# Machine Learning

## Deep Learning Introduction

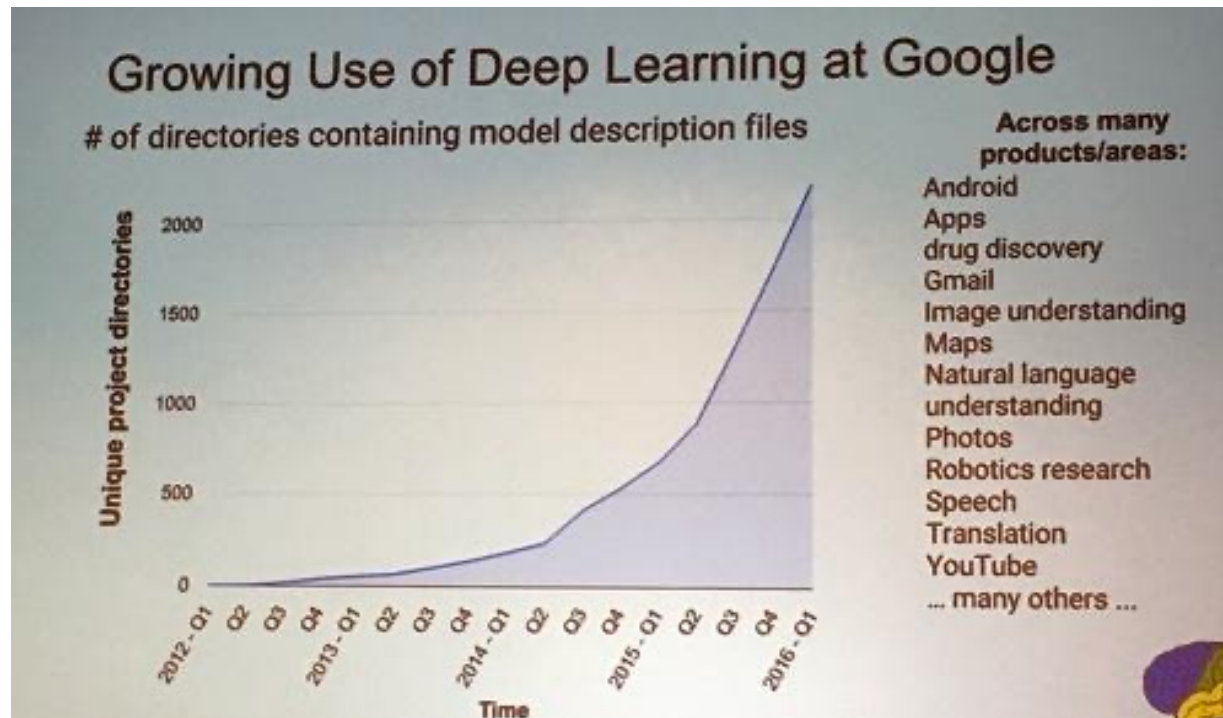Python tutorial: http://learnpython.org/

TensorFlow tutorial: https://www.tensorflow.org/tutorials/

PyTorch tutorial: https://pytorch.org/tutorials/

# Deep learning attracts lots of attention.

- I believe you have seen lots of exciting results before.



Deep learning trends at Google. Source: SIGMOD 2016/Jeff Dean
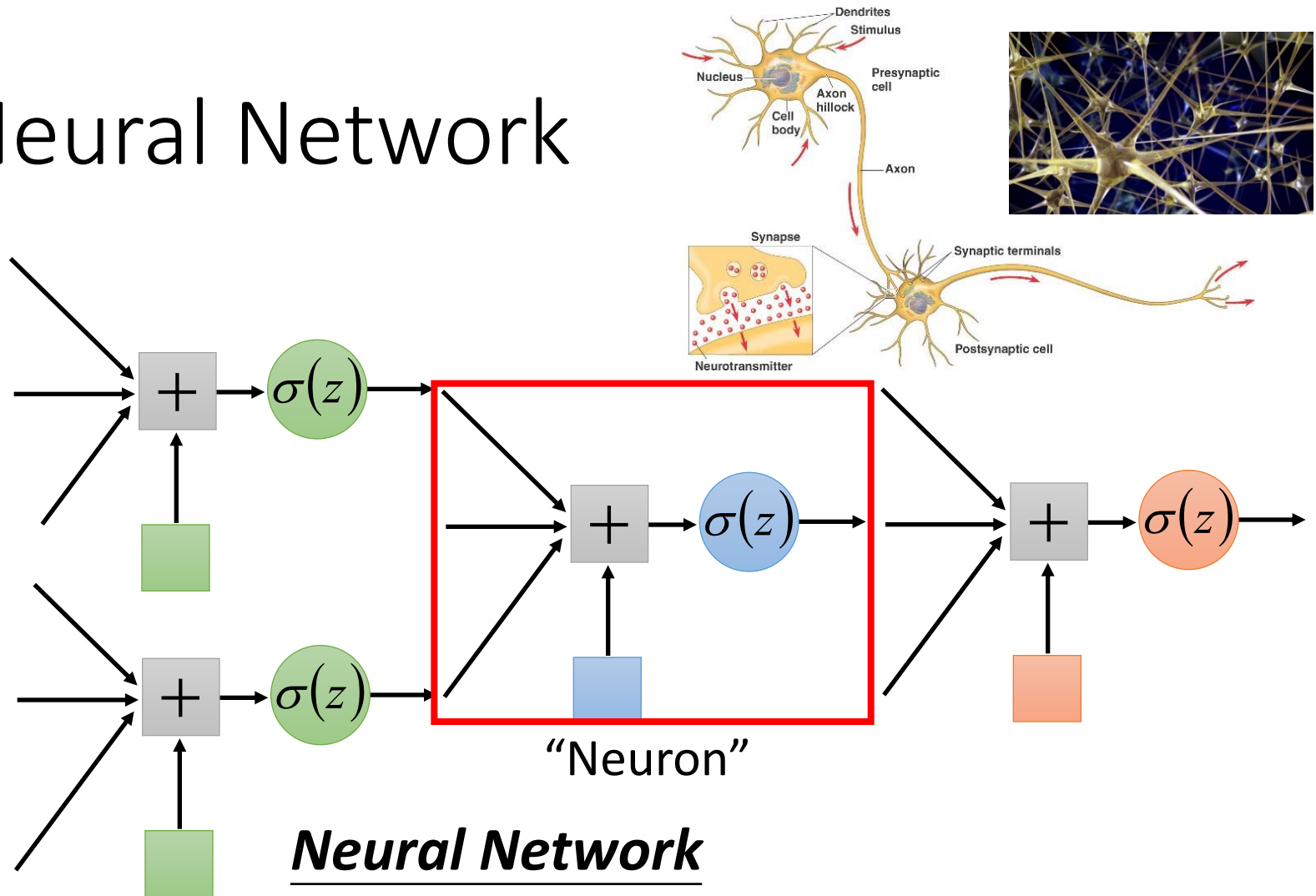
# *Ups and downs of Deep Learning*

