

---

# CAP 5516

# Medical Image Computing

# (Spring 2022)

Dr. Chen Chen

Center for Research in Computer Vision (CRCV)

University of Central Florida

Office: HEC 221

Address: 4328 Scorpius St., Orlando, FL 32816-2365

Email: [chen.chen@crcv.ucf.edu](mailto:chen.chen@crcv.ucf.edu)

Web: <https://www.crcv.ucf.edu/chenchen/>

---

# Lecture 5: Introduction to Deep Learning (1)

---

## ARTIFICIAL INTELLIGENCE

A program that can sense, reason,  
act, and adapt

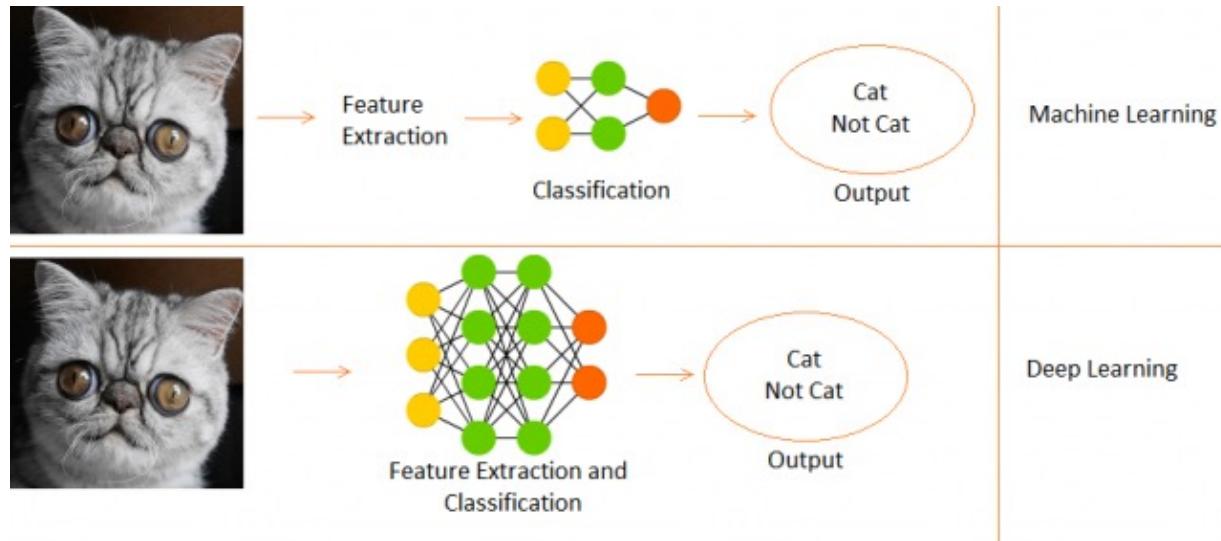
## MACHINE LEARNING

Algorithms whose performance improve  
as they are exposed to more data over time

## DEEP LEARNING

Subset of machine learning in  
which multilayered neural  
networks learn from  
vast amounts of data

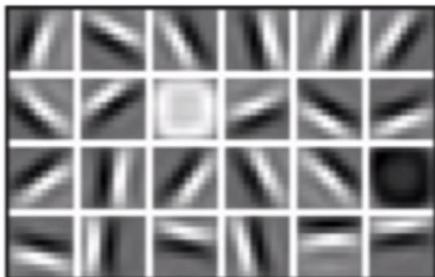
# Why Deep Learning?



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



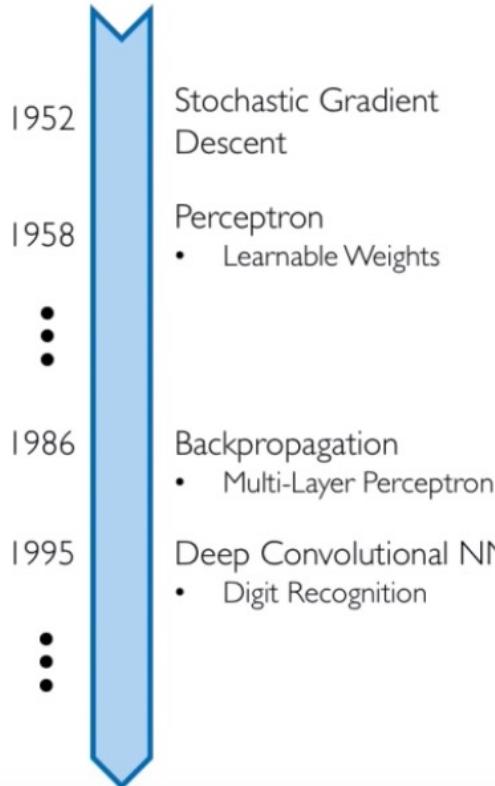
Eyes & Nose & Ears

High Level Features



Facial Structure

# Why Now?



Neural Networks date back decades, so why the resurgence?

## I. Big Data

- Larger Datasets
- Easier Collection & Storage



## 2. Hardware

- Graphics Processing Units (GPUs)
- Massively Parallelizable



## 3. Software

- Improved Techniques
- New Models
- Toolboxes



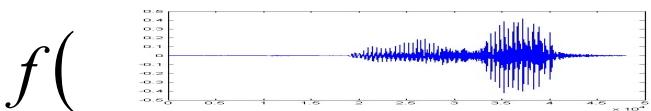
## 4. Open resource (GitHub)

Credit: Lex Fridman

# Machine Learning ≈ Looking for a Function

---

- Speech Recognition



) = “How are you”

- Image Recognition

$$f\left( \begin{array}{c} \text{[Image icon]} \\ \vdots \\ \text{[Image icon]} \end{array} \right)$$



) = “Cat”

- Playing Go

$$f\left( \begin{array}{c} \text{[Image icon]} \\ \vdots \\ \text{[Image icon]} \end{array} \right)$$



) = “5-5” (next move)

- Dialogue System

$$f\left( \begin{array}{c} \text{[Text icon]} \\ \vdots \\ \text{[Text icon]} \end{array} \right)$$

“Hi”

(what the user said)

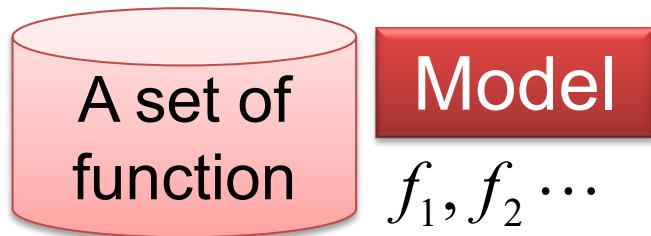
) = “Hello”

(system response)

# Framework

## Image Recognition:

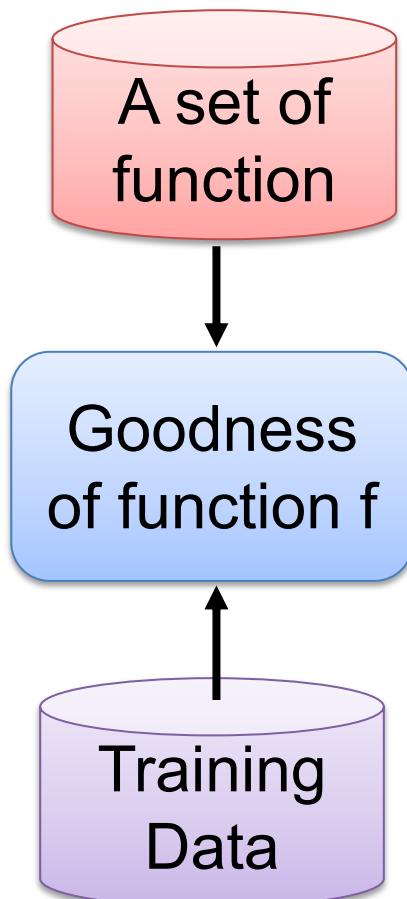
$$f(\text{) = "cat"}$$



$$f_1(\text{}) = "cat" \quad f_2(\text{}) = "money"$$

$$f_1(\text{}) = "dog" \quad f_2(\text{}) = "snake"$$

# Framework



## Image Recognition:

$f($    $) = \text{"cat"}$

$f_1($    $) = \text{"cat"}$        $f_2($    $) = \text{"money"}$   
**Better!**  
 $f_1($    $) = \text{"dog"}$        $f_2($    $) = \text{"snake"}$

## Supervised Learning

function input:

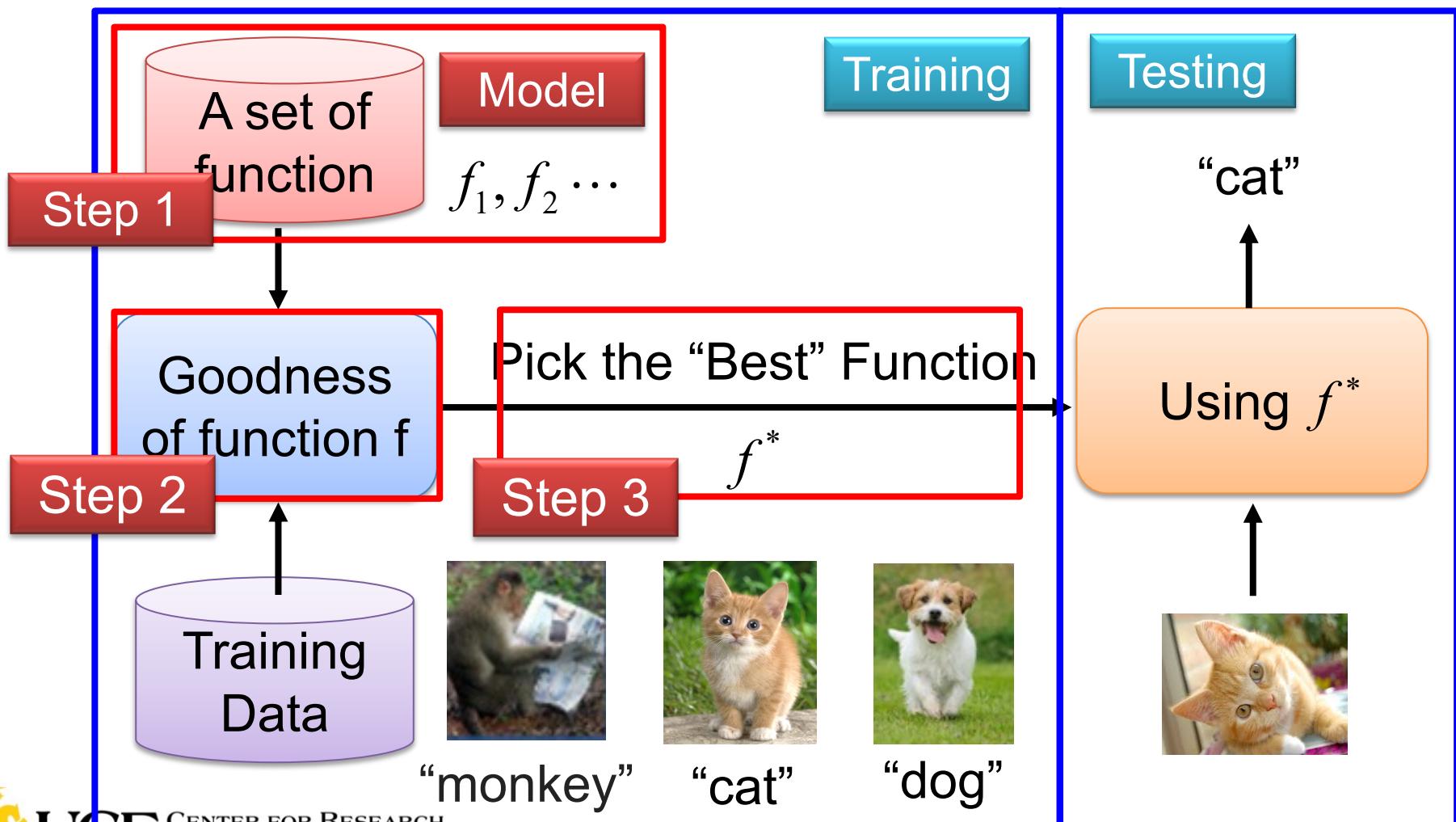


function output: "monkey"    "cat"    "dog"

