# CSCE 636 Neural Networks (Deep Learning)

Lecture 8: Deep Learning for Text and Sequences

Anxiao (Andrew) Jiang

Based on the interesting lecture of Prof. Hung-yi Lee "Recurrent Neural Network"
https://www.youtube.com/watch?v=xCGidAeyS4M&list=PLJV_el3uVTsPy9oCRY30oBPNLCo89yu49&index=30

# Recurrent Neural Network (RNN)

# Example Application

- Slot Filling
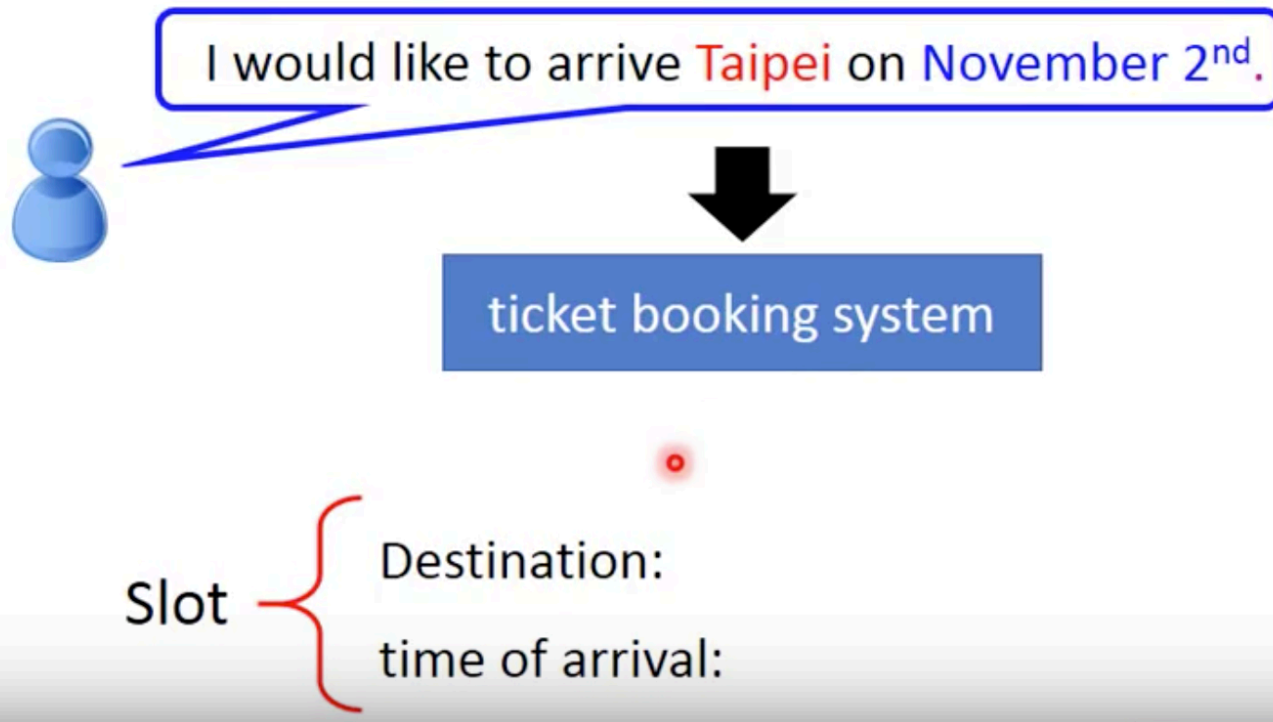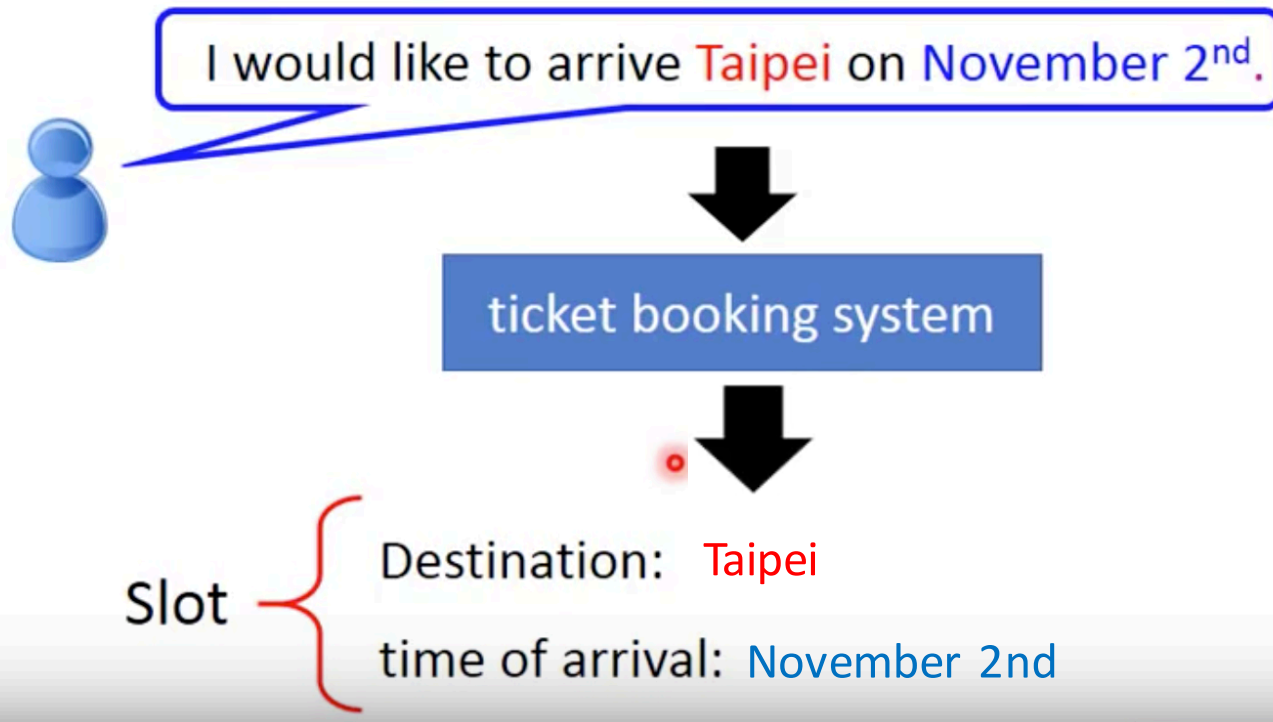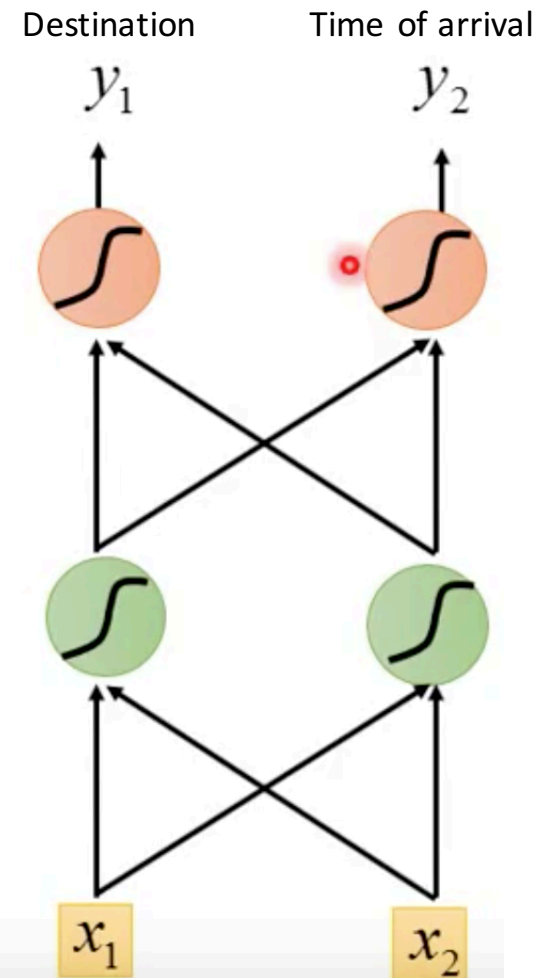
# Example Application