- 1958: Perceptron (linear model)

- 1969: Perceptron has limitation

- 1980s: Multi-layer perceptron
    - Do not have significant difference from DNN today

- 1986: Backpropagation
    - Usually more than 3 hidden layers is not helpful

- 1989: 1 hidden layer is "good enough", why deep?

- 2006: RBM initialization

- 2009: GPU

- 2011: Start to be popular in speech recognition

- 2012: win ILSVRC image competition

- 2015.2: Image recognition surpassing human-level performance

- 2016.3: Alpha GO beats Lee Sedol

- 2016.10: Speech recognition system as good as humans

# Three Steps for Deep Learning



Deep Learning is so simple ……

# Neural Network



"Neuron"
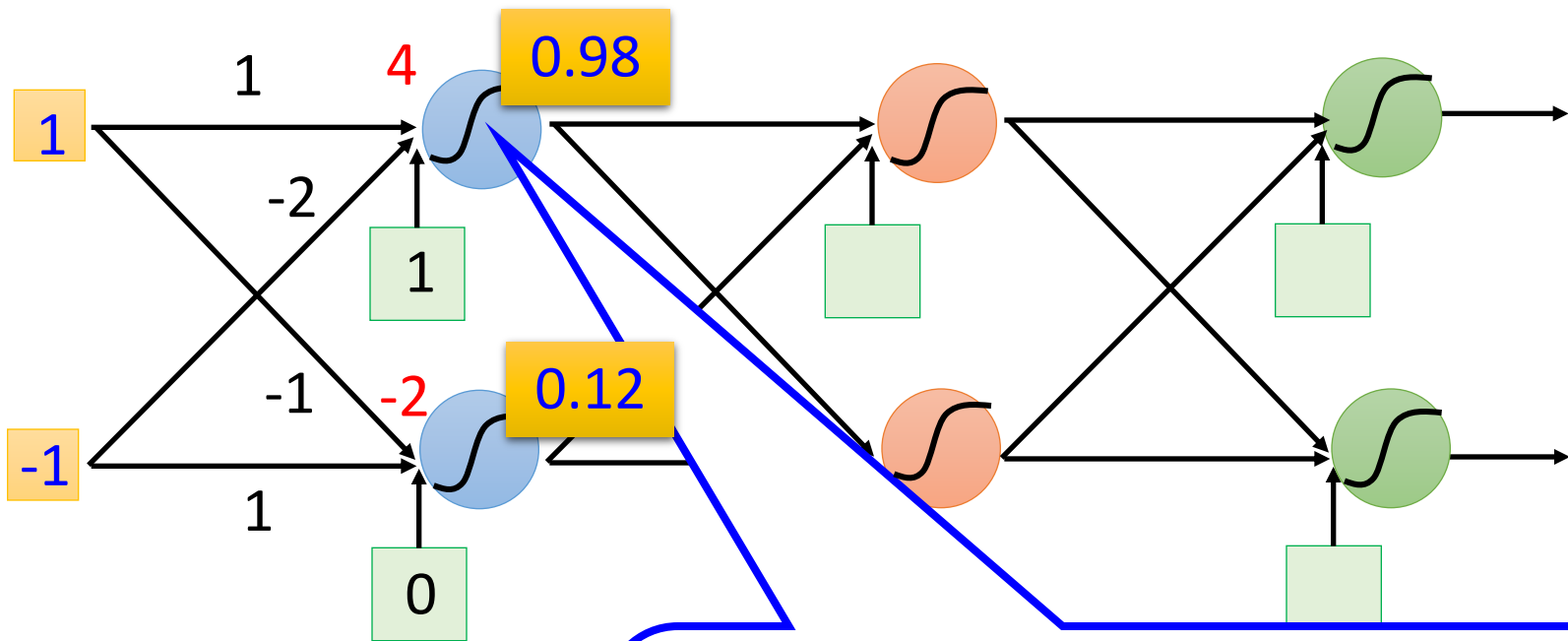
## *Neural Network*

Different connection leads to different network structures

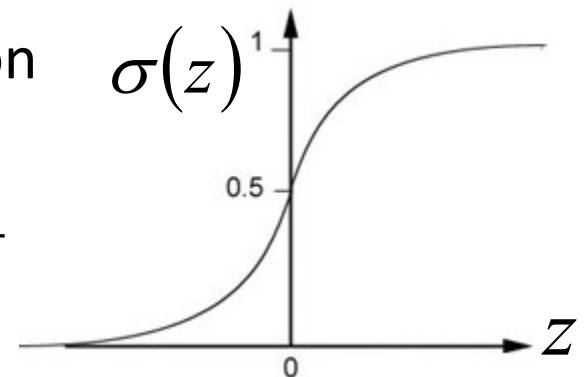Network parameter $\theta$: all the weights and biases in the "neurons"
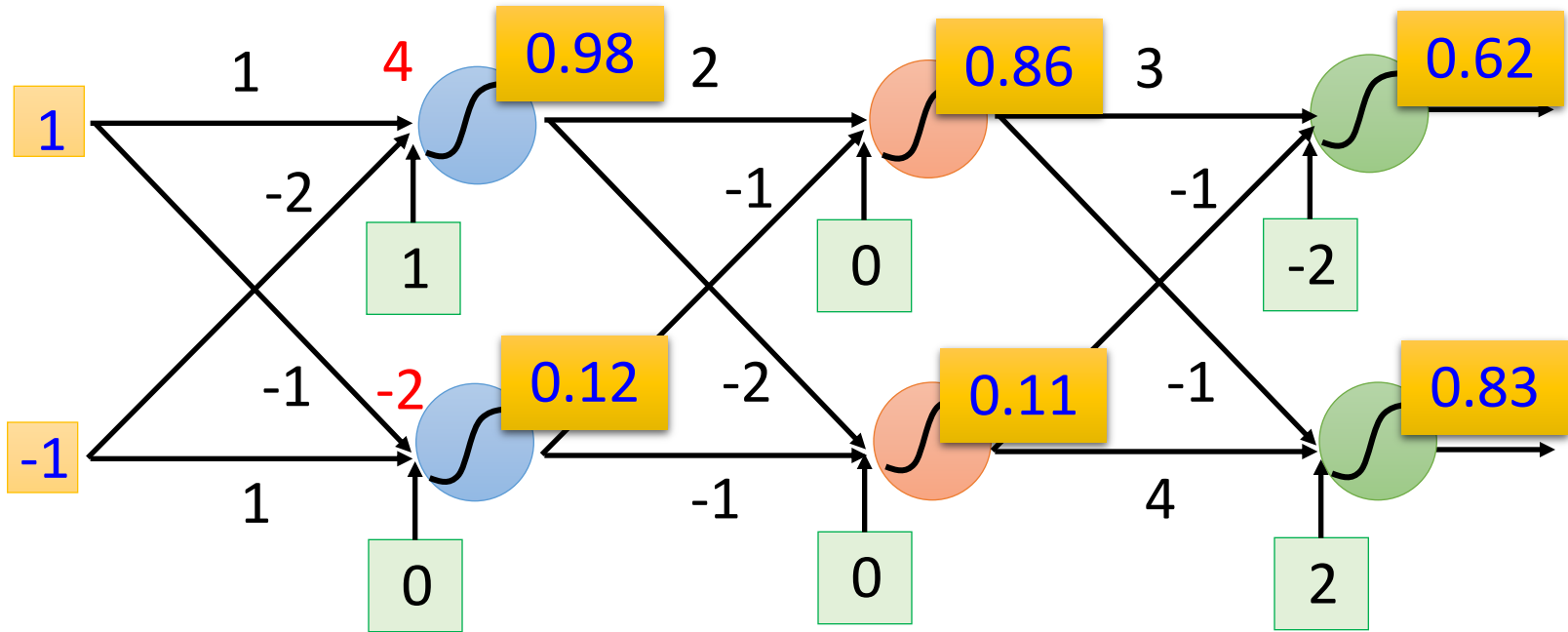
# Fully Connect Feedforward Network
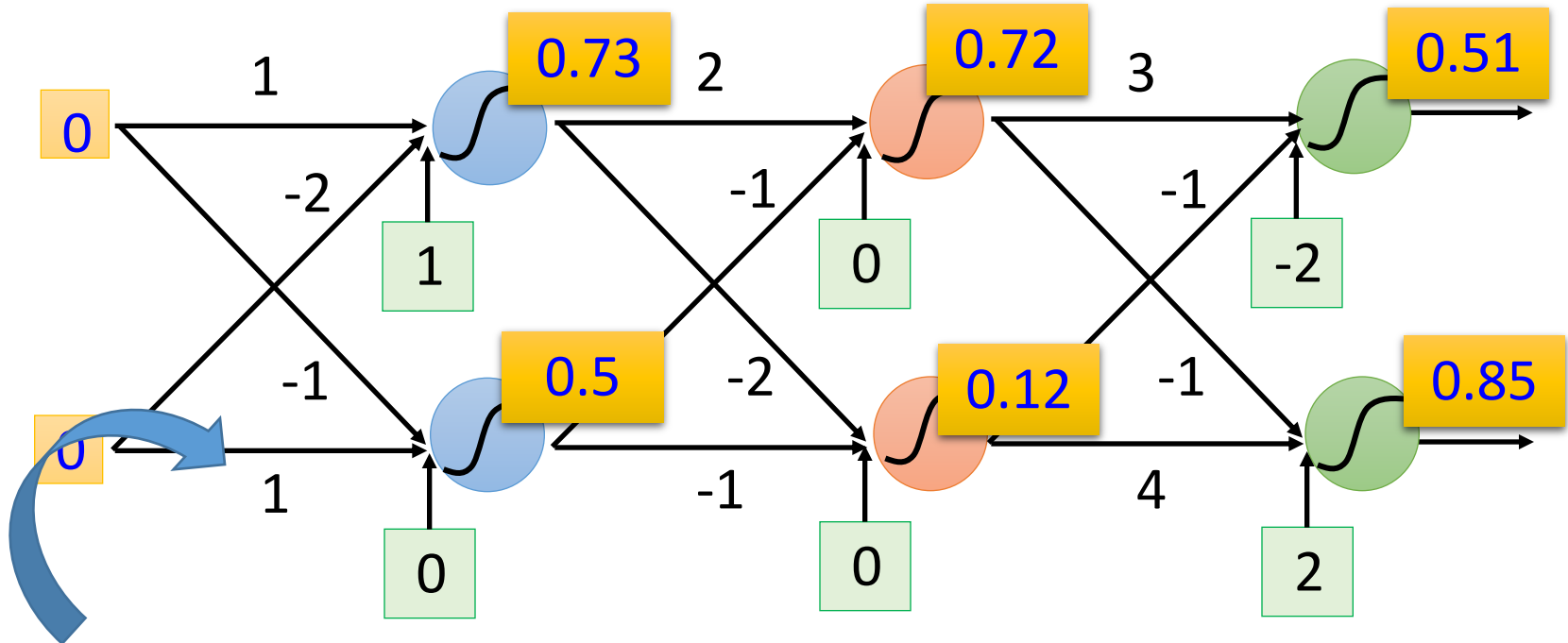