# Framework

## Image Recognition:

$$f(\text{cat image}) = \text{"cat"}$$



# Three Steps for Deep Learning

---

Step 1: define a set of function

Neural Network



Step 2: goodness of function



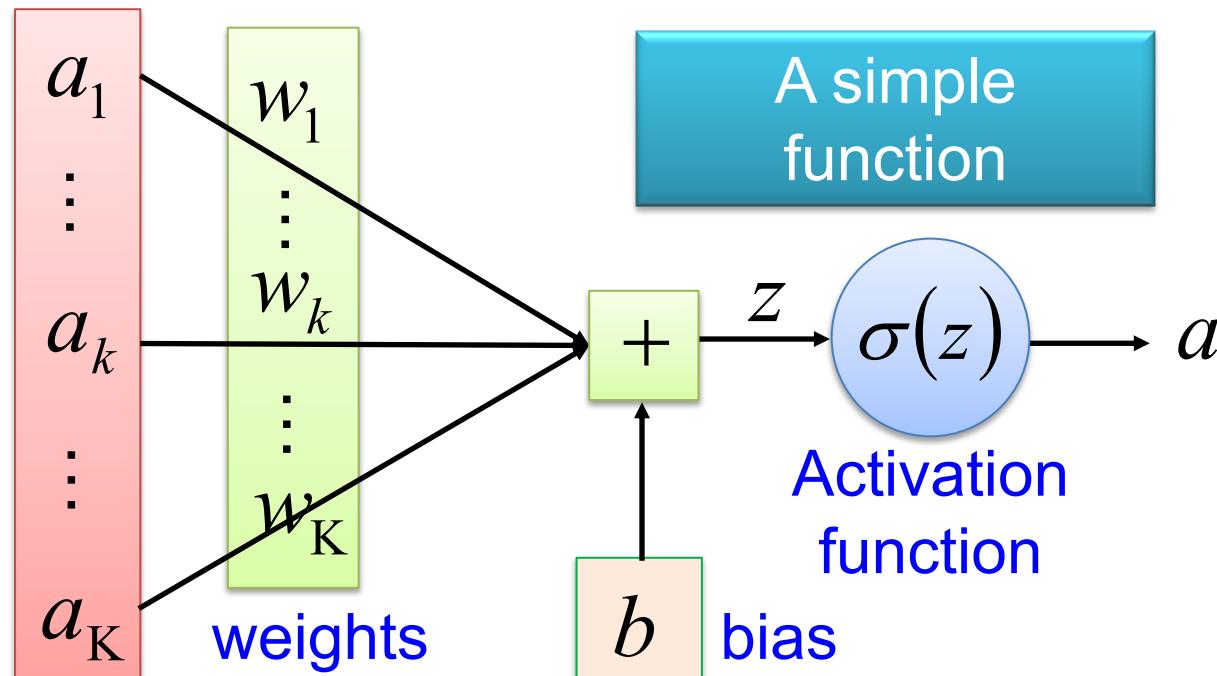
Step 3: pick the best function

# Neural Network

---

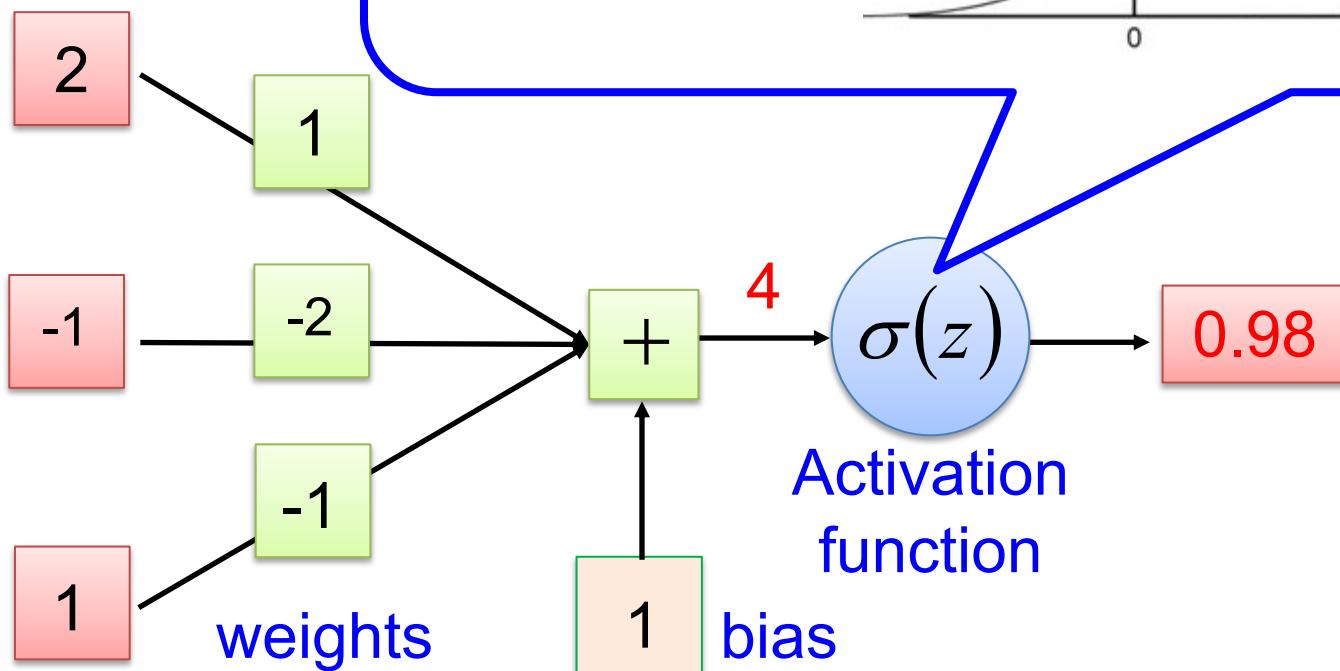
## Neuron

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



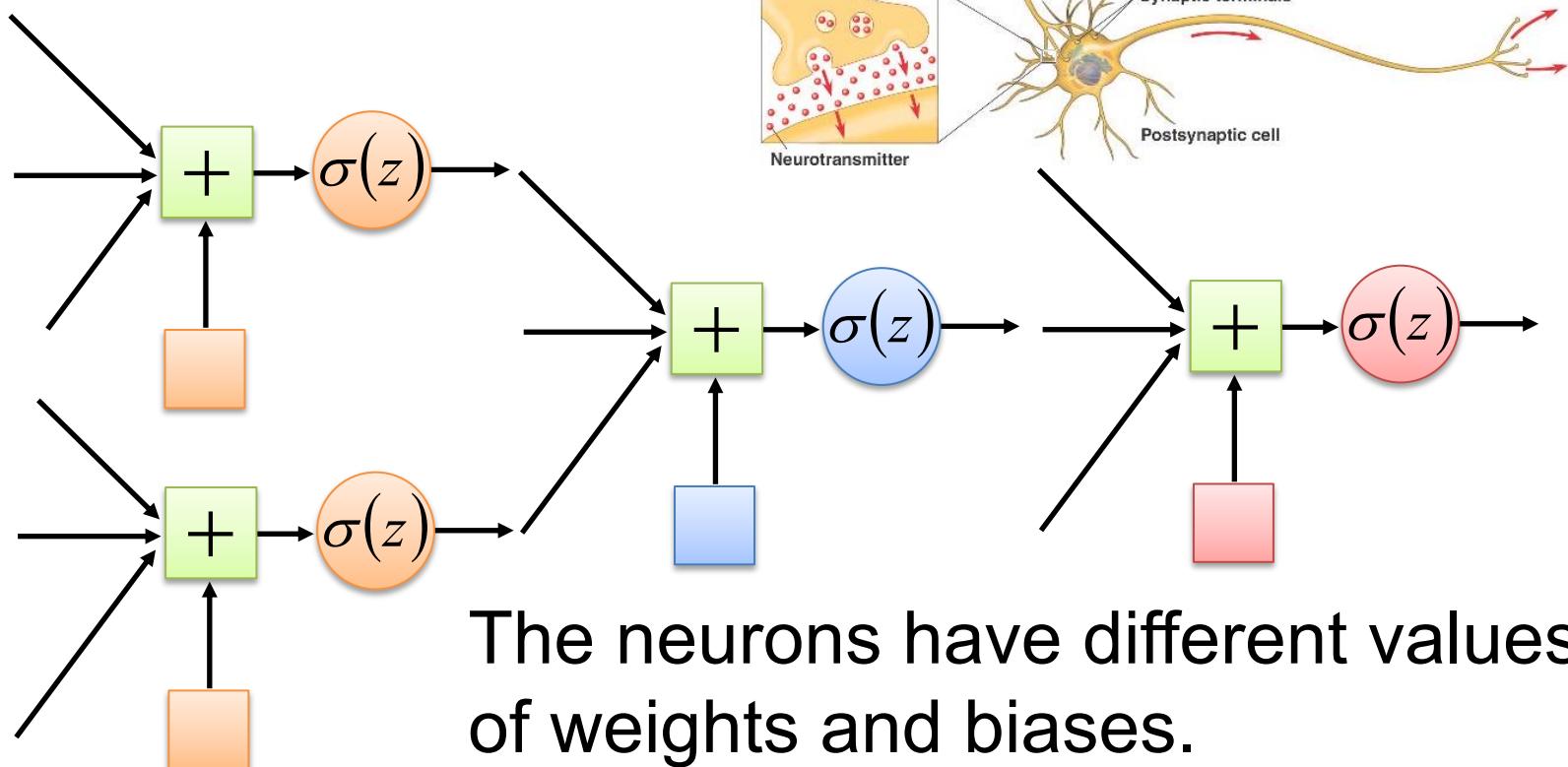
# Neural Network

## Neuron

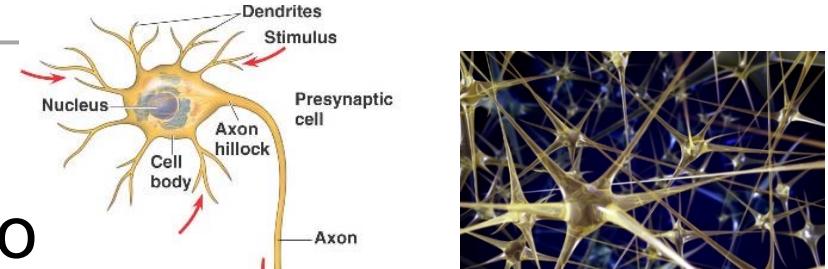


# Neural Network

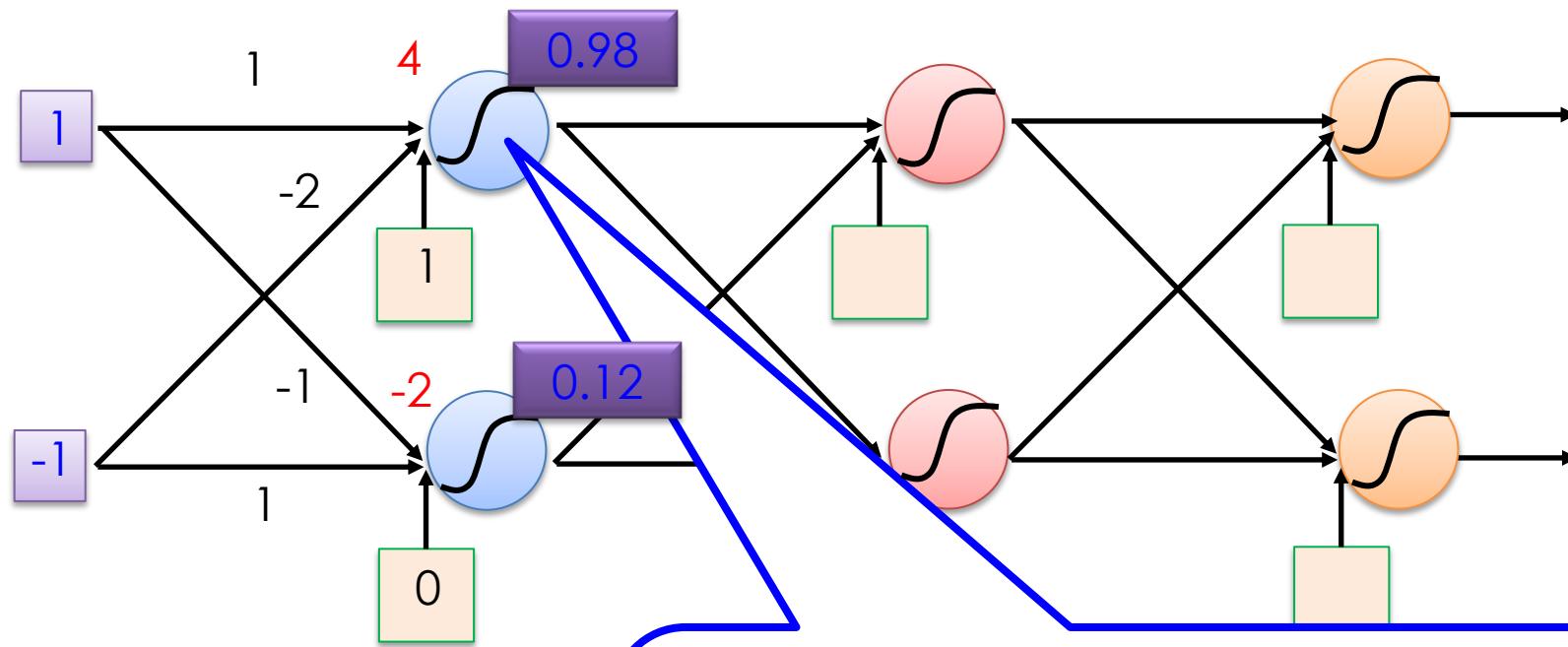
Different connections lead to different network structures



