

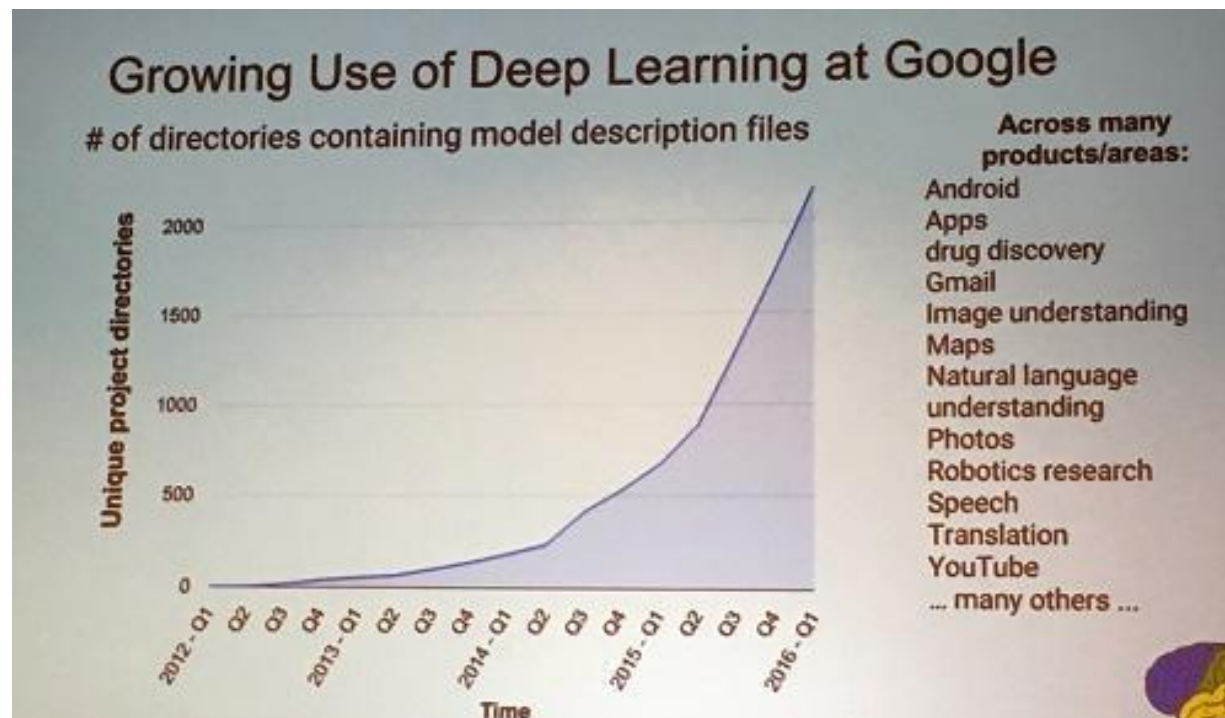
# Deep Learning

Hung-yi Lee

李宏毅

# 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) 少了sigmoid的概念
- 1969: Perceptron has limitation
- 1980s: Multi-layer perceptron Neural Network
  - 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 現在比較少用
- 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

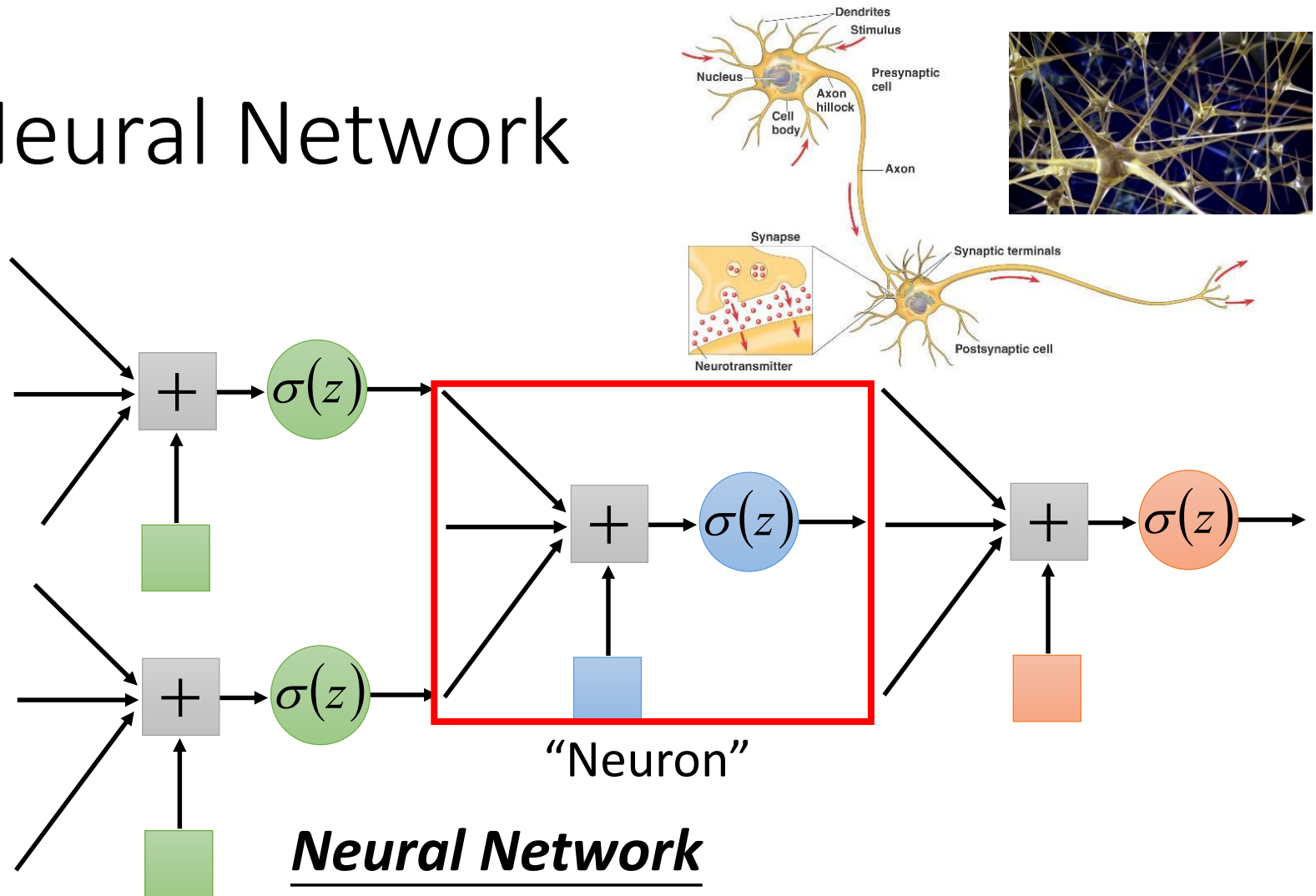
# Three Steps for Deep Learning



Deep Learning is so simple .....