Sigmoid Function

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

# Fully Connect Feedforward Network

# Fully Connect Feedforward Network
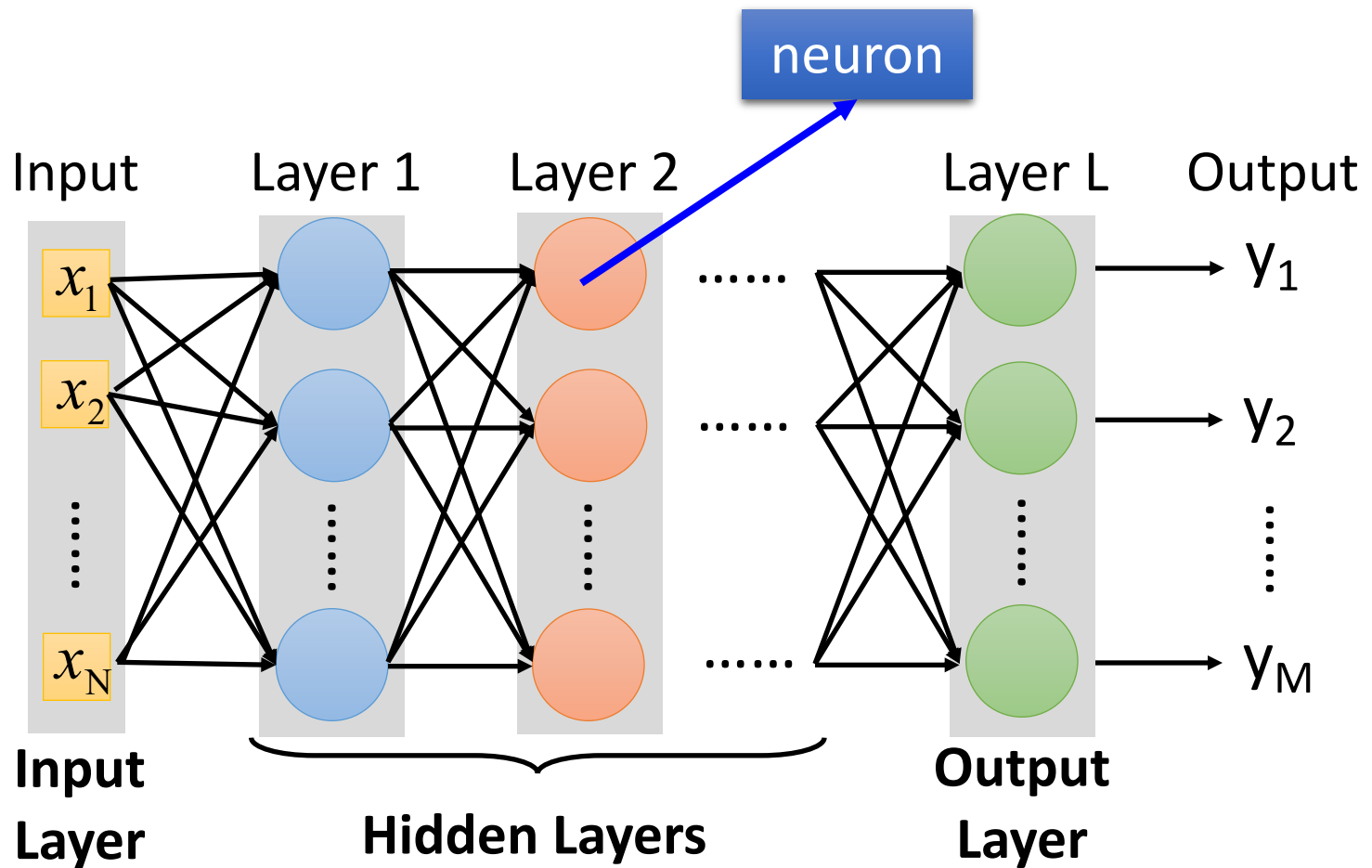


This is a function.
Input vector, output vector

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

Given network structure, define **_a function set_**
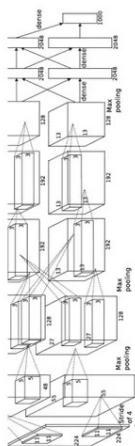
# Fully Connect Feedforward Network

