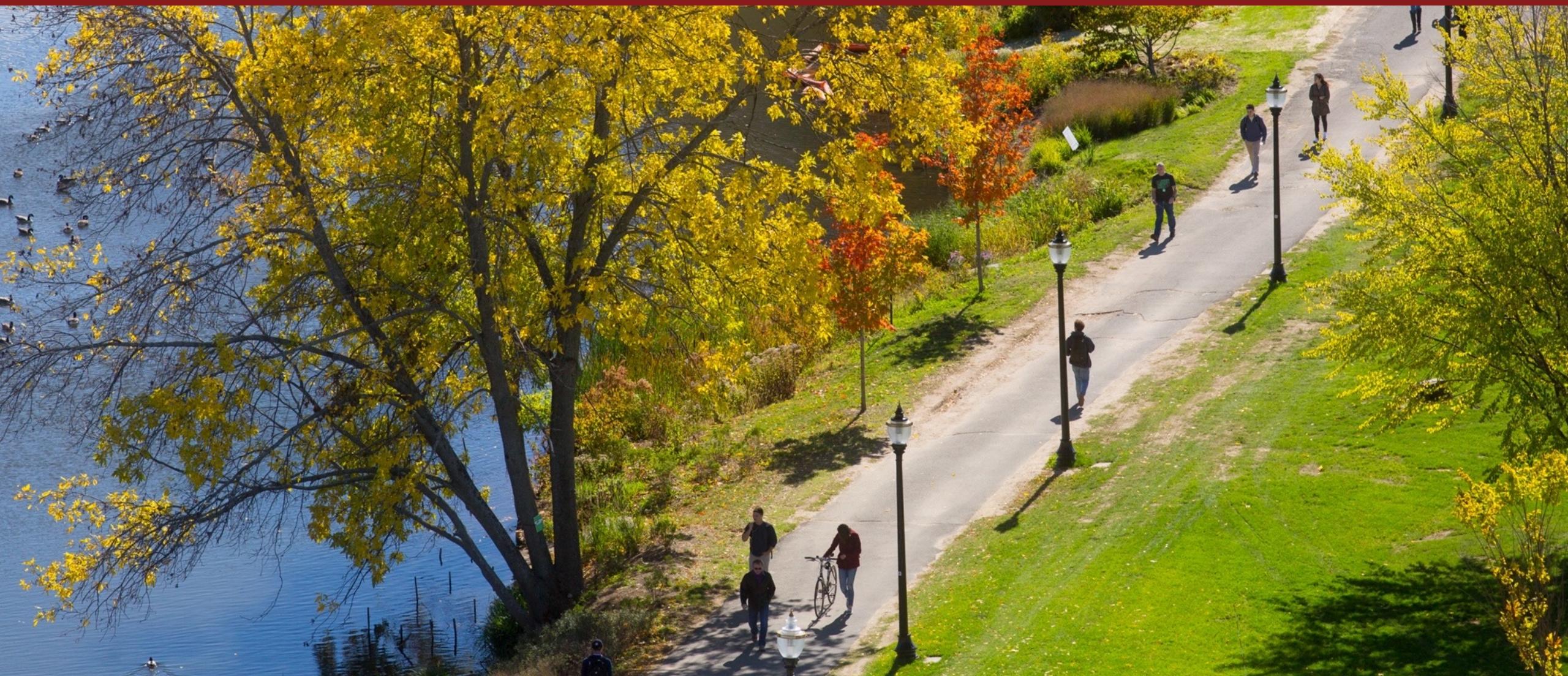


Digital Image Processing ECE 566

Ahmad Ghasemi

Department of Electrical and Computer Engineering



Segmentation Evaluation

Can be considered to consist of two components:

(1) Theoretical

Study mathematical equivalence among algorithms.

(2) Empirical

Study practical performance of algorithms in specific application domains.

Segmentation Evaluation: Theoretical

Attributes commonly used by segmentation methods:

- (1) Connectedness
- (2) Texture
- (3) Smoothness of boundary
- (4) Gradient / homogeneity
- (5) Shape information about object
- (6) Noise handling
- (7) Optimization employed
- (8) Orientedness of boundary

Segmentation Evaluation: Theoretical

	Connected	Gradient	Texture	Smooth	Shape	Noise	Optimize
Chan-Vese	No	No	Yes	Yes	No	No	Yes
Mum-Shah	No	No	Yes	Yes	No	Yes	Yes
Live wire	Boundary	Yes	Yes	Yes	User	No	Yes
Act. shape	Yes	No	No	No	Yes	No	Yes
Act. app	Yes	No	Yes	No	Yes	No	Yes
Graph cut	Usly not	Yes	Possible	No	No	No	Yes
Clustering	No	No	Yes	No	No	No	Yes

Segmentation Evaluation: Theoretical

	Connected	Gradient	Texture	Smooth	Shape	Noise	Optimize
Chan-Vese	No	No	Yes	Yes	No	No	Yes
Mum-Shah	No	No	Yes	Yes	No	Yes	Yes
Live wire	Boundary	Yes	Yes	Yes	User	No	Yes
Act. shape	Yes	No	No	No	Yes	No	Yes
Act. app	Yes	No	Yes	No	Yes	No	Yes
Graph cut	Usly not	Yes	Possible	No	No	No	Yes
Clustering	No	No	Yes	No	No	No	Yes
Deep Learning	Yes	Yes	Yes	Yes	Yes	Yes	Yes

Segmentation Evaluation: Empirical

Precision : Repeatability taking into account all subjective actions influencing the result.

Accuracy : Degree to which the result agrees with truth.

Efficiency : Practical viability of the method.

Segmentation Evaluation: Empirical

Precision (Reliability)

A measure of precision for method M in a trial that produces $C_M^{O_1}$ and $C_M^{O_2}$

$$PR_M^{T_i} = \frac{|C_M^{O_1} \cap C_M^{O_2}|}{|C_M^{O_1} \cup C_M^{O_2}|},$$

Segmentation Evaluation: Empirical

Accuracy (Validity)

The degree to which segmentations agree with true segmentation.

Surrogates of truth are needed.

For any image C acquired for application domain

C_M^O - segmentation of O in C by method M ,

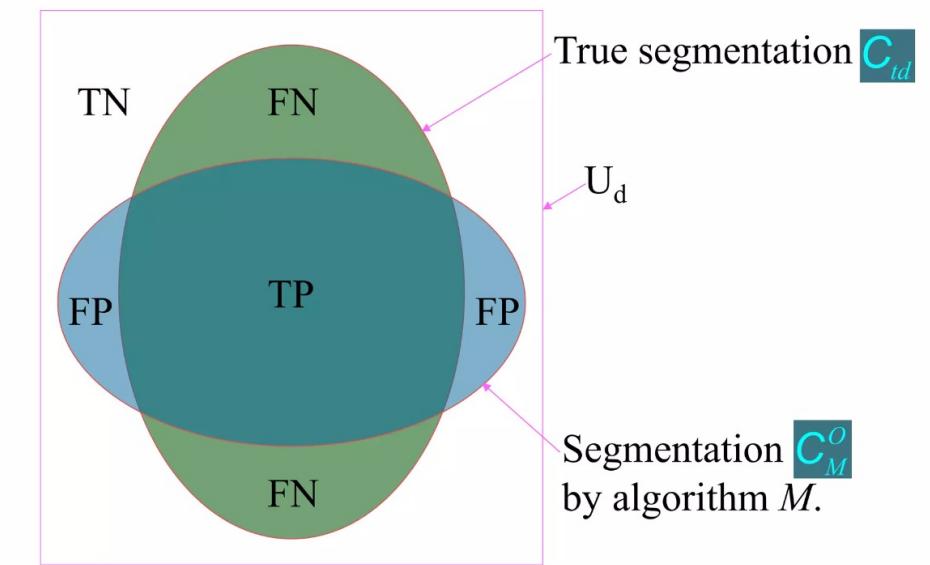
C_{td} - surrogate of true delineation of O in C .

Segmentation Evaluation: Empirical

Accuracy (Validity)

$$FNVF_M^d = \frac{|C_{td} - C_M^O|}{|C_{td}|}, \quad TPVF_M^d = \frac{|C_{td} \cap C_M^O|}{|C_{td}|}$$

$$FPVF_M^d = \frac{|C_M^O - C_{td}|}{|U_d - C_{td}|}, \quad TNVF_M^d = \frac{|U_d - C_M^O - C_{td}|}{|U_d - C_{td}|}$$



U_d : A binary scene representing a reference super set
(for example, this may be the body region that is imaged).

$FNVF_M^d$: Amount of tissue truly in O that is missed by M .

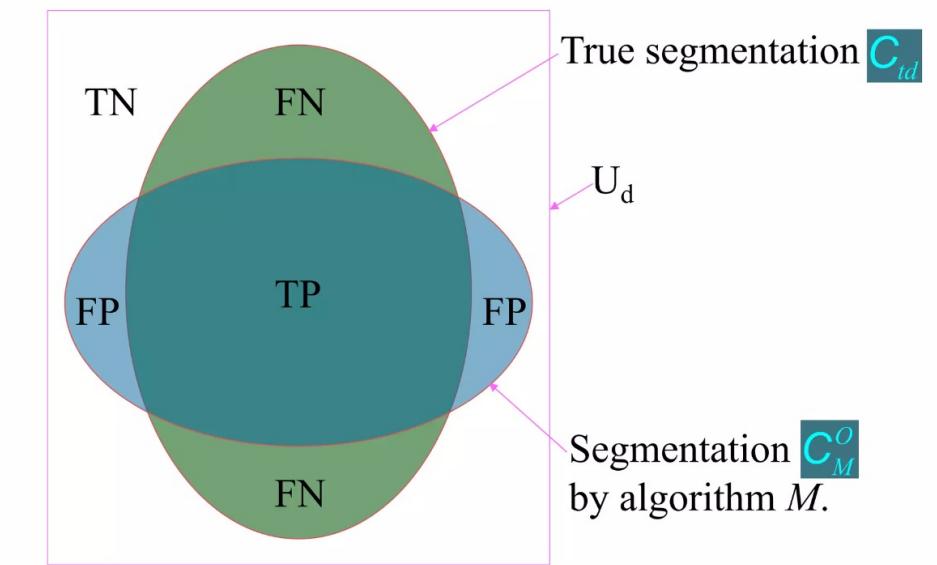
$FPVF_M^d$: Amount of tissue falsely delineated by M .

Segmentation Evaluation: Empirical

Accuracy (Validity)

$$FNVF_M^d = \frac{|C_{td} - C_M^O|}{|C_{td}|}, \quad TPVF_M^d = \frac{|C_{td} \cap C_M^O|}{|C_{td}|}$$

$$FPVF_M^d = \frac{|C_M^O - C_{td}|}{|U_d - C_{td}|}, \quad TNVF_M^d = \frac{|U_d - C_M^O - C_{td}|}{|U_d - C_{td}|}$$



U_d : A binary scene representing a reference super set
(for example, this may be the body region that is imaged).

$FNVF_M^d$: Amount of tissue truly in O that is missed by M .

$FPVF_M^d$: Amount of tissue falsely delineated by M .

$$FNVF_M^d = 1 - TPVF_M^d$$

$$FPVF_M^d = 1 - TNVF_M^d$$

Segmentation Evaluation: Empirical

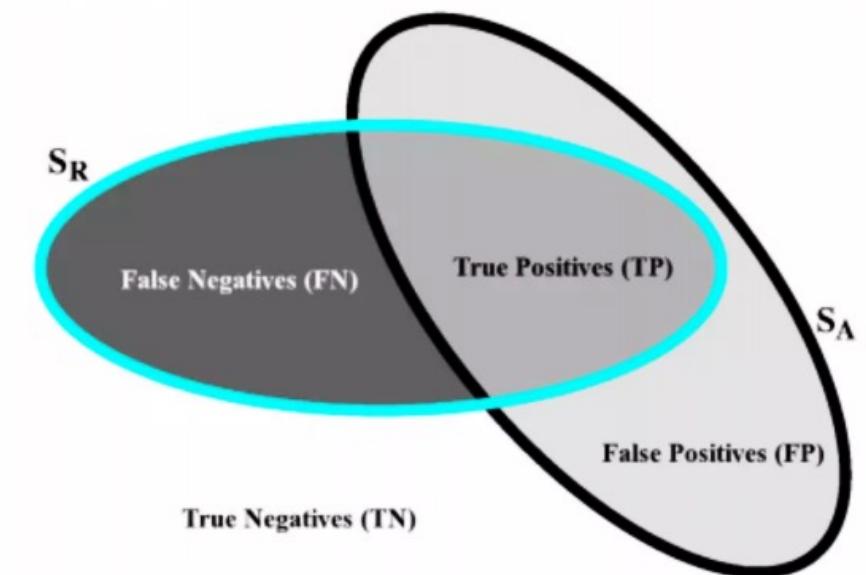
Overlap Measures

- Sensitivity: rate of detection of the structure

$$\text{Sens} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

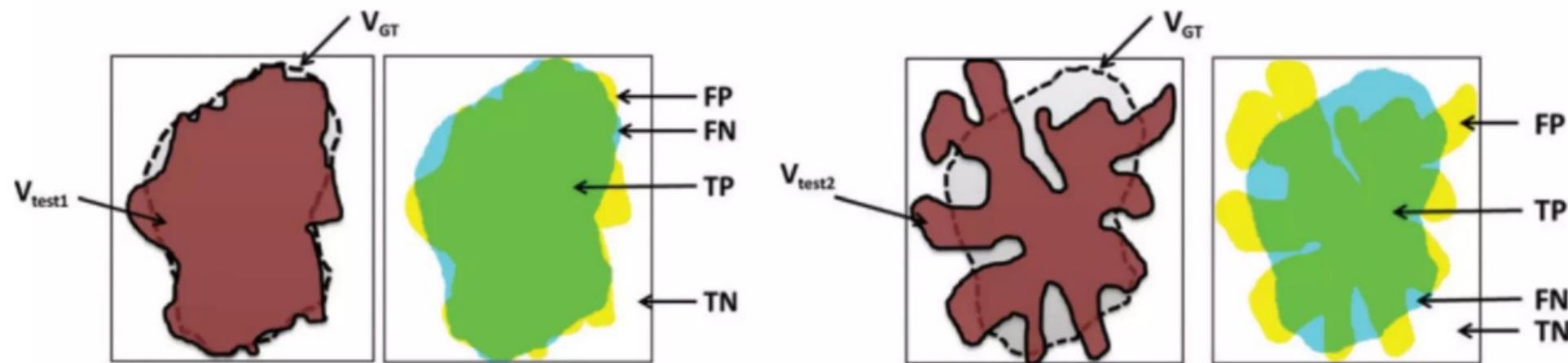
- Specificity: rate of detection of the structure

$$\text{Spec} = \frac{\text{TN}}{\text{TN} + \text{FP}}$$



Segmentation Evaluation: Empirical

Overlap Measures: example



Sensitivity=94.69%
Specificity=94.19%

Sensitivity=72.99%
Specificity=78.16%

Segmentation Evaluation: Empirical

Efficiency (Viability)

Describes practical viability of a method.

Four factors should be considered:

- (1) Computational time – for one time training of M $\left(t_M^{c_1} \right)$
- (2) Computational time – for segmenting each scene $\left(t_M^{c_2} \right)$
- (3) Human time – for one-time training of M $\left(t_M^{h_1} \right)$
- (4) Human time – for segmenting each scene $\left(t_M^{h_2} \right)$

(2) and (4) are crucial. (4) determines the degree of automation of M .

Segmentation Evaluation: Empirical

Remark

Precision, accuracy, efficiency are interdependent.

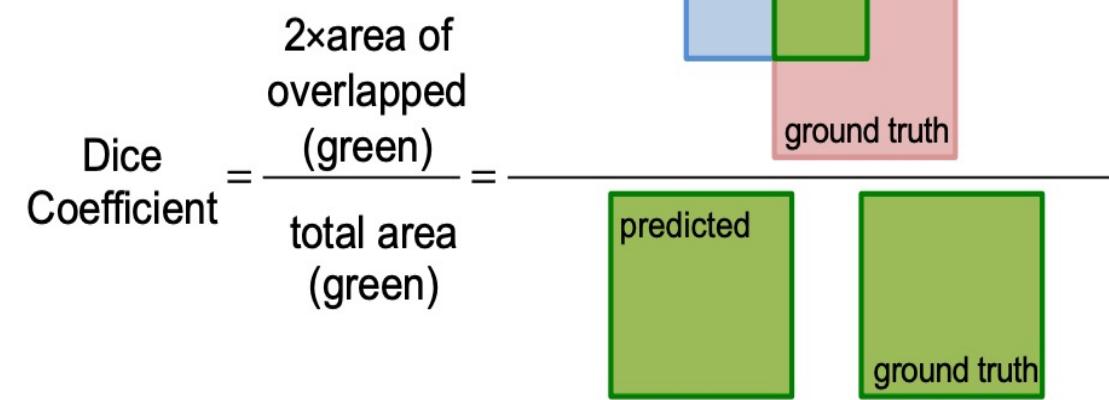
\uparrow accuracy \rightarrow \downarrow efficiency.

\uparrow precision and \uparrow accuracy \rightarrow difficult.

Segmentation Evaluation: Empirical

Dice index

$$Dice(A, B) = 2 \frac{|A \cap B|}{|A| + |B|}$$



Jaccard index(IoU)

$$Jaccard(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

Image Segmentation

**Traditional
Methods**

- Threshold-based
- Edge-based
- Region-based
- Energy-based
- Graph-based
- ...