Weights and biases are network parameters  $\theta$

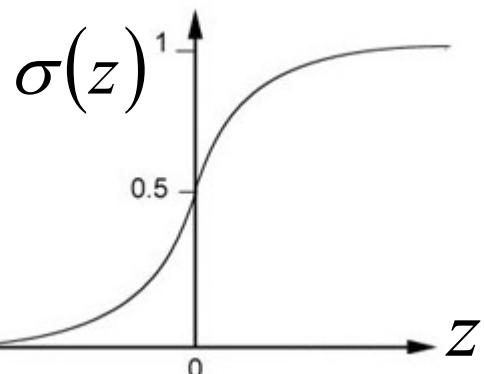


# Fully Connected Feedforward Network

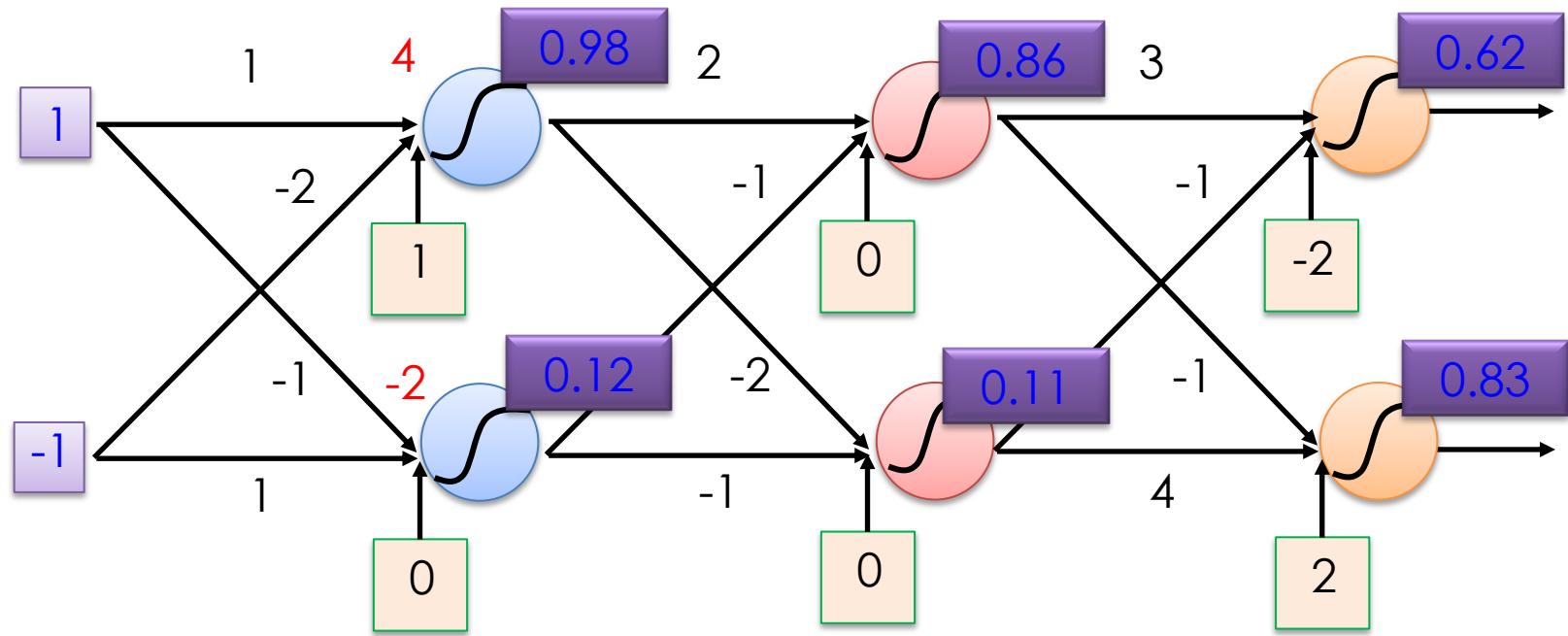


Sigmoid Function

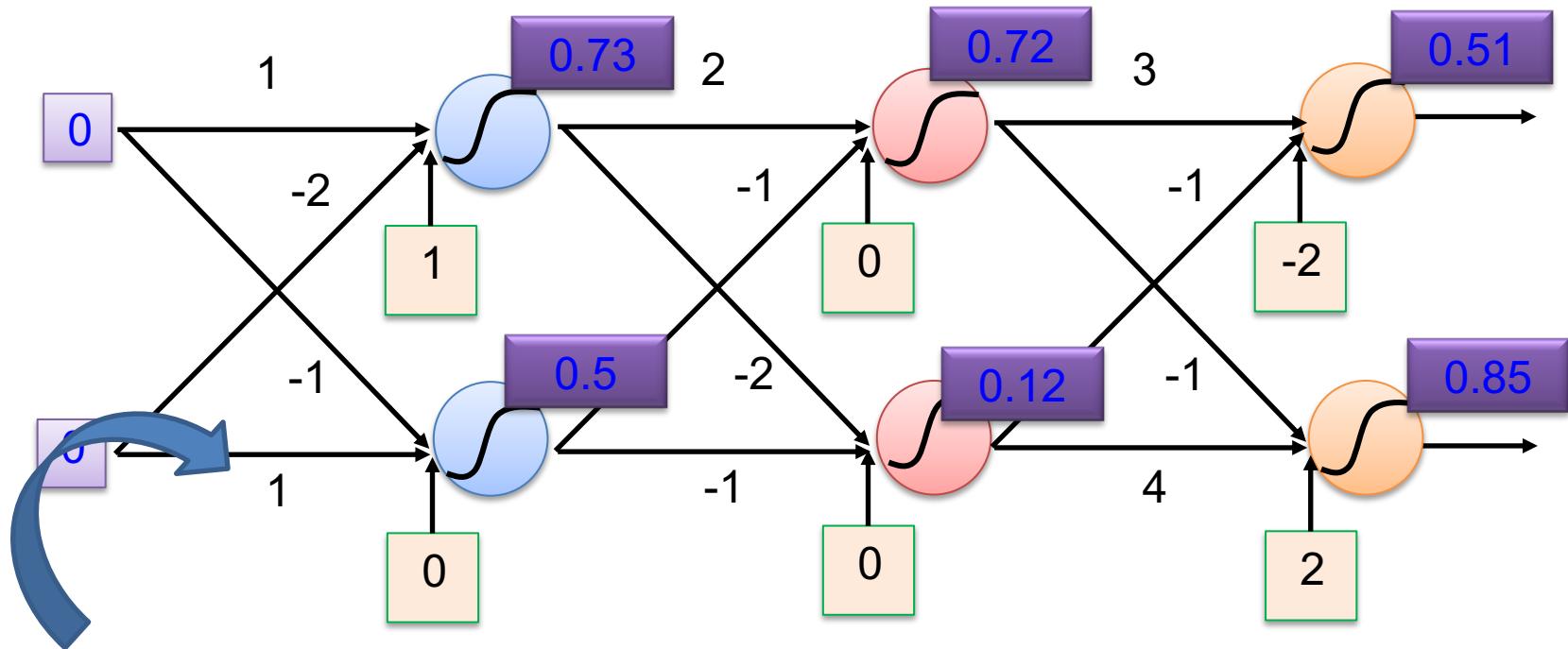
$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



# Fully Connected Feedforward Network



# Fully Connected Feedforward Network



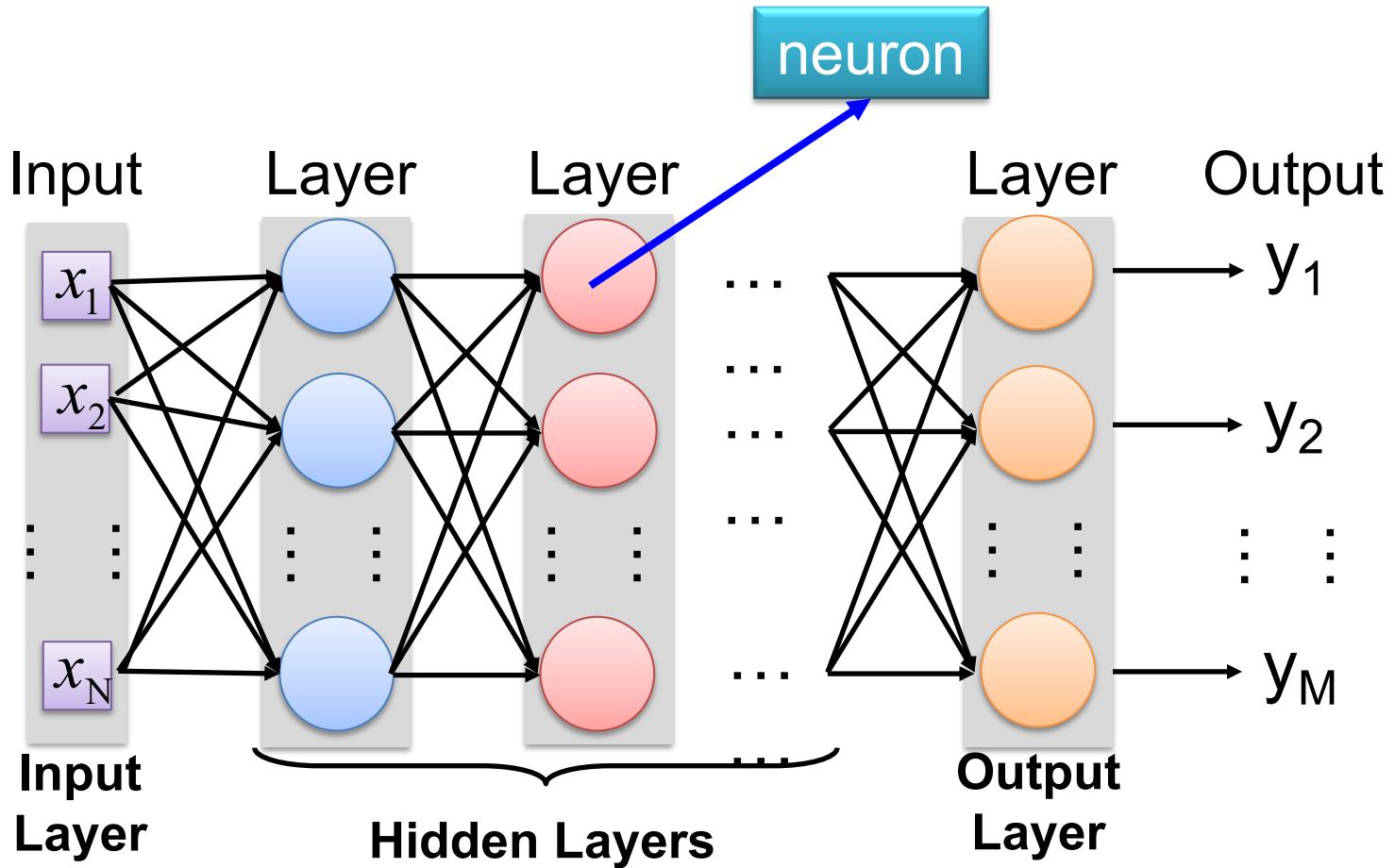
This is a function.  
Input vector, output vector

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

Given parameters  $\theta$ , define a function

Given network structure, define a function set

# Fully Connected Feedforward Network

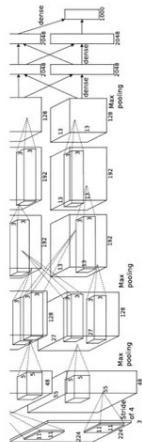


Deep means many hidden layers

# Deep = Many hidden layers

[http://cs231n.stanford.edu/slides/winter1516\\_lecture8.pdf](http://cs231n.stanford.edu/slides/winter1516_lecture8.pdf)

8 layers

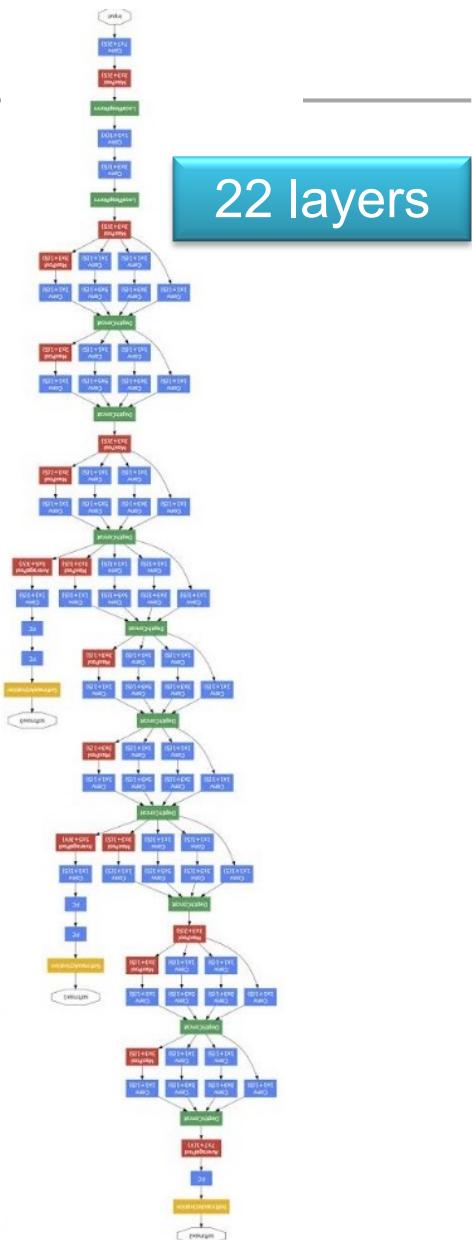


**AlexNet (2012)**  
UCF CENTER FOR RESEARCH  
IN COMPUTER VISION



**VGG (2014)**

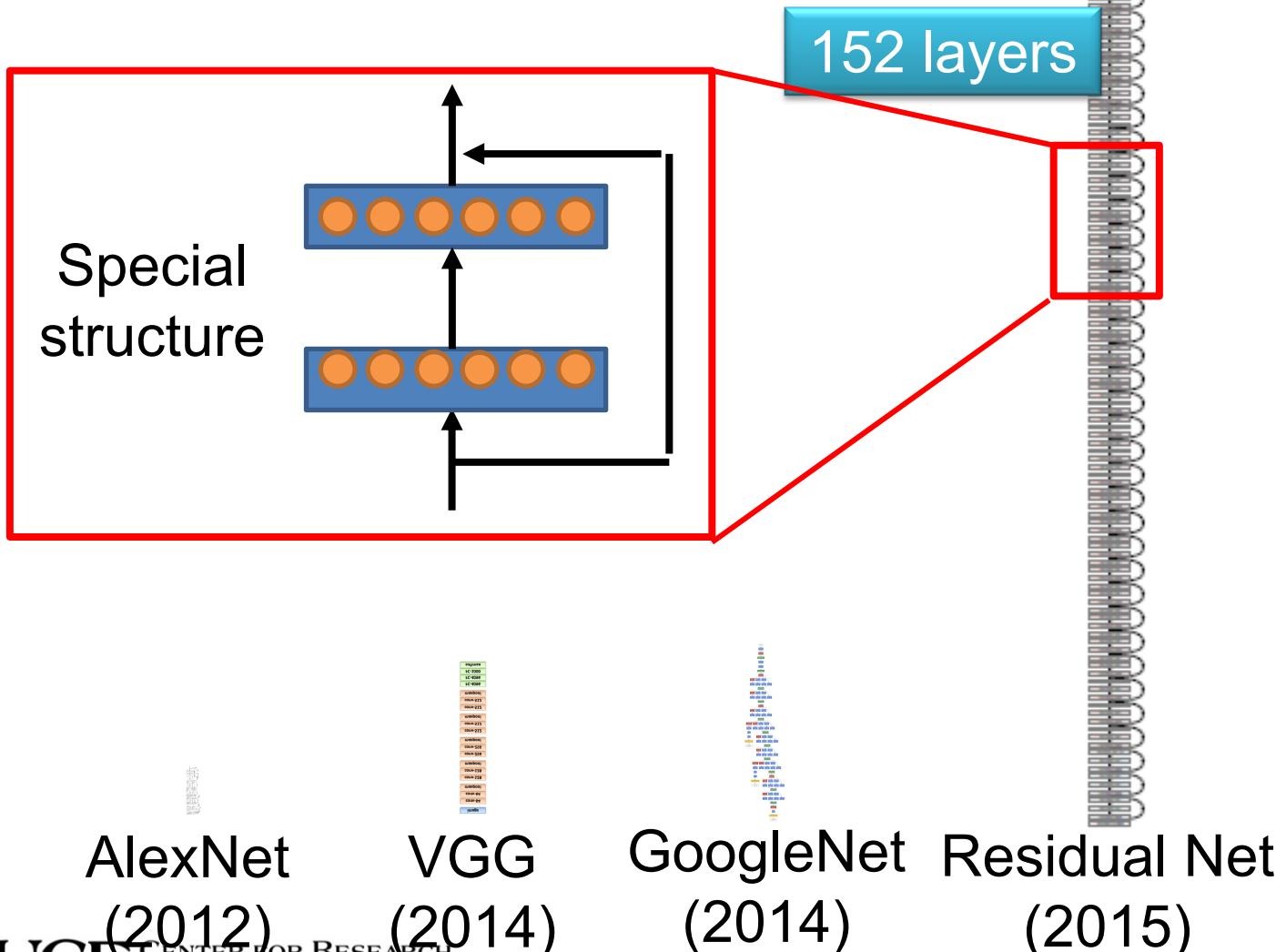
19 layers



**GoogleNet (2014)**

22 layers

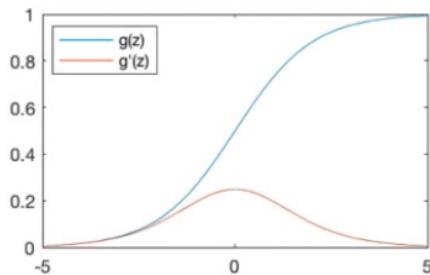
# Deep = Many hidden layers



# Activation Function

## Common Activation Functions

Sigmoid Function

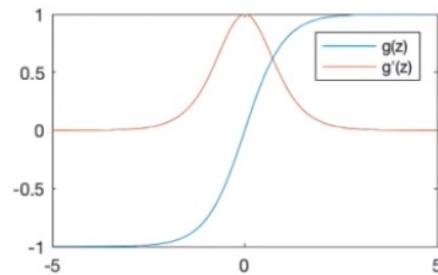


$$g(z) = \frac{1}{1 + e^{-z}}$$

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

`tf.nn.sigmoid(z)`

Hyperbolic Tangent

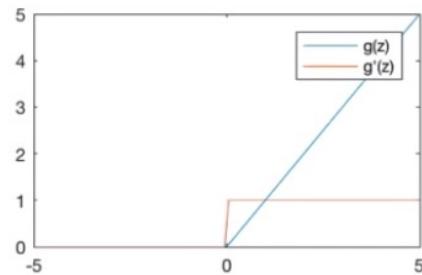


$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

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

`tf.nn.tanh(z)`

Rectified Linear Unit (ReLU)



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

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

`tf.nn.relu(z)`

NOTE: All activation functions are non-linear

[https://ml-cheatsheet.readthedocs.io/en/latest/activation\\_functions.html](https://ml-cheatsheet.readthedocs.io/en/latest/activation_functions.html)

More details about the activation functions: <https://www.v7labs.com/blog/neural-networks-activation-functions>

# Activation Function

---

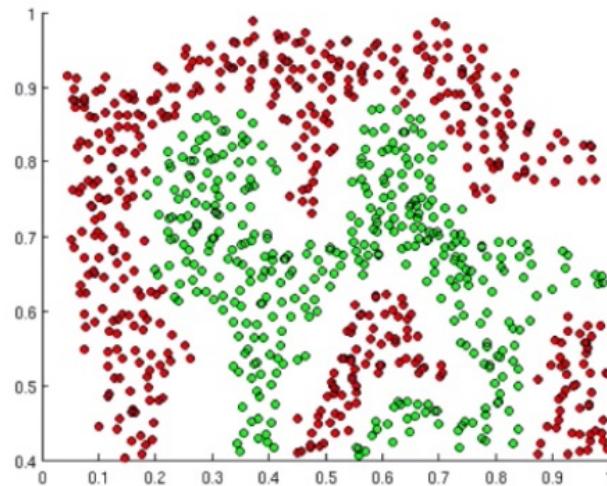
- Why non-linear activation function?