# Deep = Many hidden layers

http://cs231n.stanford.edu/slides/winter1516_lecture8.pdf

22 layers

19 layers

8 layers

16.4%

7.3%

6.7%

AlexNet (2012)

VGG (2014)

GoogleNet (2014)

# Deep = Many hidden layers



Special structure

152 layers

101 layers

Ref: https://www.youtube.com/watch?v=dxB6299gpvI

3.57%

16.4%
AlexNet (2012)

7.3%
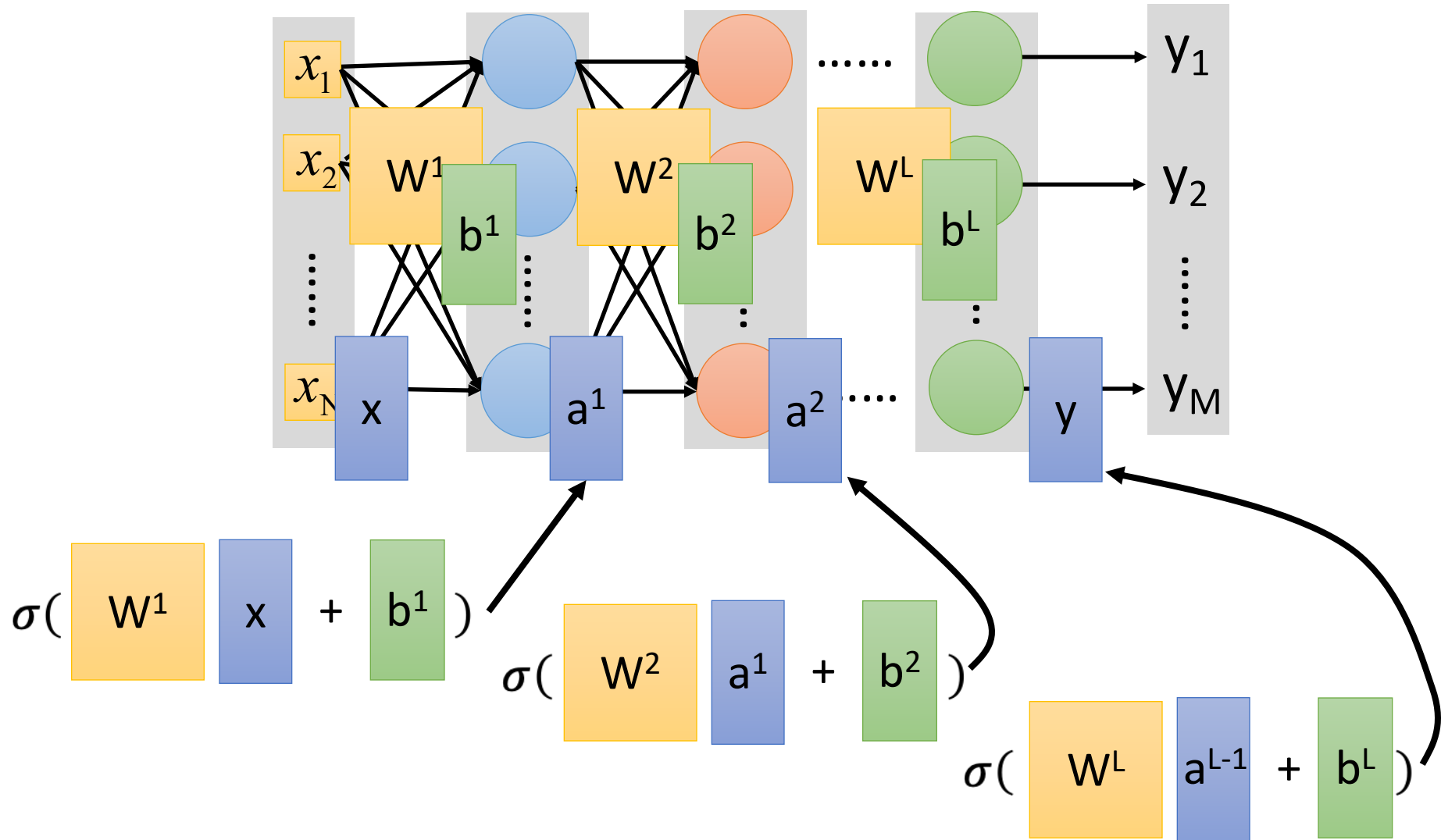VGG (2014)

6.7%
GoogleNet (2014)

Residual Net (2015)