Deep Learning based Image Segmentation

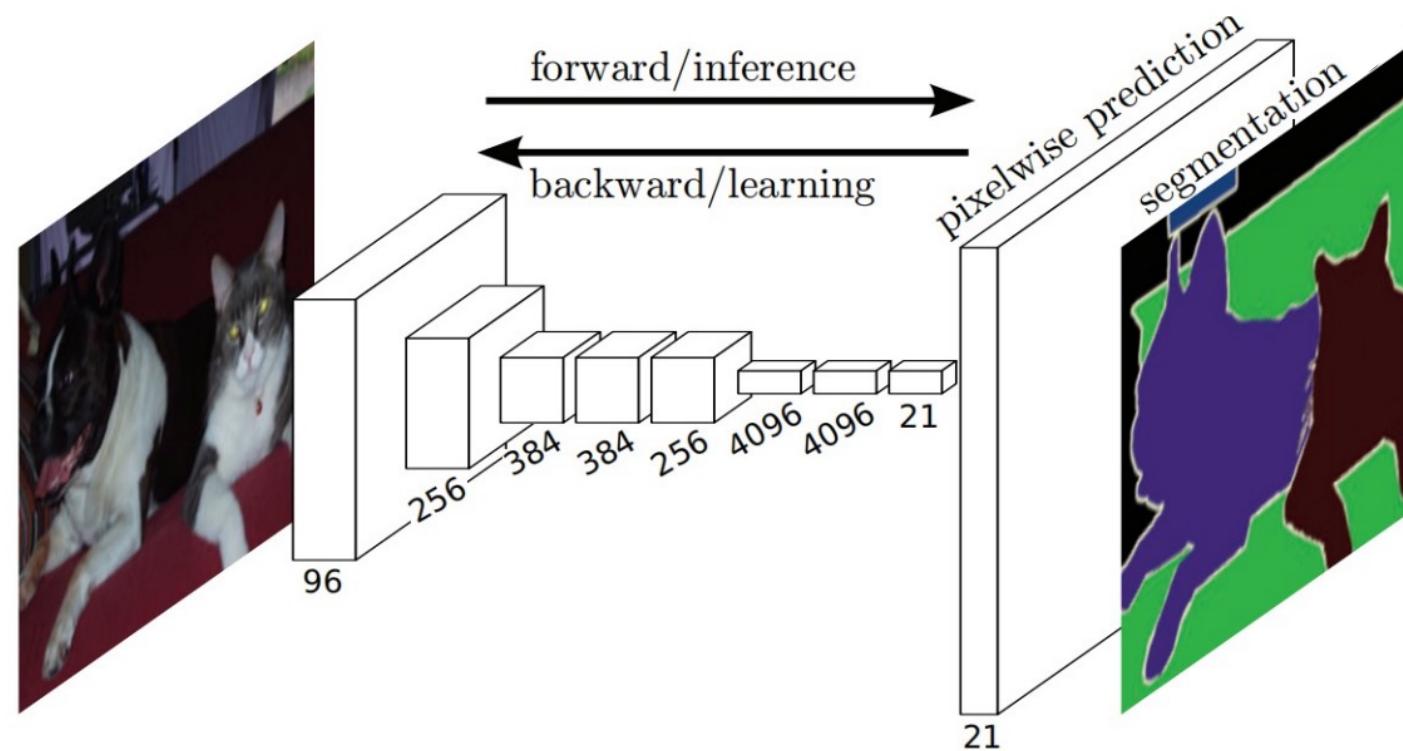
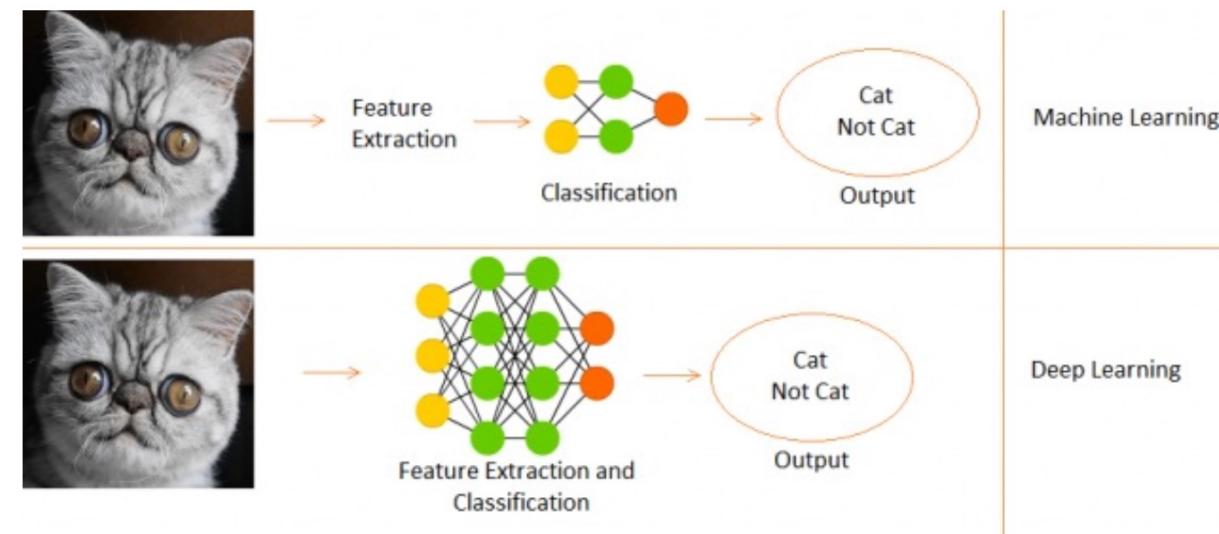


Figure 1. Fully convolutional networks can efficiently learn to make dense predictions for per-pixel tasks like semantic segmentation.

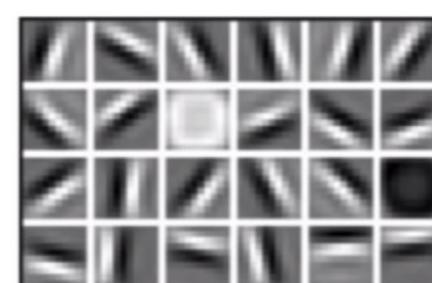
Why Deep Learning?



Hand engineered features are time consuming, brittle and not scalable in practice

Can we learn the **underlying features** directly from data?

Low Level Features



Lines & Edges

Mid Level Features



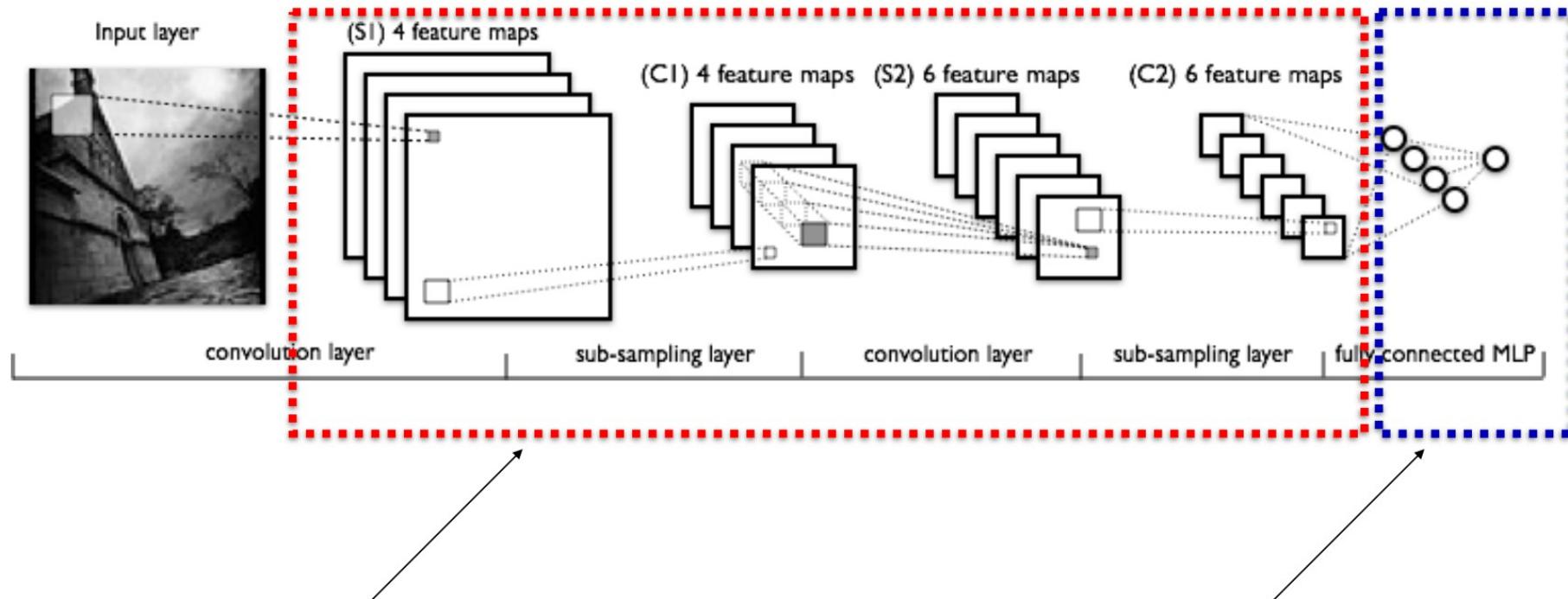
Eyes & Nose & Ears

High Level Features



Facial Structure

Deep Convolutional Neural Networks (DCNN)



**Feature extraction using
convolutional layers**

**Classification using
multilayer perceptron (MLP)**

Machine Learning ≈ Looking for a Function

- Speech Recognition

$$f\left(\text{[sound波形图]} \right) = \text{“How are you”}$$

- Image Recognition

$$f\left(\text{[猫的照片]} \right) = \text{“Cat”}$$

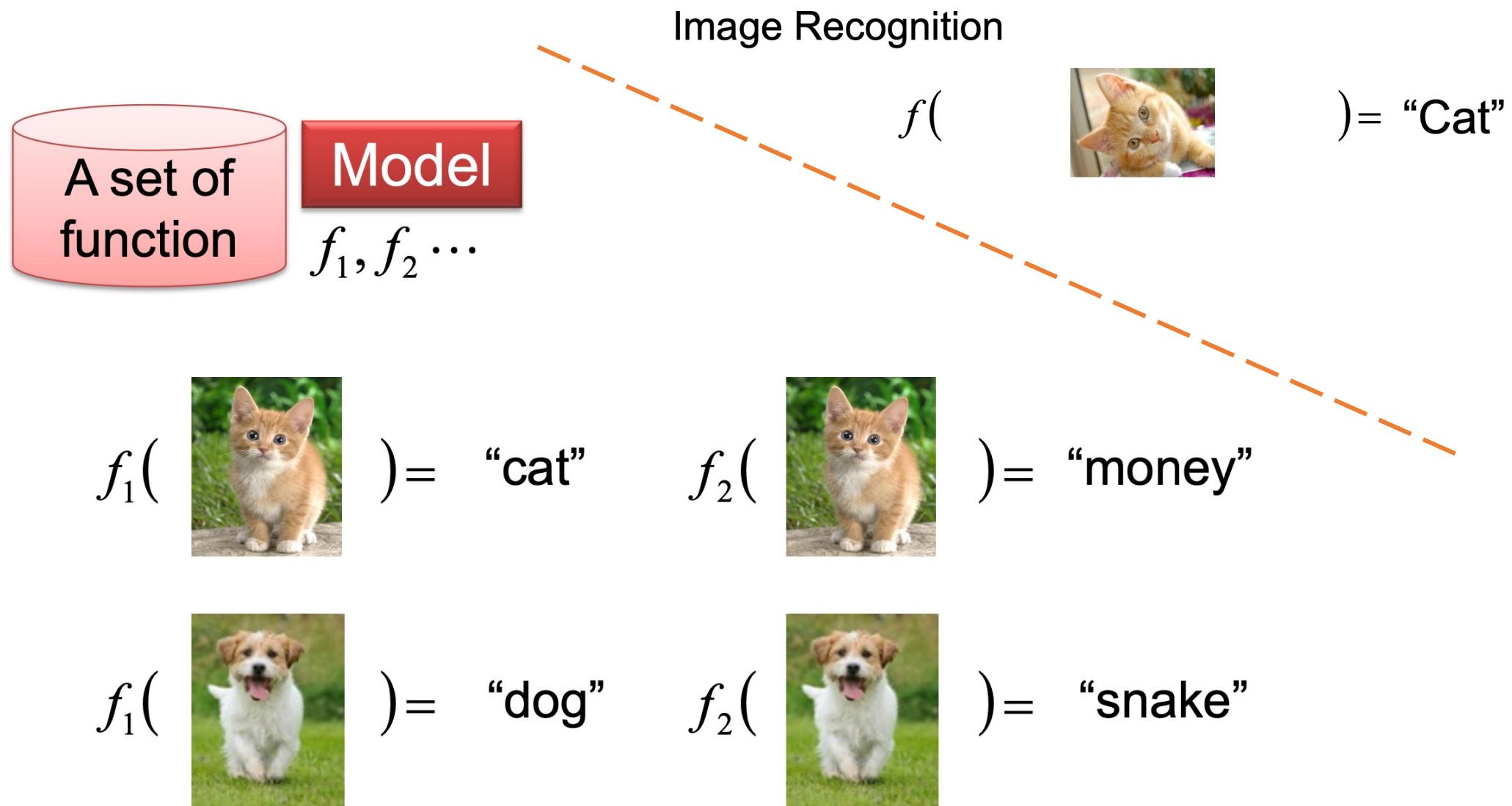
- Playing Go

$$f\left(\text{[围棋棋盘]} \right) = \text{“5-5” (next move)}$$

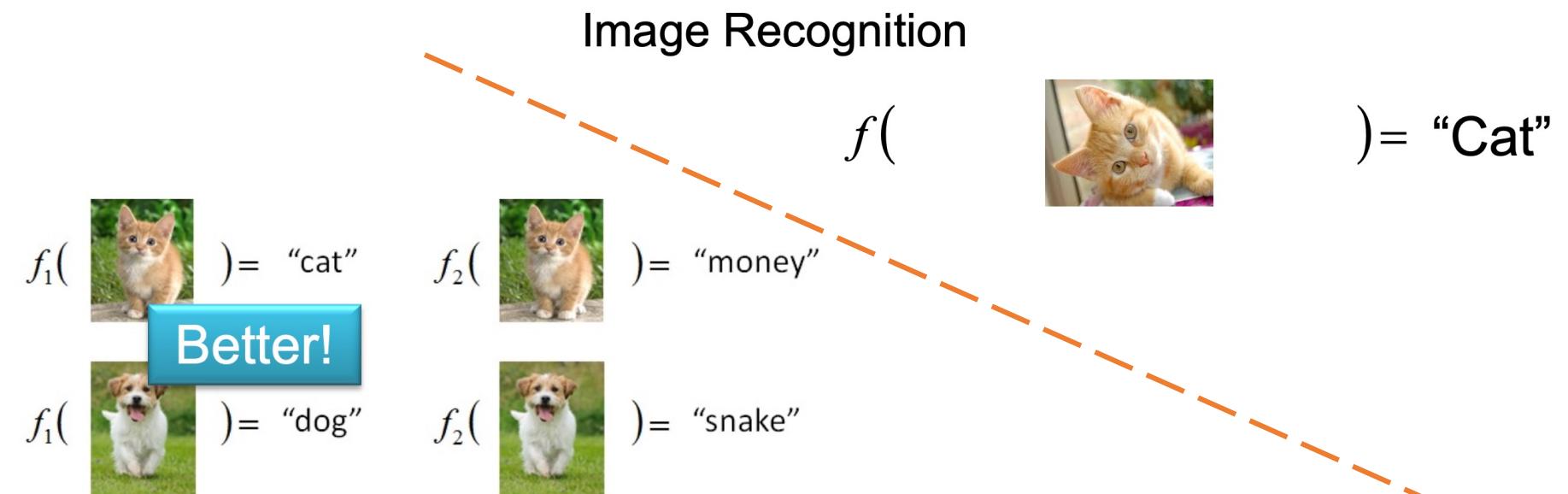
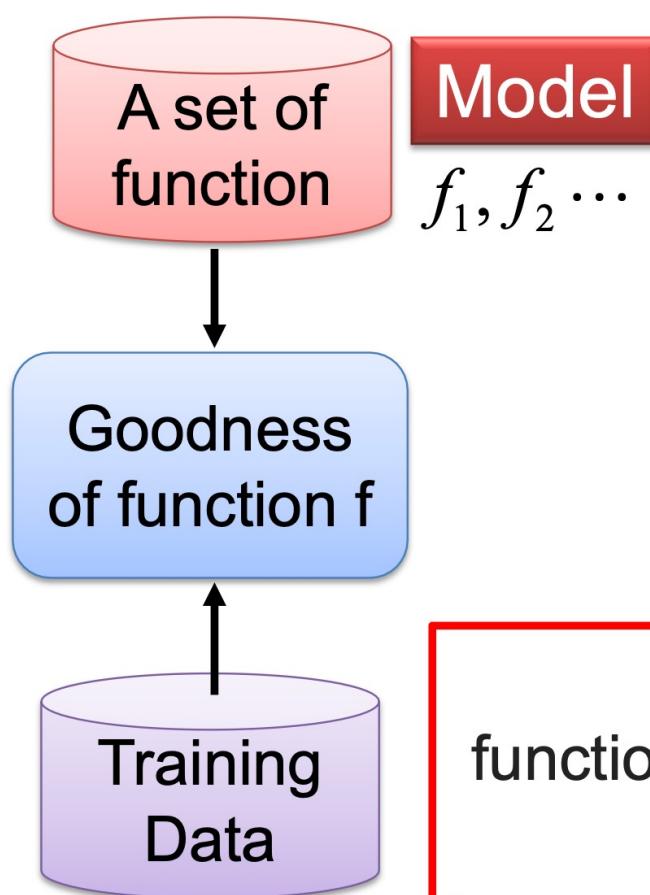
- Dialogue System

$$f\left(\begin{array}{c} \text{“Hi”} \\ \text{(what the user said)} \end{array} \right) = \begin{array}{c} \text{“Hello”} \\ \text{(system response)} \end{array}$$