# Activation Function

---

- Why non-linear activation function?

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



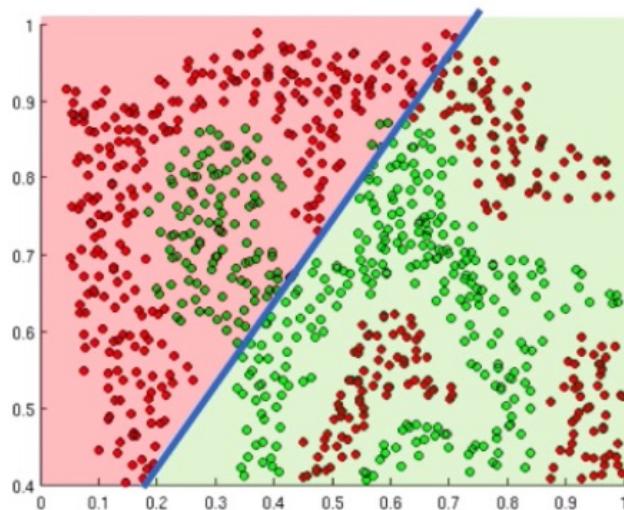
What if we wanted to build a Neural Network to  
distinguish green vs red points?

# Activation Function

---

- Why non-linear activation function?

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

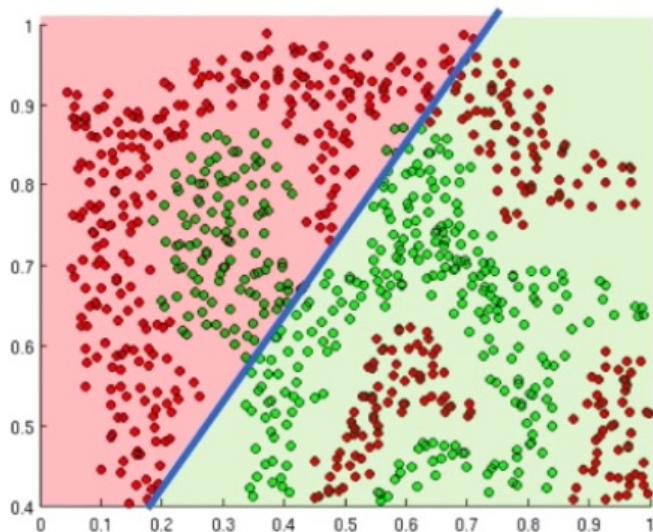


Linear Activation functions produce linear decisions no matter the network size

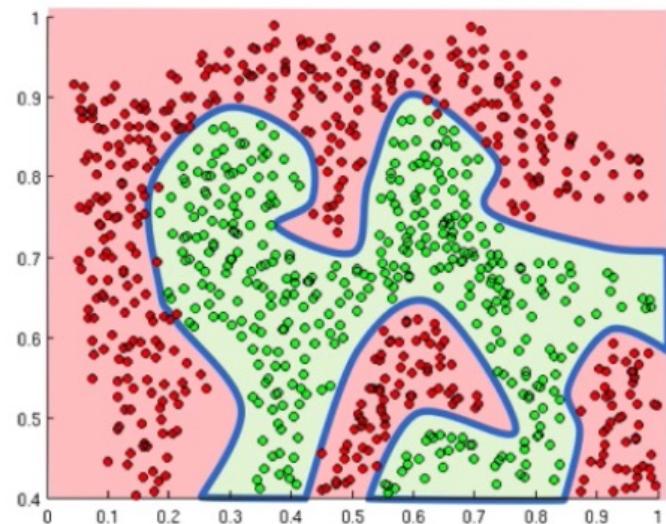
# Activation Function

- 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

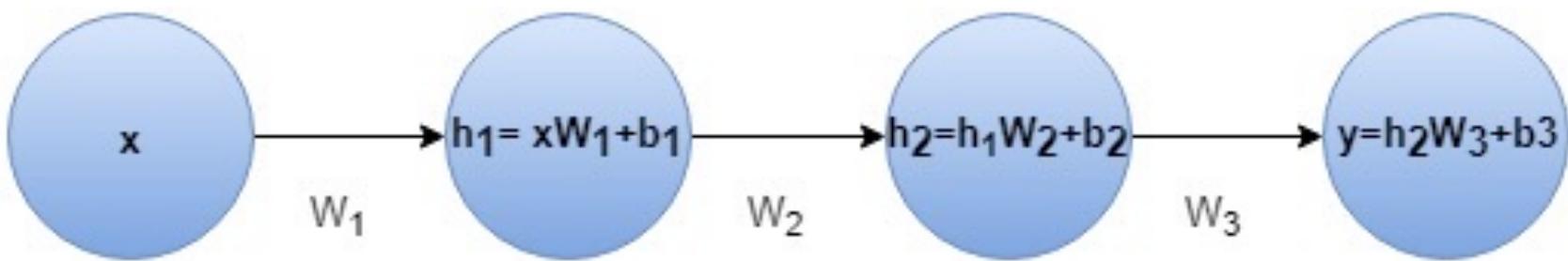
# Activation Function

---

- Why non-linear activation function?
- **All layers of the neural network collapse into one**—with linear activation functions, no matter how many layers in the neural network, the last layer will be a linear function of the first layer (because a linear combination of linear functions is still a linear function). So a linear activation function turns the neural network into just one layer.
- A neural network 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.

# Activation Function

- Why non-linear activation function?



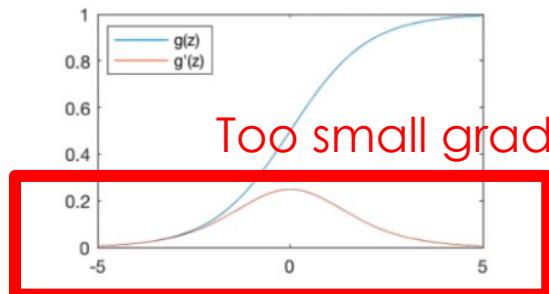
$$\begin{aligned} y &= h2 * w3 + b3 \\ &= (h1 * w2 + b2) * w3 + b3 \\ &= h1 * w2 * w3 + b2 * w3 + b3 \\ &= (x * w1 + b1) * w2 * w3 + b2 * w3 + b3 \\ &= x * w1 * w2 * w3 + b1 * w2 * w3 + b2 * w3 + b3 \\ &= x * w' + b' \end{aligned}$$

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

# Activation Function

## Common Activation Functions

Sigmoid Function

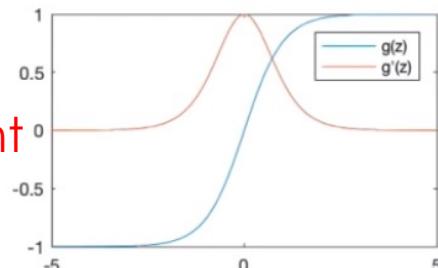


$$g(z) = \frac{1}{1 + e^{-z}}$$

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

 `tf.nn.sigmoid(z)`

Hyperbolic Tangent

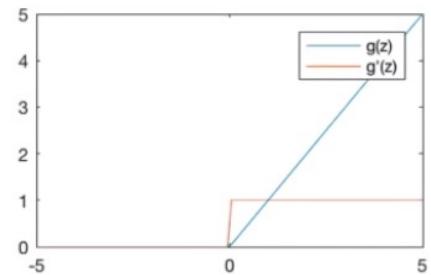


$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

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

 `tf.nn.tanh(z)`

Rectified Linear Unit (ReLU)



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

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

 `tf.nn.relu(z)`

NOTE: All activation functions are non-linear

# Activation Function

- Activation function ReLU (rectified linear unit)
  - $\text{ReLU}(z) = \max\{z, 0\}$

Gradient 1

Gradient 0

The Rectified Linear Activation Function

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

0

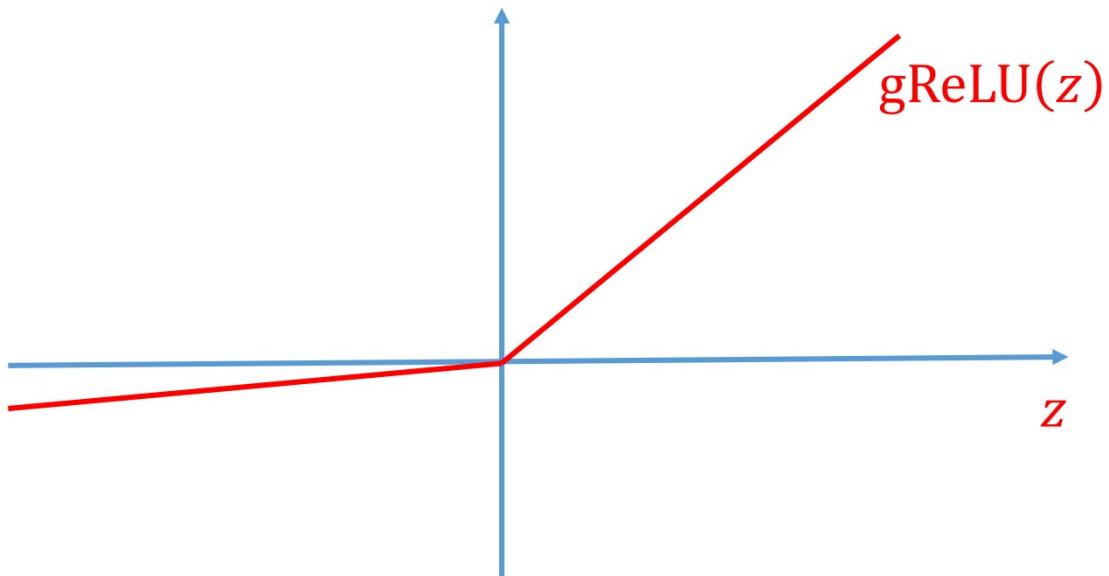
0

$z$

# Activation Function

---

- Generalizations of ReLU  $g\text{ReLU}(z) = \max\{z, 0\} + \alpha \min\{z, 0\}$ 
  - Leaky-ReLU( $z$ ) =  $\max\{z, 0\} + 0.01 \min\{z, 0\}$
  - Parametric-ReLU( $z$ ):  $\alpha$  learnable



# Output Layer

---

- Softmax layer as the output layer

## Ordinary Layer

$$z_1 \rightarrow \sigma \rightarrow y_1 = \sigma(z_1)$$

$$z_2 \rightarrow \sigma \rightarrow y_2 = \sigma(z_2)$$

$$z_3 \rightarrow \sigma \rightarrow y_3 = \sigma(z_3)$$

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

May not be easy to interpret

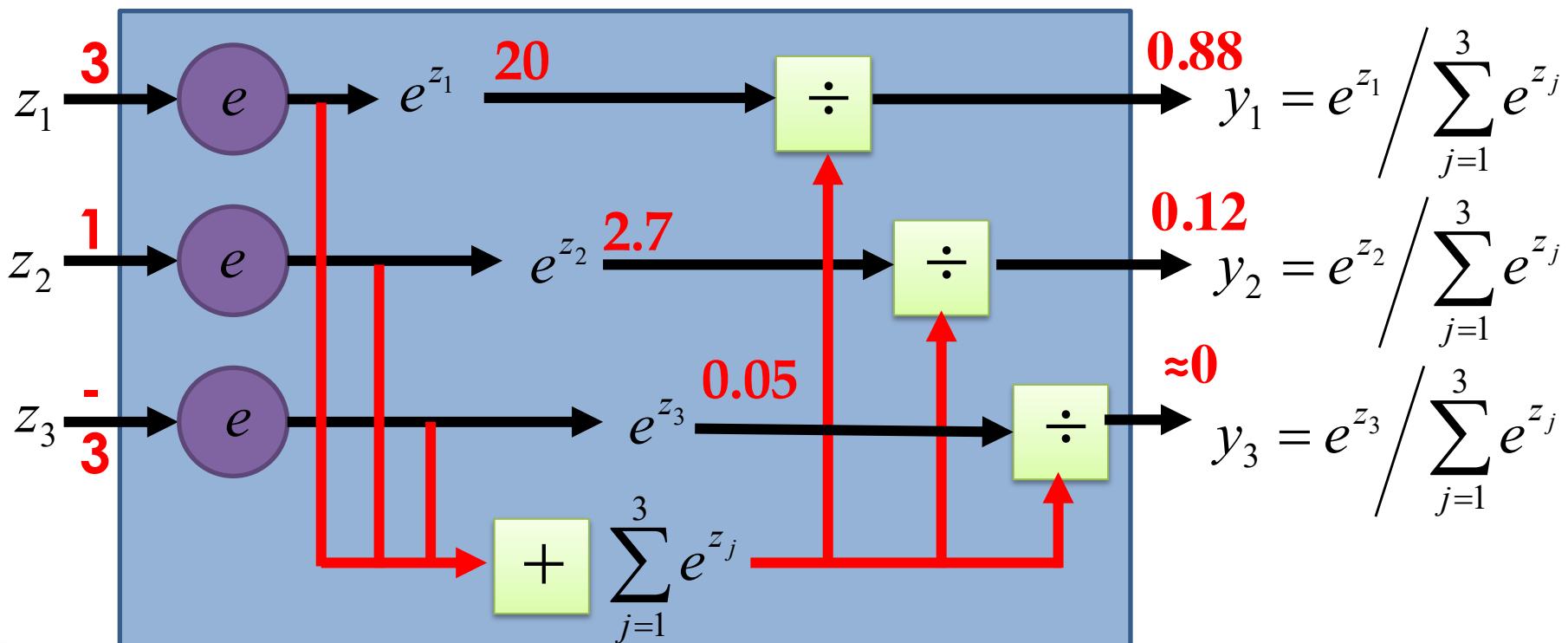
# Output Layer

- Softmax layer as the output layer

Probability:

