4 | |
| | | |

Convolved
Feature
$I * K$

# Convolution of images

$$(I * K)(i,j) = \sum_m \sum_n I(i+m, j+n) K(m,n)$$

Image $I$



hand crafted filter

NN: learn filters

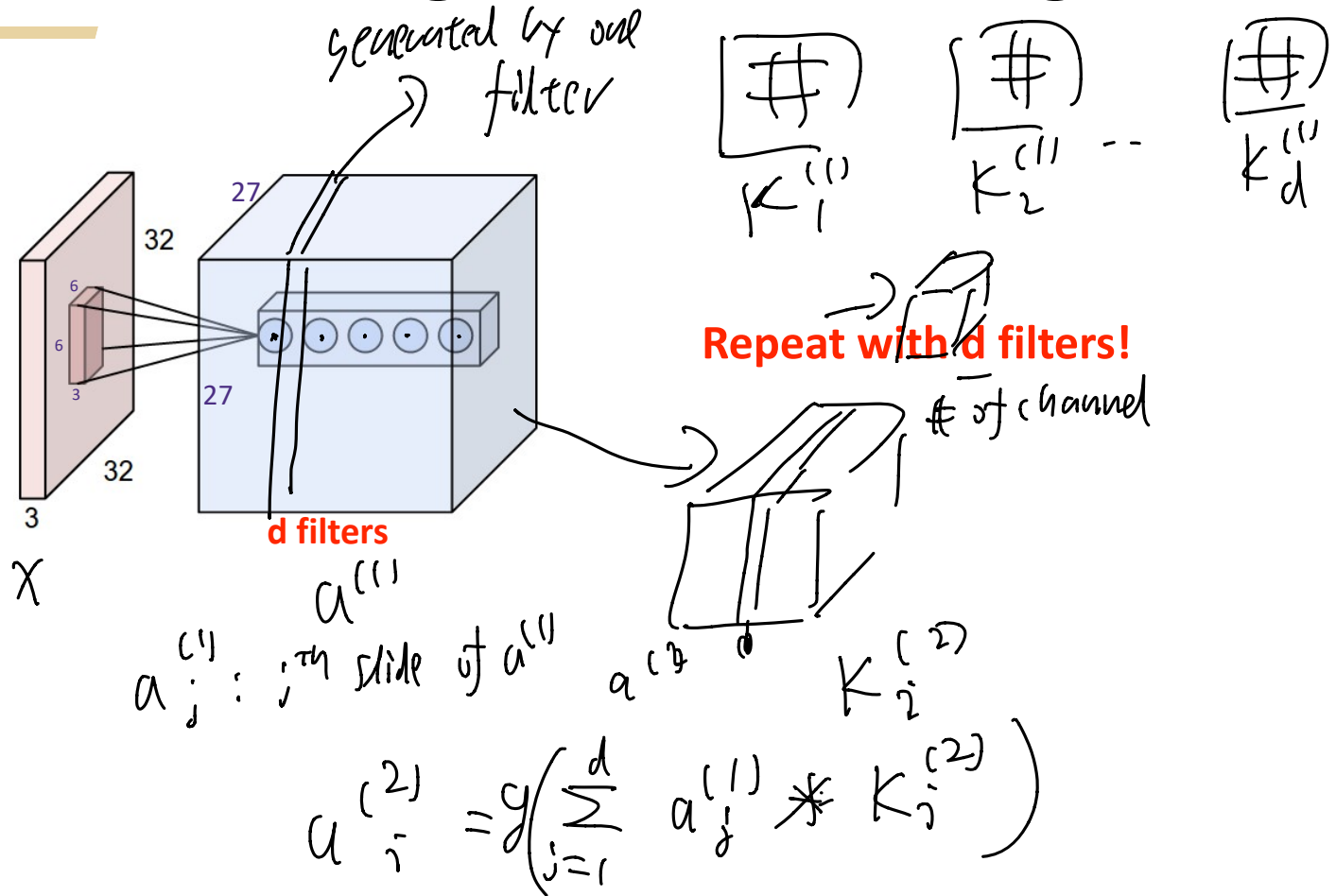| Operation | Filter $K$ | Convolved Image $I * K$ |
|---|---|---|
| Edge detection | $\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix}$ |  |
| | $\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$ |  |
| | $\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$ |  |
| Sharpen | $\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$ |  |
| Box blur (normalized) | $\frac{1}{9}\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$ |  |
| Gaussian blur (approximation) | $\frac{1}{16}\begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$ |  |