# Neural Network



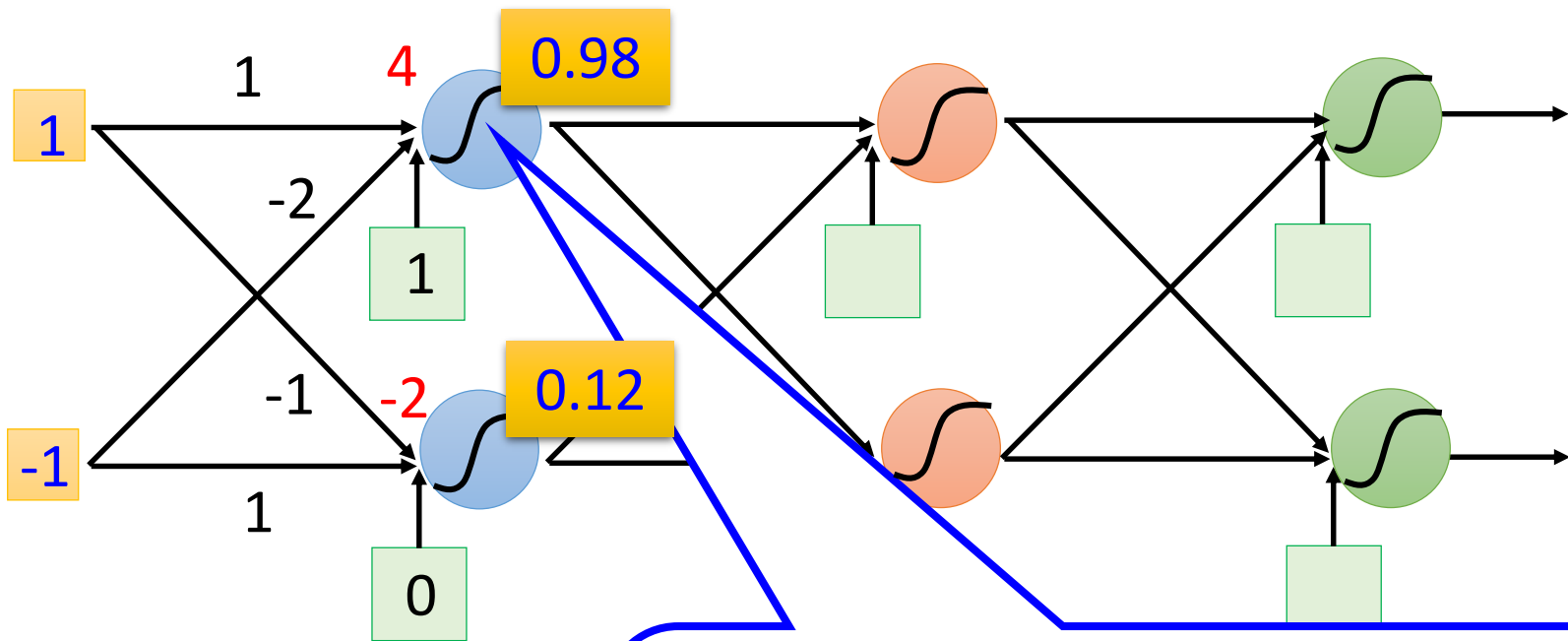
## Neural Network

Different connection leads to different network structures

每一個logistic regression的weight and bias結合起來

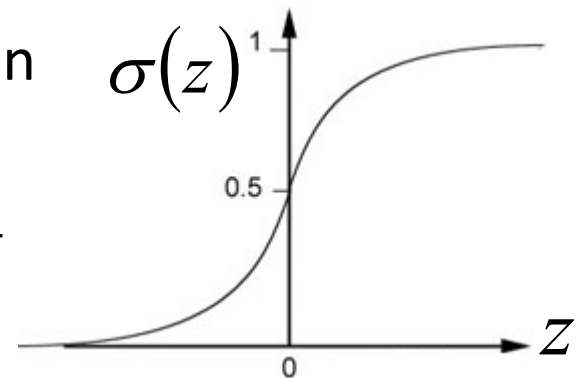
Network parameter  $\theta$ : all the weights and biases in the "neurons"

# Fully Connect Feedforward Network

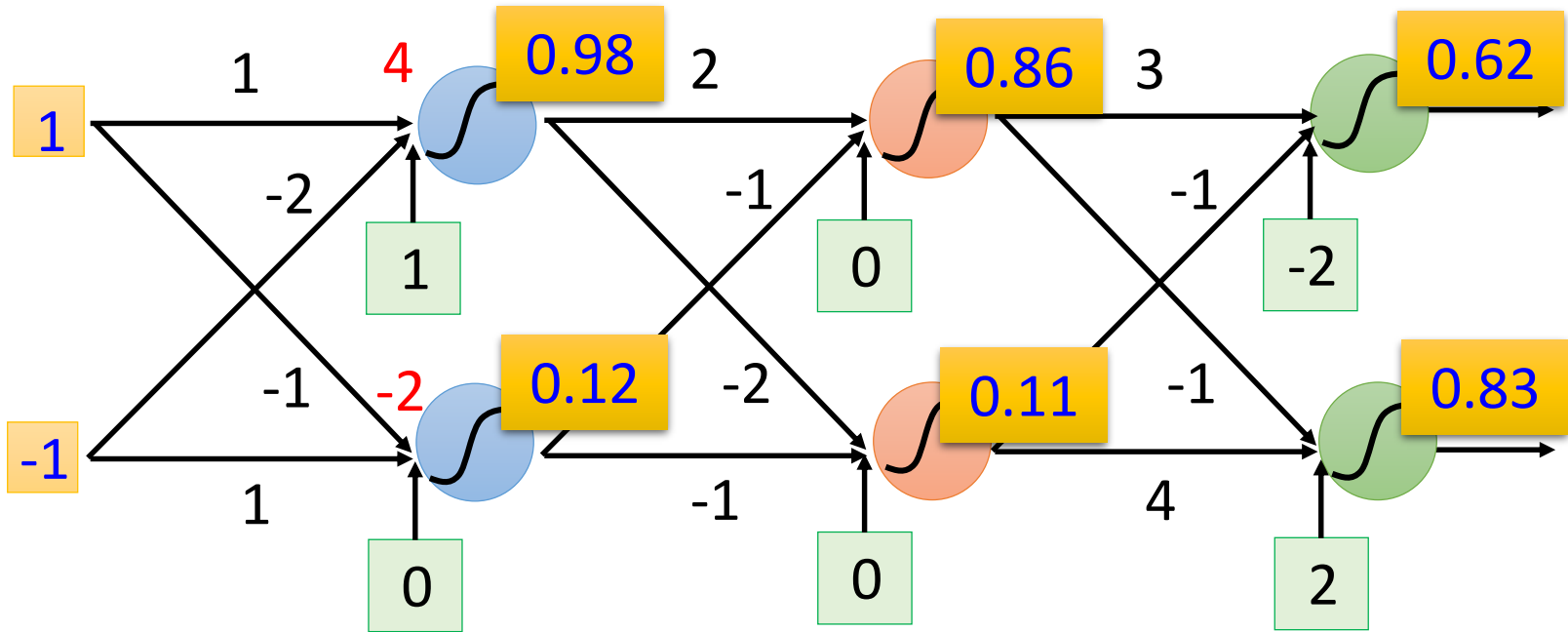


Sigmoid Function

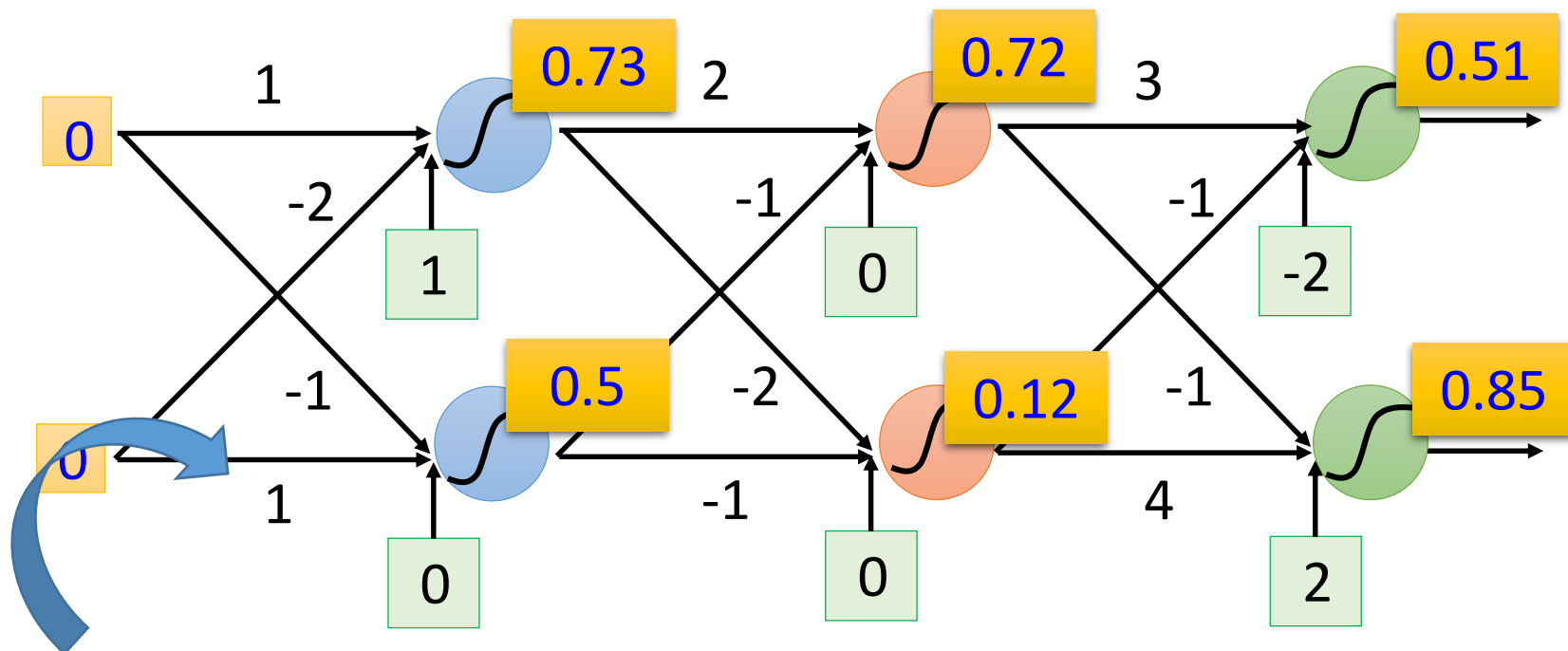
$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



# Fully Connect Feedforward Network



# Fully Connect Feedforward Network



This is a function.