Taipei 101

# Matrix Operation


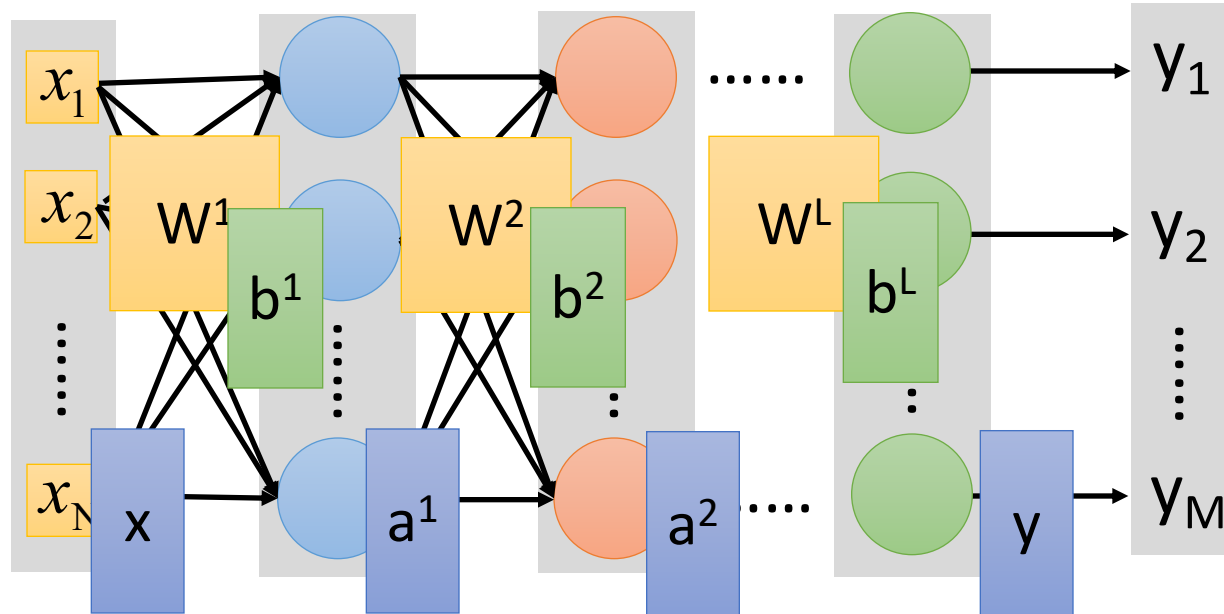
$$\sigma\left( \begin{bmatrix} 1 & -2 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \end{bmatrix} \right) = \begin{bmatrix} 0.98 \\ 0.12 \end{bmatrix}$$

$$\begin{bmatrix} 4 \\ -2 \end{bmatrix}$$

# Neural Network



$$\sigma(\ W^1\ x\ +\ b^1\ )$$

$$\sigma(\ W^2\ a^1\ +\ b^2\ )$$

$$\sigma(\ W^L\ a^{L-1}\ +\ b^L\ )$$

# Neural Network

# Output Layer
# as Multi-Class Classifier

# Example Application


Machine → "2"

## Input



16 x 16 = 256

Ink → 1
No ink → 0

$x_1$
$x_2$
$x_{256}$

## Output

| 0.1 | is 1 |
| 0.7 | is 2 |
| 0.2 | is 0 |

The image is "2"

Each dimension represents the confidence of a digit.

# Example Application

- Handwriting Digit Recognition



Input:
256-dim vector

Neural Network

What is needed is a function ……

output:
10-dim vector

# Example Application



Input | Layer 1 | Layer 2 | Layer L | Output

$x_1$, $x_2$, $x_N$

**Input Layer**

**Hidden Layers**

**Output Layer**

$y_1$, $y_2$, $y_{10}$

is 1, is 2, is 0

A function set containing the candidates for
Handwriting Digit Recognition
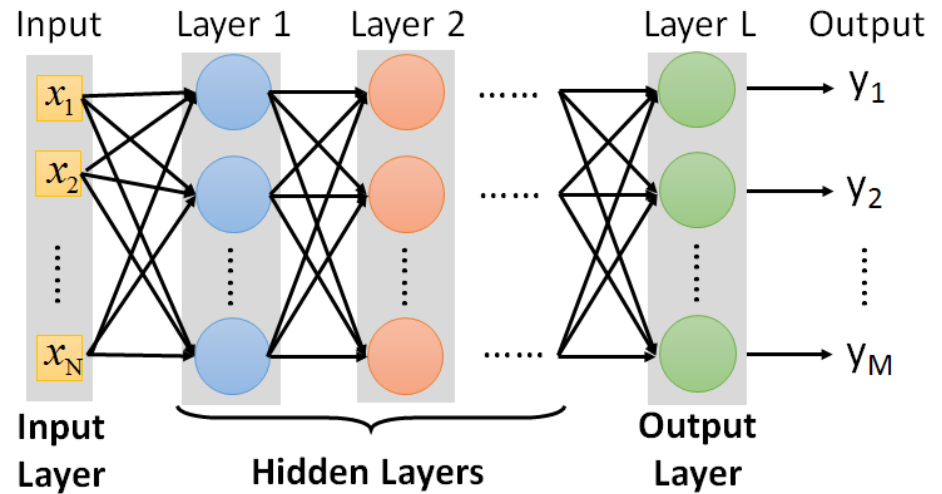
You need to decide the network structure to let a good function in your function set.

# FAQ



- Q: How many layers? How many neurons for each layer?