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

## Softmax Layer



# Three Steps for Deep Learning

---

Step 1: define a set of function



Step 2: goodness of function

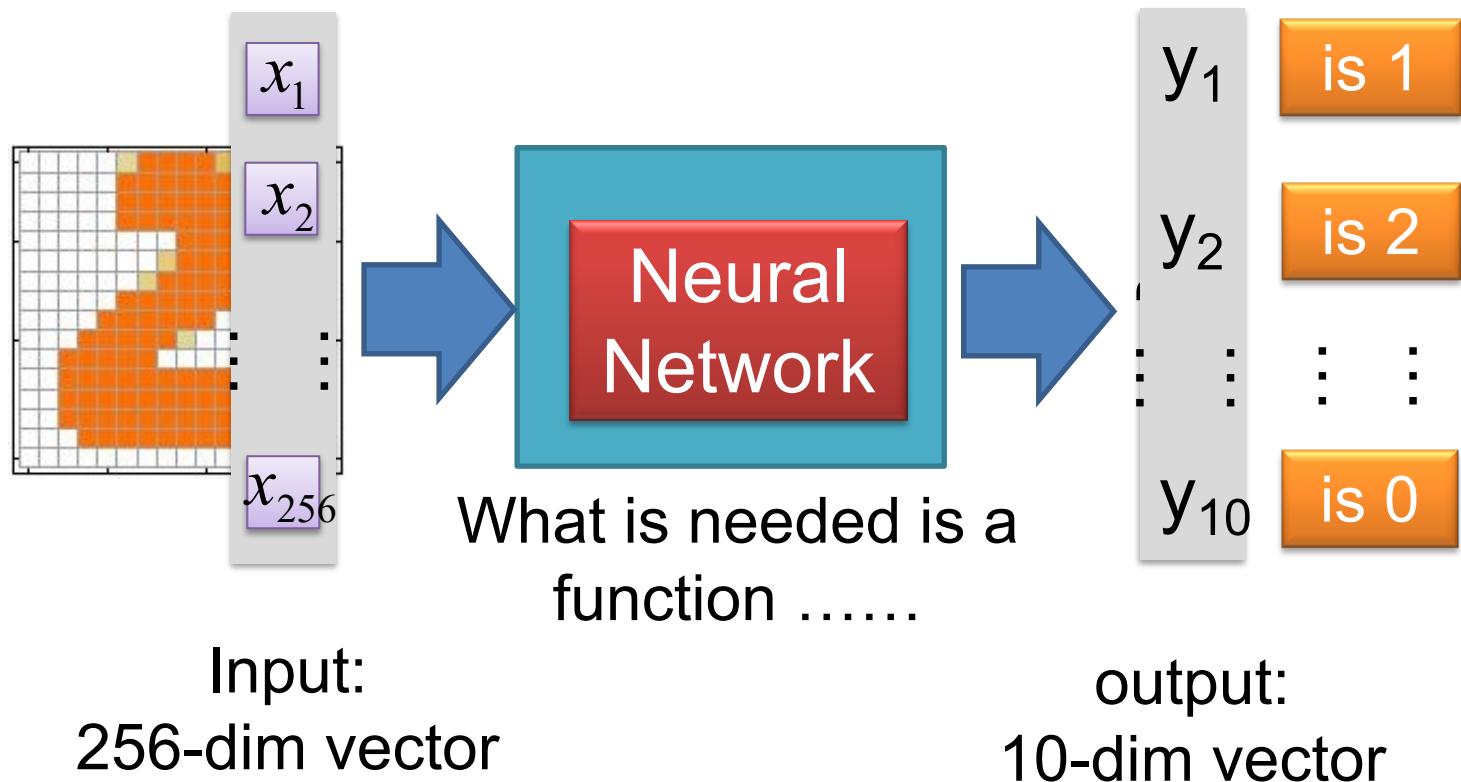


Step 3: pick the best function

# Example Application

---

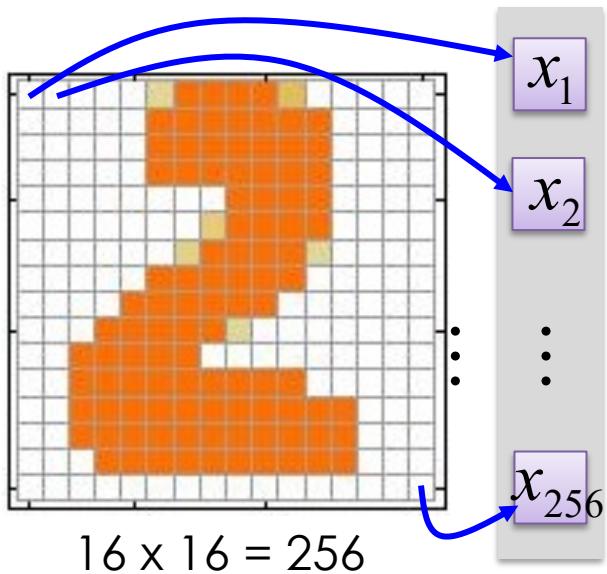
- Handwriting Digit Recognition



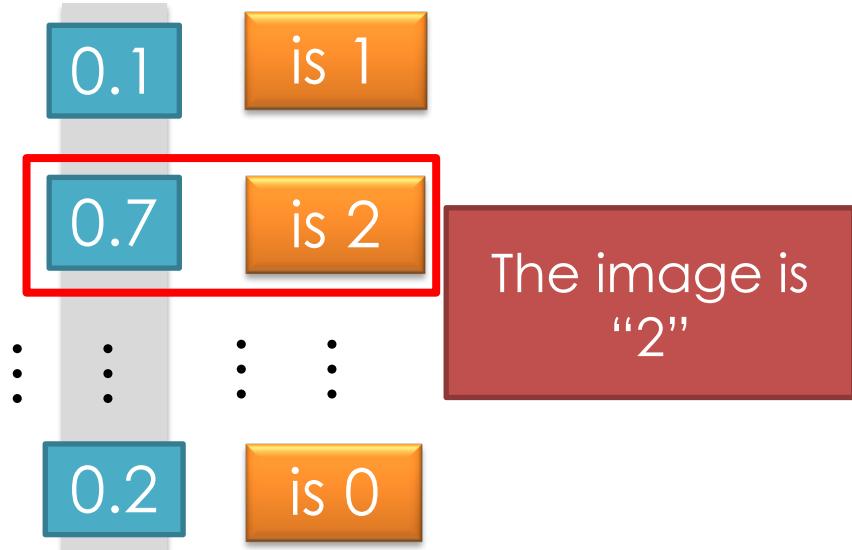
# Example Application



## Input



## Output



Each dimension represents the confidence of a digit.

# Training Data

---

- Preparing training data: images and their labels



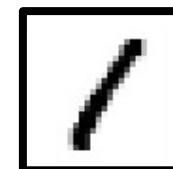
“5”



“0”



“4”



“1”



“9”



“2”



“1”

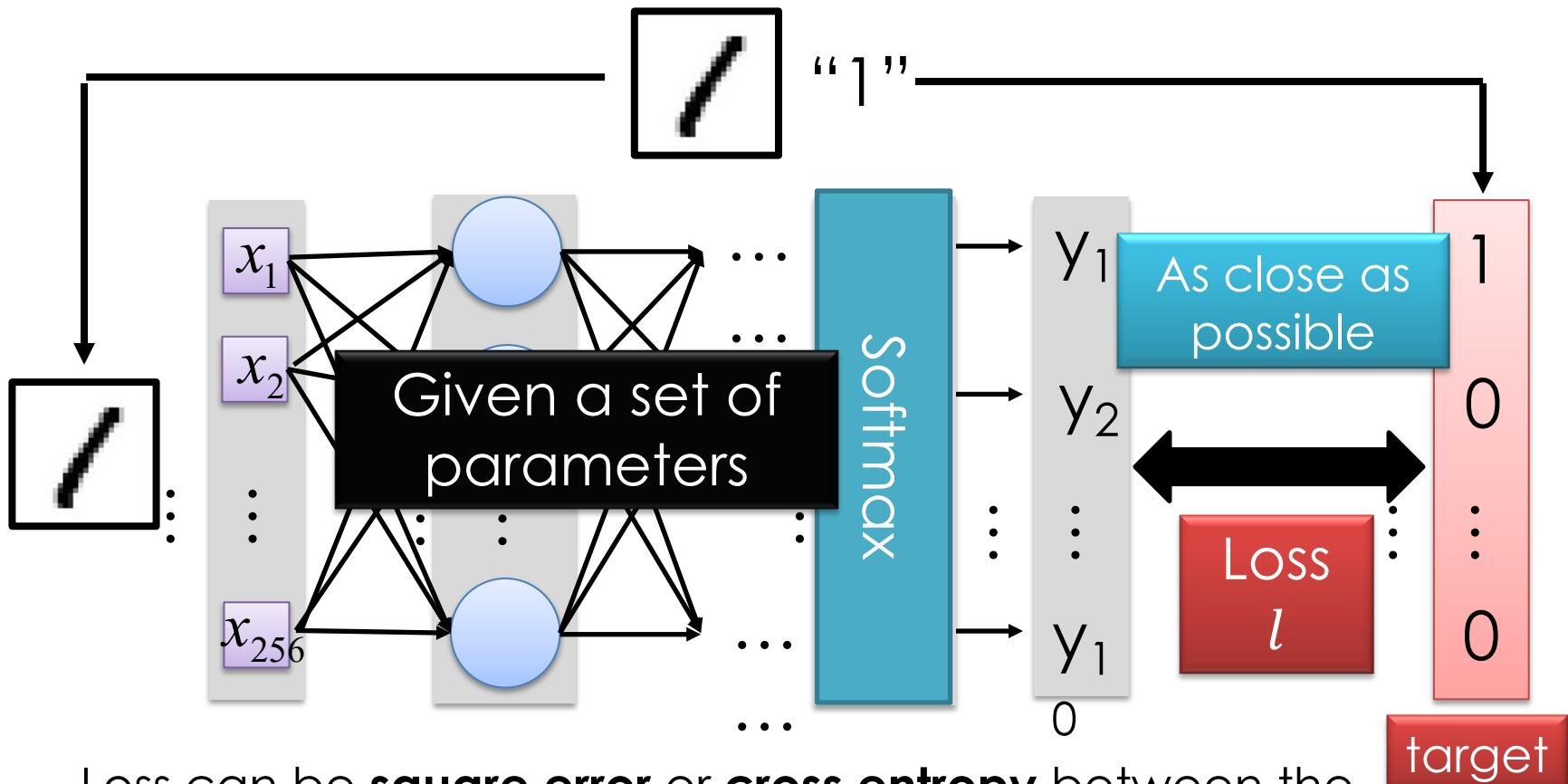


“3”

The learning target is defined on the  
training data.

# Loss

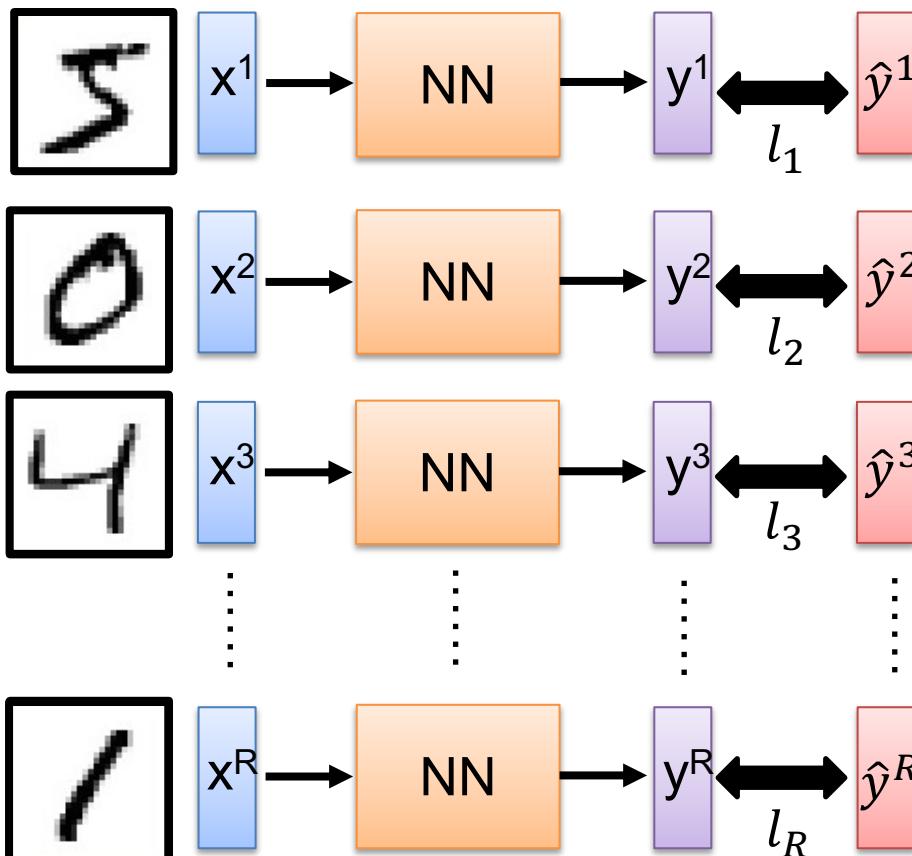
A good function should make the loss of all examples as small as possible.



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  $\theta^*$  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 Image

- Some patterns are much smaller than the whole image

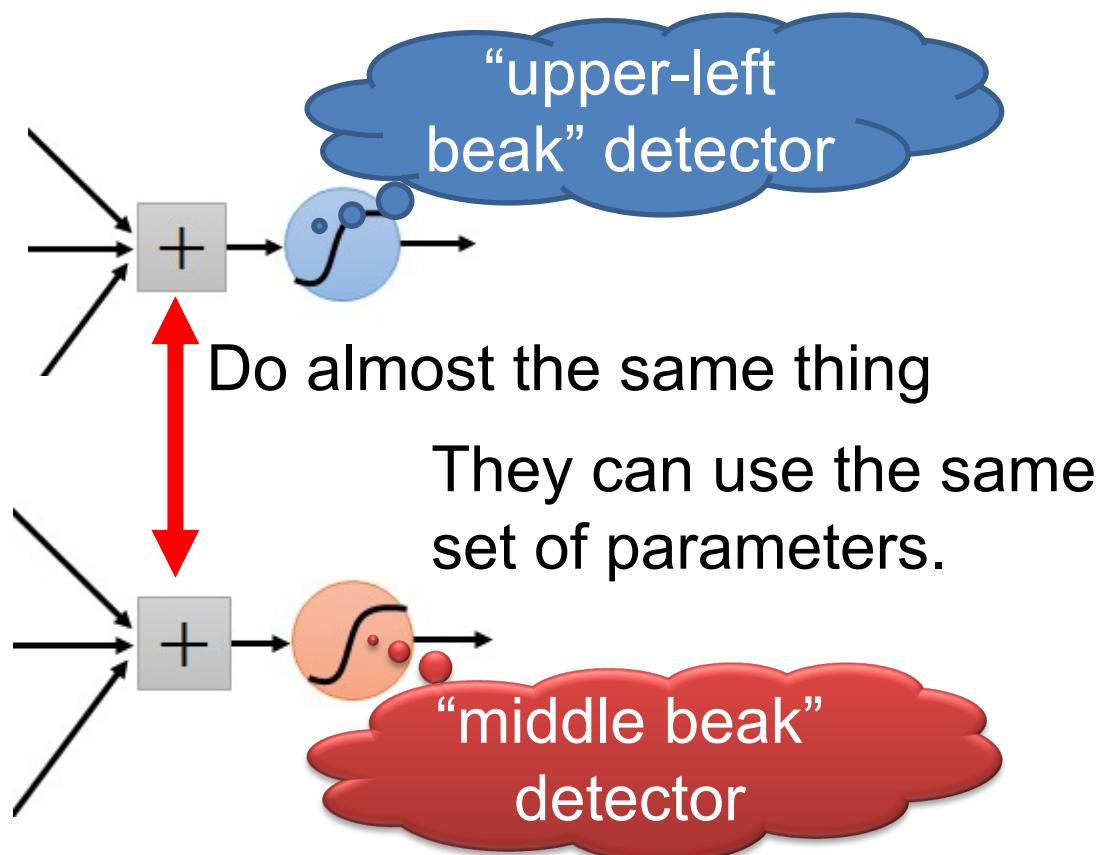
A neuron does not have to see the whole image to discover the pattern.