• Slot Filling

I would like to arrive Taipei on November 2nd.

ticket booking system

# Example Application

- Slot Filling

# Example Application

- Slot Filling

I would like to arrive **Taipei** on **November 2nd**.

ticket booking system

Slot
- Destination: **Taipei**
- time of arrival: **November 2nd**

# Example Application

Solving slot filling by
Feedforward network?

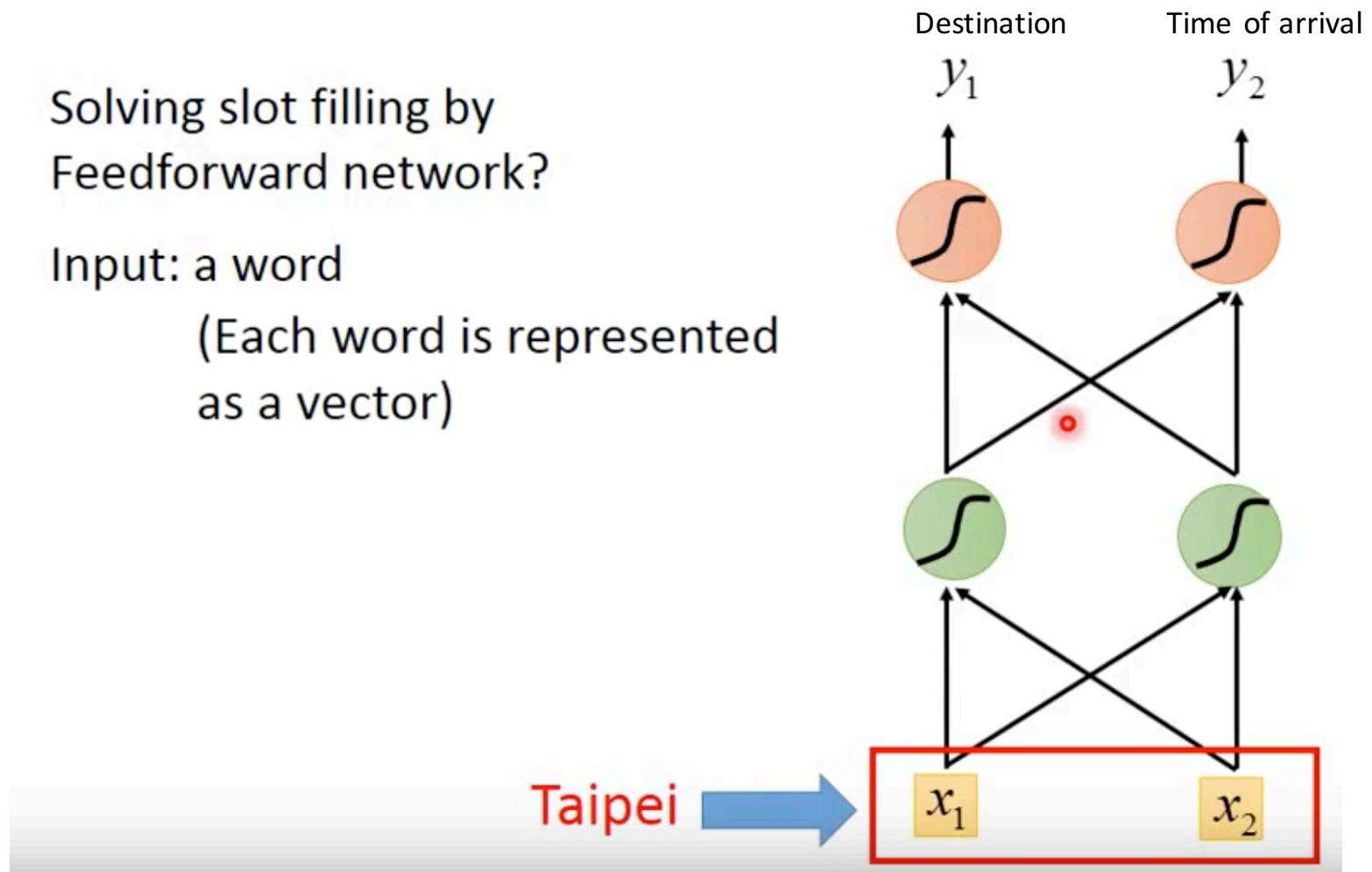Destination      Time of arrival

$y_1$         $y_2$

$x_1$         $x_2$

# Example Application

Solving slot filling by
Feedforward network?

Input: a word

(Each word is represented
as a vector)

Destination · Time of arrival

$y_1$ · $y_2$

Taipei · $x_1$ · $x_2$

# 1-of-N encoding (that is, one-hot encoding)

How to represent each word as a vector?

**1-of-N Encoding**    lexicon = {apple, bag, cat, dog, elephant}

The vector is lexicon size.

Each dimension corresponds to a word in the lexicon

The dimension for the word is 1, and others are 0

$$apple = [1 \ 0 \ 0 \ 0 \ 0]$$