Framework

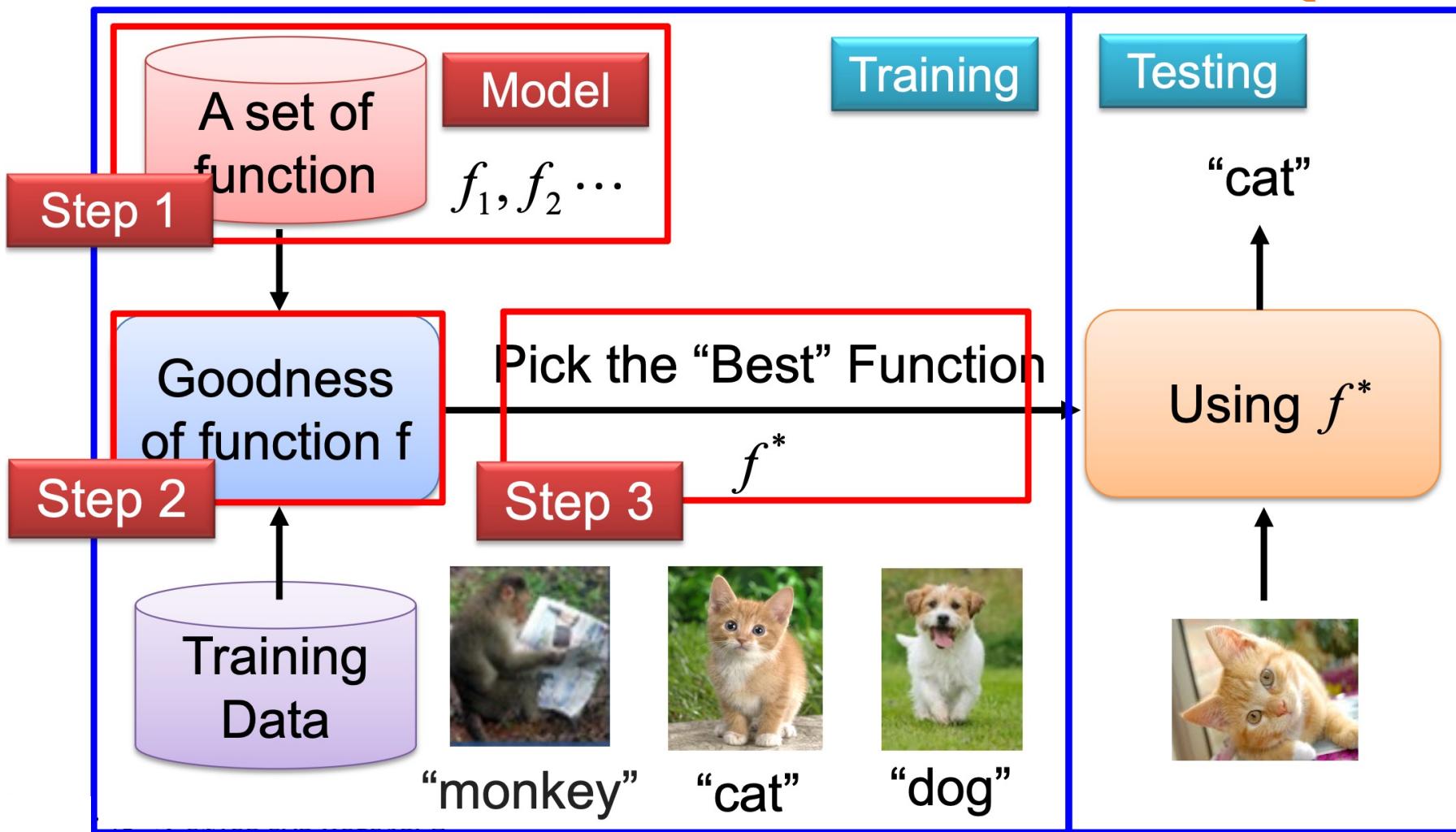


Framework



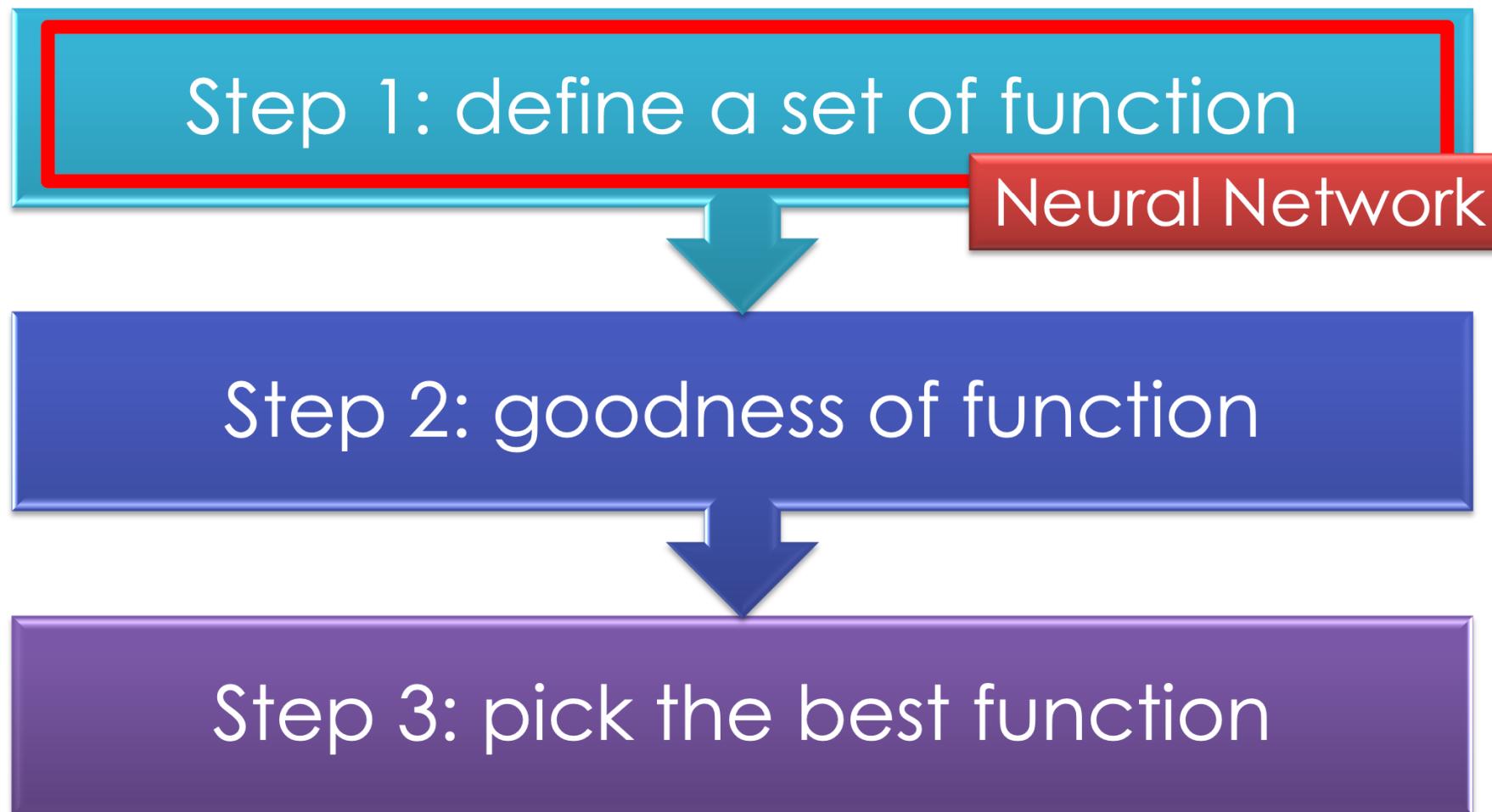
Framework

Image Recognition



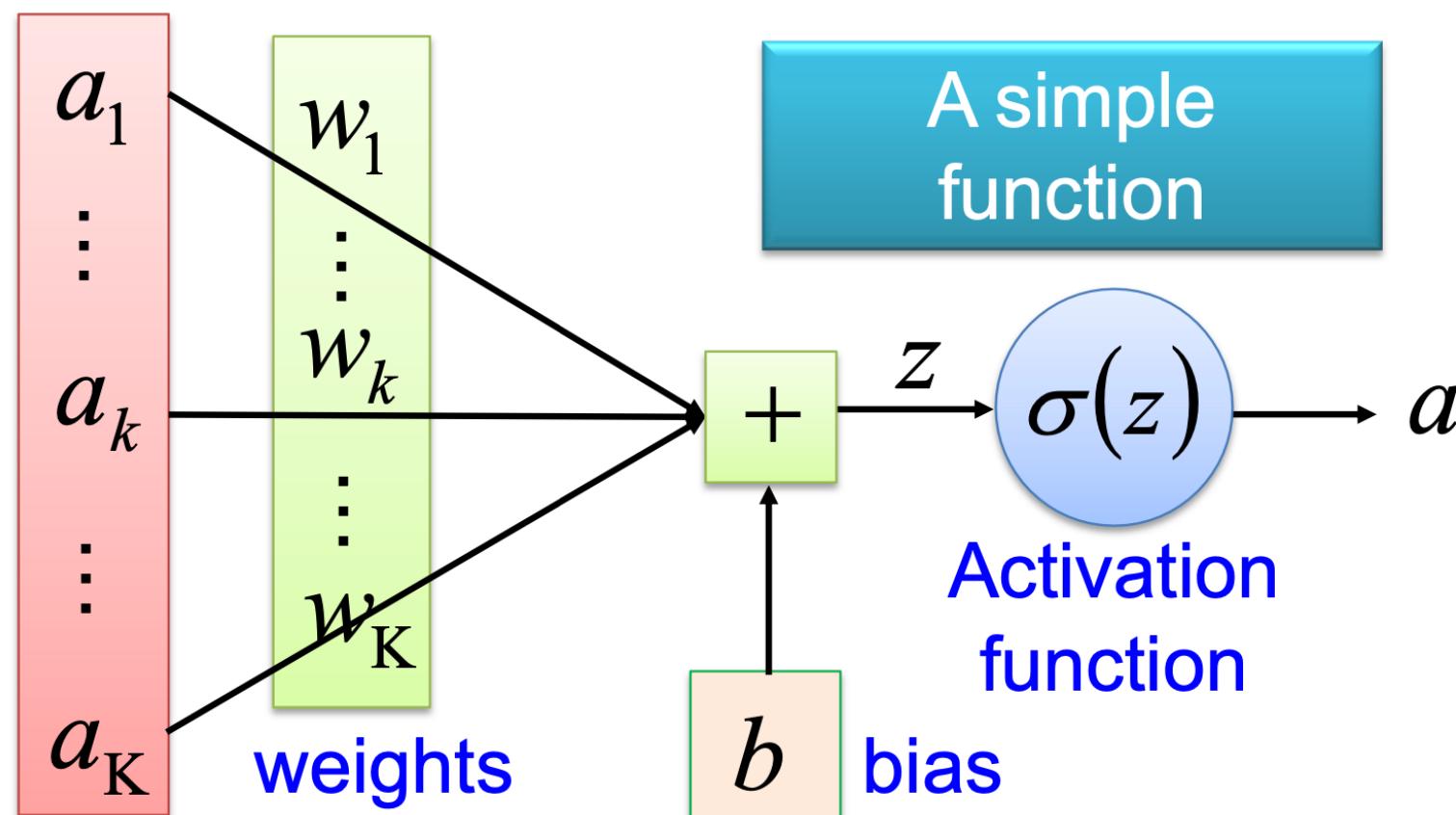
) = “Cat”

Three steps for Deep Learning

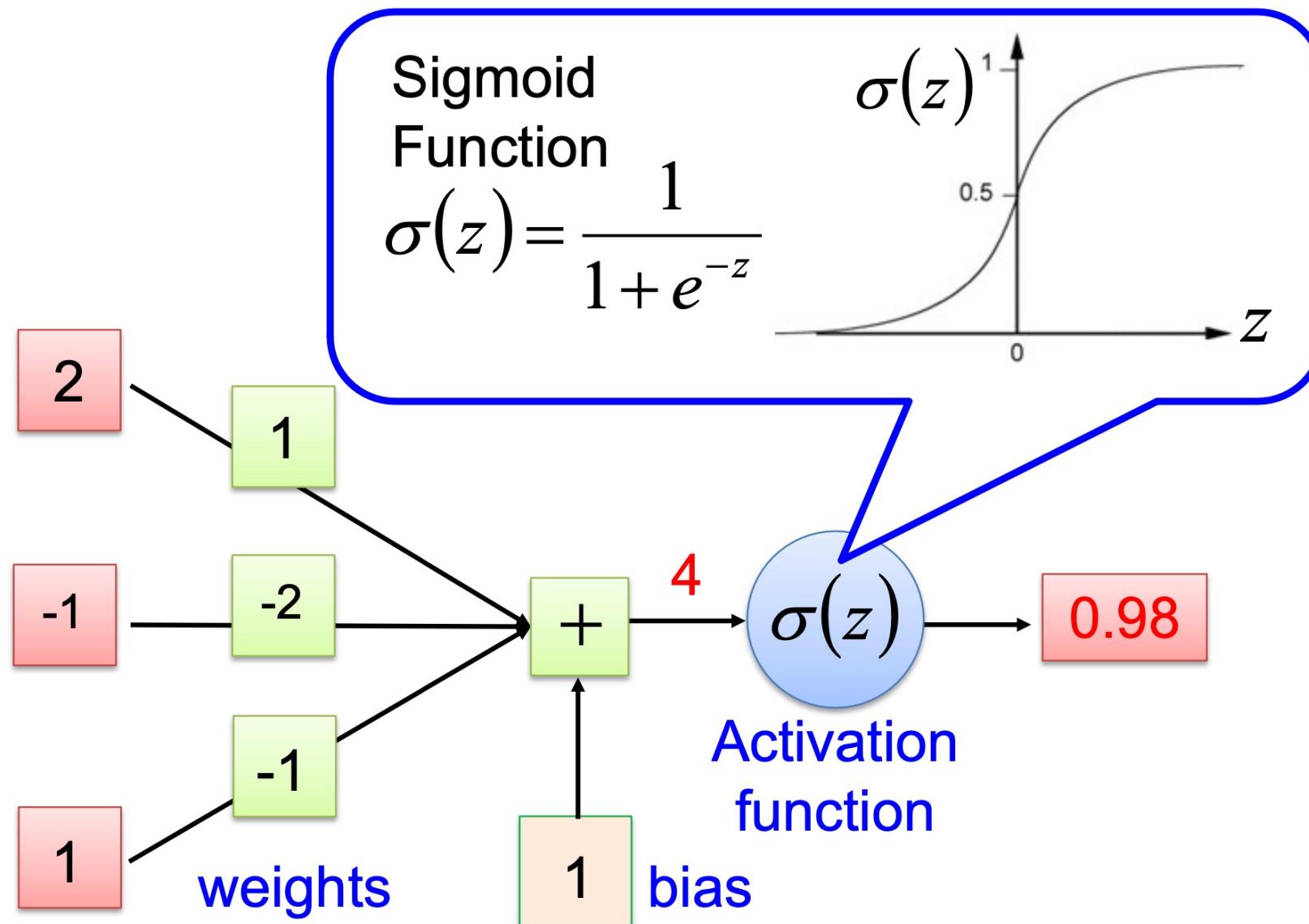


Neural Network: Neuron

$$z = a_1 w_1 + \cdots + a_k w_k + \cdots + a_K w_K + b$$

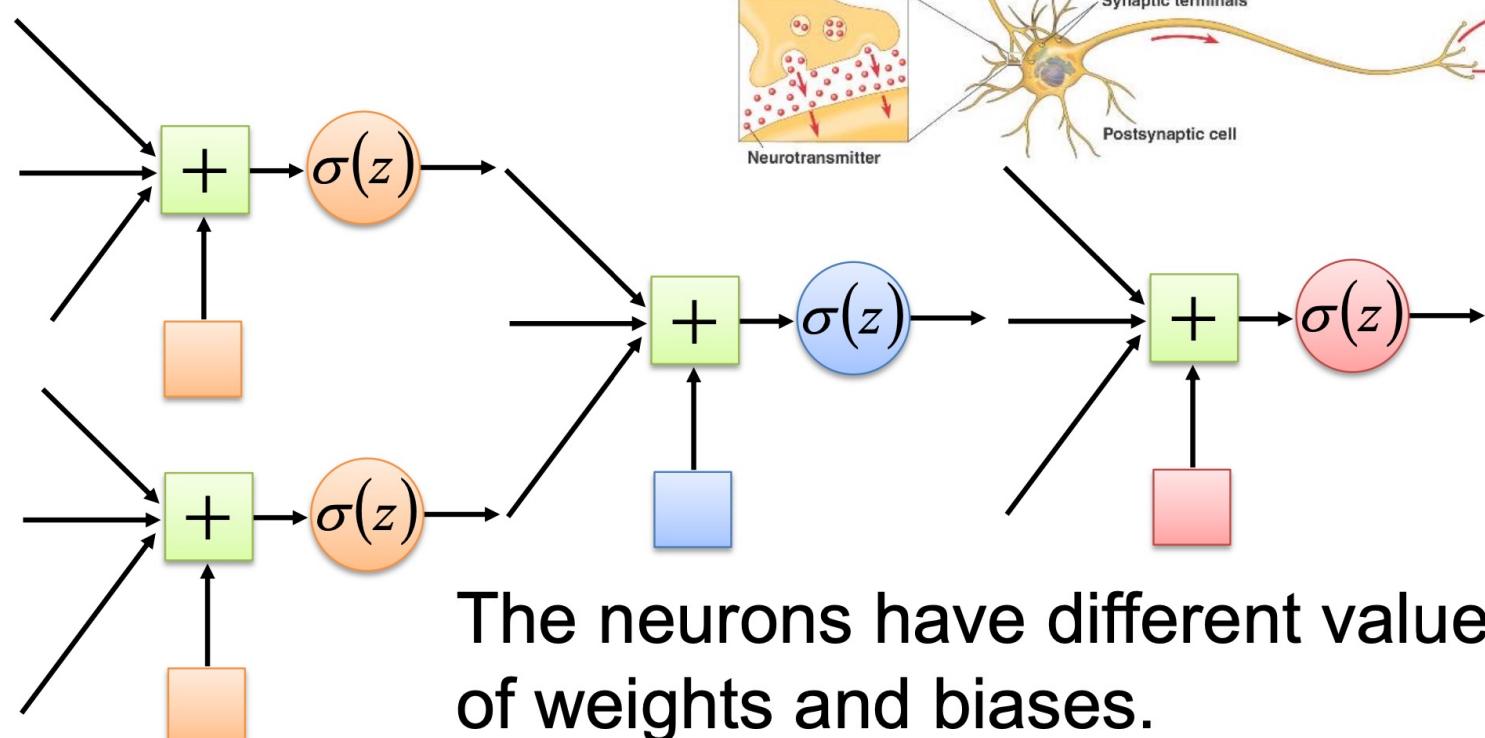
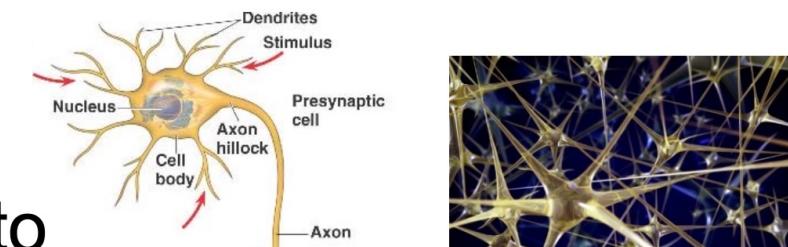