Input vector, output vector

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

一個neural network可以把它看做一個function (with parameters)

Given network structure, define a function set

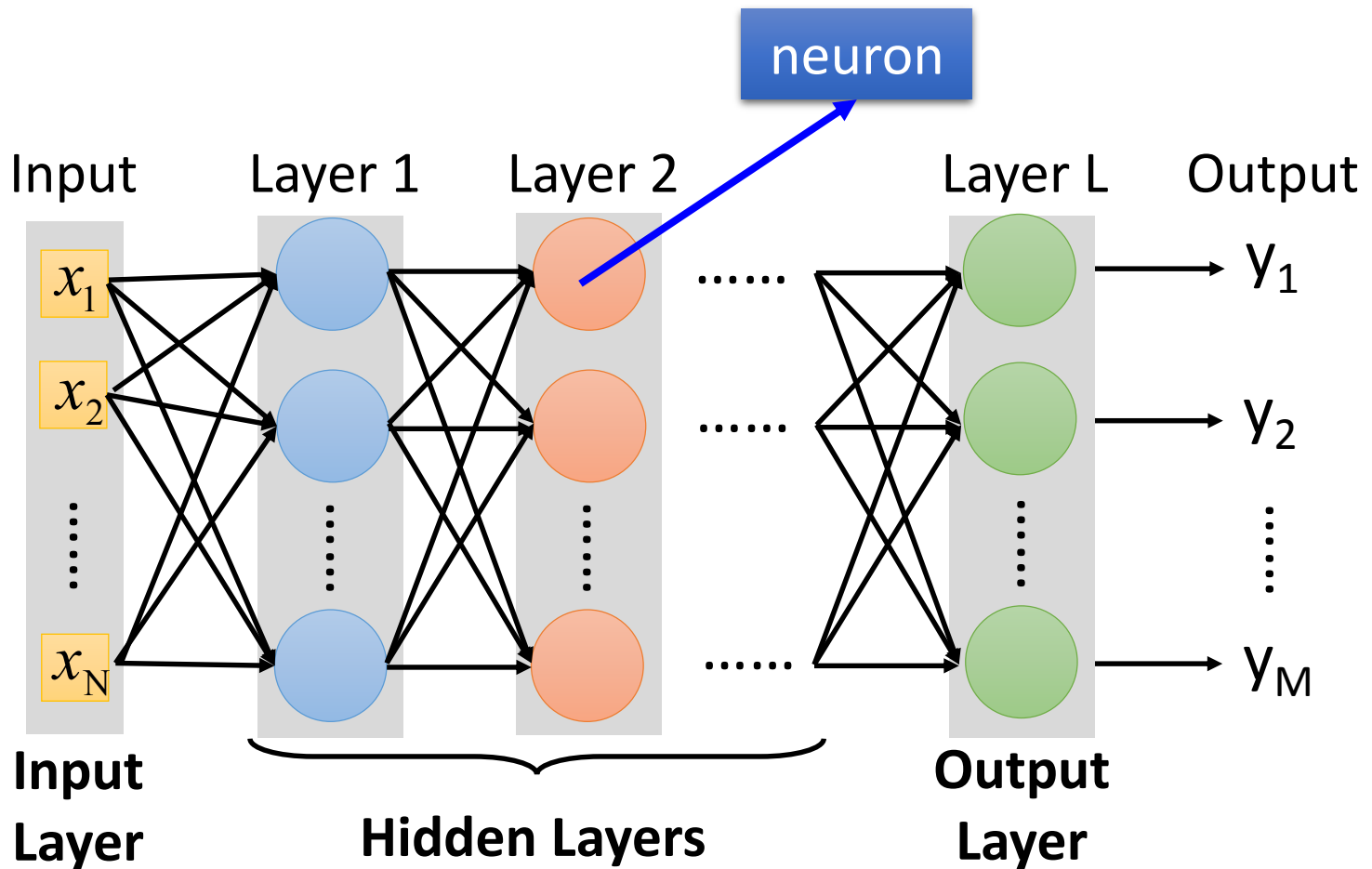
一個neural network的連接方式(structure) 就相當於 function set (without parameters)



所有layer都兩兩連接

# Fully Connect Feedforward Network

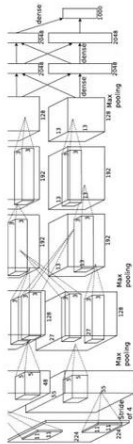
1 -2 ->3 ...



# Deep = Many hidden layers

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

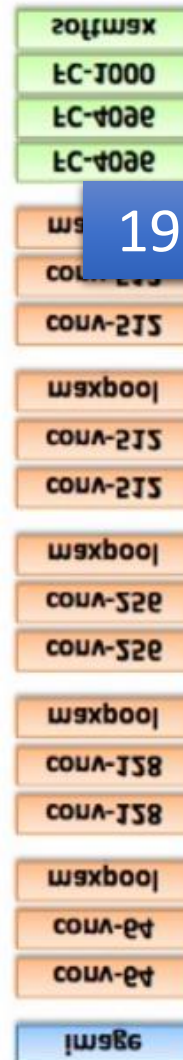
8 layers



16.4%

# AlexNet (2012)

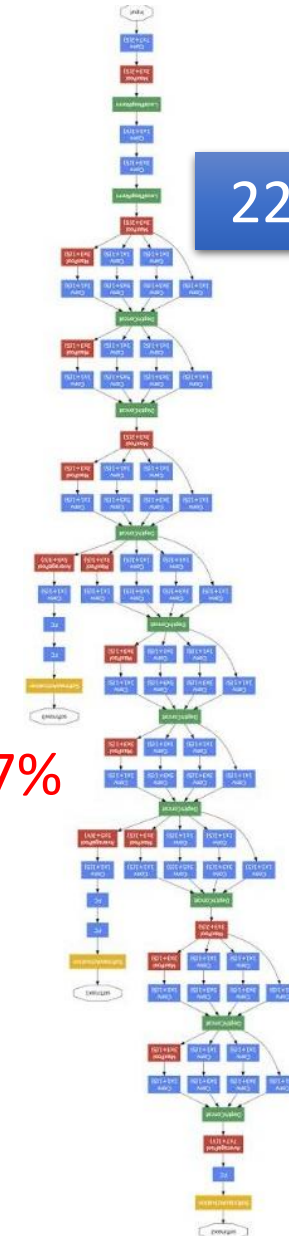
7.3%



19 layers

## VGG (2014)

22 layers



6.7%

## GoogleNet (2014)

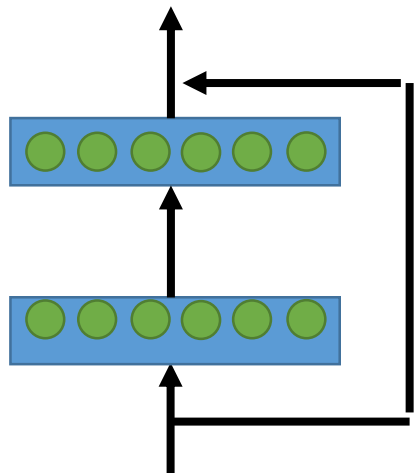
# Deep = Many hidden layers

不是一般的fully connected  
feedforward network

152 layers

101 layers

Special  
structure



Ref:  
<https://www.youtube.com/watch?v=dxB6299gpvl>

3.57%

16.4%

7.3%

6.7%

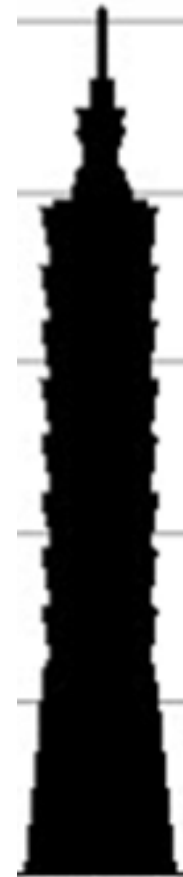
AlexNet  
(2012)

VGG  
(2014)