$$bag \ = [0 \ 1 \ 0 \ 0 \ 0]$$
$$cat \ = [0 \ 0 \ 1 \ 0 \ 0]$$
$$dog \ = [0 \ 0 \ 0 \ 1 \ 0]$$
$$elephant \ = [0 \ 0 \ 0 \ 0 \ 1]$$

# Beyond 1-of-N encoding

## Dimension for "Other"



**Dense word embedding**

Word2Vec (pre-trained)

GloVe (pre-trained)

Train your own embedding

# Example Application

Destination      Time of arrival

$y_1$          $y_2$

Solving slot filling by Feedforward network?

Input: a word
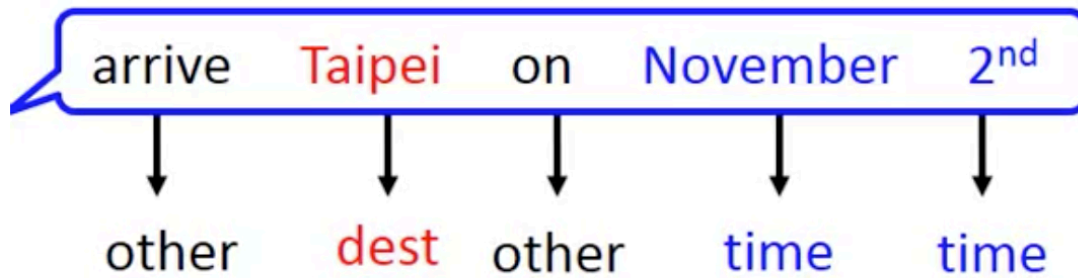
(Each word is represented as a vector)

Output:

Probability distribution that the input word belonging to the slots
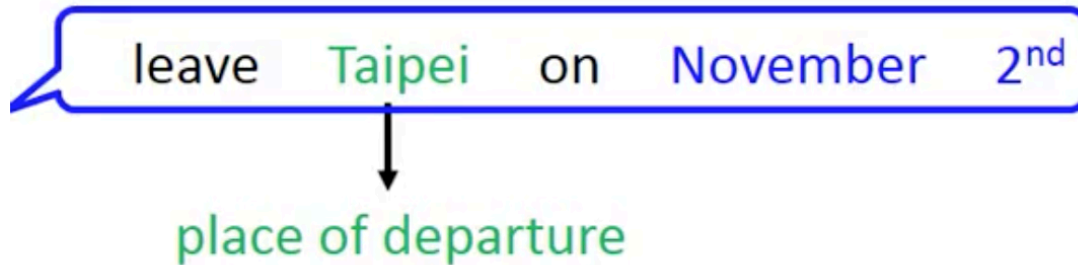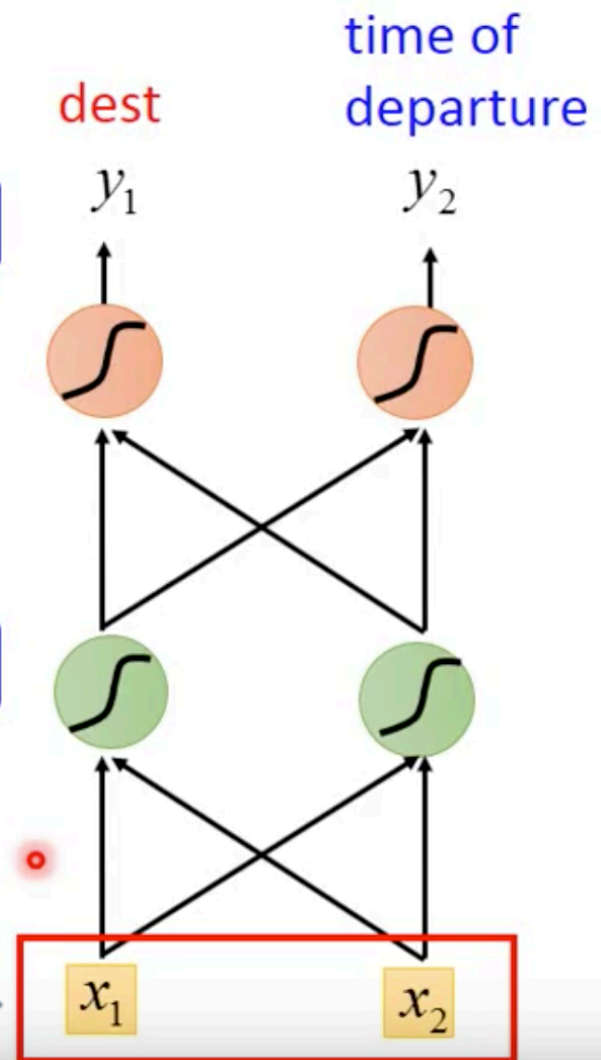
Taipei ➡ $x_1$      $x_2$

# Example Application

arrive    Taipei    on    November    2nd

other    dest    other    time    time

**Problem?**

leave    Taipei    on    November    2nd

place of departure

dest

time of departure

$y_1$      $y_2$

Taipei $\Rightarrow$ $x_1$    $x_2$

# Recurrent Neural Network (RNN)

The output of hidden layer are stored in the memory.