Neural Network: Neuron



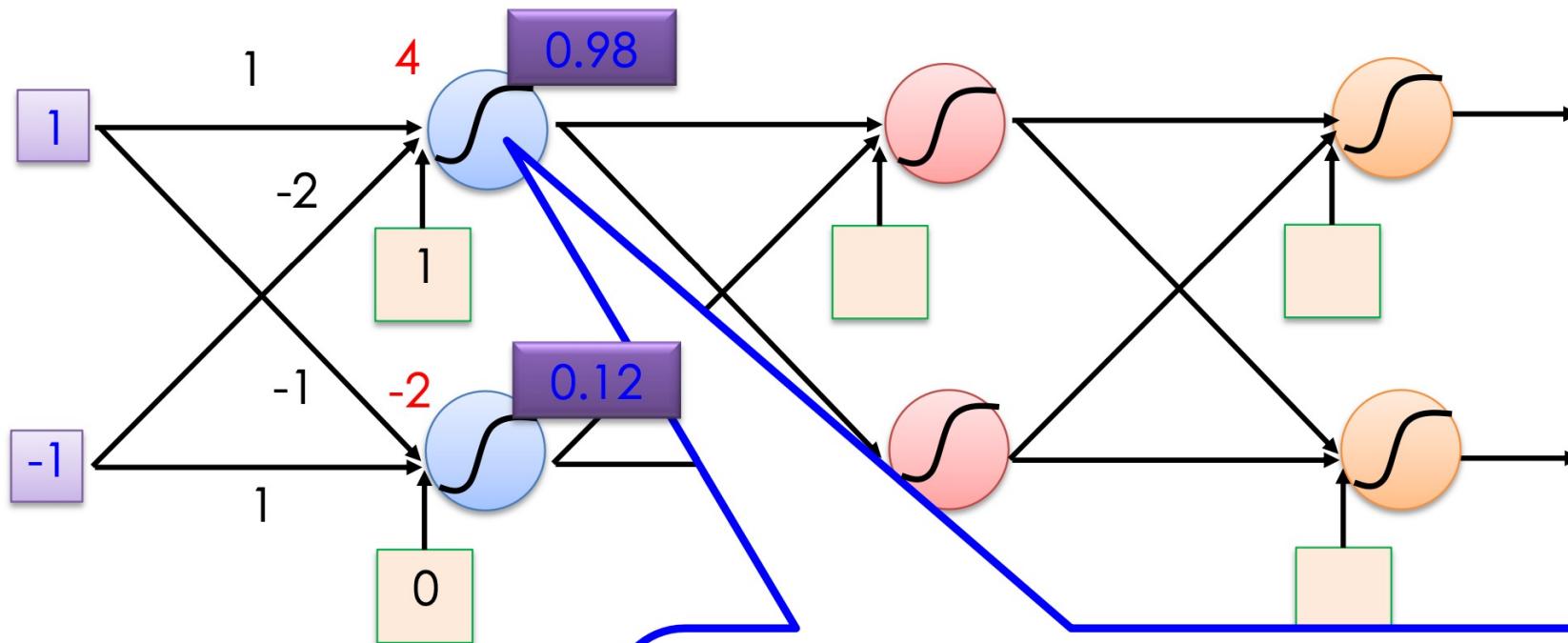
Neural Network: Neuron

Different connections lead to different network structures



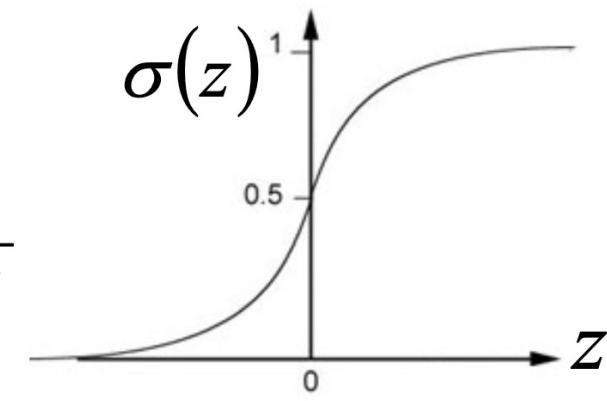
Weights and biases are network parameters θ

Fully Connected Feedforward Network

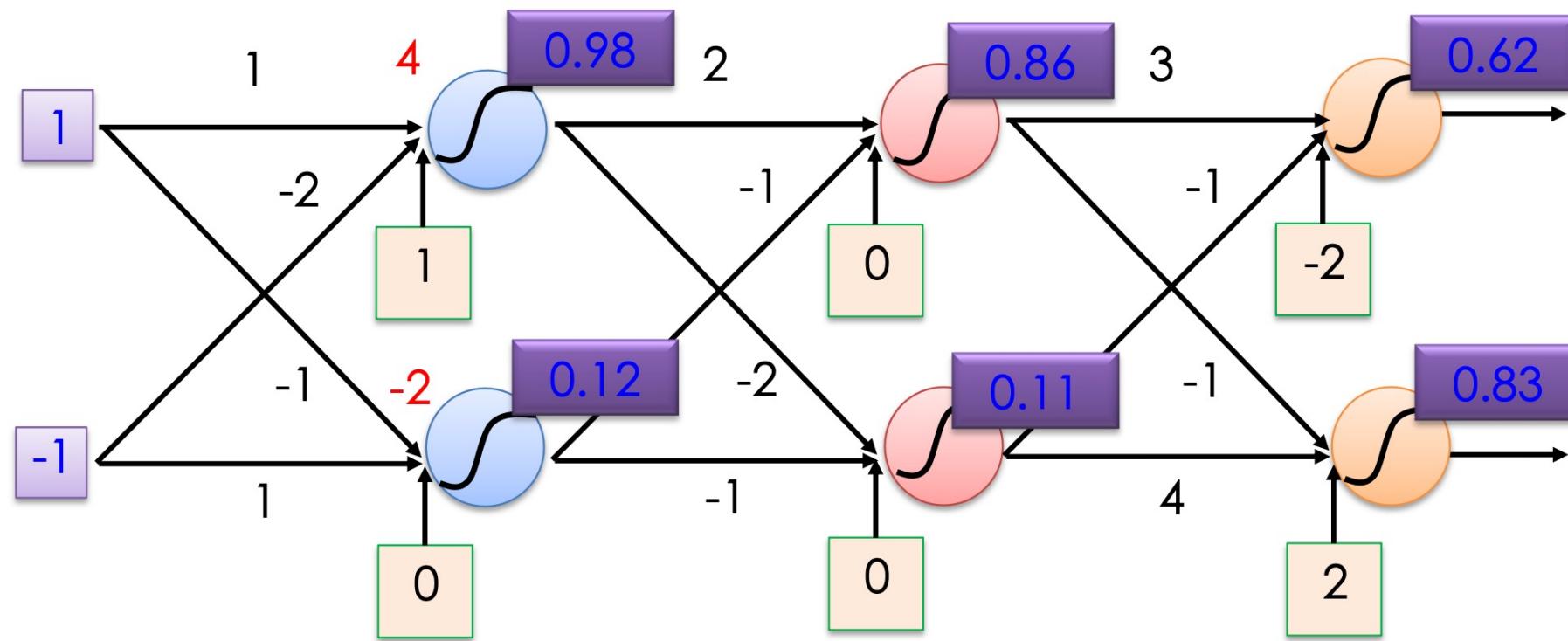


Sigmoid Function

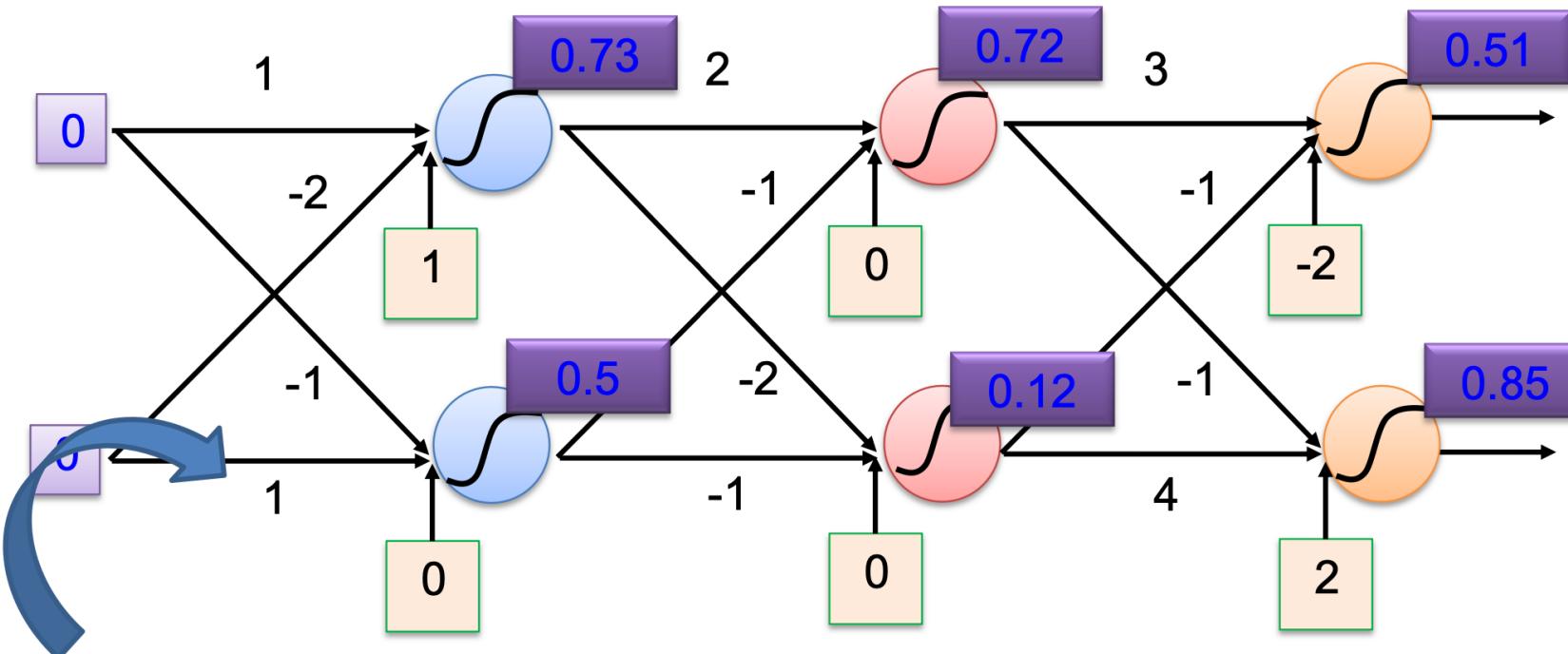
$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



Fully Connected Feedforward Network



Fully Connected Feedforward Network



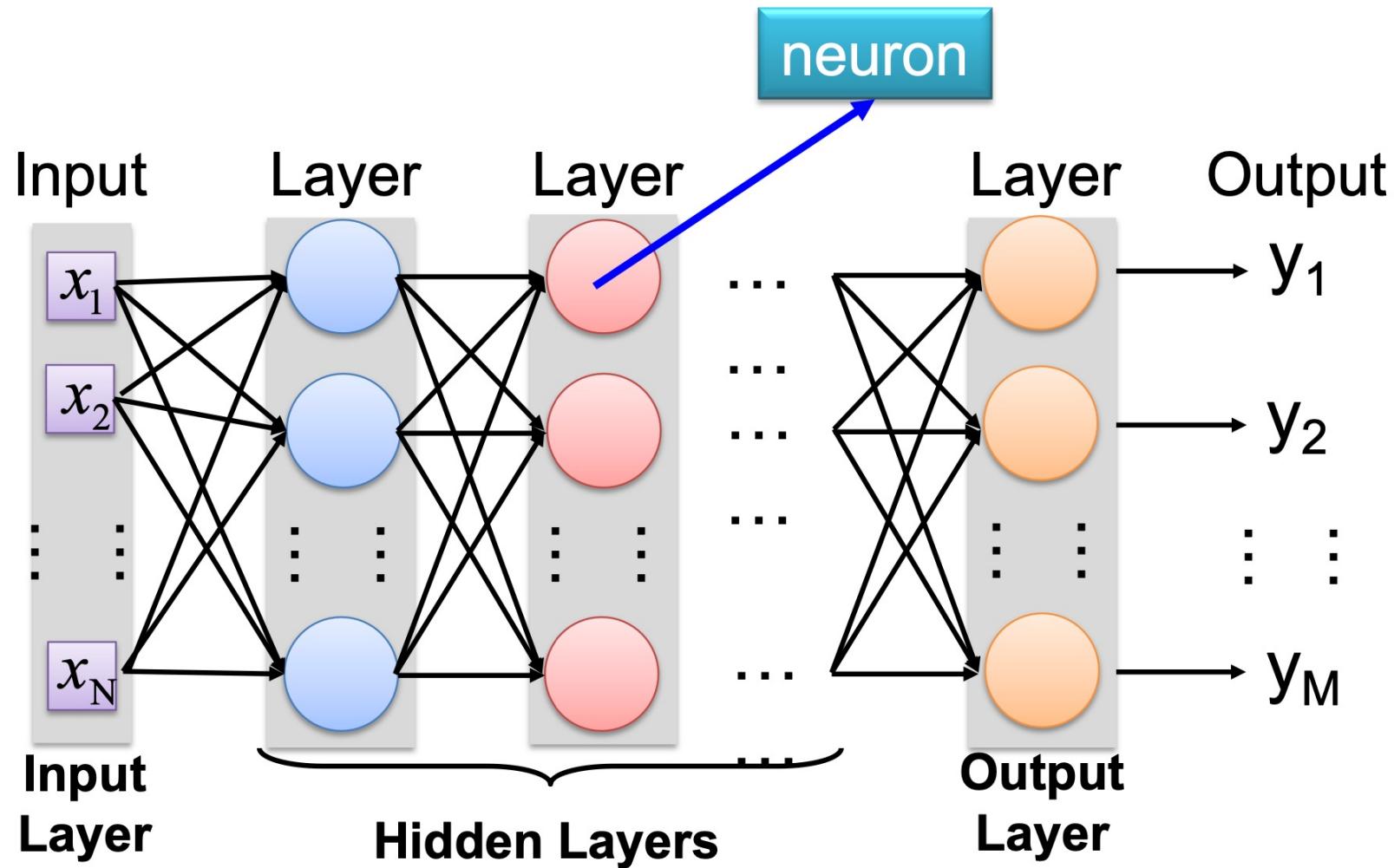
This is a function.
Input vector, output vector

$$f \left(\begin{bmatrix} 1 \\ -1 \end{bmatrix} \right) = \begin{bmatrix} 0.62 \\ 0.83 \end{bmatrix} \quad f \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix} \right) = \begin{bmatrix} 0.51 \\ 0.85 \end{bmatrix}$$

Given parameters θ , define a function

Given network structure, define a function set

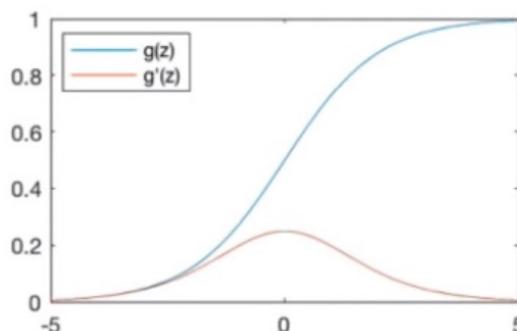
Fully Connected Feedforward Network



Deep means many hidden layers

Activation Function

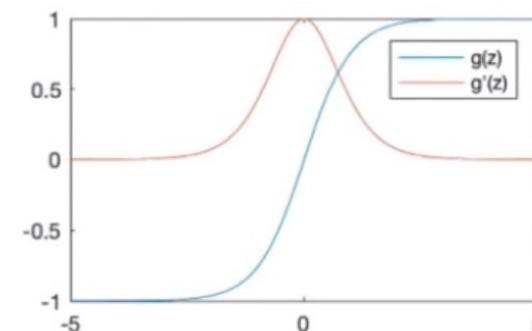
Sigmoid Function



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

$$g'(z) = g(z)(1 - g(z))$$

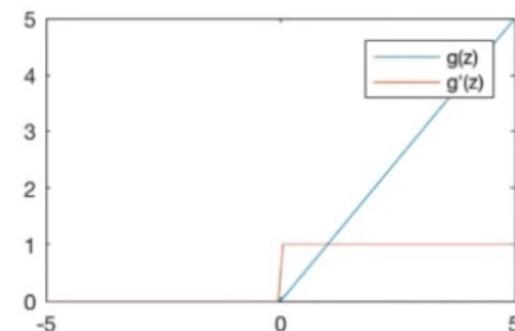
Hyperbolic Tangent



$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$g'(z) = 1 - g(z)^2$$

Rectified Linear Unit (ReLU)



$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$

NOTE: All activation functions are non-linear

https://ml-cheatsheet.readthedocs.io/en/latest/activation_functions.html

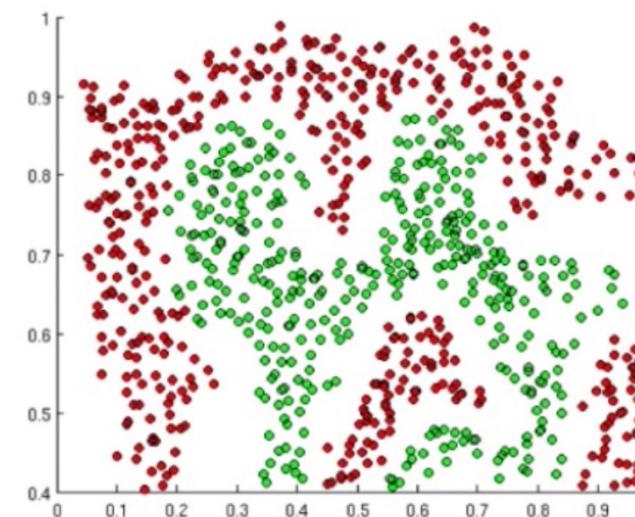
Activation Function

- Why non-linear activation function?

Activation Function

- Why non-linear activation function?

*The purpose of activation functions is to **introduce non-linearities** into the network*

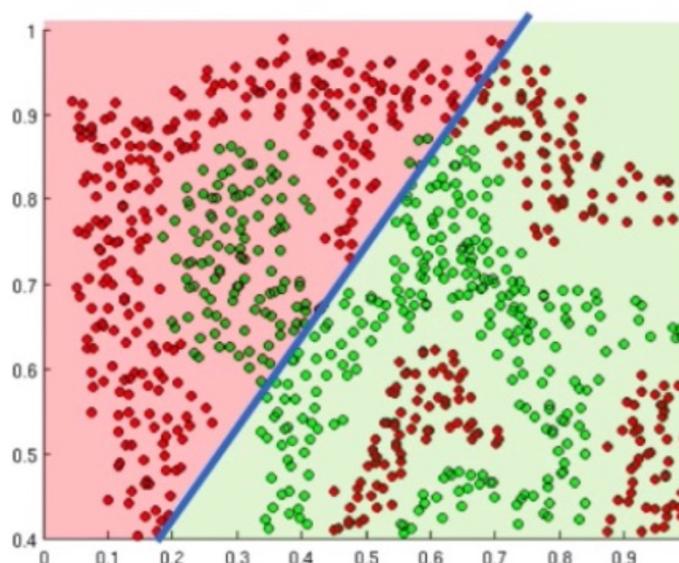


What if we wanted to build a Neural Network to
distinguish green vs red points?