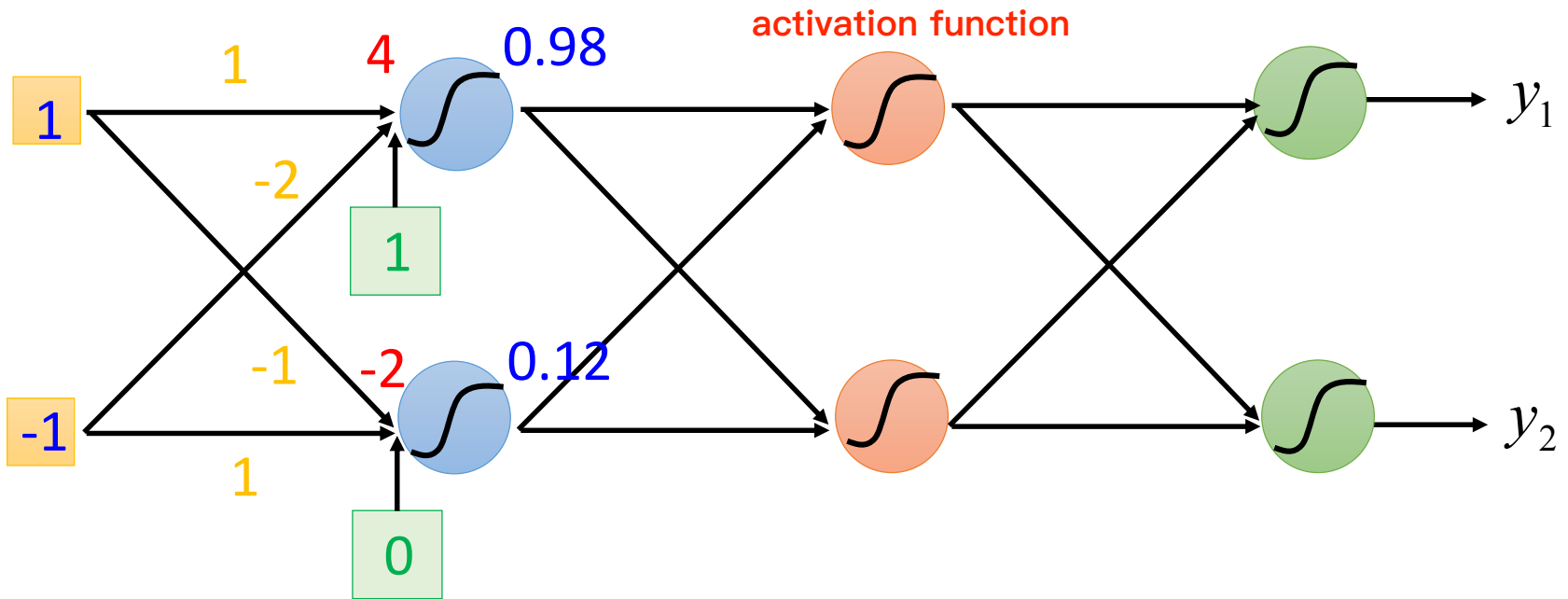
GoogleNet  
(2014)

Residual Net  
(2015)

Taipei  
101

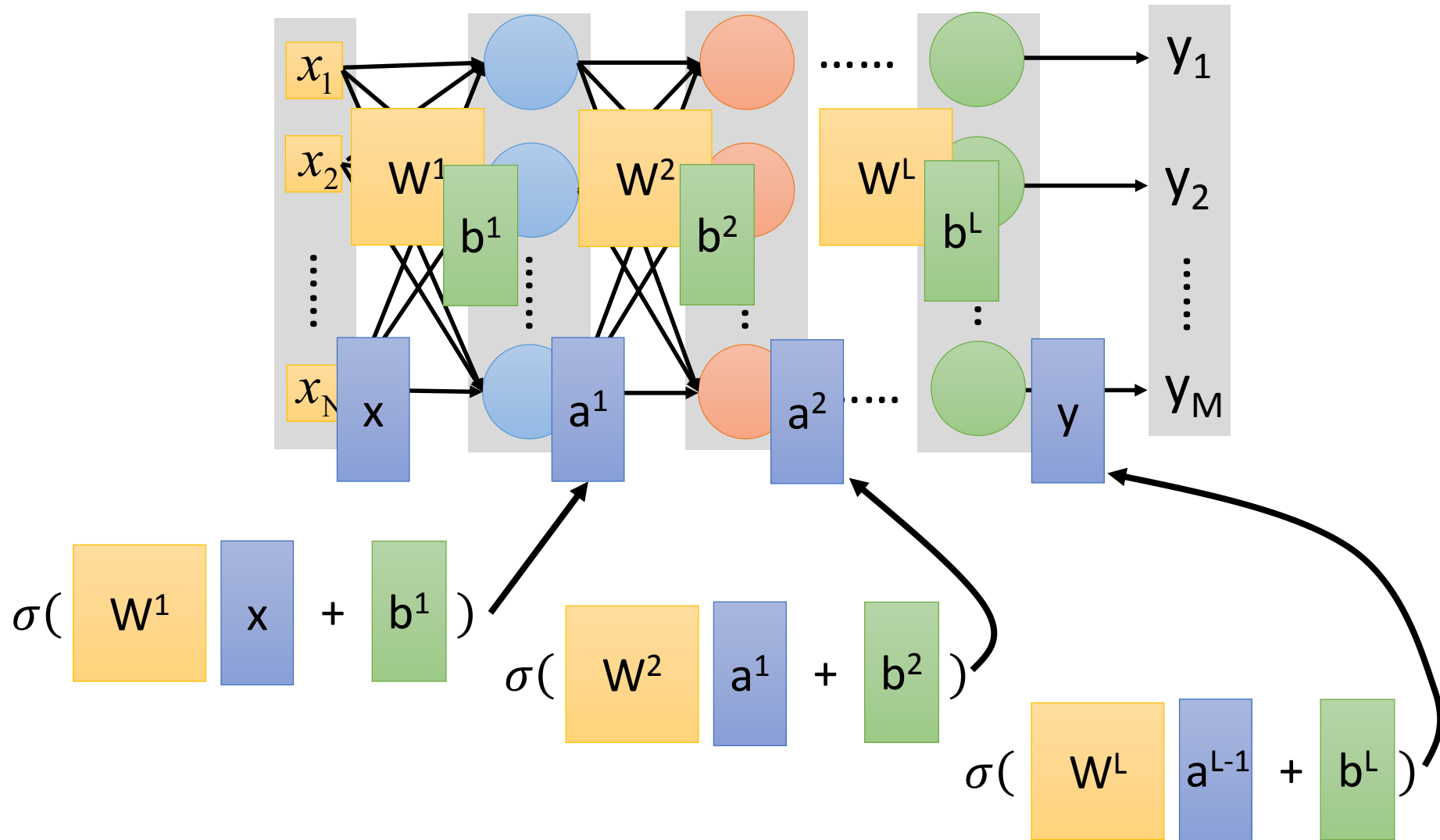


# Matrix Operation

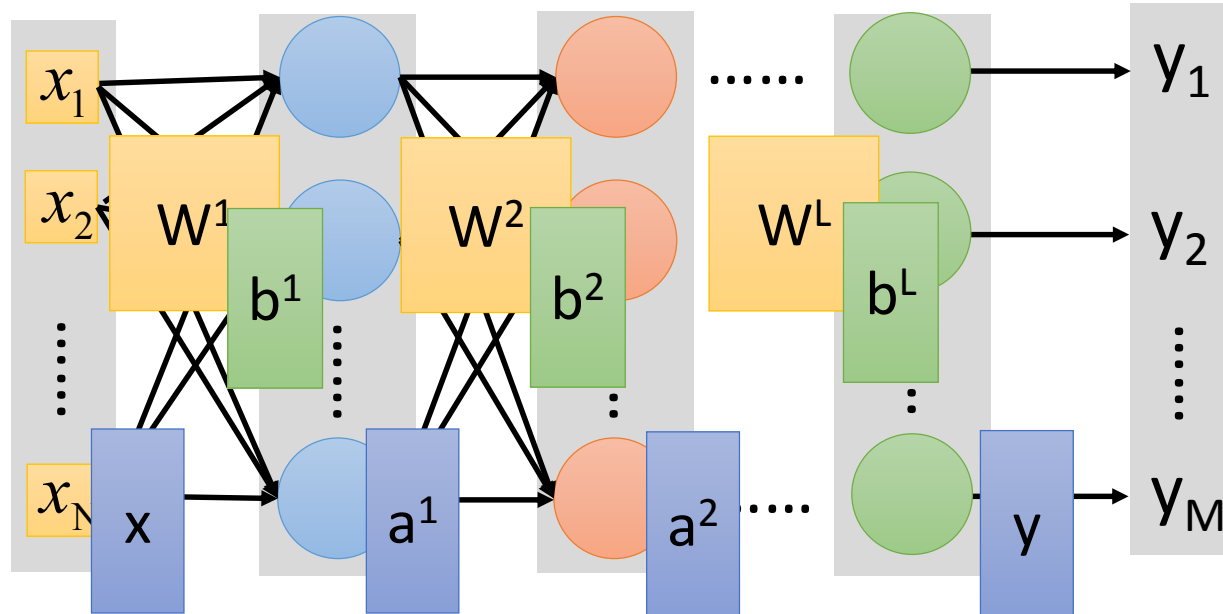


$$\sigma\left(\underbrace{\begin{bmatrix} 1 & -2 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \end{bmatrix}}_{\begin{bmatrix} 4 \\ -2 \end{bmatrix}}\right) = \begin{bmatrix} 0.98 \\ 0.12 \end{bmatrix}$$