$y_1$

$y_2$

$a_1$

$a_2$

$x_1$

$x_2$

# Recurrent Neural Network (RNN)

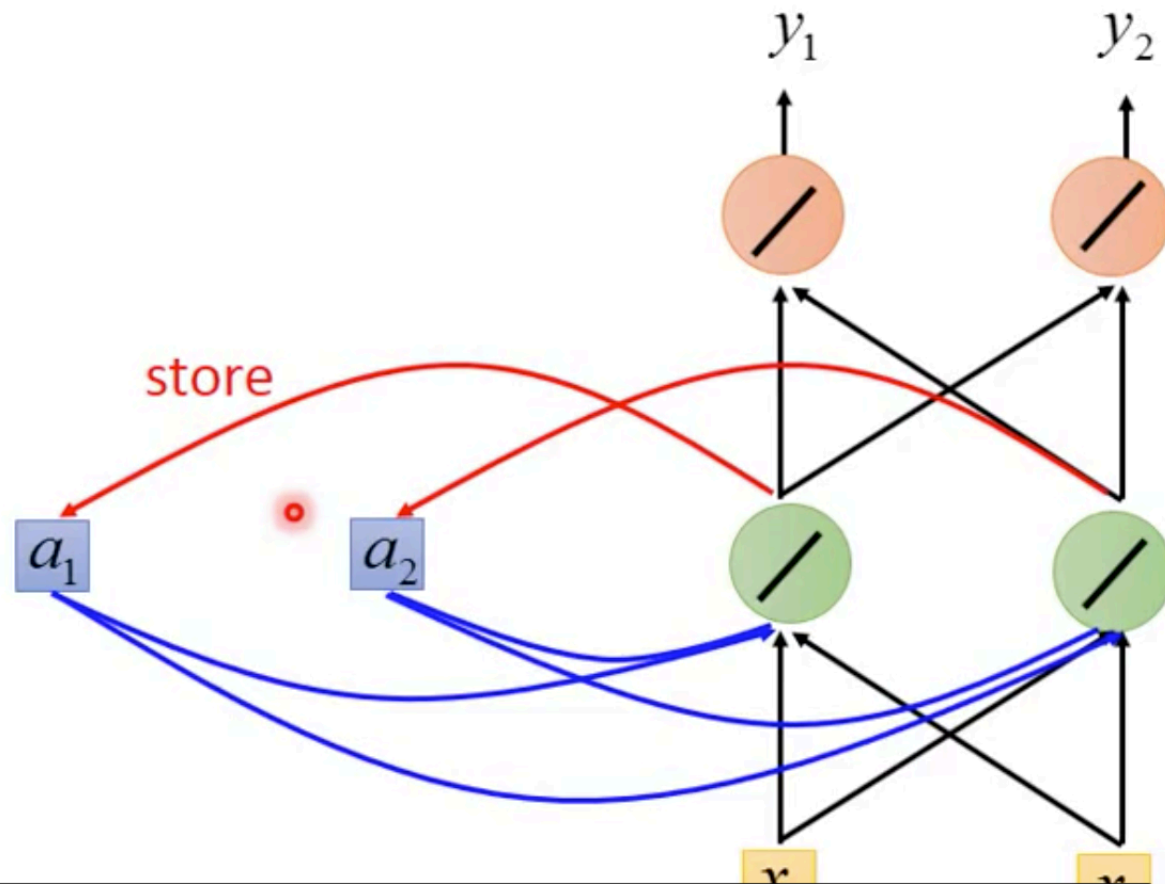The output of hidden layer are stored in the memory.