Activation Function

- Why non-linear activation function?

*The purpose of activation functions is to **introduce non-linearities** into the network*

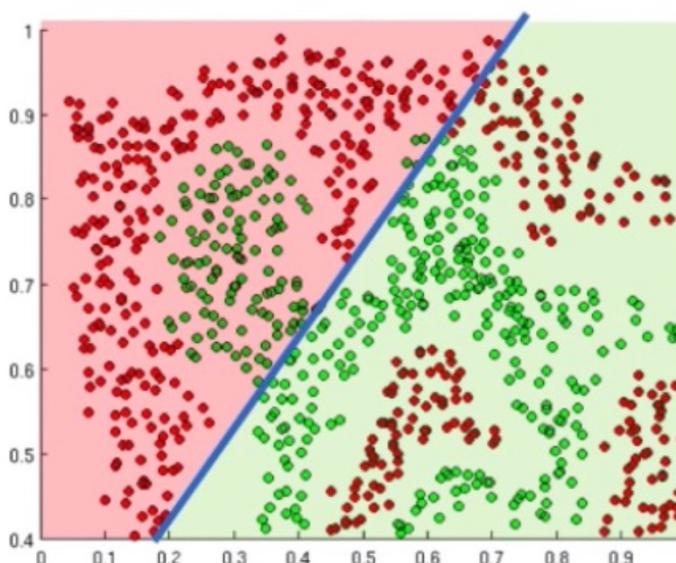


Linear Activation functions produce linear decisions no matter the network size

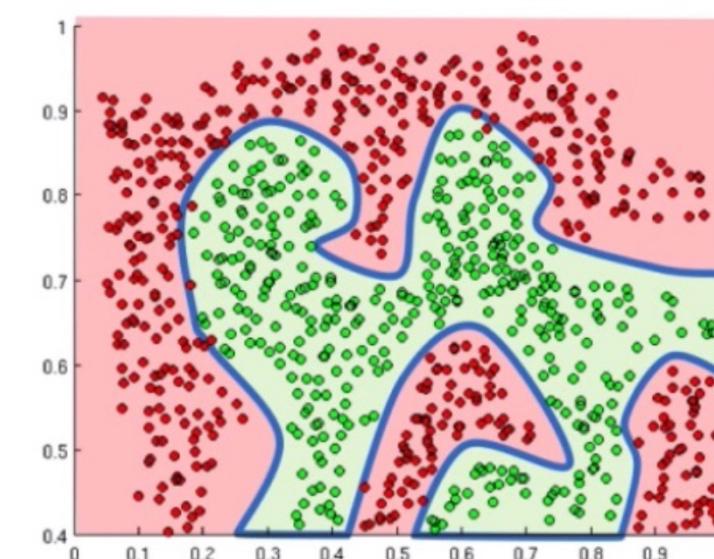
Activation Function

- Why non-linear activation function?

*The purpose of activation functions is to **introduce non-linearities** into the network*



Linear Activation functions produce linear decisions no matter the network size

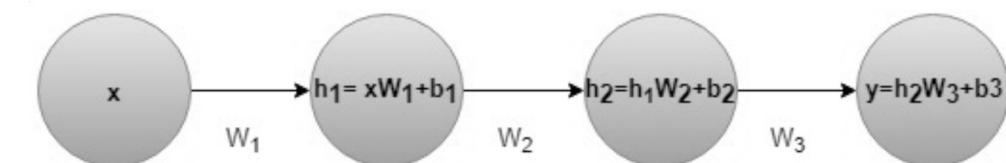


Non-linearities allow us to approximate arbitrarily complex functions

Activation Function

- Why non-linear activation function?

➤ With **linear** activation functions, no matter how many layers in the NN, the last layer will be a linear function of the first layer (because a linear combination of linear functions is still a linear function) → a linear activation function turns the NN into just one layer.



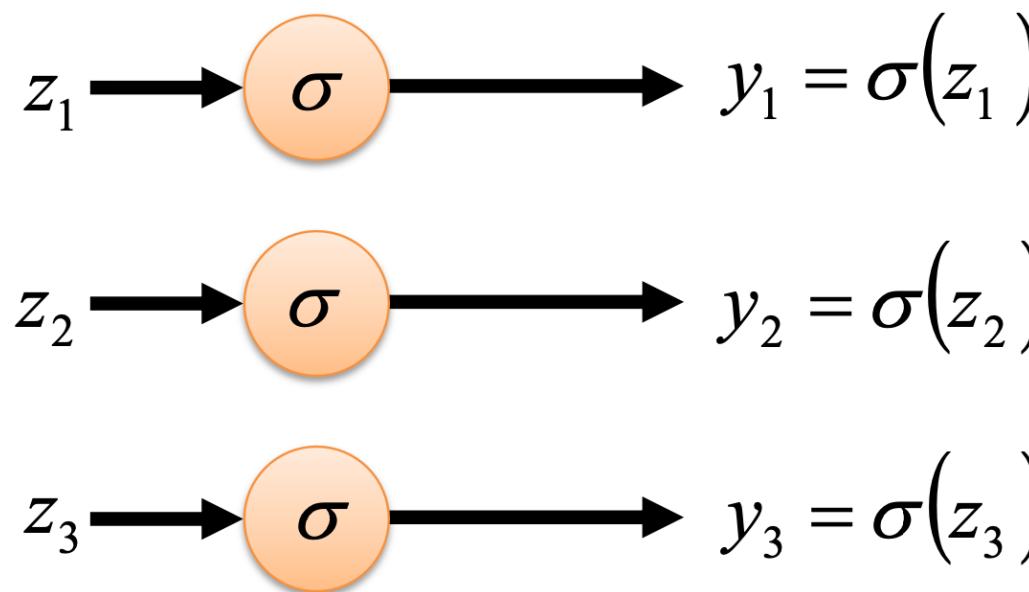
$$\begin{aligned} y &= h2 * w3 + b3 \\ &= (h1 * w2 + b2) * w3 + b3 \\ &= h1 * w2 * w3 + b2 * w3 + b3 \\ &= (x * w1 + b1) * w2 * w3 + b2 * w3 + b3 \\ &= x * w1 * w2 * w3 + b1 * w2 * w3 + b2 * w3 + b3 \\ &= x * w' + b' \end{aligned}$$

A linear activation function turns the neural network into just one layer.

➤ A NN with a linear activation function is simply a linear regression model. It has limited **power** and **ability** to handle complexity varying parameters of input data.

Output Layer: Softmax Function

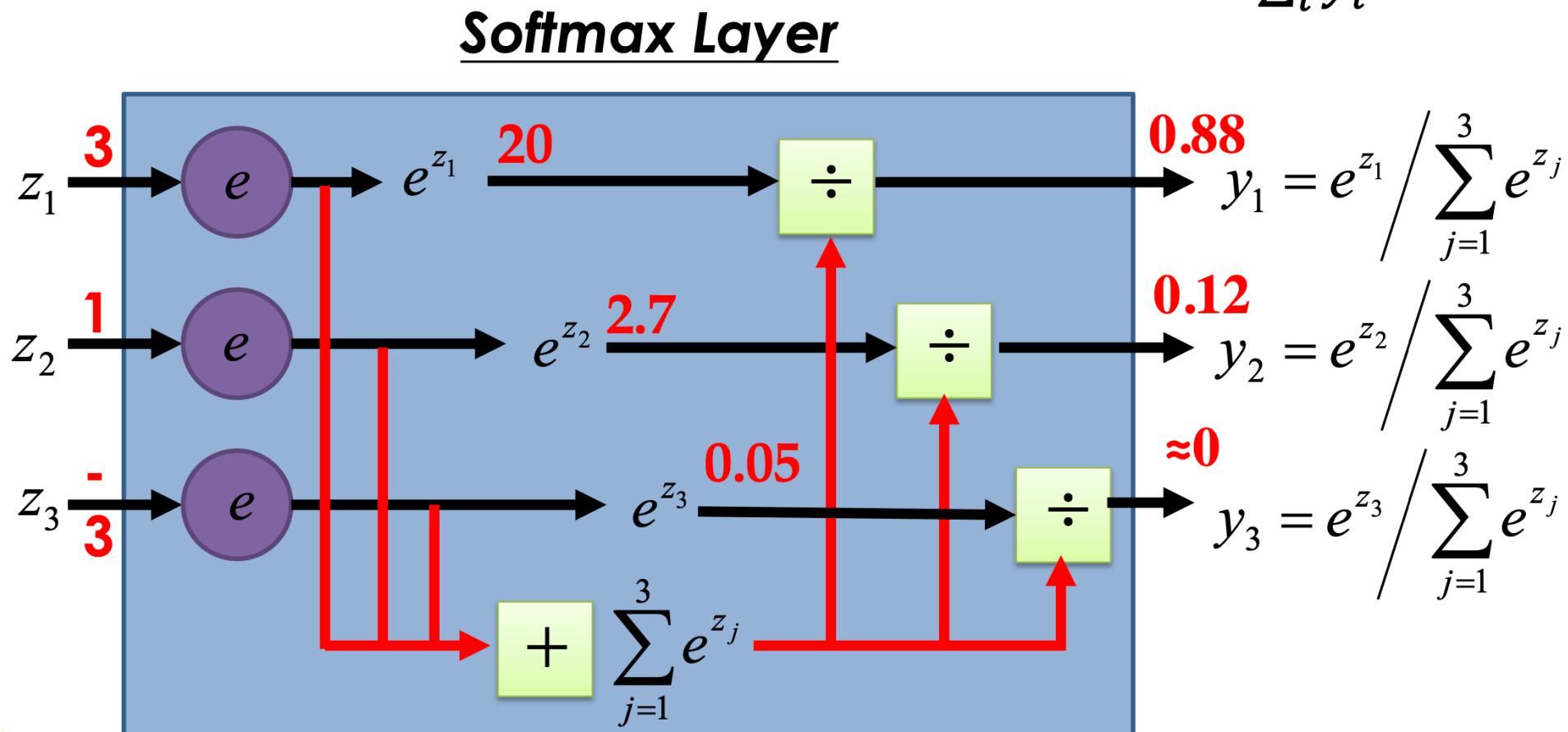
Ordinary Layer



In general, the output of network can be any value.

May not be easy to interpret

Output Layer: Softmax Function



Probability:

- $1 > y_i > 0$
- $\sum_i y_i = 1$

Three steps for Deep Learning

Step 1: define a set of function

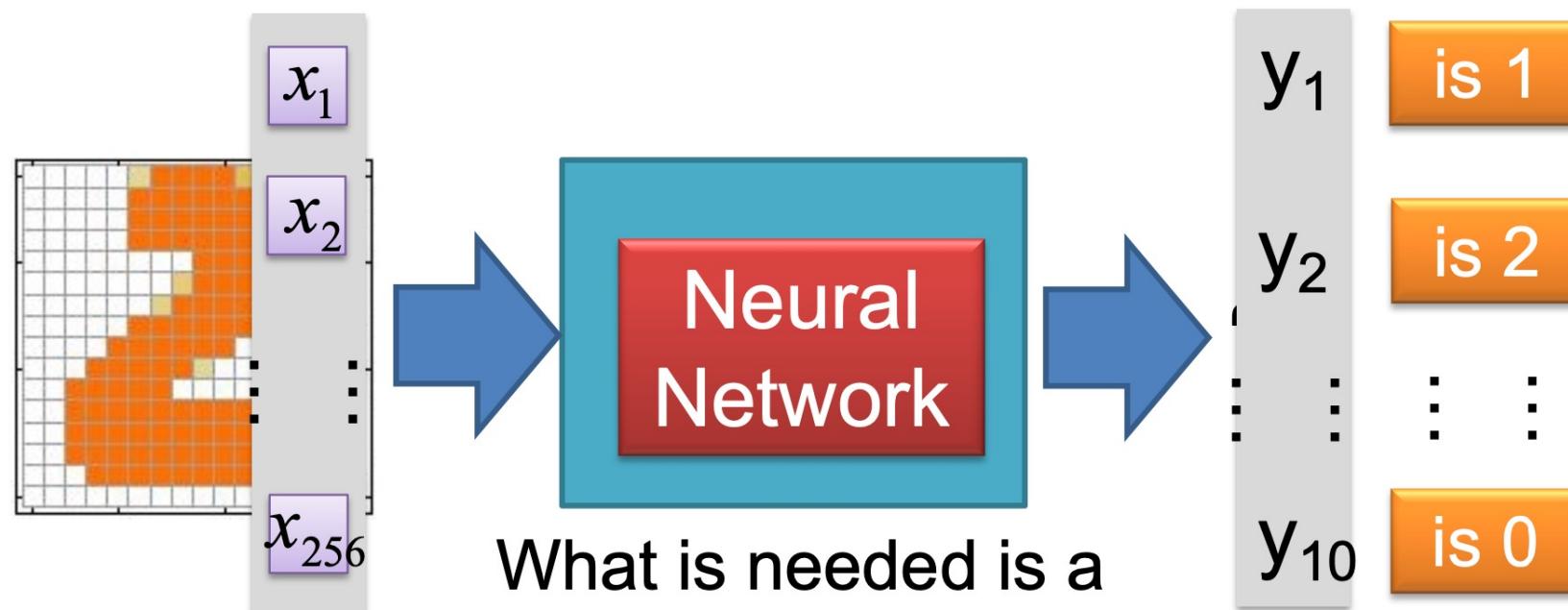


Step 2: goodness of function



Step 3: pick the best function

Application Example: Handwriting Digit Recognition



What is needed is a
function

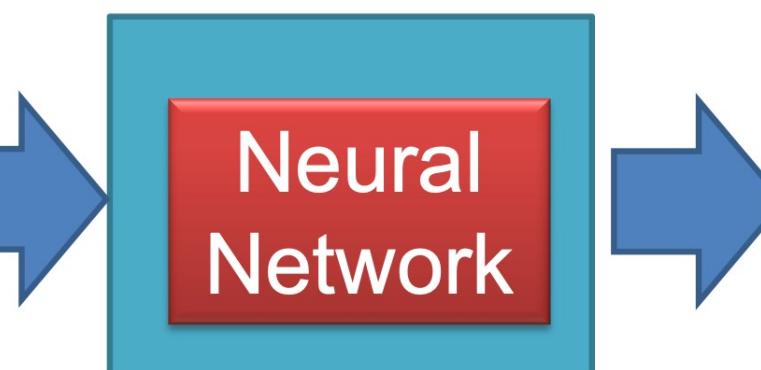
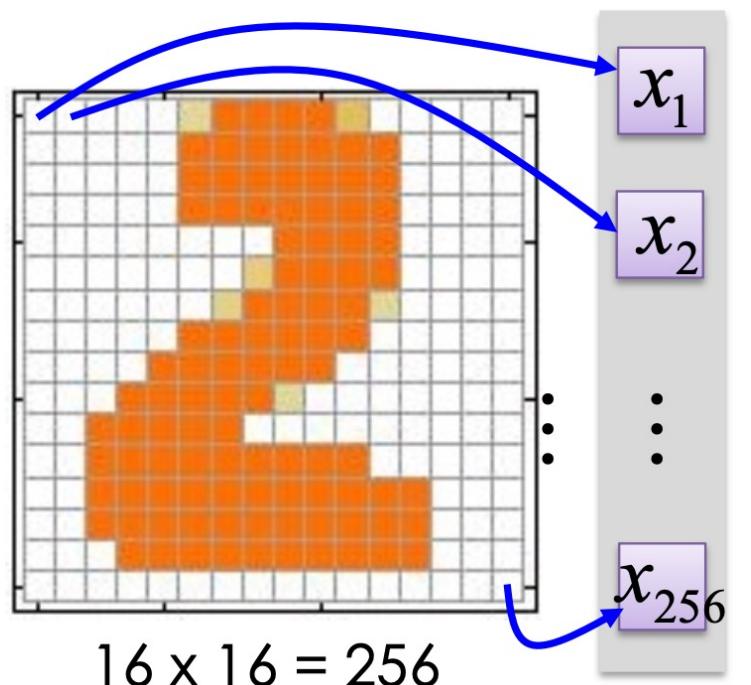
Input:
256-dim vector

output:
10-dim vector

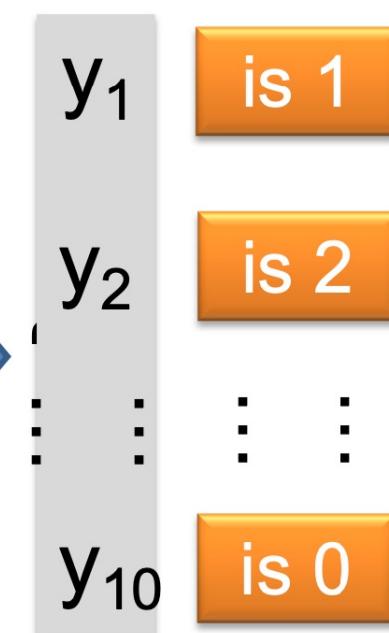
Application Example: Handwriting Digit Recognition