# Stacking convolved images



$$K \qquad m = 3$$

channel

$$x \in \mathbb{R}^{n \times n \times r}$$

R G B

$$a = \sum_{\alpha=1}^{r} x(:,:,\alpha) * K$$

# Stacking convolved images



generated by one filter

$K_1^{(1)}$   $K_2^{(1)}$ -- $K_d^{(1)}$

27
32
6
6
3
27
32
3

$X$

d filters

**Repeat with d filters!**

# of channel

$a^{(1)}$

$a_j^{(1)}$ : $j^{th}$ slide of $a^{(1)}$

$a^{(2)}$

$K_i^{(2)}$

$$a_i^{(2)} = g\left(\sum_{j=1}^{d} a_j^{(1)} * K_j^{(2)}\right)$$

# Pooling

Pooling reduces the dimension and can be interpreted as "This filter had a high response in this general region"

**Single depth slice**
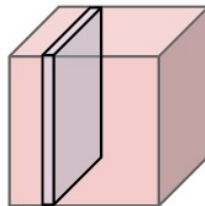
| 1 | 1 | 2 | 4 |
|---|---|---|---|
| 5 | 6 | 7 | 8 |
| 3 | 2 | 1 | 0 |
| 1 | 2 | 3 | 4 |

x

4

4   y

max pool with 2x2 filters and stride 2

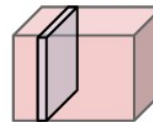| 6 | 8 |
|---|---|
| 3 | 4 |

only use most important info from each part

27x27x64

pool

14x14x64

# Pooling Convolution layer



27

32

6

6

3

32

3

27

64 filters

14x14x64

Convolve
with 64 6x6x3 filters

MaxPool with
2x2 filters and
stride 2

# Flattening



32

27

14x14x64

6

6

3

32

27

3

64 filters

MaxPool with
2x2 filters and
stride 2

Convolve
with 64 6x6x3 filters

Flatten into a single
vector of size
14*14*64=12544

# Training Convolutional Networks



CONV hidden layer

reshape

pool

14x14x64

FC hidden layer

output layer

32

27

6

6

3

27

32

3

to learn: filla

to learn

Recall: Convolutional neural networks (CNN) are just regular fully connected (FC) neural networks with some connections removed.
**Train with SGD!**

input layer

hidden layer 1   hidden layer 2

output layer

# Training Convolutional Networks



CONV hidden layer

reshape

FC hidden layer

pool

14x14x64

output layer

32

27

6

6

3

27

32

3

Real example network: LeNet

Convolution + ReLU   Pooling   Convolution + ReLU   Pooling   Fully Connected   Fully Connected   Output Predictions

dog (0.01)
cat (0.04)
boat (0.94)
bird (0.02)

Real example network: LeNet



Convolution + ReLU   Pooling   Convolution + ReLU   Pooling   Fully Connected   Fully Connected   Output Predictions

dog (0.01)
cat (0.04)
boat (0.94)
bird (0.02)

# Remarks

- Convolution is a fundamental operation in signal processing. Instead of hand-engineering the filters (e.g., Fourier, Wavelets, etc.) **Deep Learning *learns* the filters and CONV layers with back-propagation**, replacing fully connected (FC) layers with convolutional (CONV) layers
- **Pooling** is a dimensionality reduction operation that summarizes the output of convolving the input with a filter

- Typically the last few layers are **Fully Connected (FC)**, with the interpretation that the CONV layers are feature extractors, preparing input for the final FC layers. Can replace last layers and retrain on different dataset+task. *transfer learning*

- Just as hard to train as regular neural networks.
- More exotic network architectures for specific tasks

# Real networks

Modern networks have
dozens of parameters to tune.

Data augmentation?
Batch norm?

RELU leakiness
slope

Learning rate schedule

$\eta_1$
$\eta_2$
$\eta_3$
$\eta_4$

$t_1$  $t_3$  $t_2$

batchsize

$J v \ layer$

| 3x3 conv, 64 |
| 3x3 conv, 64 |

n0 layers of f0 filters

Reduce spatial
dimension

| 3x3 conv, 128, /2 |
| 3x3 conv, 128 |
| 3x3 conv, 128 |

n1 layers of f1 filters

Reduce spatial
dimension

| 3x3 conv, 128 |
| 3x3 conv, 256, /2 |
| 3x3 conv, 256 |
| 3x3 conv, 256 |

n2 layers of f2 filters

Reduce spatial
dimension

| 3x3 conv, 256 |
| 3x3 conv, 512, /2 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |

n3 layers of f3 filters

| 3x3 conv, 512 |