這邊呢，用藍色的方塊來表示 memory

# Recurrent Neural Network (RNN)

The output of hidden layer are stored in the memory.
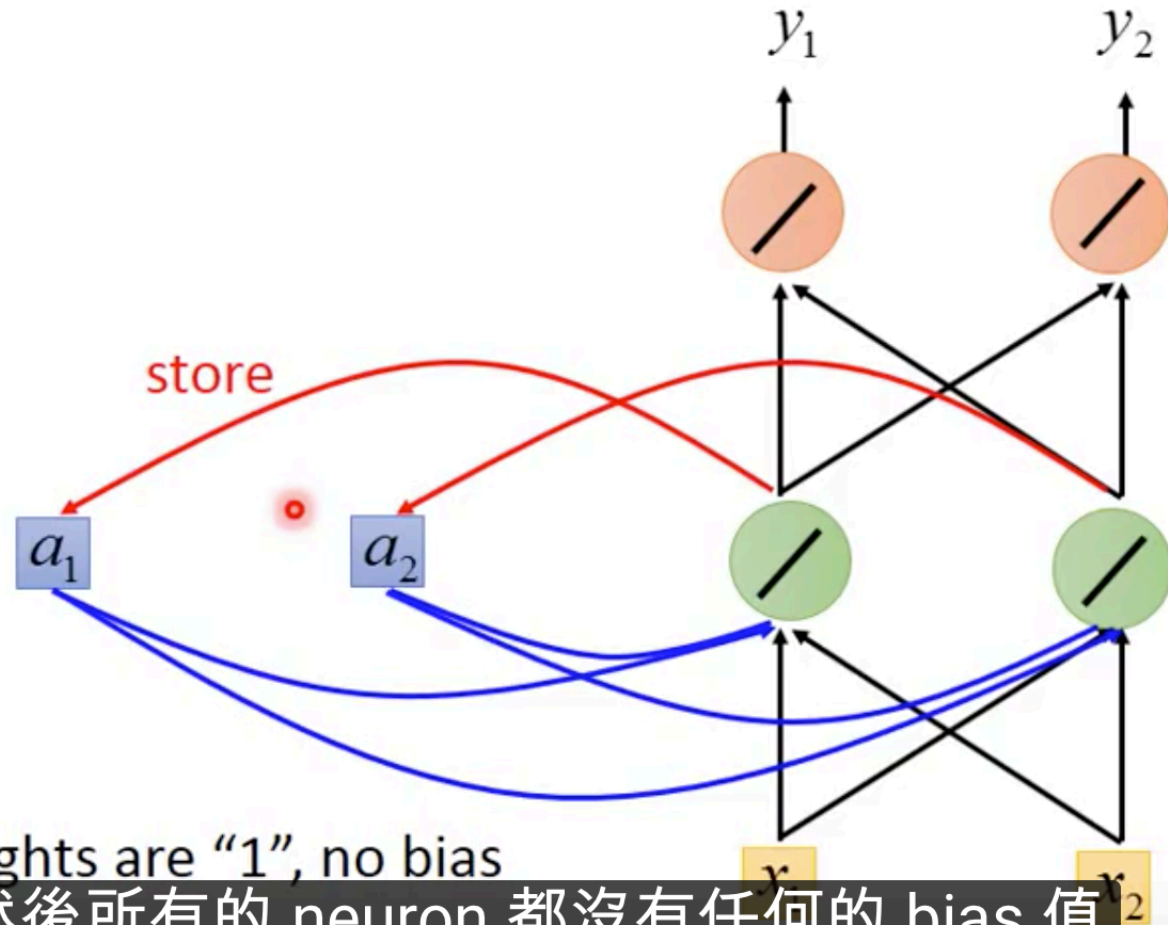
store

$a_1$ $a_2$

Memory can be considered as another input.