Connecting to small region with less parameters



# Why CNN for Image

- The same patterns appear in different regions.



# Why CNN for Image

---

- 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

# CNN for Image

---

## Property 1

- Some patterns are much smaller than the whole

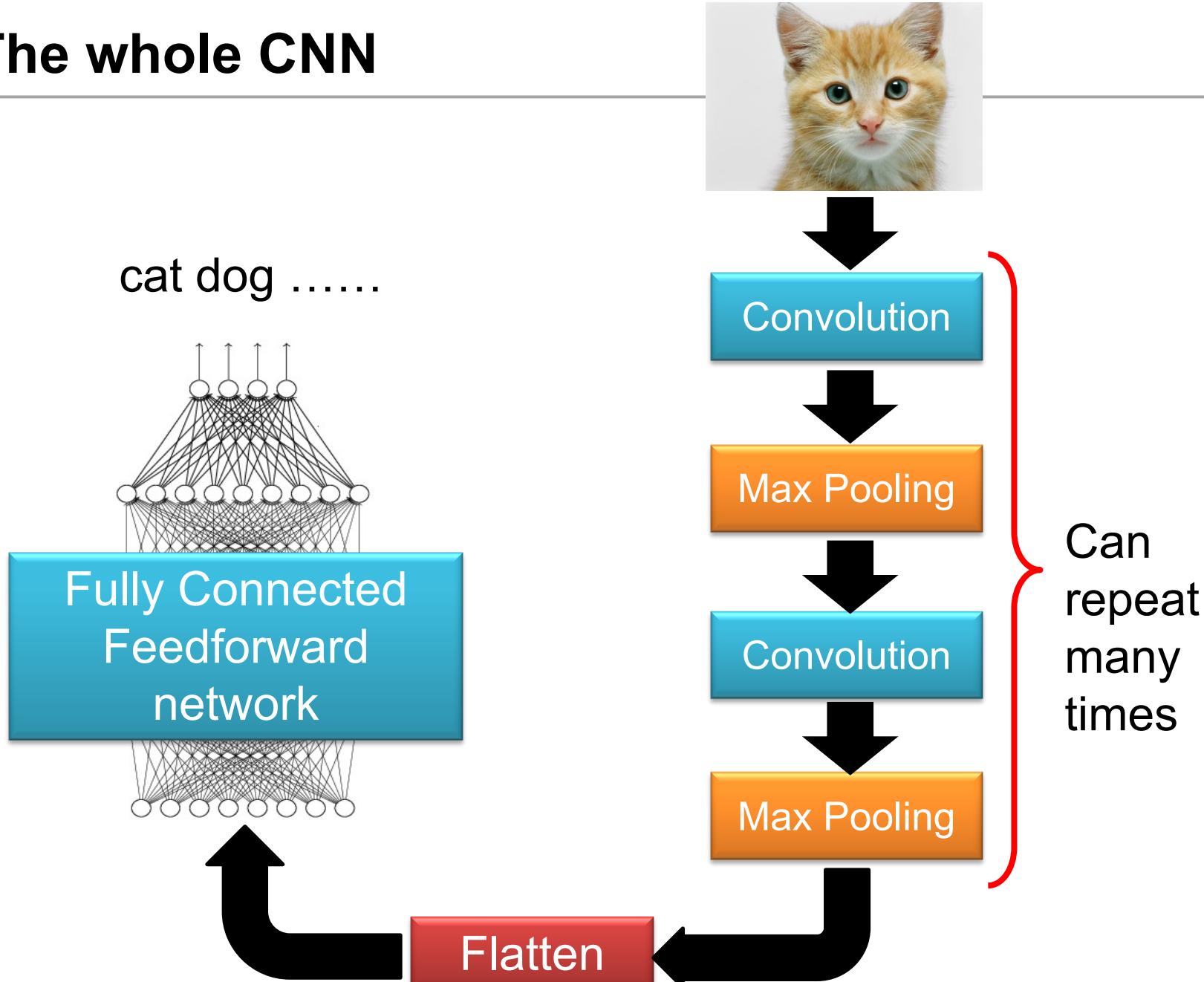
## Property 2

- The same patterns appear in different regions.

## Property 3

- Subsampling the pixels will not change the object

# The whole CNN



# The whole CNN



## Property 1

- Some patterns are much smaller than the whole image

## Property 2

- The same patterns appear in different regions.

## Property 3

- Subsampling the pixels will not change the object

Convolution

Max Pooling

Convolution

Max Pooling

Flatten

Can  
repeat  
many  
times



# CNN – Convolution

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

Those are the network parameters to be learned.

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).

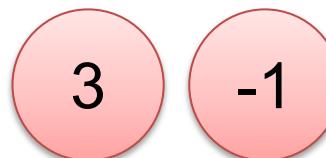
# CNN – Convolution

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

Filter/Kernel 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



6 x 6 image

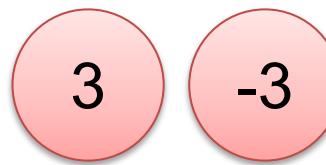
# CNN – Convolution

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

Filter 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

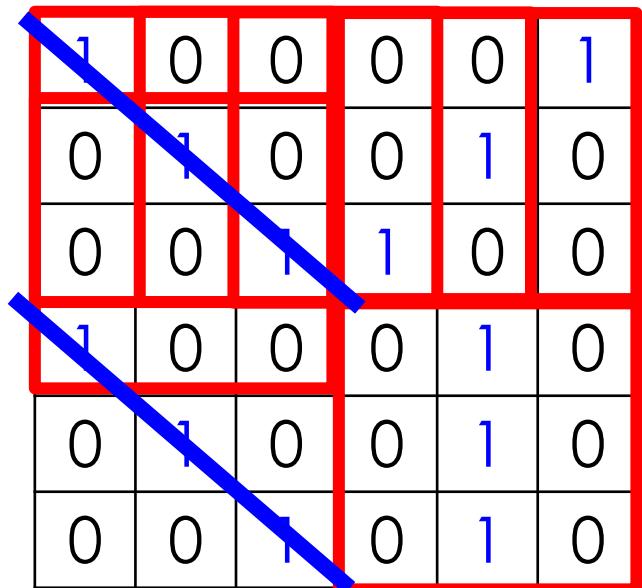


We set stride=1 below

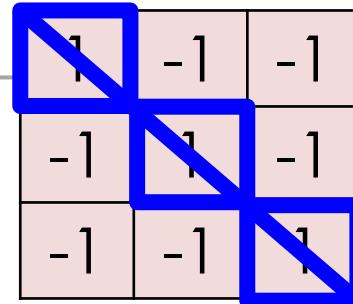
6 x 6 image

# CNN – Convolution

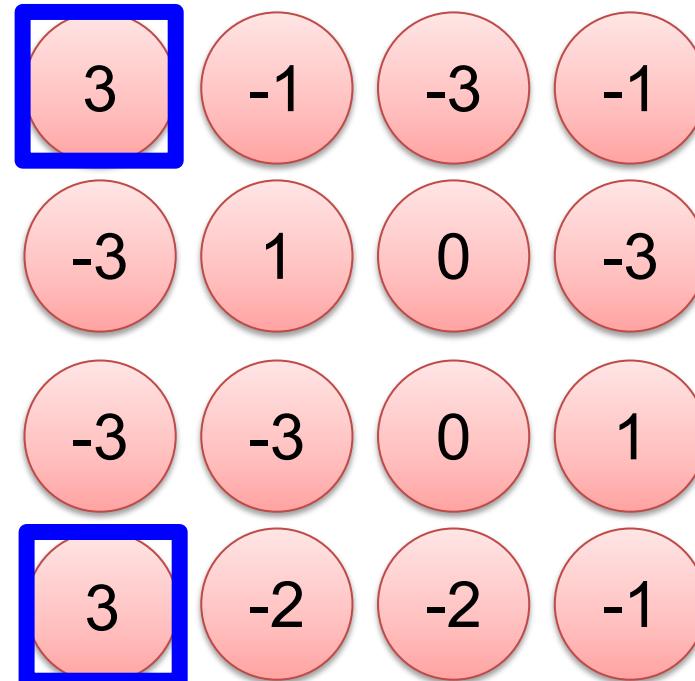
stride=1



6 x 6 image



Filter 1



Property 2

# CNN – Convolution

stride=1

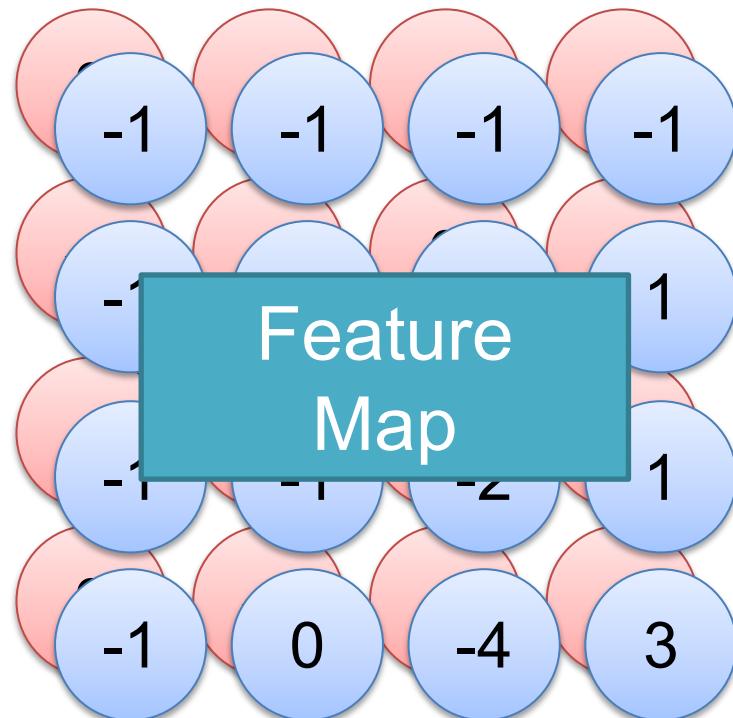
1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image

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

Filter 2

Do the same process  
for every filter



Feature  
Map

4 x 4 image (reduced size)

# Demo

1	0	1
0	1	0
1	0	1

Filter/kernel

1 <small><math>\times 1</math></small>	1 <small><math>\times 0</math></small>	1 <small><math>\times 1</math></small>	0	0
0 <small><math>\times 0</math></small>	1 <small><math>\times 1</math></small>	1 <small><math>\times 0</math></small>	1	0
0 <small><math>\times 1</math></small>	0 <small><math>\times 0</math></small>	1 <small><math>\times 1</math></small>	1	1
0	0	1	1	0
0	1	1	0	0