# Neural Network



# Neural Network



$$y = f(x)$$

Using parallel computing techniques  
to speed up matrix operation

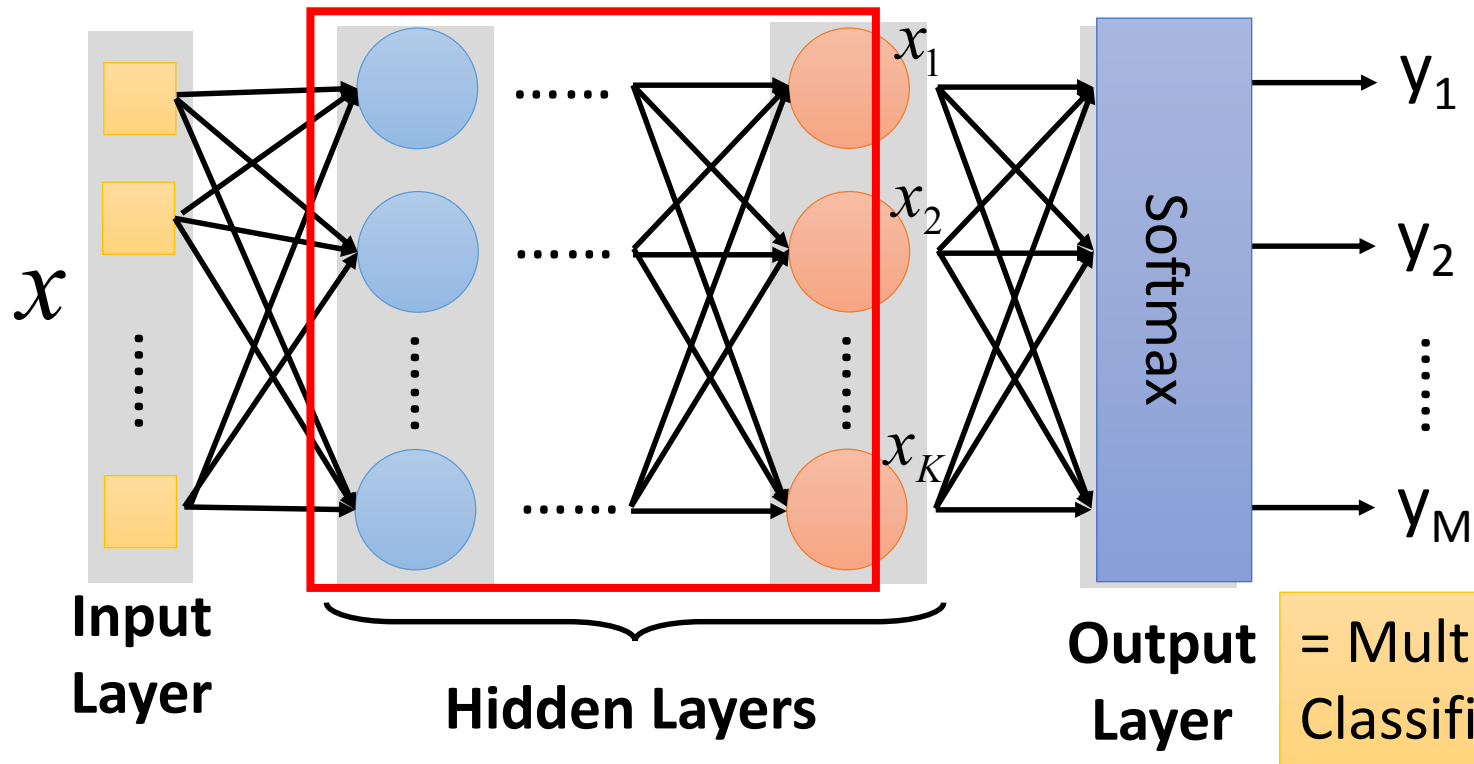
使用matrix operation利用GPU執行

$$= \sigma(W^L \dots \sigma(W^2 \sigma(W^1 x + b^1) + b^2) \dots + b^L)$$

# Output Layer as Multi-Class Classifier

Feature extractor replacing  
feature engineering

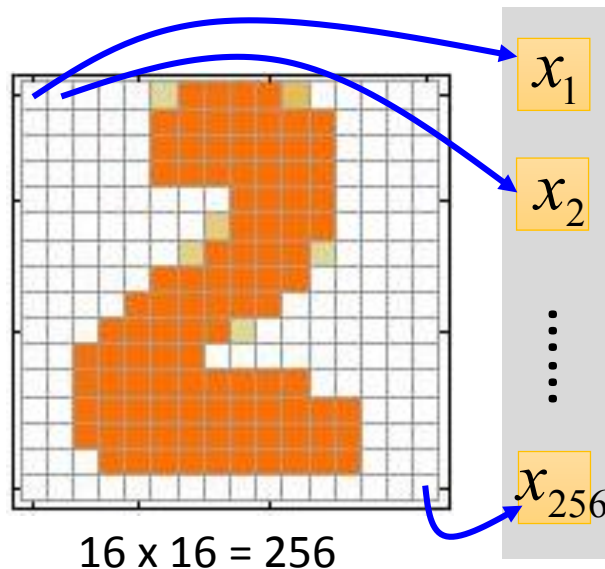
一般在output layer做multi-class  
classification時我們會採用  
softmax 作為activation function