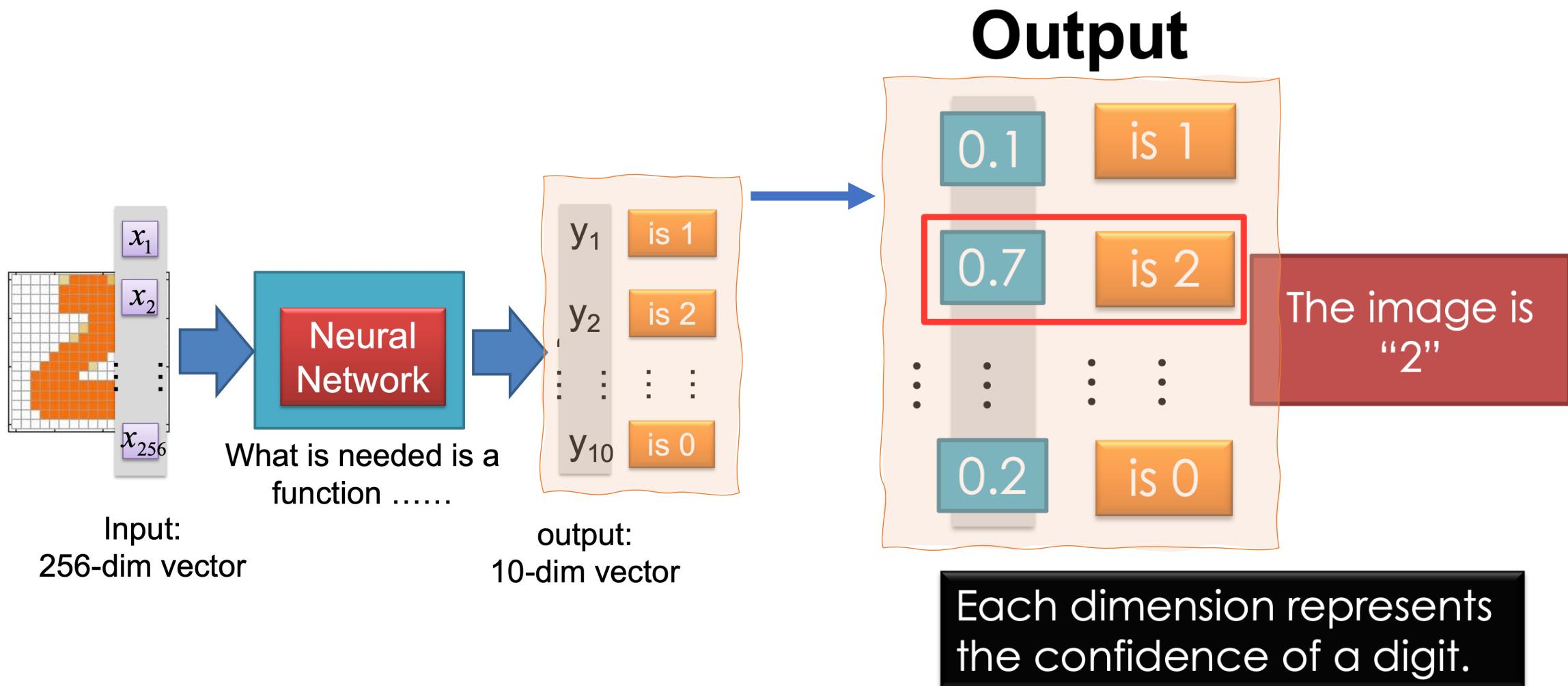
Input



What is needed is a function



Application Example: Handwriting Digit Recognition



Preparing Training Data: Images and their labels



“5”



“0”



“4”



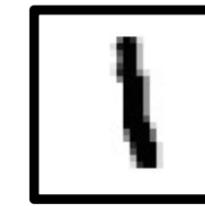
“1”



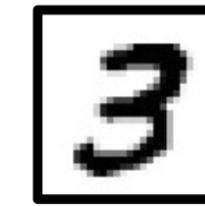
“9”



“2”



“1”

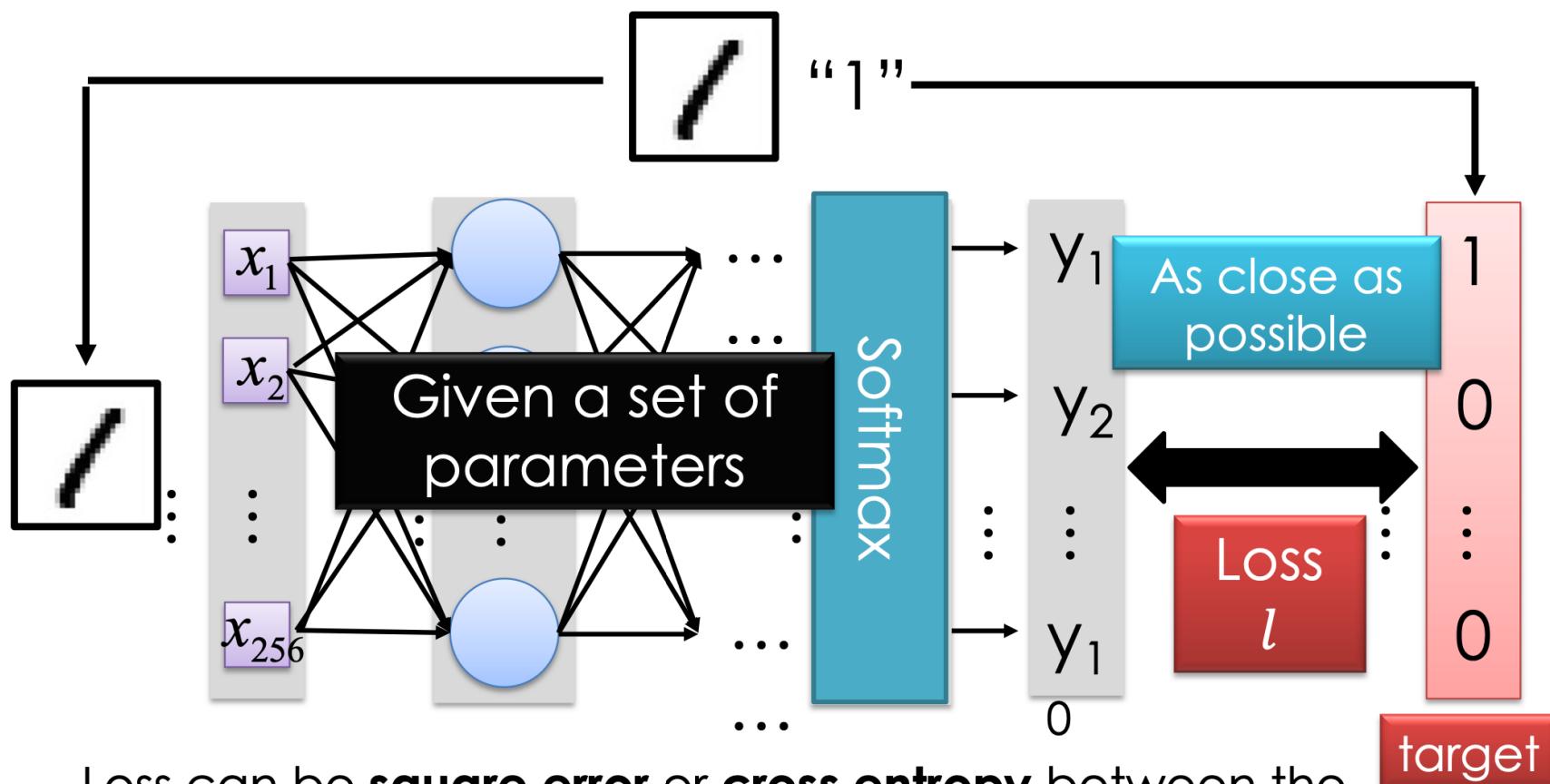


“3”

The learning target is defined on the
training data.

Loss

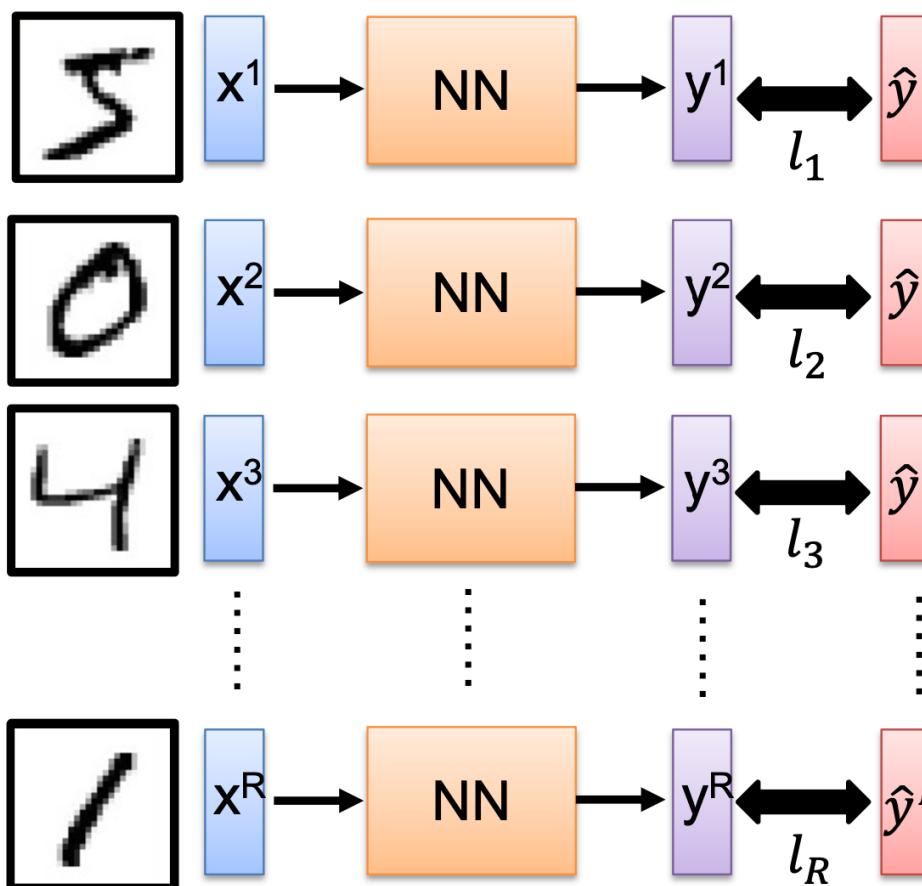
A good function should make the loss of all examples as small as possible.



Loss can be **square error** or **cross entropy** between the network output and target

Total Loss

For all training data ...



Total Loss:

$$L = \sum_{r=1}^R l_r$$

As small as possible

Find a function in function set that minimizes total loss L

Find the network parameters θ^* that minimize total loss L

Three steps for Deep Learning

Step 1: define a set of function



Step 2: goodness of function



Step 3: pick the best function

Testing or Inference

Convolutional Neural Network (CNN)

Widely used in image
processing and computer vision

Why CNN for Images?

- Some patterns are much smaller than the whole image

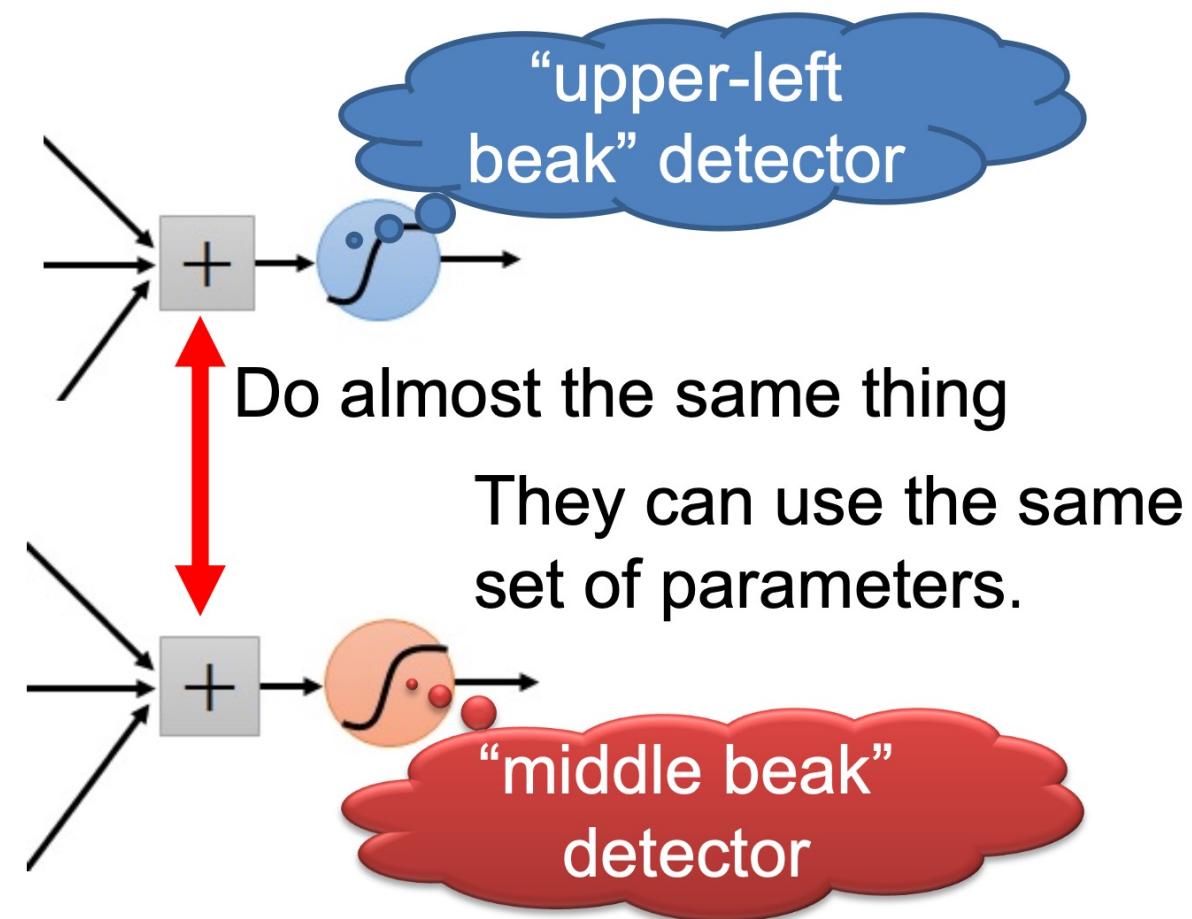
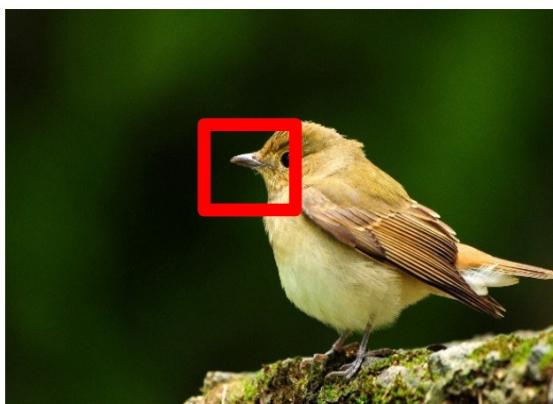
A neuron does not have to see the whole image to discover the pattern.

Connecting to small region with less parameters



Why CNN for Images?

- The same patterns appear in different regions.



Why CNN for Images?

- Subsampling the pixels will not change the object

bird



subsampling

bird



We can subsample the pixels to make image smaller



Less parameters for the network to process the image

Why CNN for Images?

Property 1

- Some patterns are much smaller than the whole

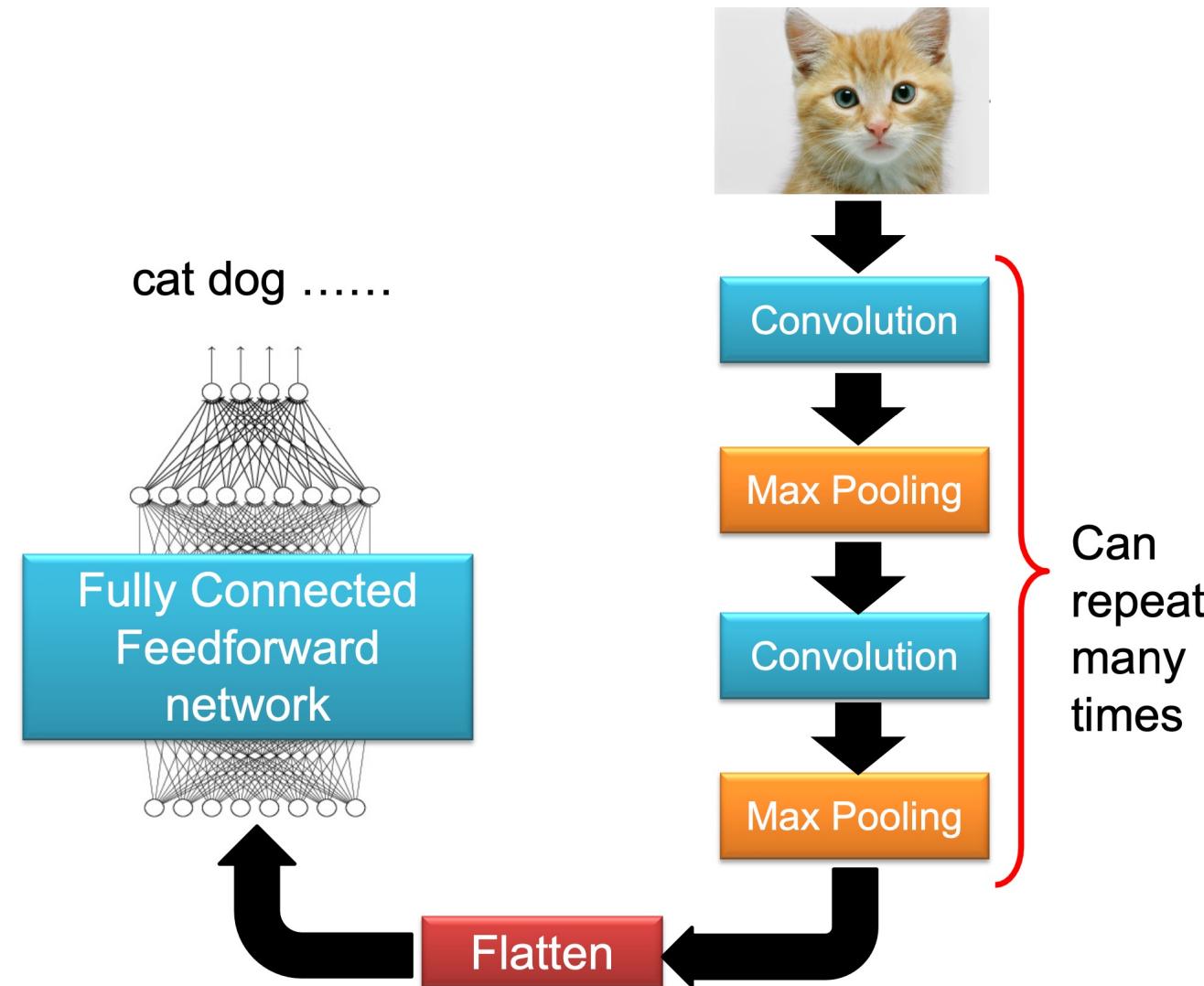
Property 2

- The same patterns appear in different regions.