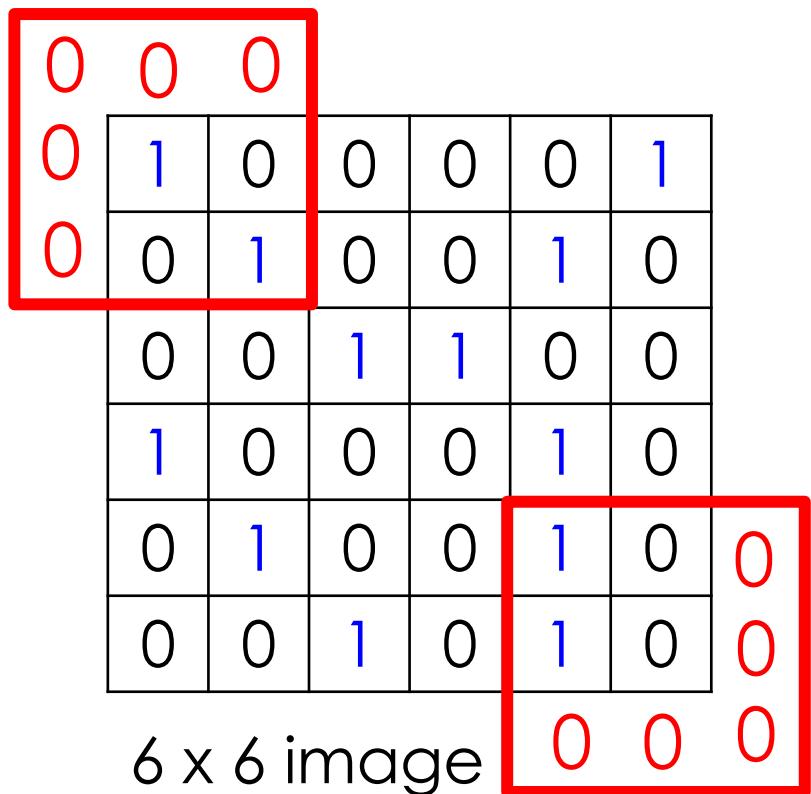
Image

4		

Convolved  
Feature

<https://ujjwalkarn.me/2016/08/11/intuitive-explanation-convnets/>

# CNN – Zero Padding



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

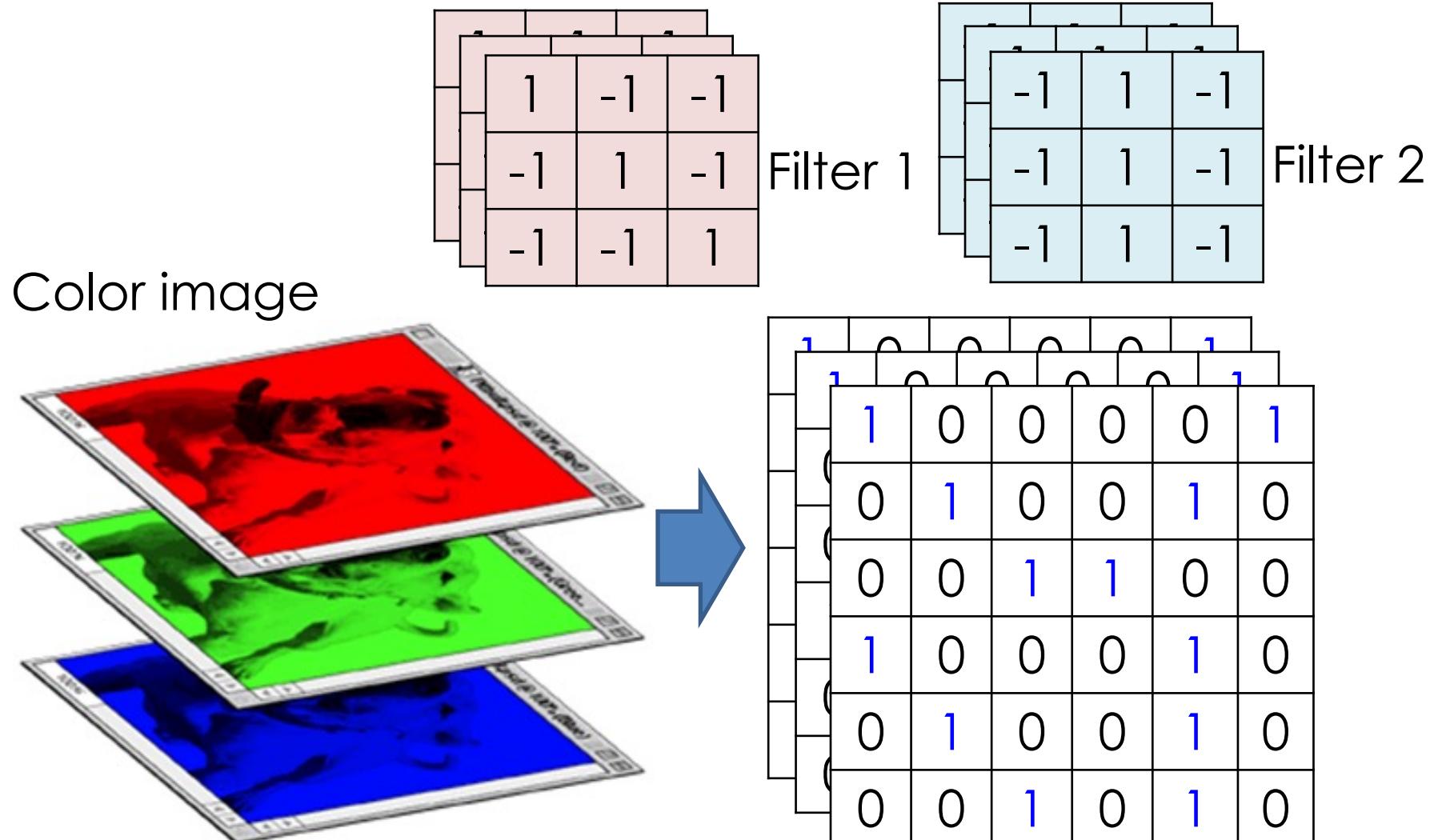
Filter 1

You will get another  $6 \times 6$  image in this way



Zero  
padding

# CNN – Color image



Input Volume (+pad 1) (7x7x3)

 $x[:, :, 0]$ 

0	0	0	0	0	0	0	0
0	2	0	0	2	1	0	
0	2	1	0	2	2	0	

0	2	1	0	2	0	0	
0	2	2	1	2	2	0	
0	0	1	1	0	1	0	
0	0	0	0	0	0	0	

0	0	0	0	0	0	0	
0	1	1	2	2	0	0	
0	1	2	1	0	1	0	

0	0	2	0	2	2	0	
0	0	2	0	0	1	0	
0	2	1	2	2	0	0	

0	0	0	0	0	0	0	
0	0	1	2	1	0	0	
0	0	0	1	1	2	0	

Filter W0 (3x3x3)

 $w0[:, :, 0]$ 

-1	-1	1
1	0	1
-1	0	-1

0	0	1
-1	-1	0
1	-1	0

0	1	1
1	1	0
1	-1	0

b0[:, :, 0]
1

Filter W1 (3x3x3)

 $w1[:, :, 0]$ 

0	-1	0
1	1	1
0	1	1

w1[:, :, 1]		
-1	-1	-1
-1	1	-1
-1	1	-1

w1[:, :, 2]		
-1	-1	1
-1	-1	1
1	-1	-1

Bias b1 (1x1x1)

 $b1[:, :, 0]$ 

0
---

Output Volume (3x3x2)

 $o[:, :, 0]$ 

-2	2	-1
3	6	-3
4	1	-2

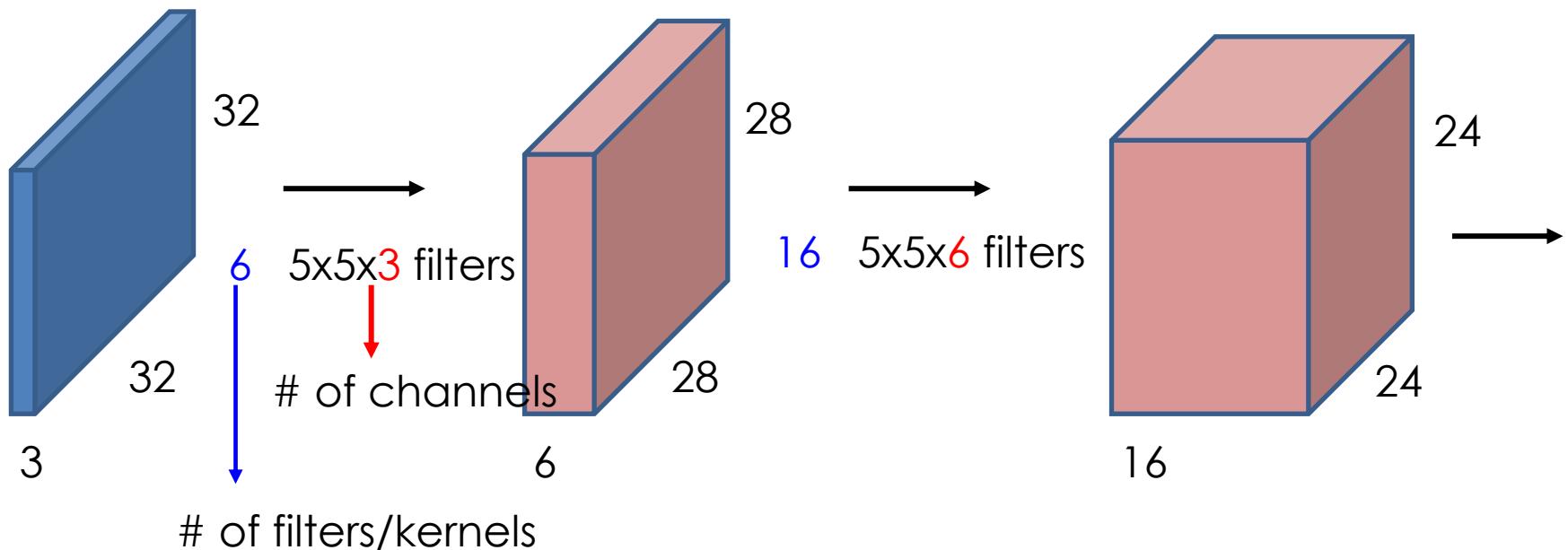
o[:, :, 1]		
4	-1	3
-4	-4	-3
-2	-3	-8

toggle movement



# Convolutional Network

- Convolution network is a sequence of these layers



# Convolution - Intuition

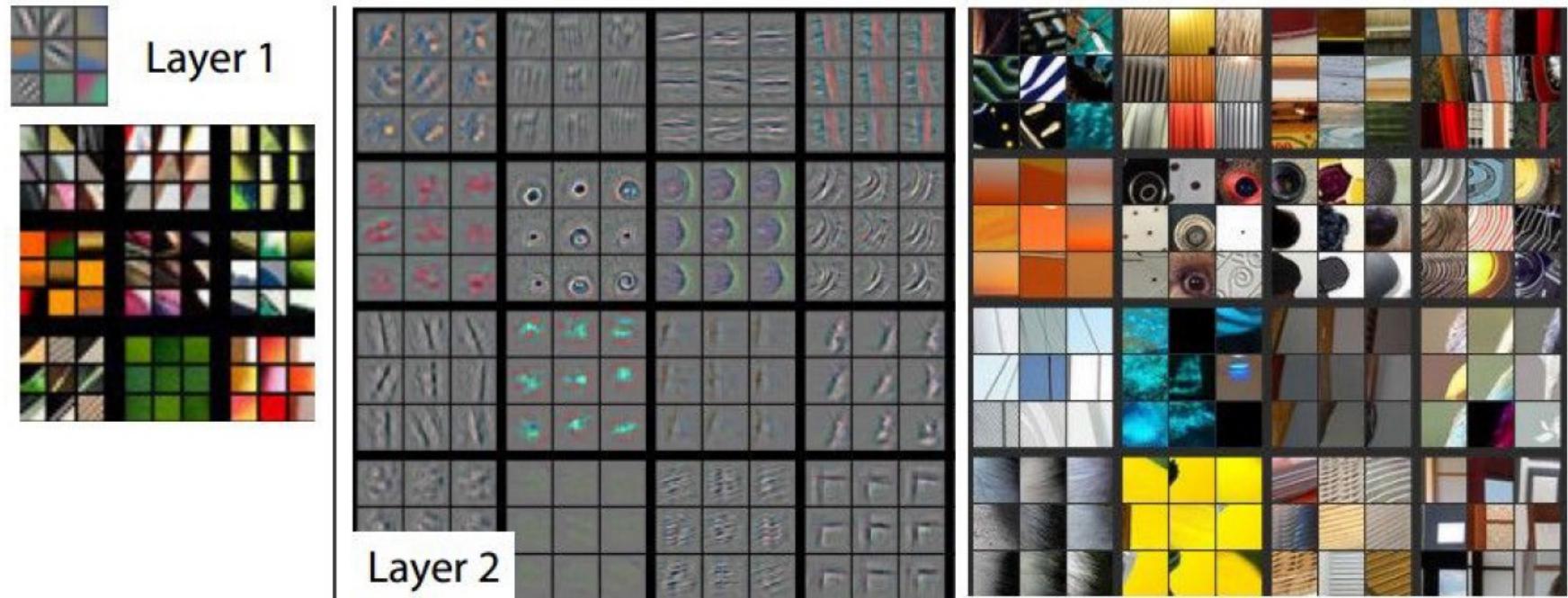
---



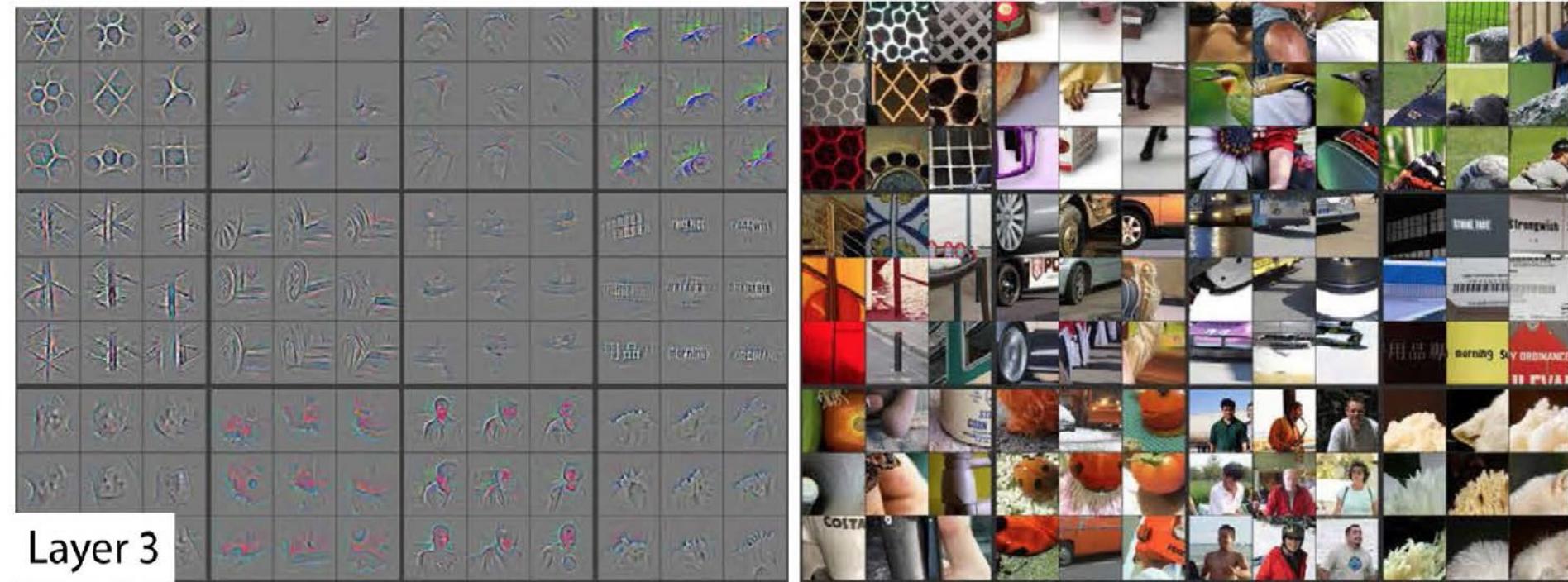
Source : [https://cs.nyu.edu/~fergus/tutorials/deep\\_learning\\_cvpr12/](https://cs.nyu.edu/~fergus/tutorials/deep_learning_cvpr12/)

# Visualizing CNN

---



Zeiler, Matthew D., and Rob Fergus. "Visualizing and understanding convolutional networks." *European conference on computer vision*. Springer, Cham, 2014.



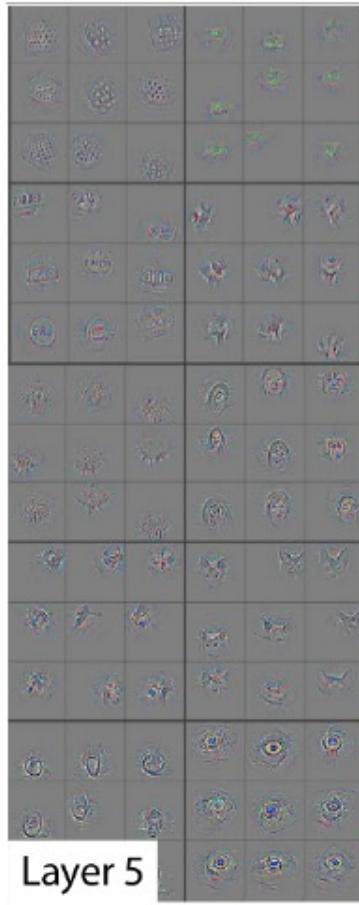
Zeiler, Matthew D., and Rob Fergus. "Visualizing and understanding convolutional networks." European conference on computer vision. Springer, Cham, 2014.