# Example Application



## Input

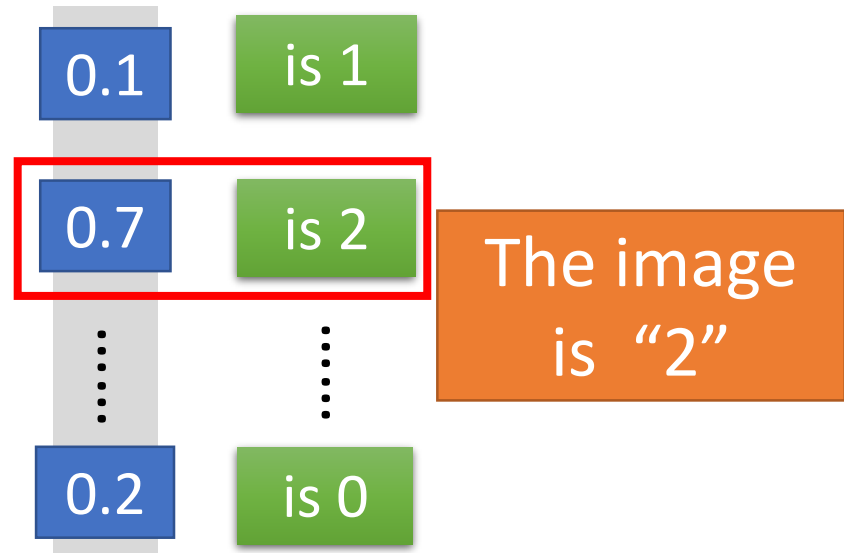


Ink  $\rightarrow$  1

No ink  $\rightarrow$  0

## Output

### Multi-class Classifier

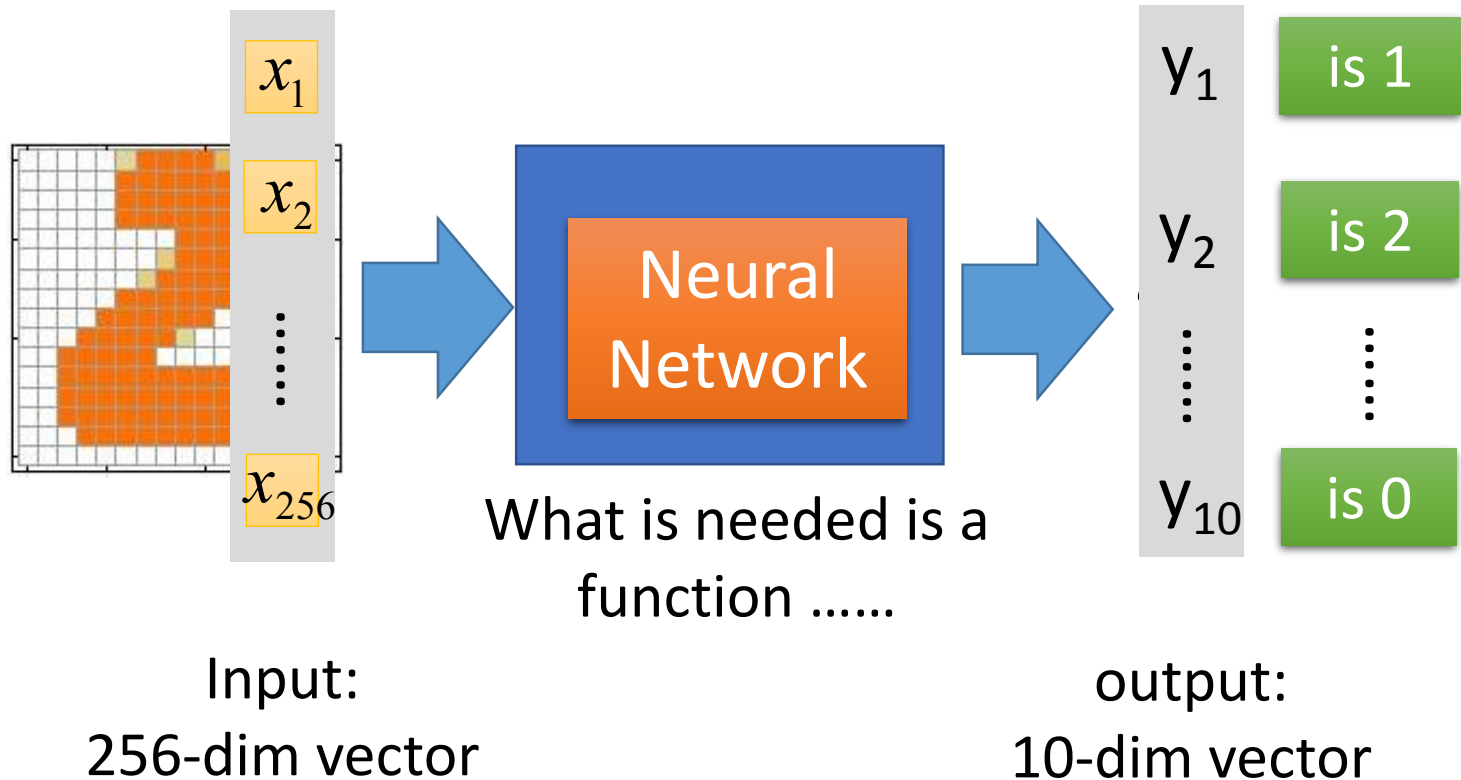


Each dimension represents the confidence of a digit.

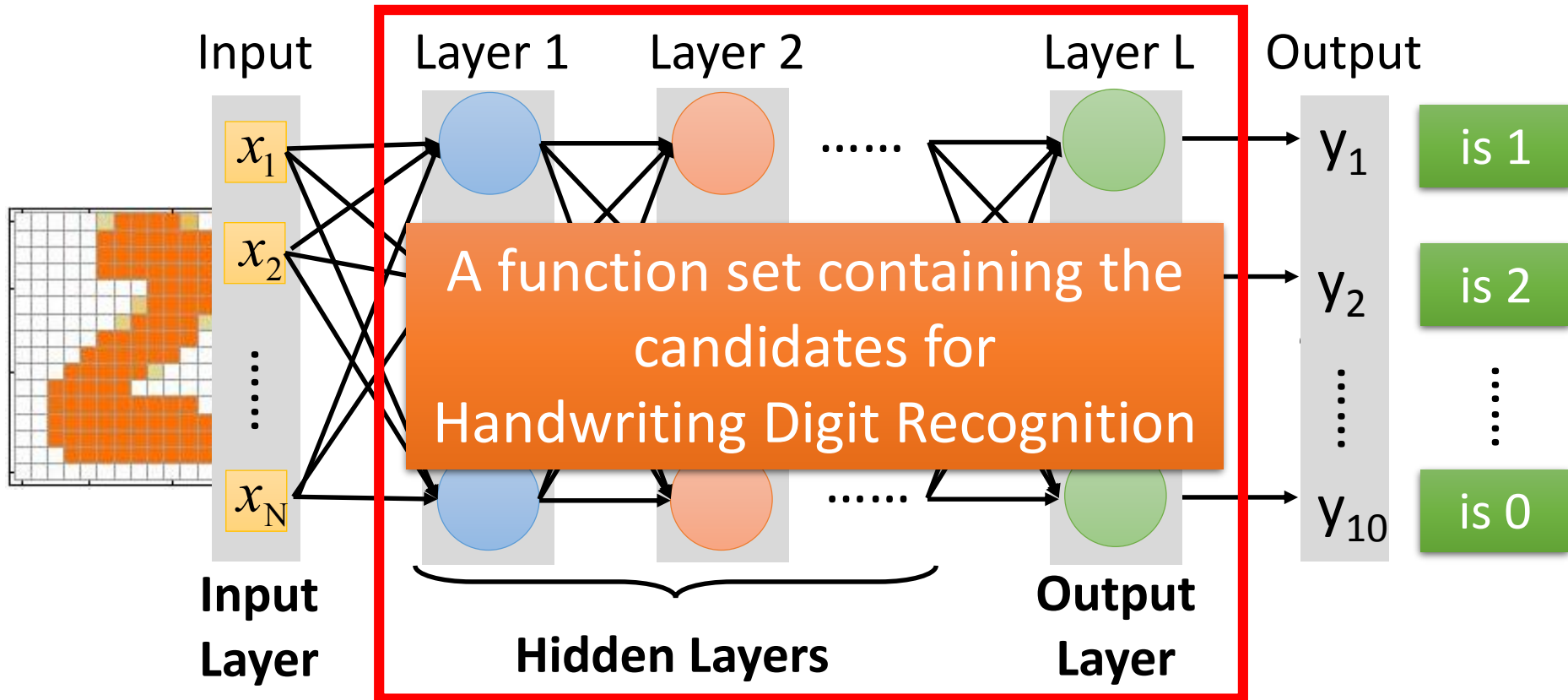


# Example Application

- Handwriting Digit Recognition

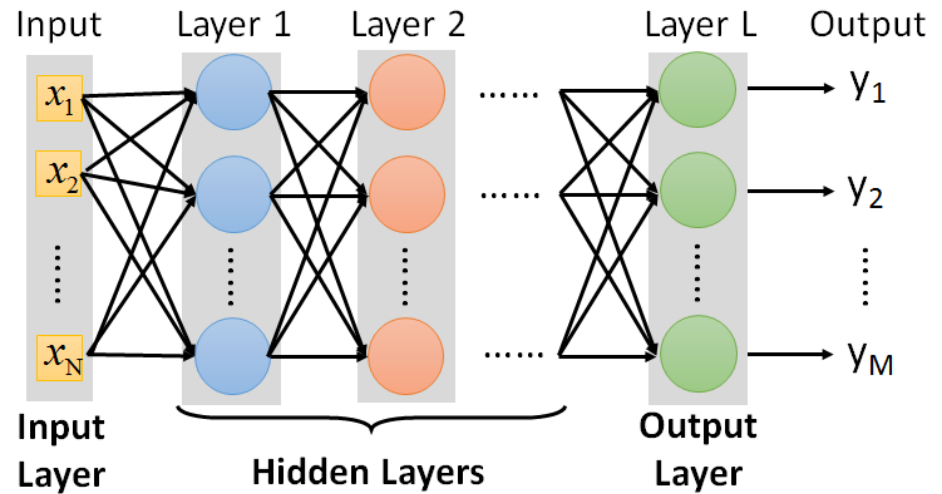


# Example Application



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?