$y_1$ $y_2$

$x_1$ $x_2$

# Example



那我想直接舉個例子，大家可能會比較清楚

# Example



$y_1$　　　$y_2$

store

$a_1$　○　$a_2$

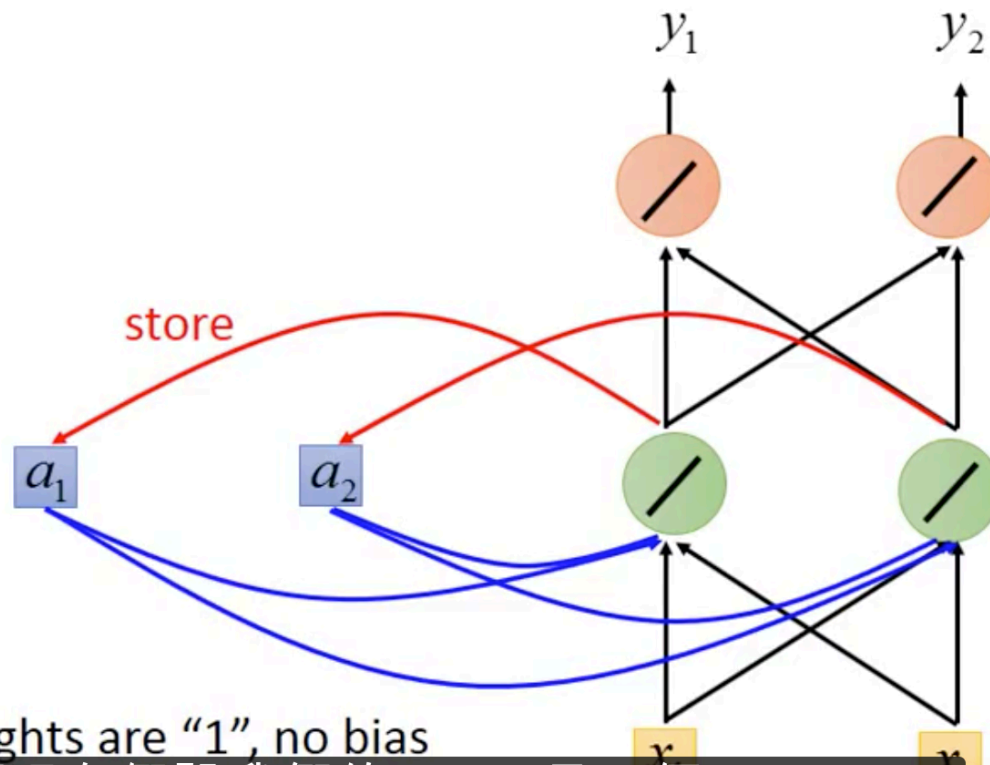All the weights are "1", no bias

然後所有的 neuron 都沒有任何的 bias 值

# Example



All the weights are "1", no bias

這樣可以不要讓計算太複雜

All activation functions are linear

# Example

Input sequence: $\begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \end{bmatrix}$ …… ……
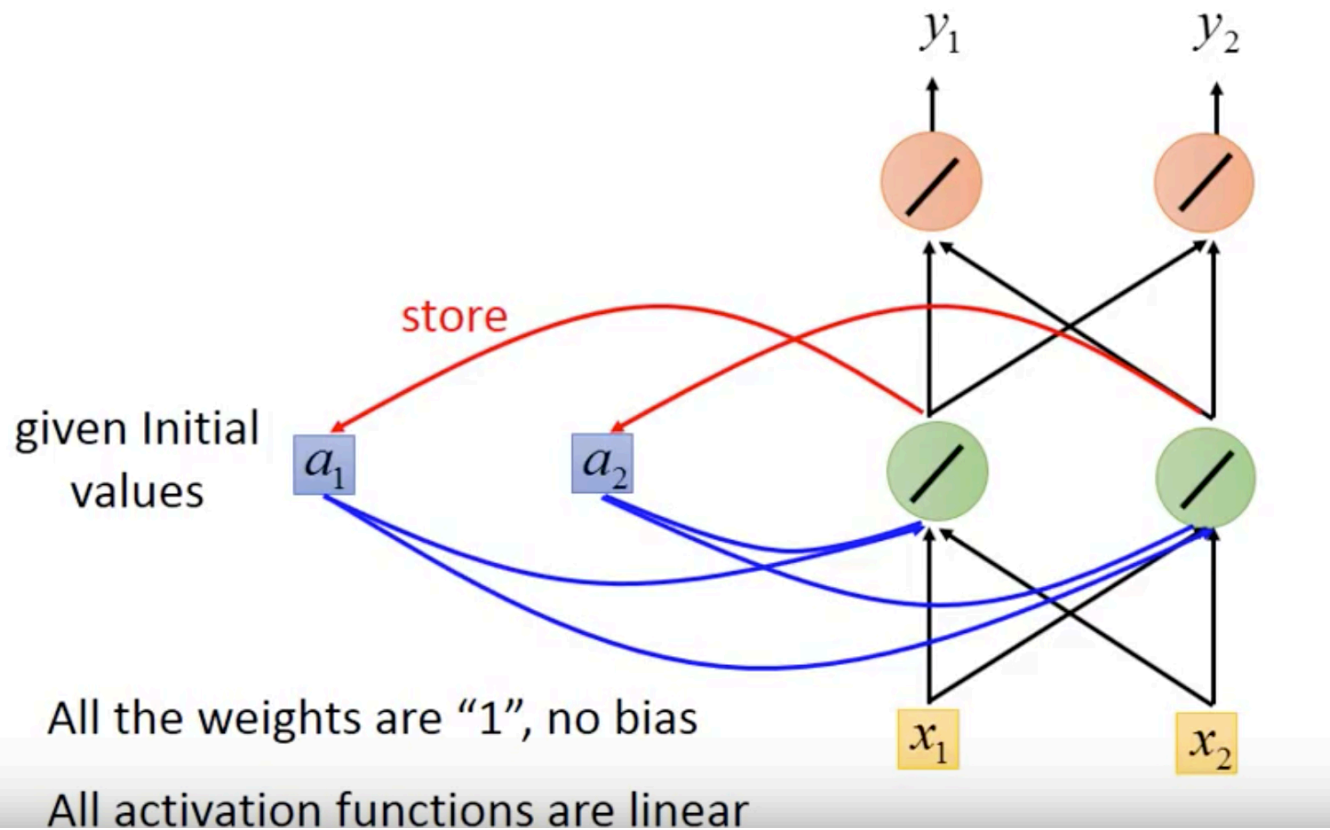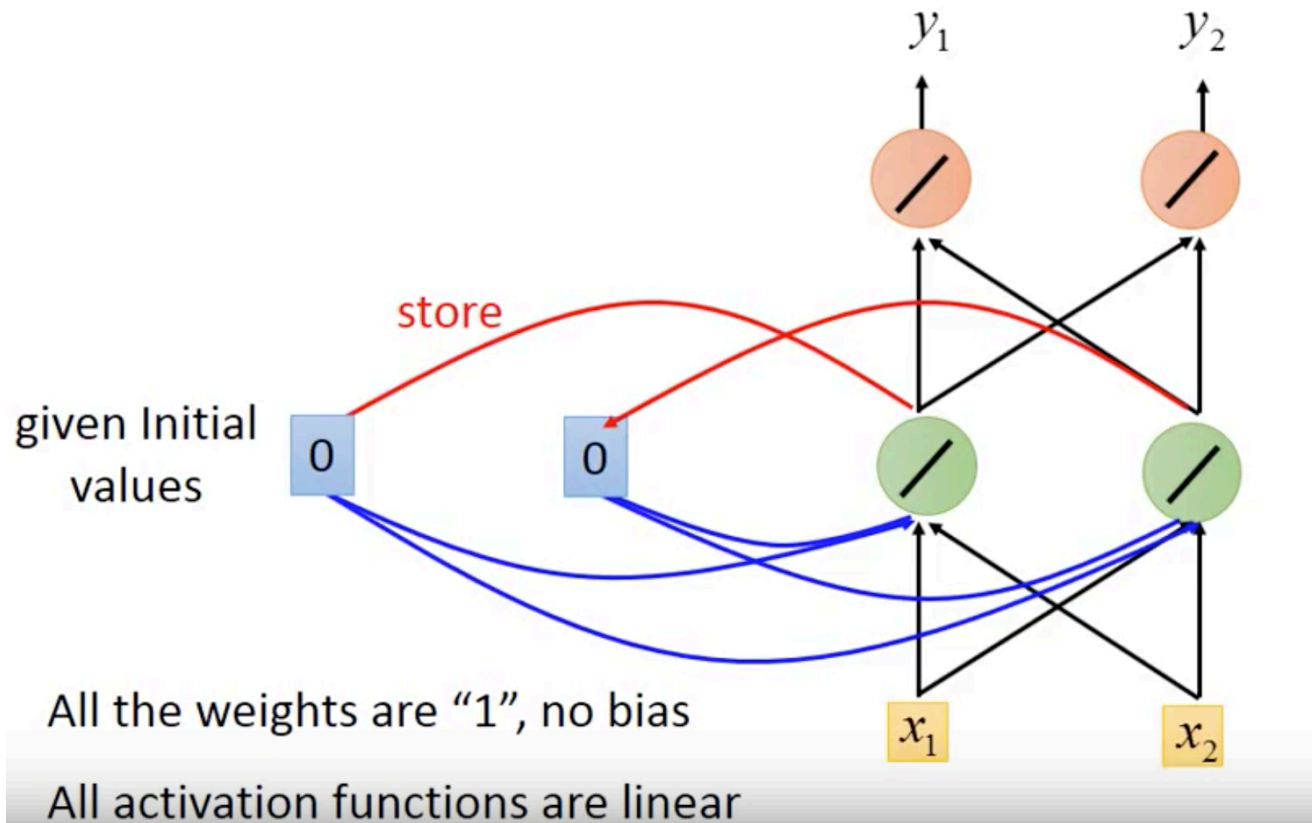