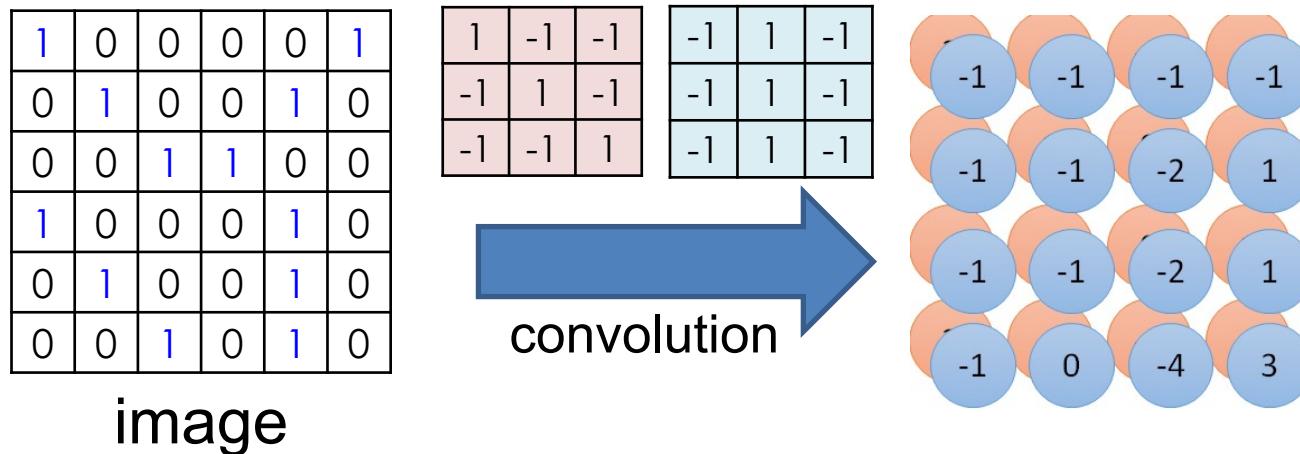
## Layer 4



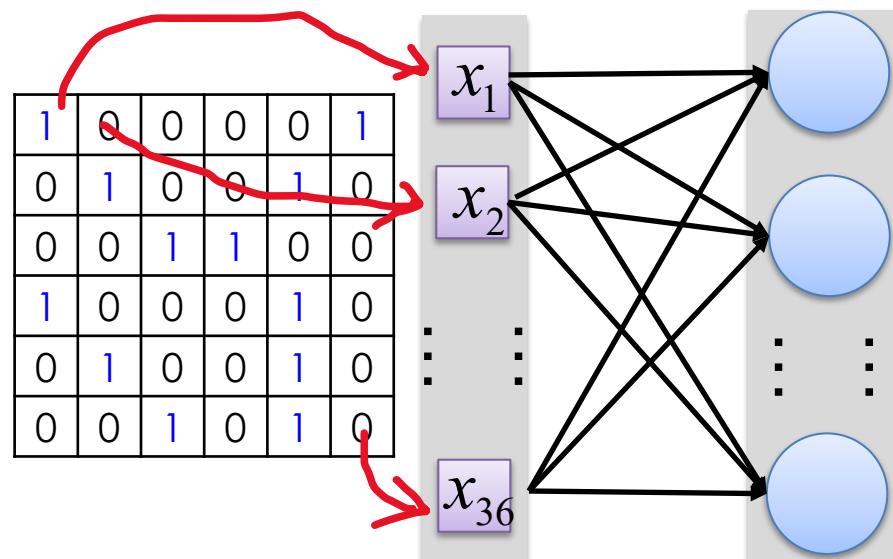
## Layer 5

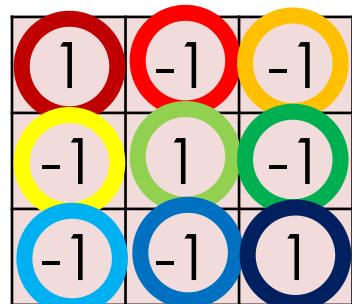


# Convolution v.s. Fully Connected

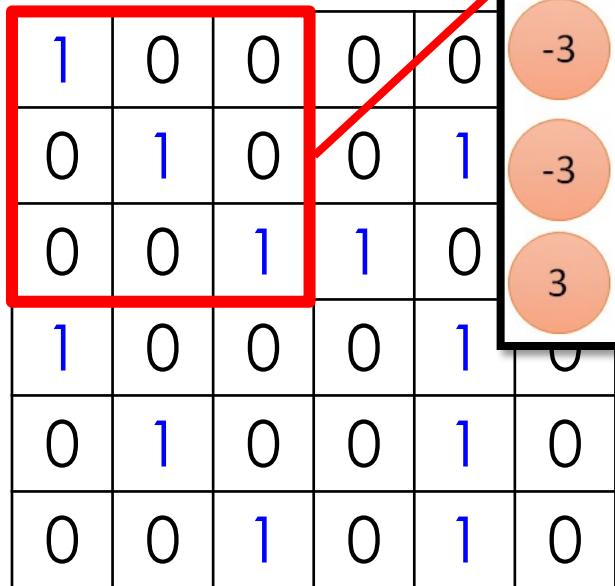


Fully-connected



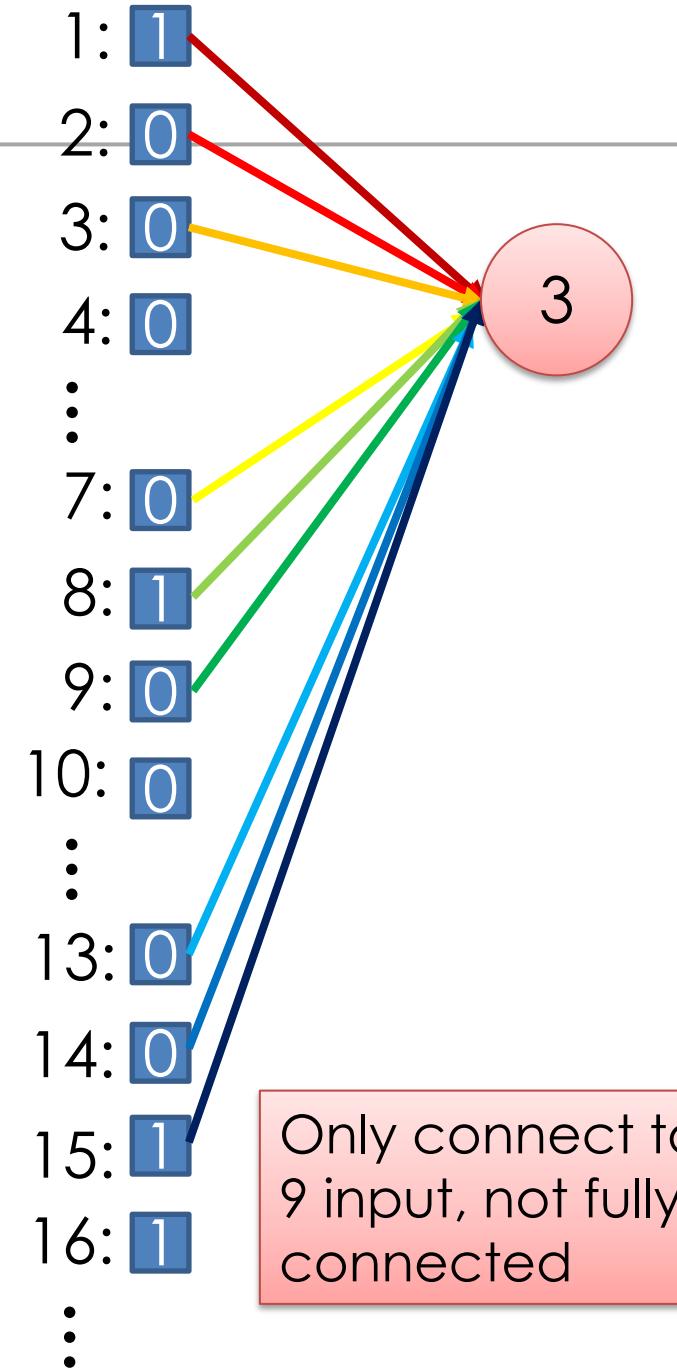
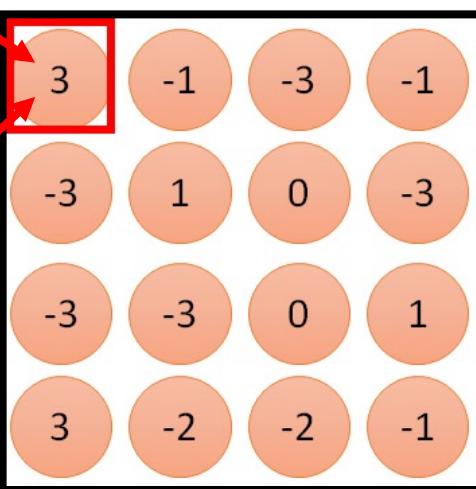


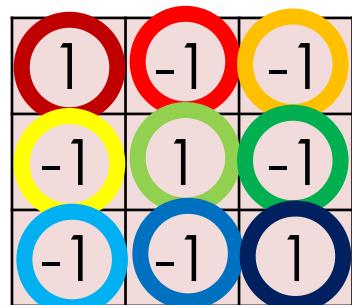
Filter 1



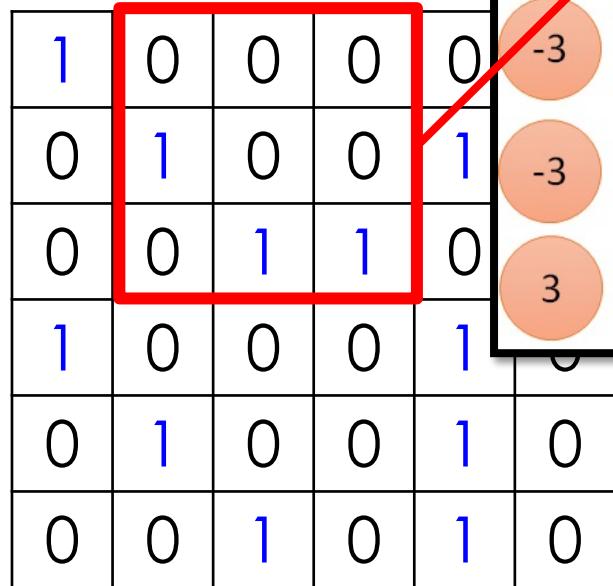
$6 \times 6$  image

Less parameters!



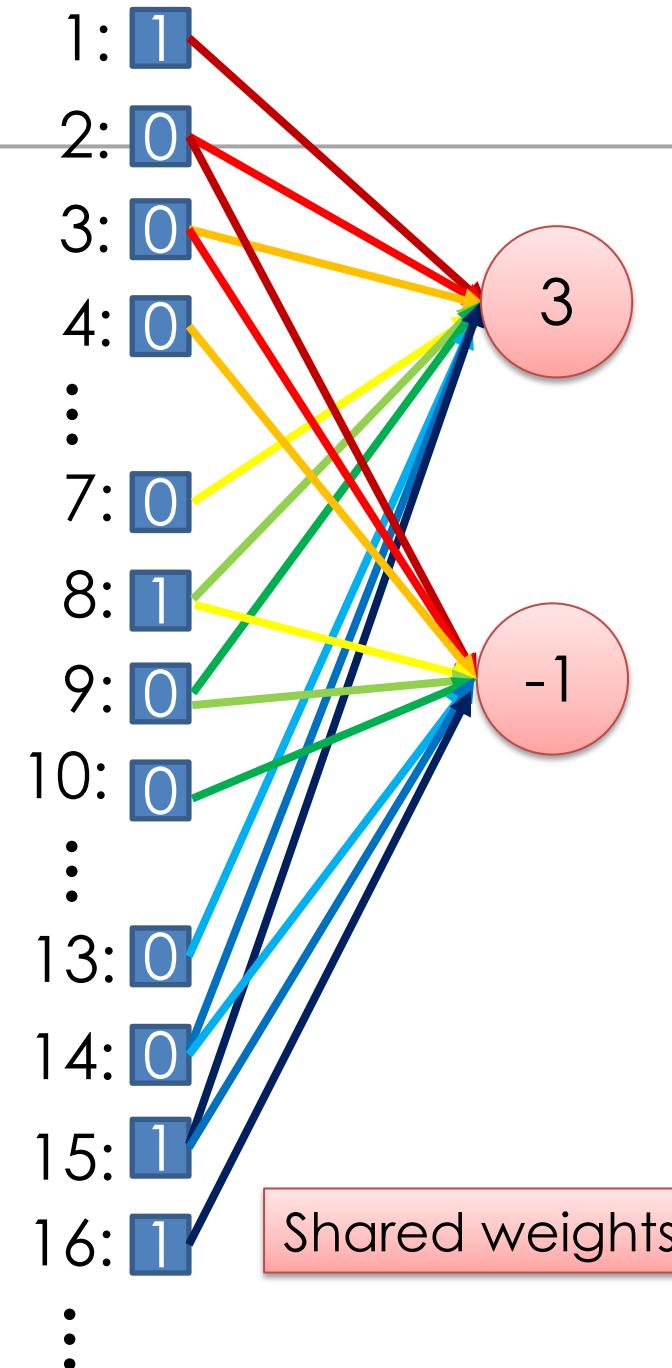
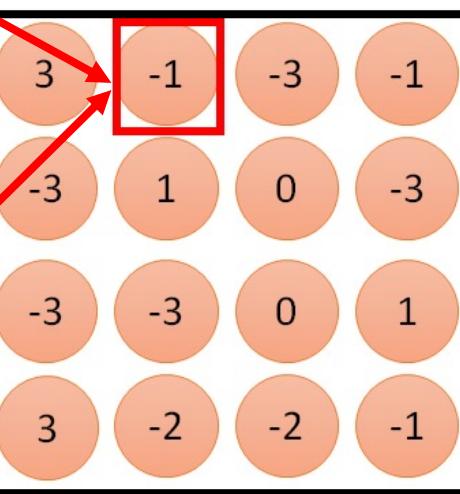


Filter 1



Less parameters!

Even less parameters!



---

Thank you!

Question?

# References and Slide Credits

---

- Many slides are adapted from the existing teaching or tutorial slides by Hung-yi Lee, Andrew Ng, Alexander Amini, Lex Fridman, and Stanford course - CS231n: Convolutional Neural Networks for Visual Recognition
- Special thanks to Dr. Hung-yi Lee for making his machine learning course slides and materials available
- Alexander Amini, MIT 6.S191 Introduction to Deep Learning:  
<http://introtodeeplearning.com/>
  - Youtube videos:  
[https://www.youtube.com/watch?v=5tvmMX8r\\_OM&list=PLtBw6njQRU-rwp5\\_7C0oIVt26ZgjG9NI&index=1](https://www.youtube.com/watch?v=5tvmMX8r_OM&list=PLtBw6njQRU-rwp5_7C0oIVt26ZgjG9NI&index=1)
- Lex Fridman, MIT Deep Learning and Artificial Intelligence Lectures: <https://deeplearning.mit.edu/>  
<https://www.youtube.com/watch?v=O5xeyoRL95U>