經驗

Trial and Error

+

直覺

Intuition

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

Convolutional Neural Network (CNN)

# Three Steps for Deep Learning

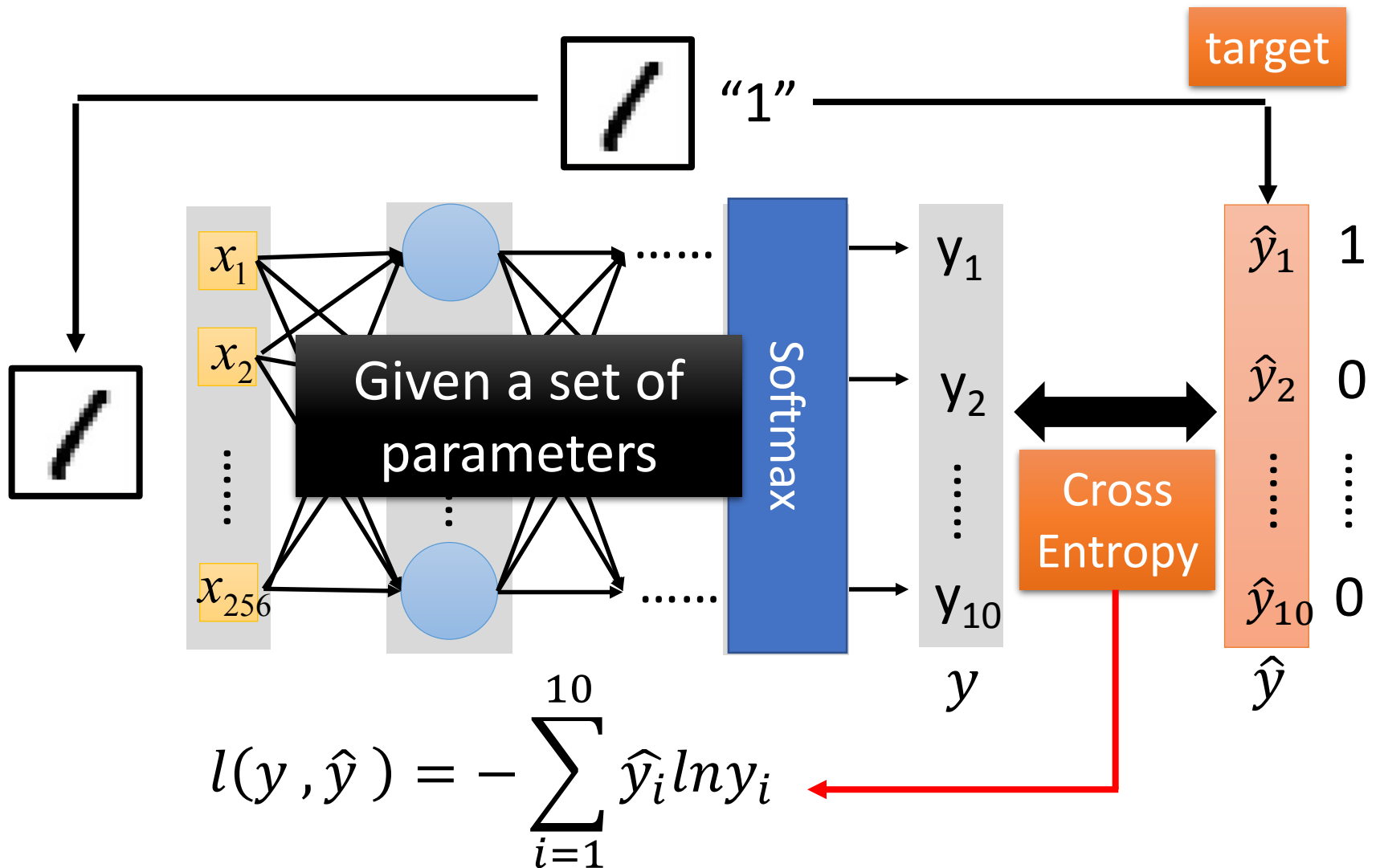


Deep Learning is so simple .....



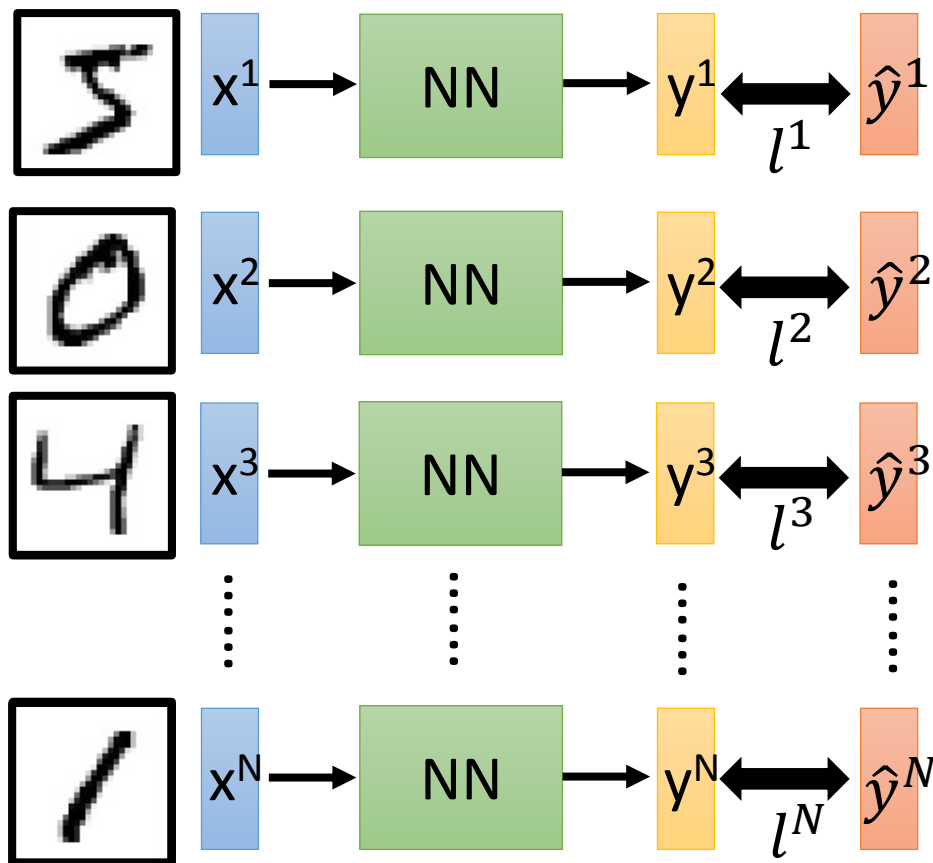
定義一組參數的好壞

# Loss for an Example



# Total Loss

For all training data ...



Total Loss:

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

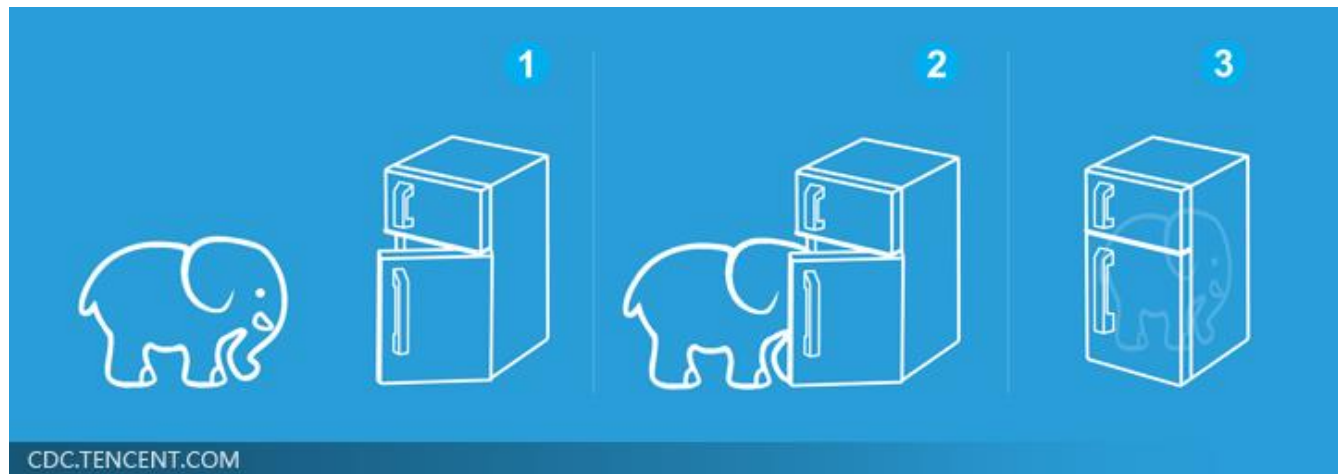
Find a function in function set that minimizes total loss  $L$