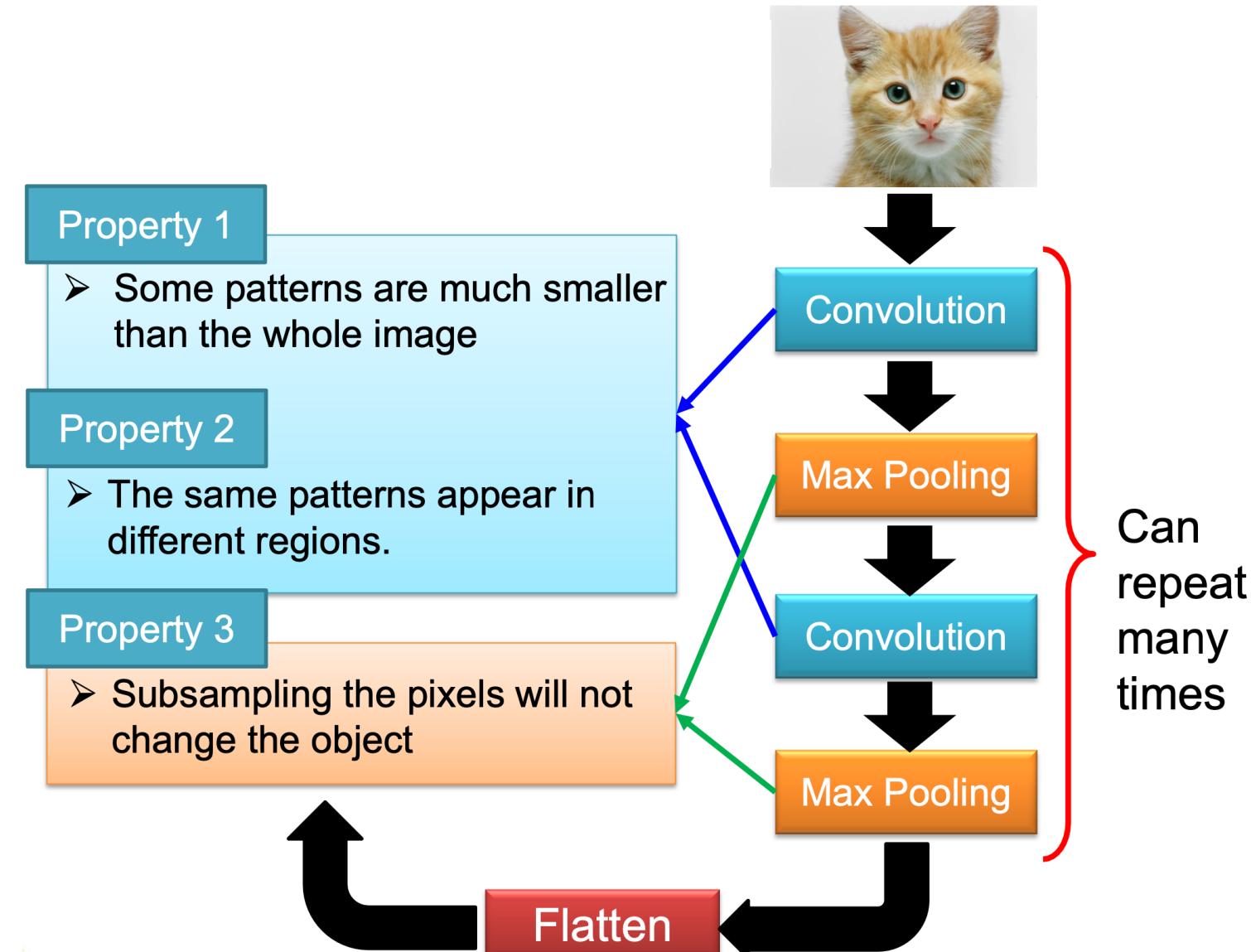
Property 3

- Subsampling the pixels will not change the object

The Whole CNN



The Whole CNN



CNN: Convolution

Those are the network
parameters to be learned.

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image

1	-1	-1
-1	1	-1
-1	-1	1

Filter/kernel 1

Matrix

-1	1	-1
-1	1	-1
-1	1	-1

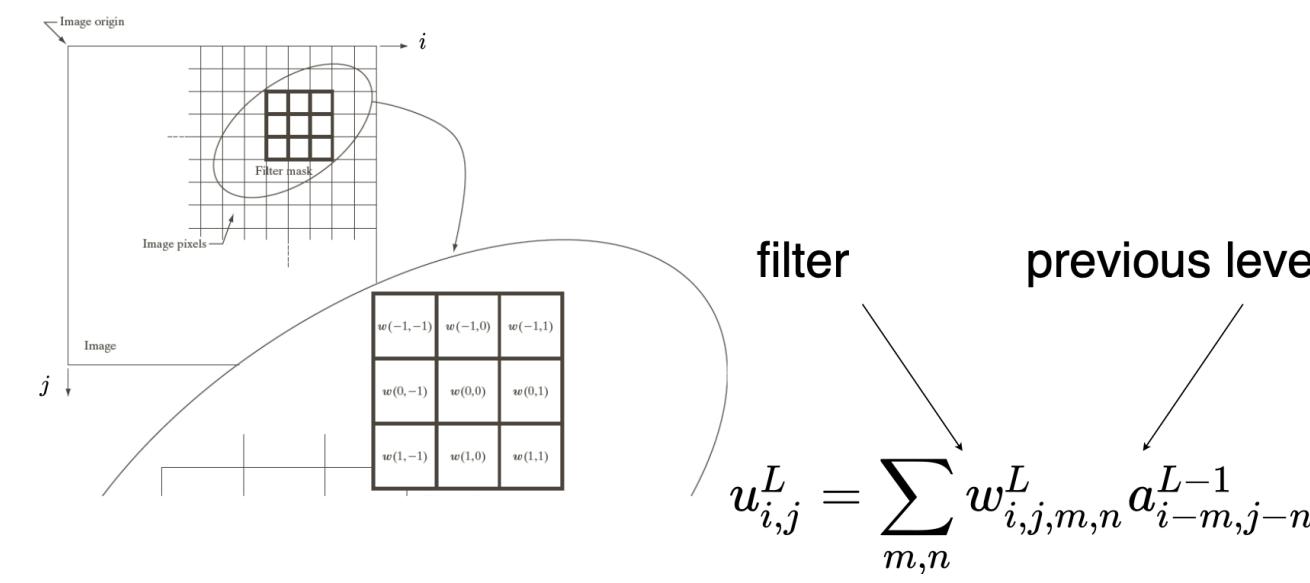
Filter 2

Matrix

⋮ ⋮

Each filter detects a
small pattern (3 x 3).