| Trial and Error | + | Intuition |

- Q: Can the structure be automatically determined?
  - E.g. Evolutionary Artificial Neural Networks
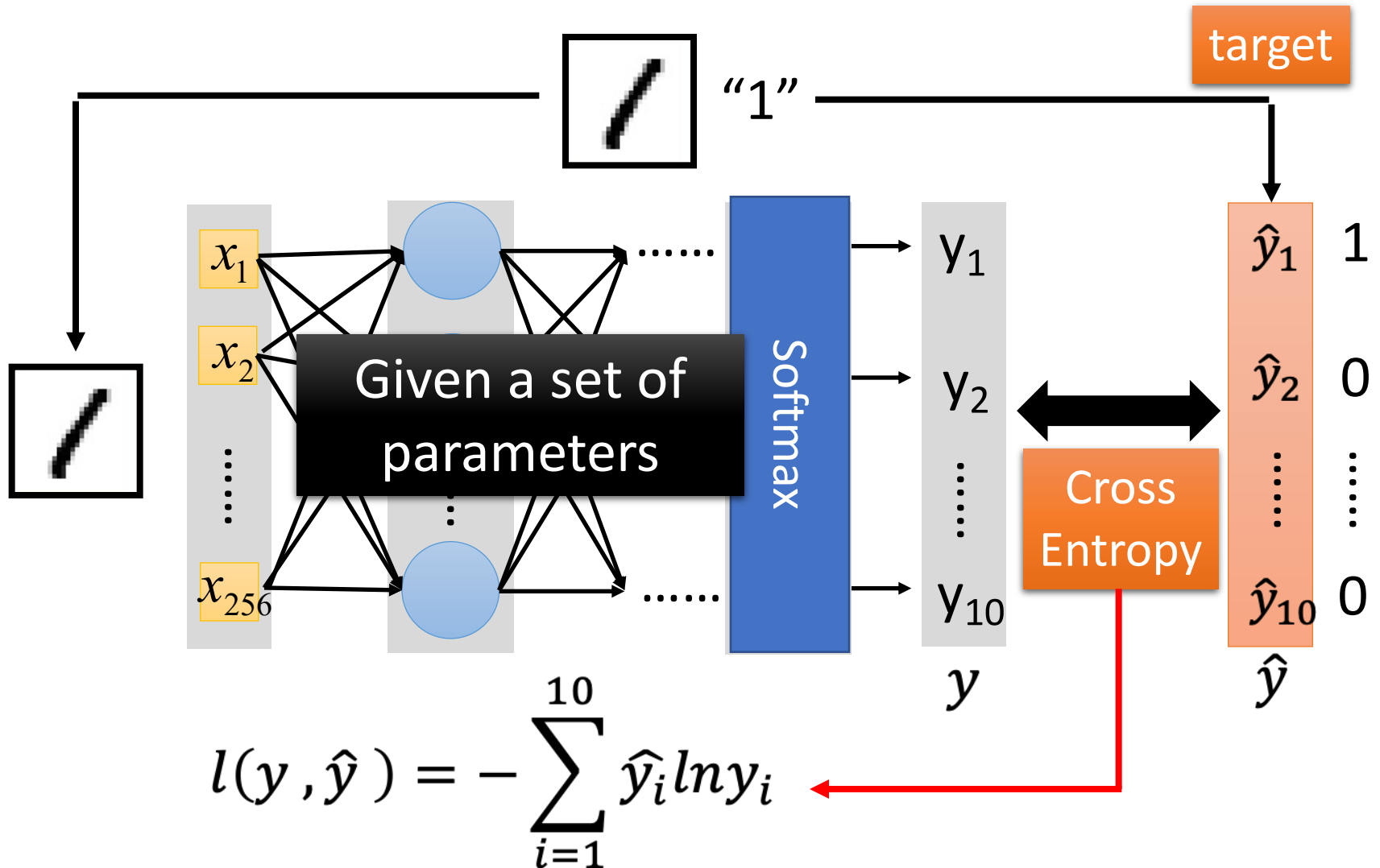- Q: Can we design the network structure?

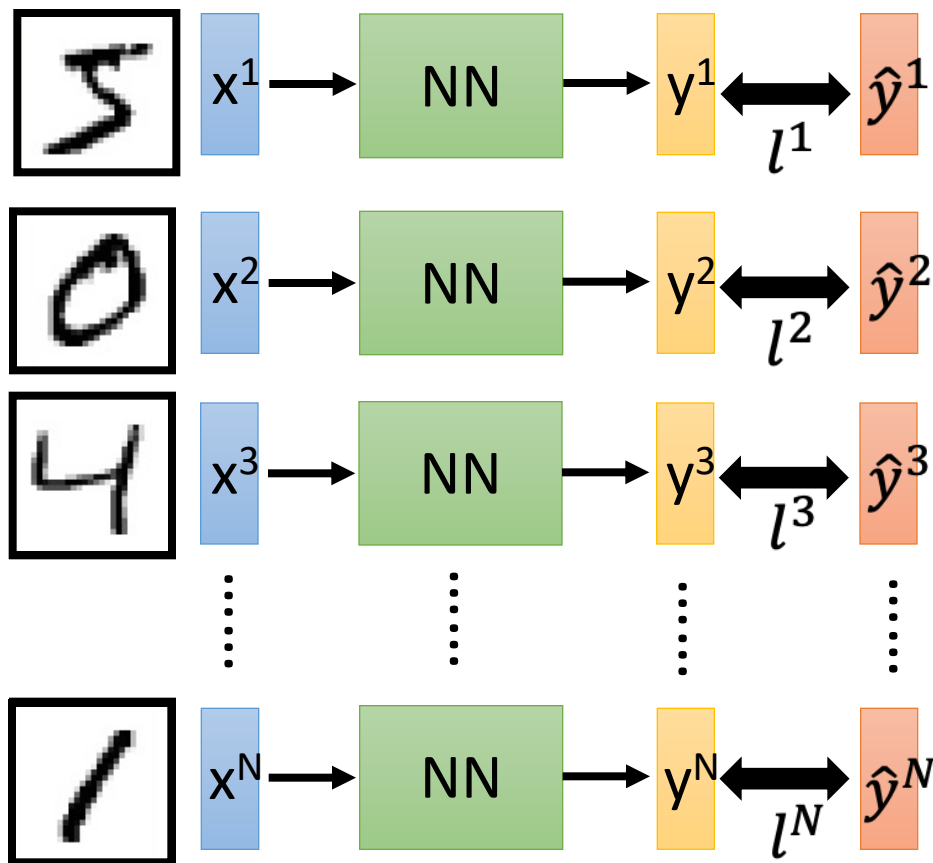Convolutional Neural Network (CNN)

# Three Steps for Deep Learning

| Step 1: Neural Network | → | Step 2: goodness of function | → | Step 3: pick the best function |