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

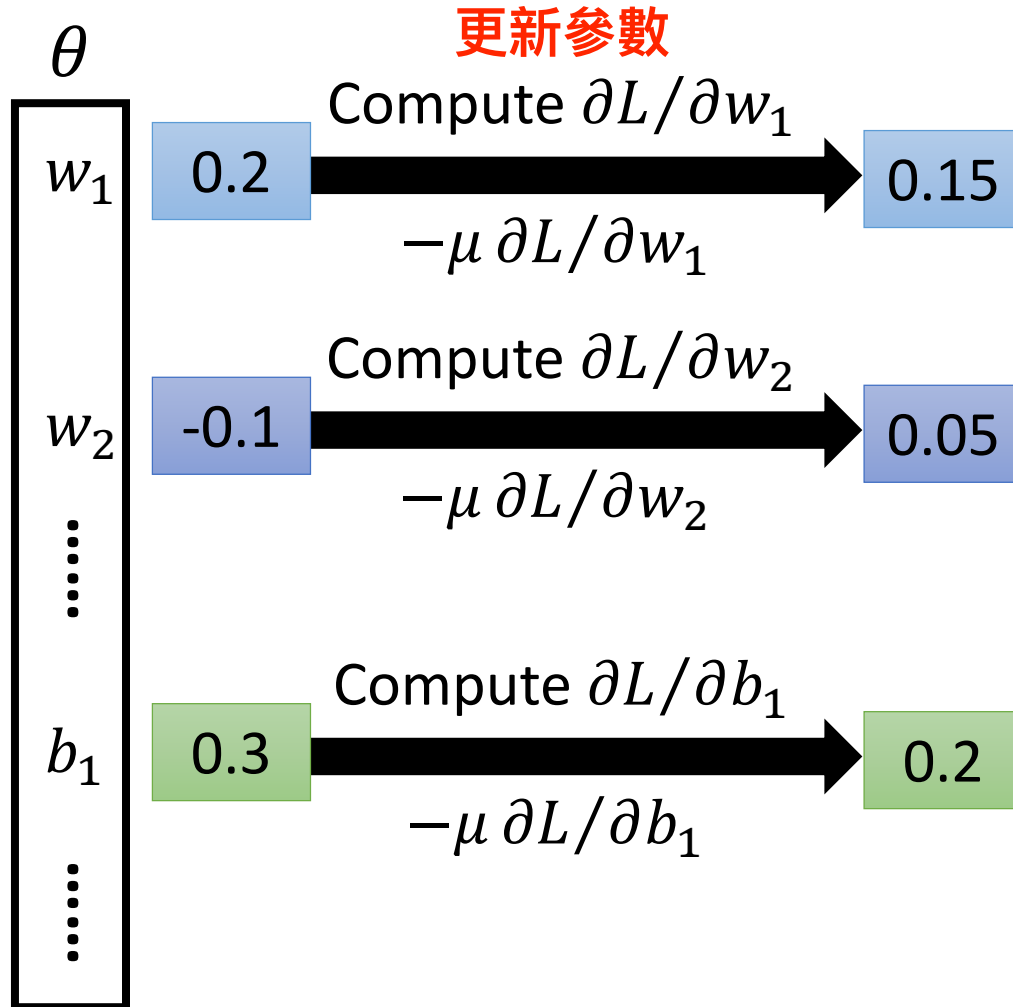
# Three Steps for Deep Learning



Deep Learning is so simple .....



# Gradient Descent

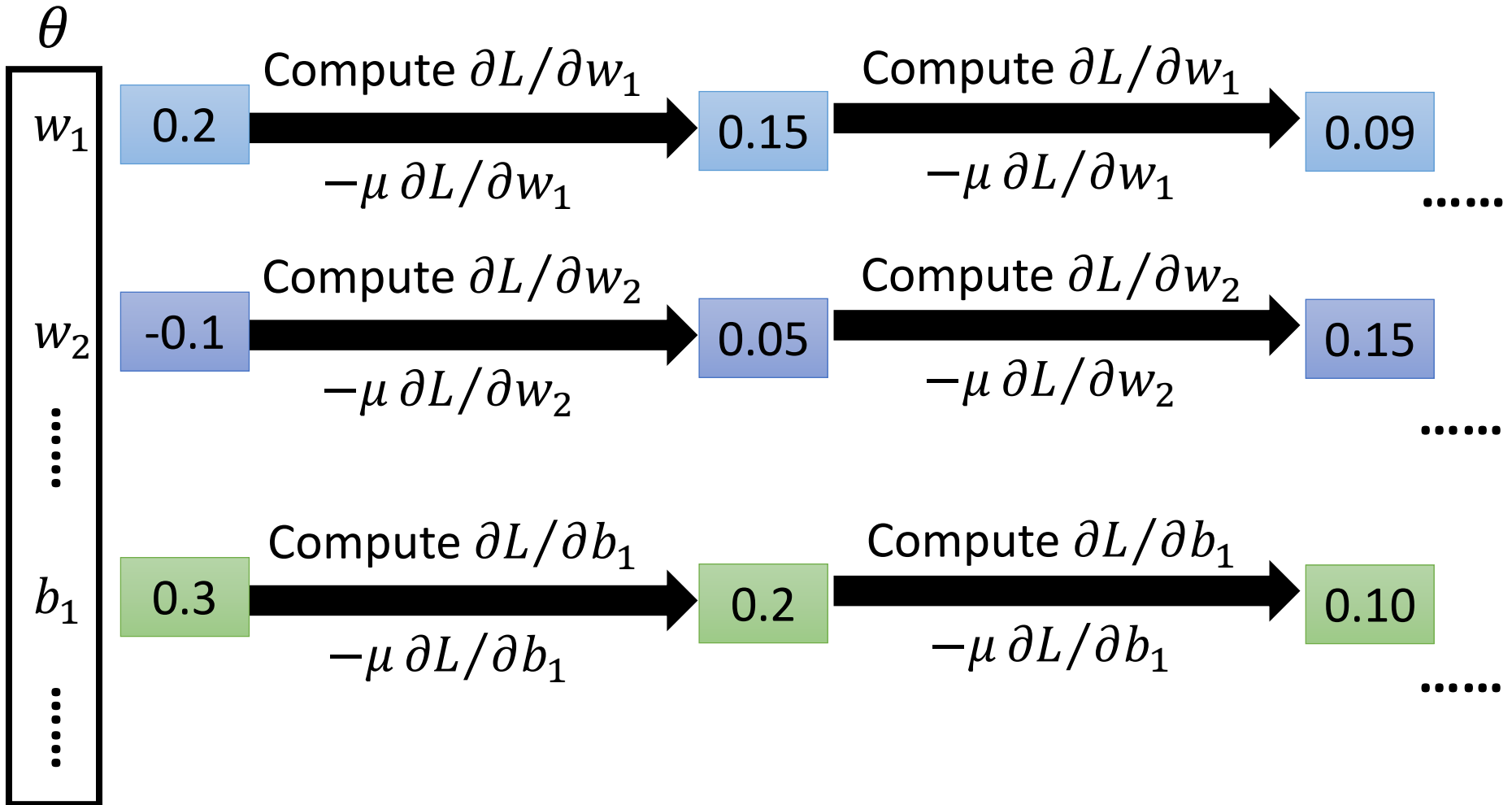


$$\nabla L = \begin{bmatrix} \frac{\partial L}{\partial w_1} \\ \frac{\partial L}{\partial w_2} \\ \vdots \\ \frac{\partial L}{\partial b_1} \\ \vdots \end{bmatrix}$$

gradient



# Gradient Descent

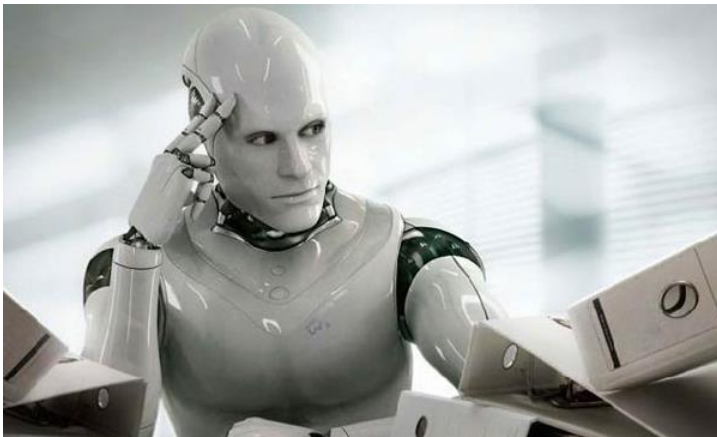


# Gradient Descent

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

比較有效率的算微分的方式

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



theano

Caffe



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

# Three Steps for Deep Learning



Deep Learning is so simple .....



# Acknowledgment

- 感謝 Victor Chen 發現投影片上的打字錯誤