CNN: Convolution



Filter/Kernel 1		
1	-1	-1
-1	1	-1
-1	-1	1

stride=1

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

3 -1

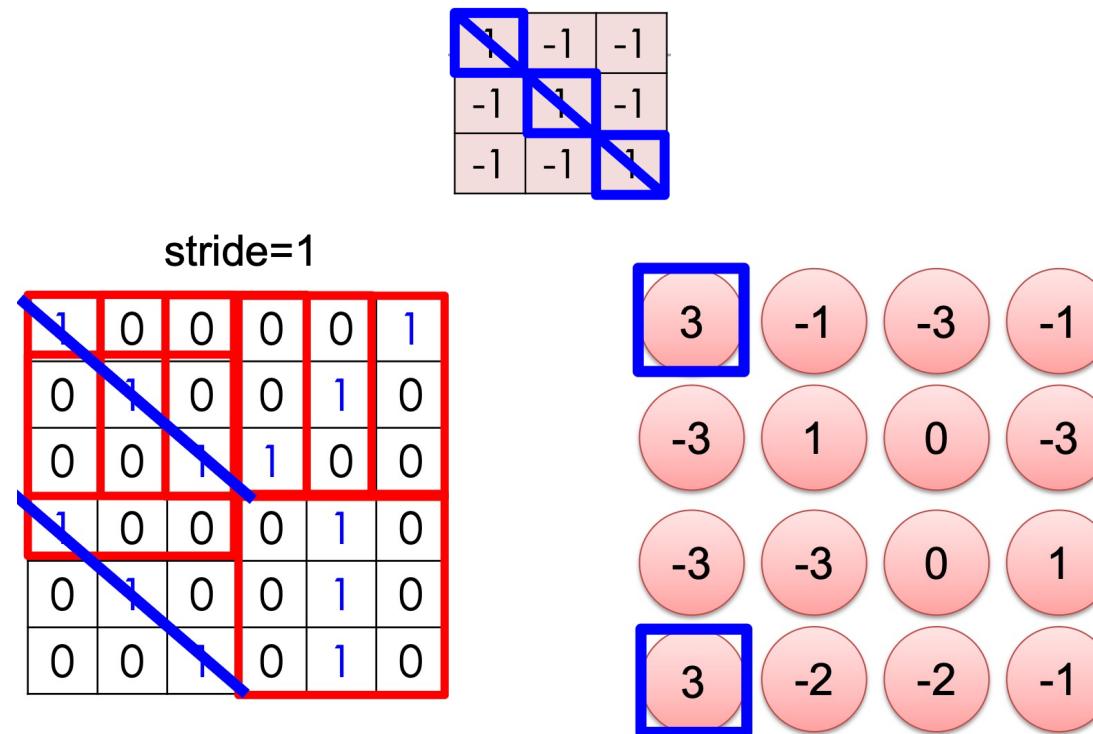
If stride=2

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

3 -3

We set stride=1 below

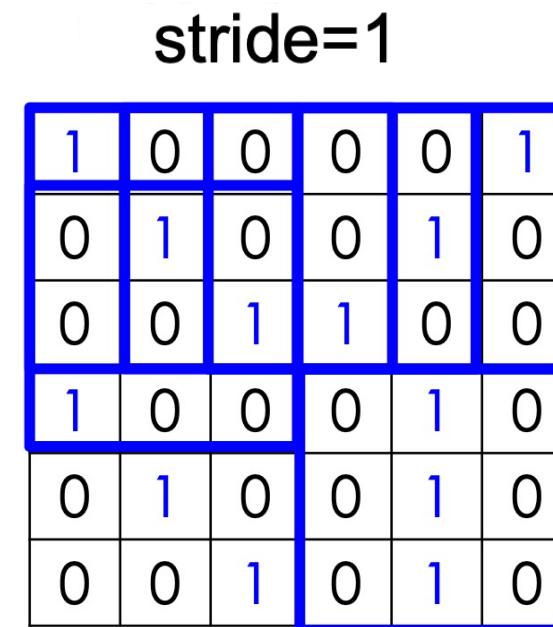
CNN: Convolution



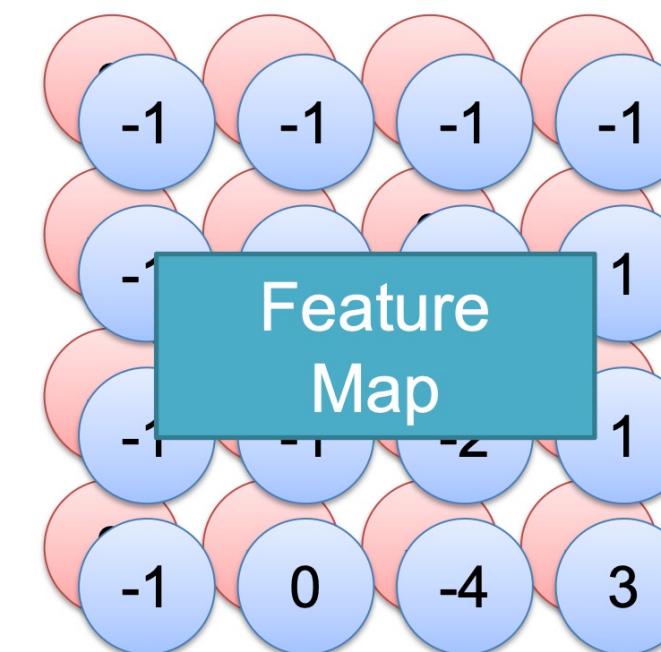
CNN: Convolution

Filter 2

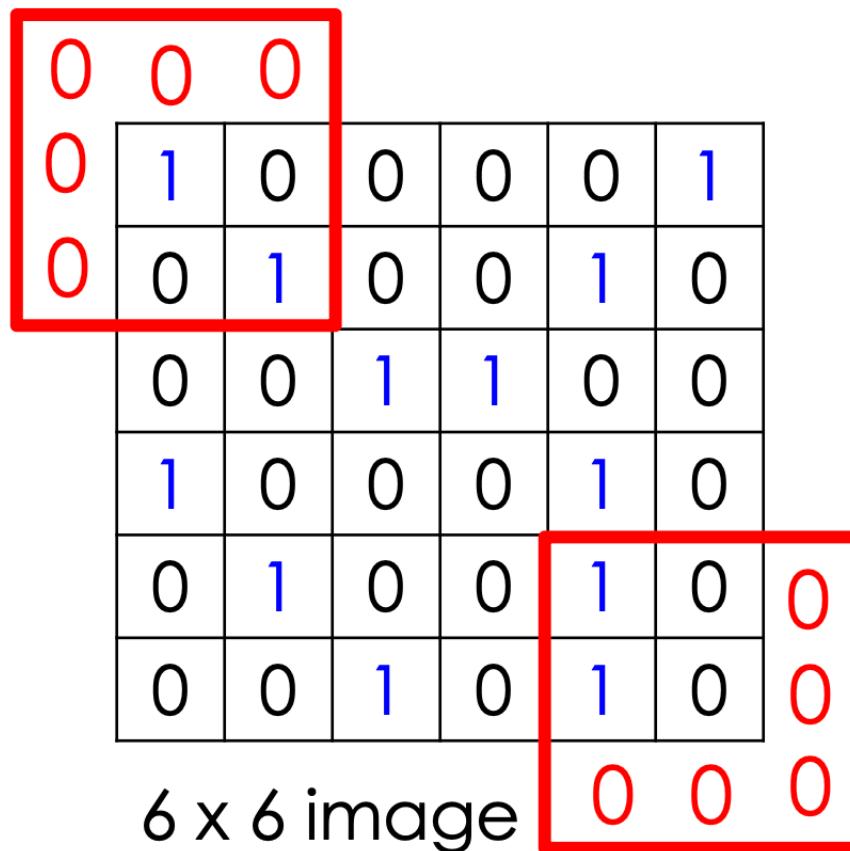
-1	1	-1
-1	1	-1
-1	1	-1



Do the same process
for every filter



CNN: Zero Padding

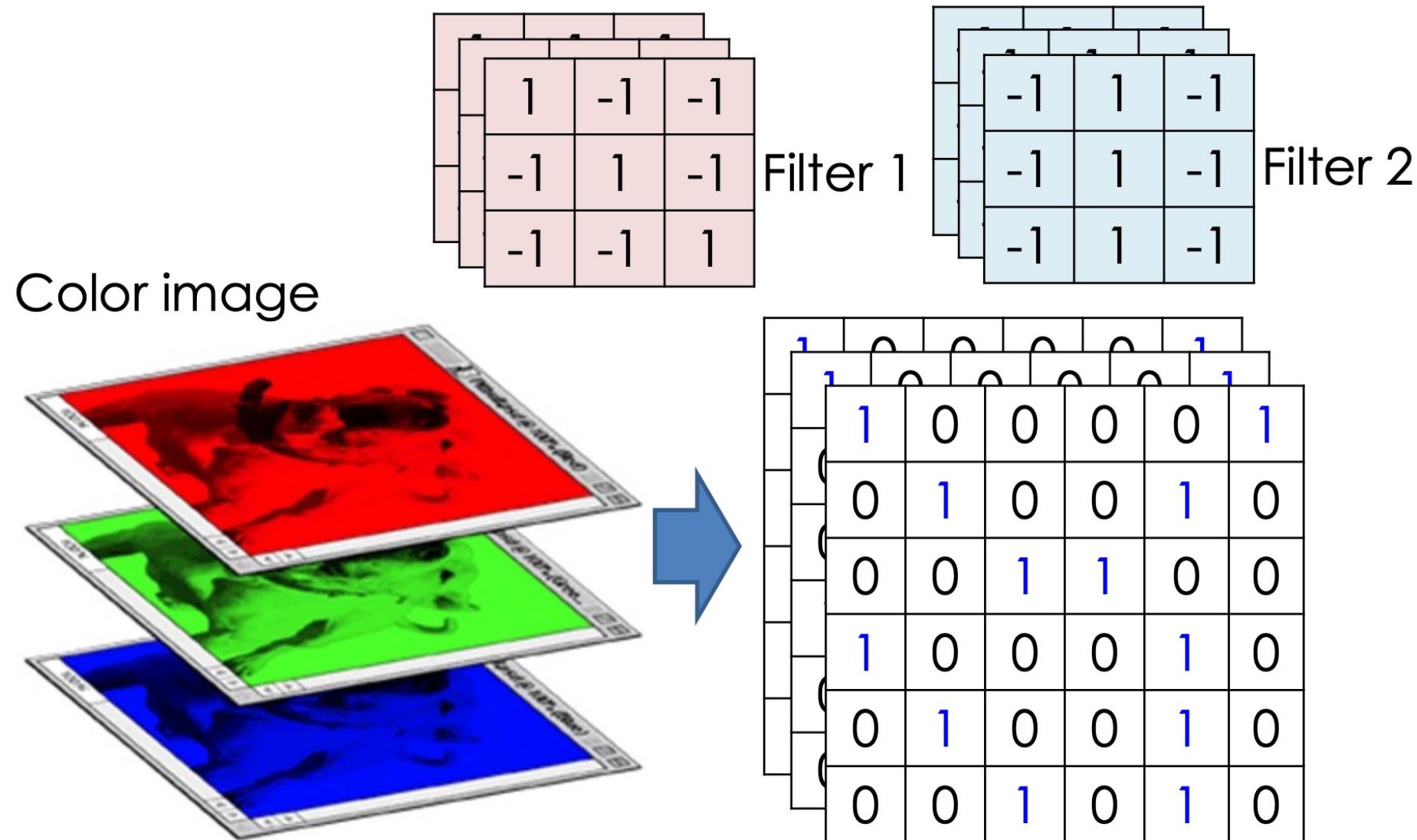


You will get another 6×6 image in this way

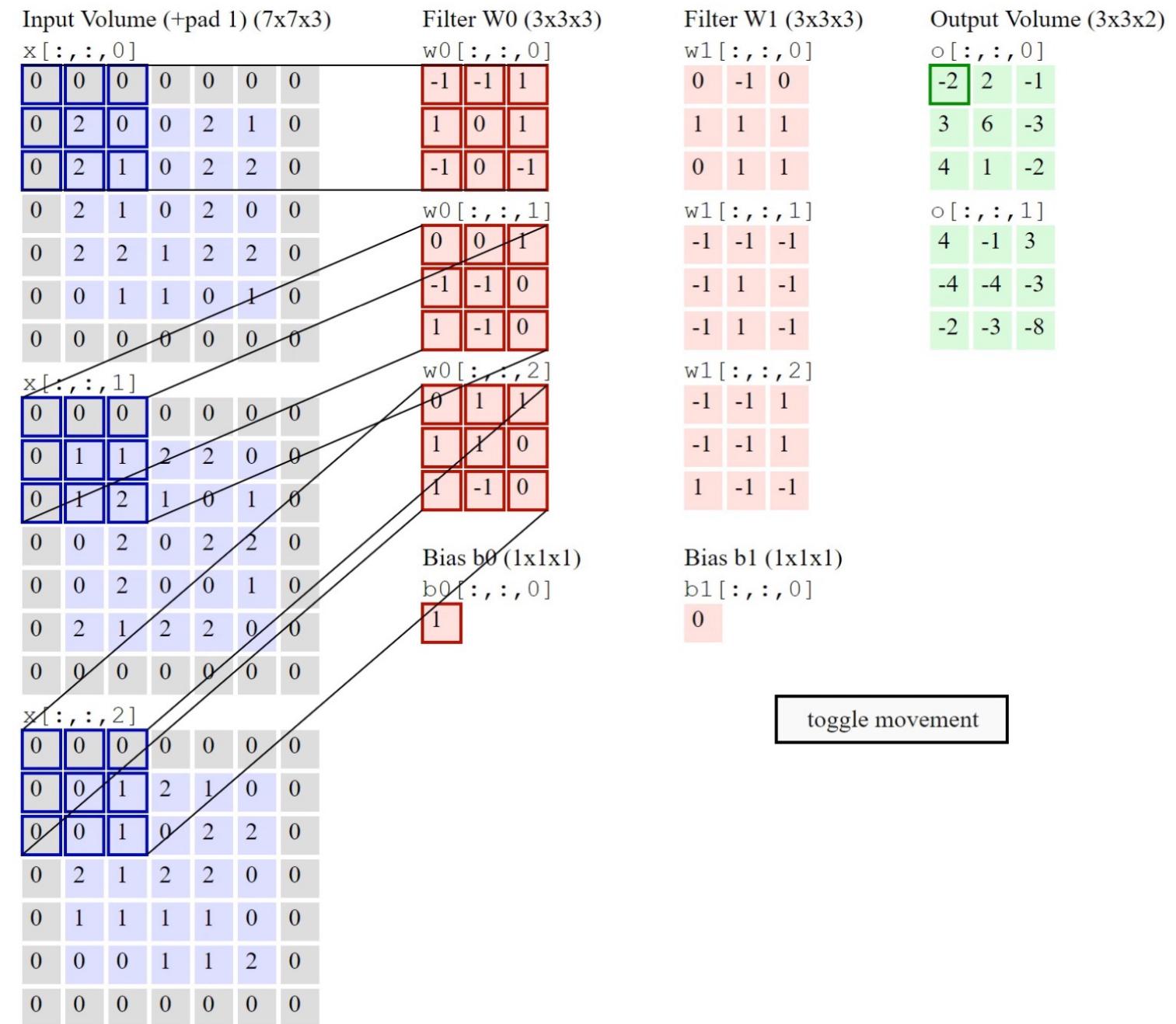


Zero
padding

CNN: Color Image

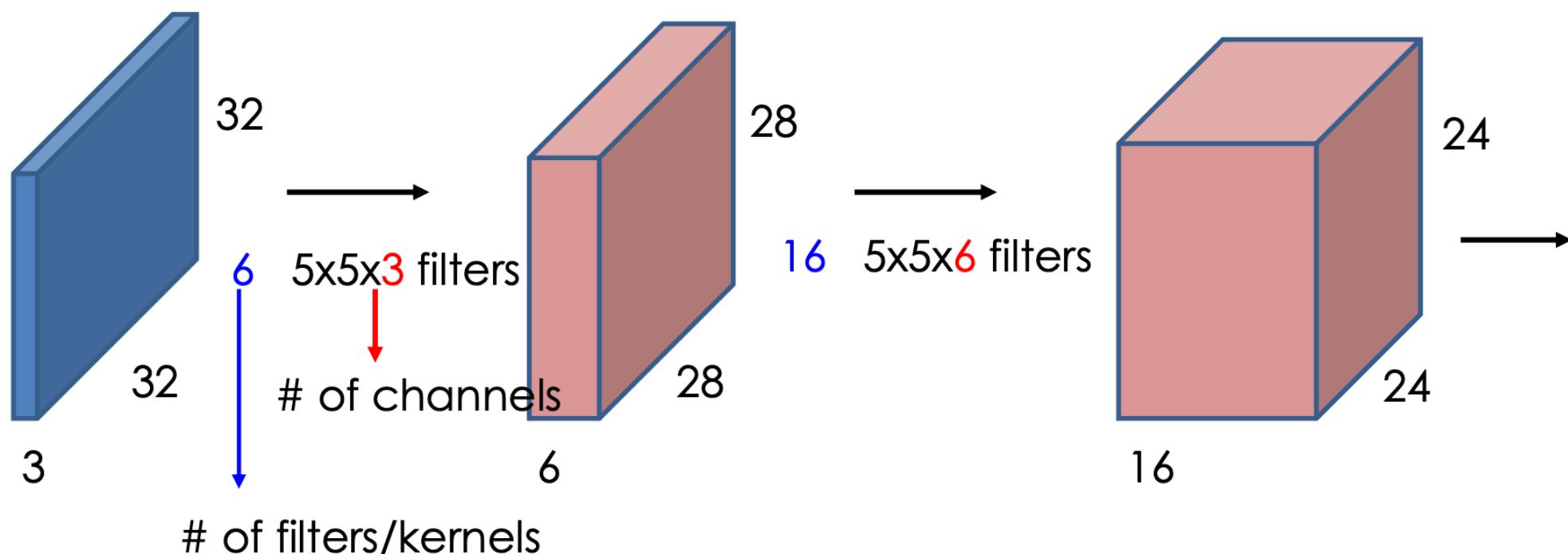


CNN: Color Image

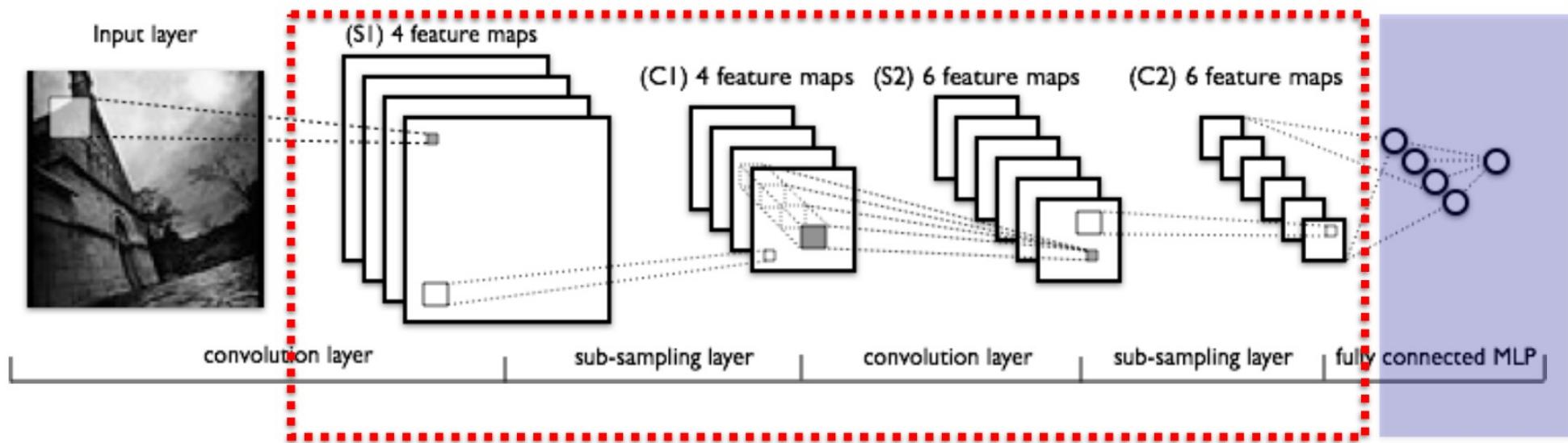


Convolutional Network

- Convolution network is a sequence of these layers

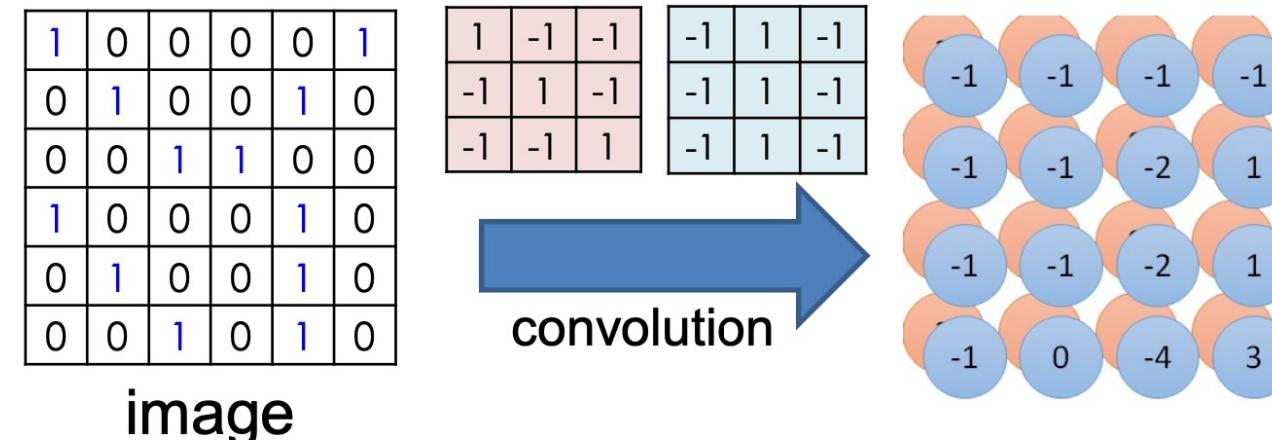


Convolutional Network

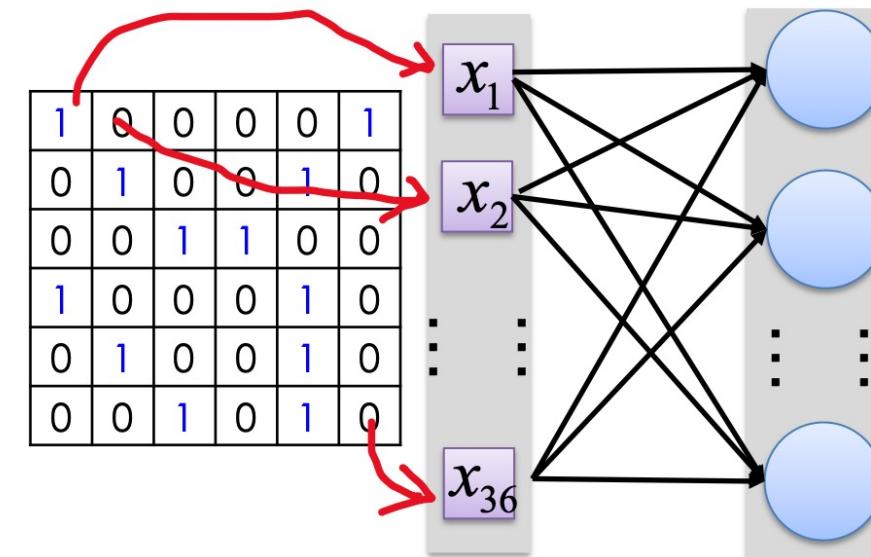


**Feature extraction using
convolutional layers**

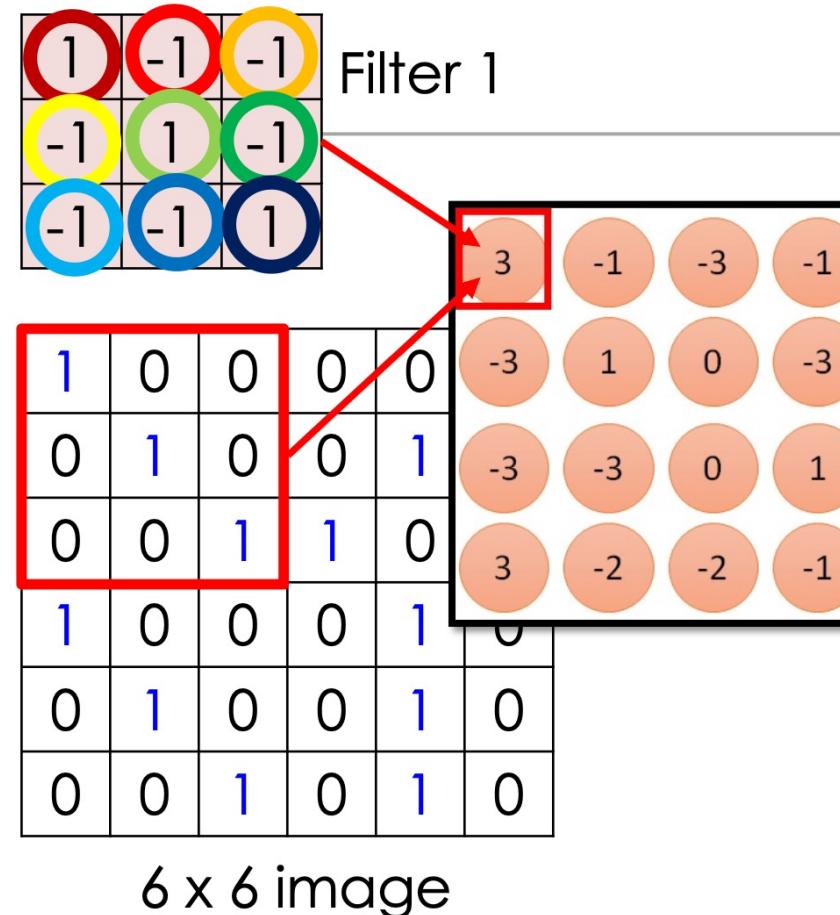
Convolution vs. Fully Connected



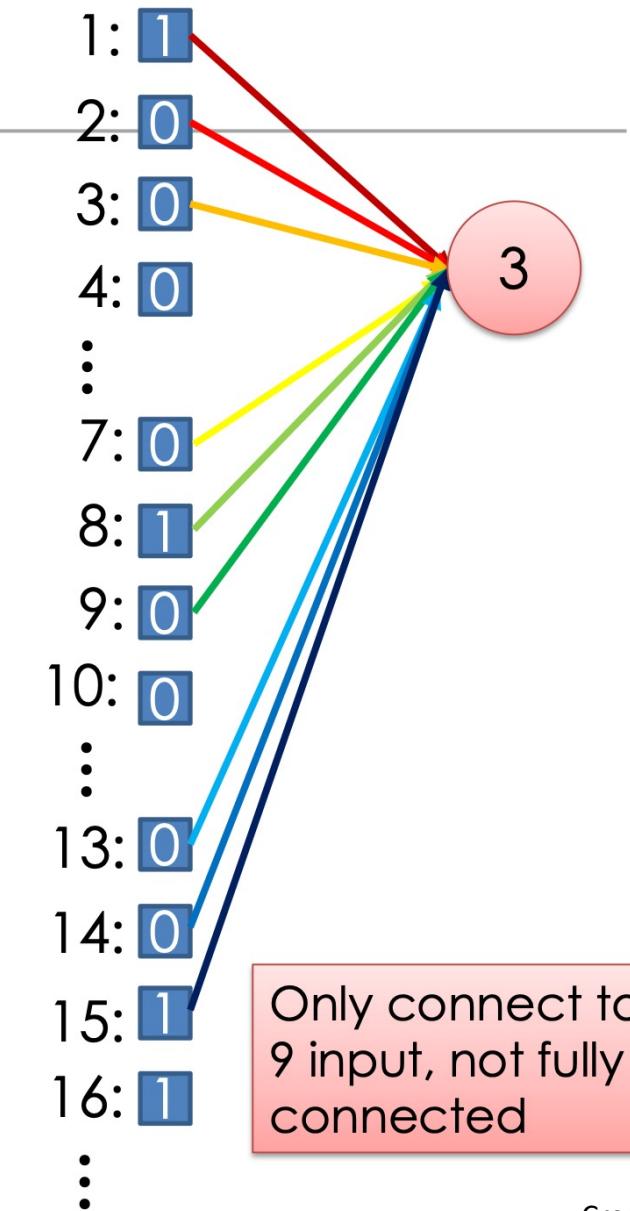
Fully-connected



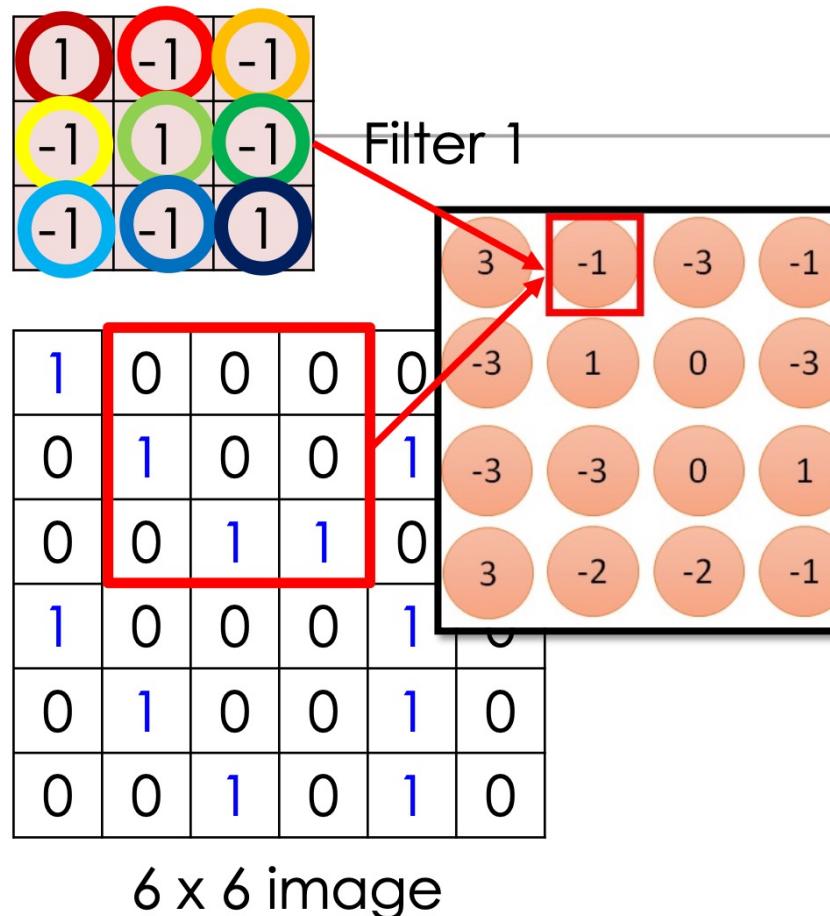
Convolution vs. Fully Connected



Less parameters!



Convolution vs. Fully Connected



Less parameters!

Even less parameters!