Deep Learning is so simple ……

# Loss for an Example



$$l(y, \hat{y}) = -\sum_{i=1}^{10} \hat{y}_i \, ln \, y_i$$

# Total Loss

Total Loss:

For all training data …



$$L = \sum_{n=1}^{N} l^n$$

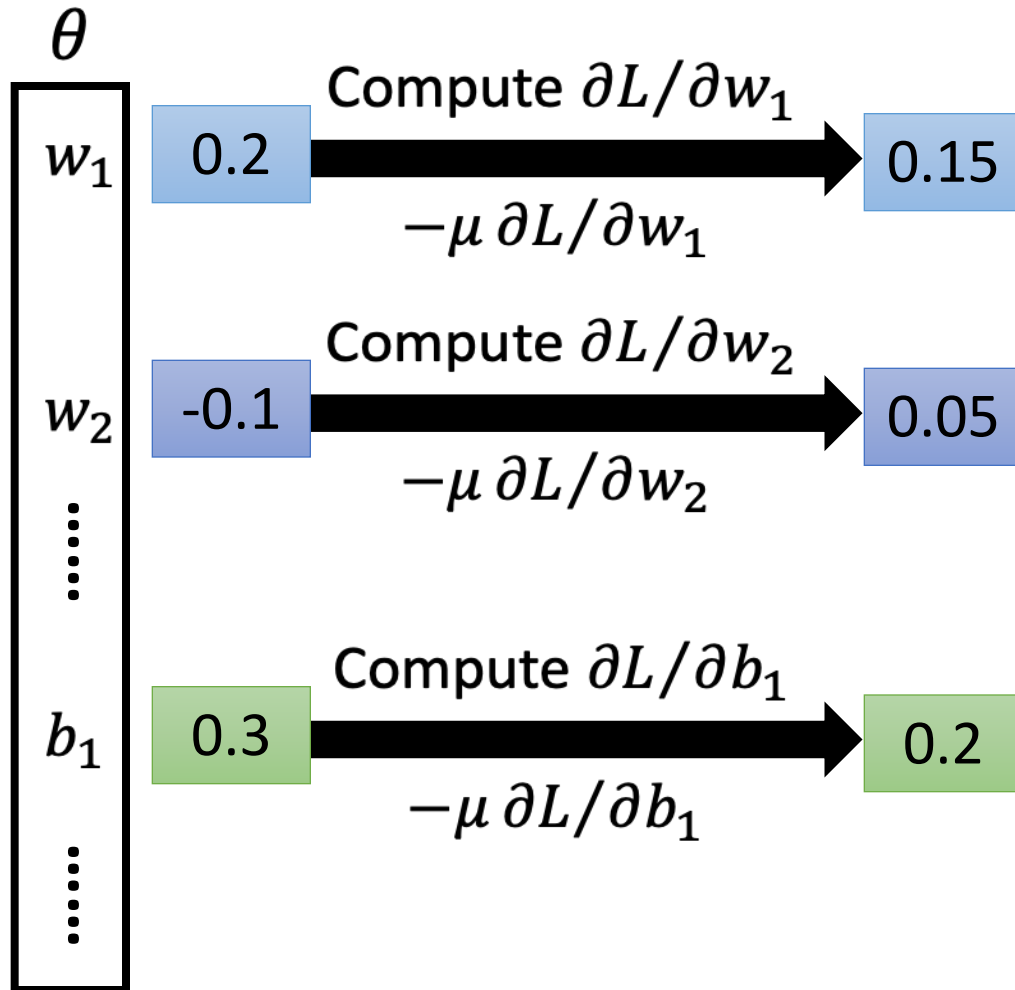Find *a function in function set* that minimizes total loss L

Find *the network parameters $\theta^*$* that minimize total loss L

# Three Steps for Deep Learning

| Step 1: Neural Network | → | Step 2: goodness of function | → | Step 3: pick the best function |
|---|---|---|---|---|

Deep Learning is so simple ……

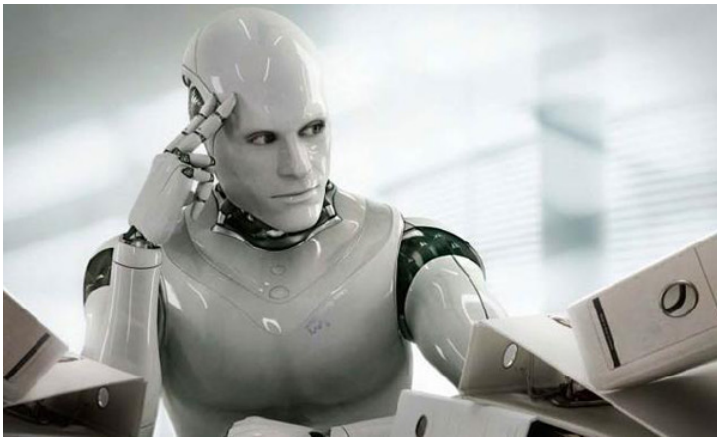# Gradient Descent

# Gradient Descent

# Gradient Descent

This is the "learning" of machines in deep learning ......

➡️ Even alpha go using this approach.

People image ......



Actually .....



I hope you are not too disappointed :p

# Backpropagation

- Backpropagation: an efficient way to compute $\partial L / \partial w$ in neural network



libdnn
台大周伯威
同學開發

Ref:
http://speech.ee.ntu.edu.tw/~tlkagk/courses/MLDS_2015_2/Lecture/DNN%20backprop.ecm.mp4/index.html

# Three Steps for Deep Learning

Step 1: Neural Network

Step 2: goodness of function

Step 3: pick the best function

Deep Learning is so simple ……