All the weights are "1", no bias

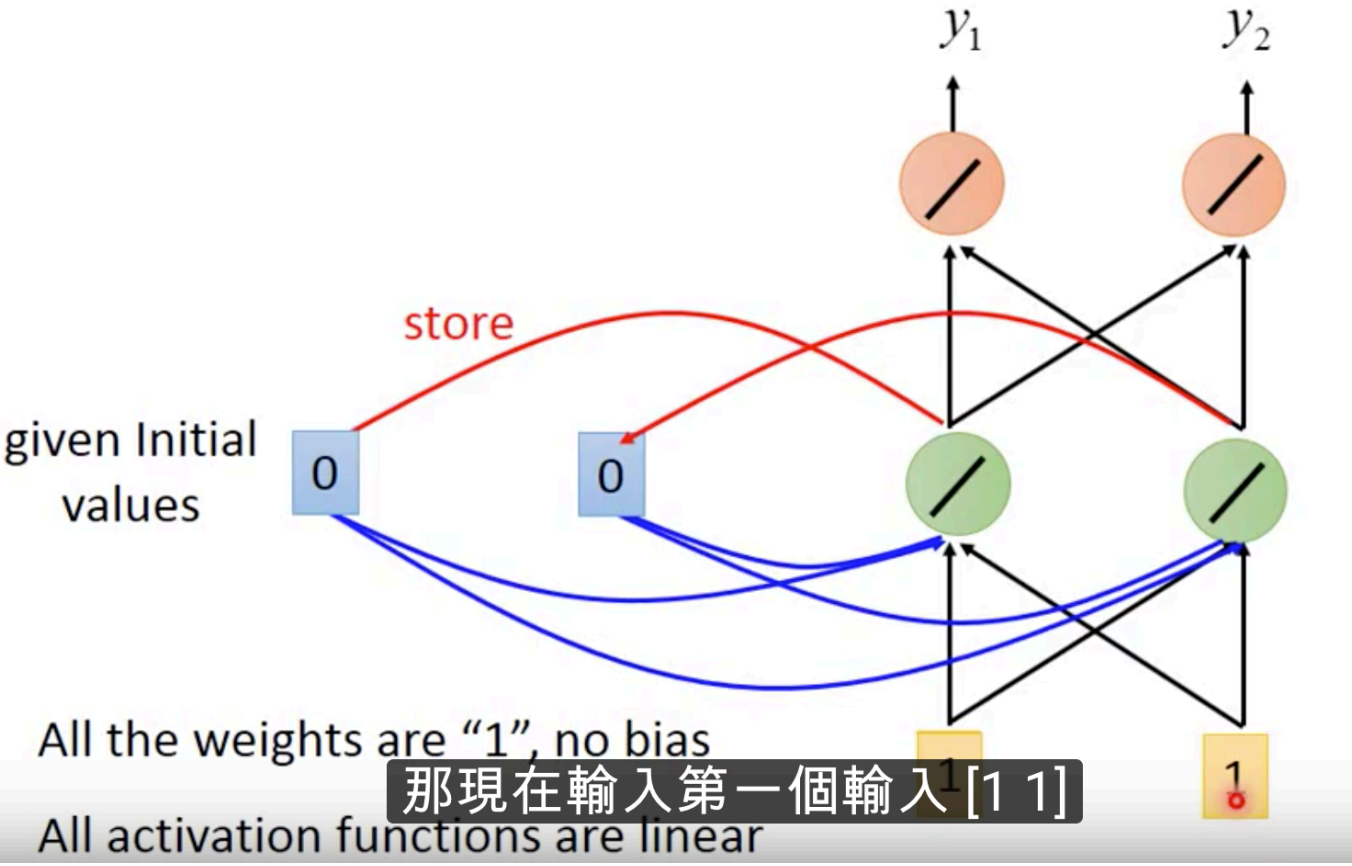All activation functions are linear

那現在假設我們的 input 是一個 sequence

Input sequence: $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$ $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$ $\begin{bmatrix} 2 \\ 2 \end{bmatrix}$ ... ...
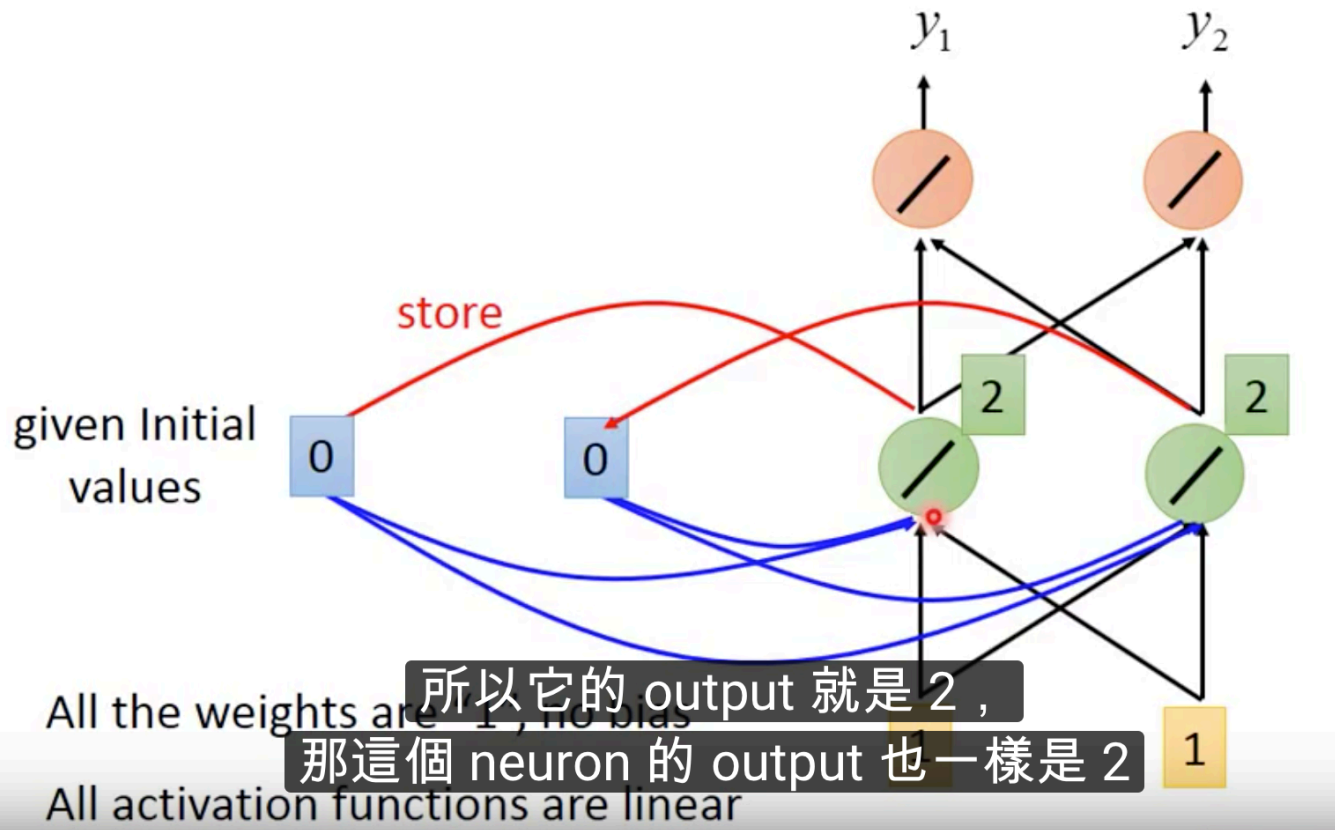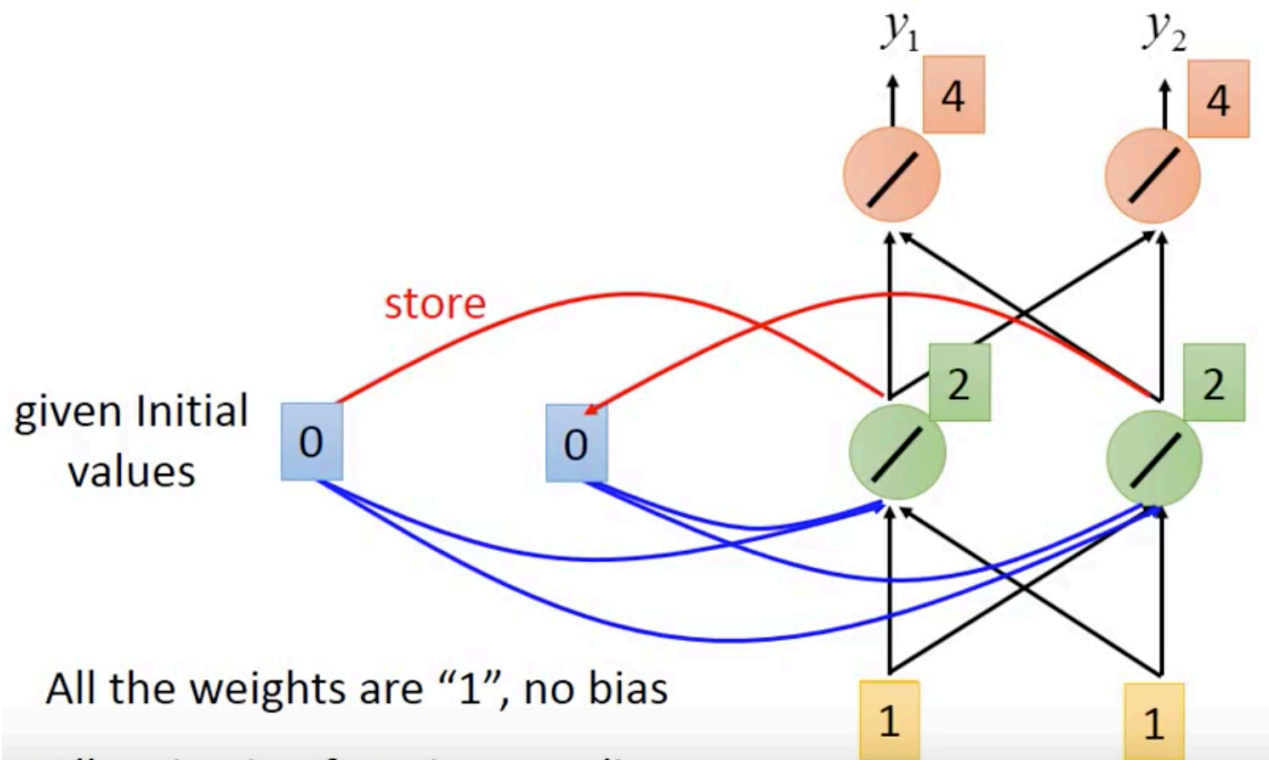
# Example



given Initial values
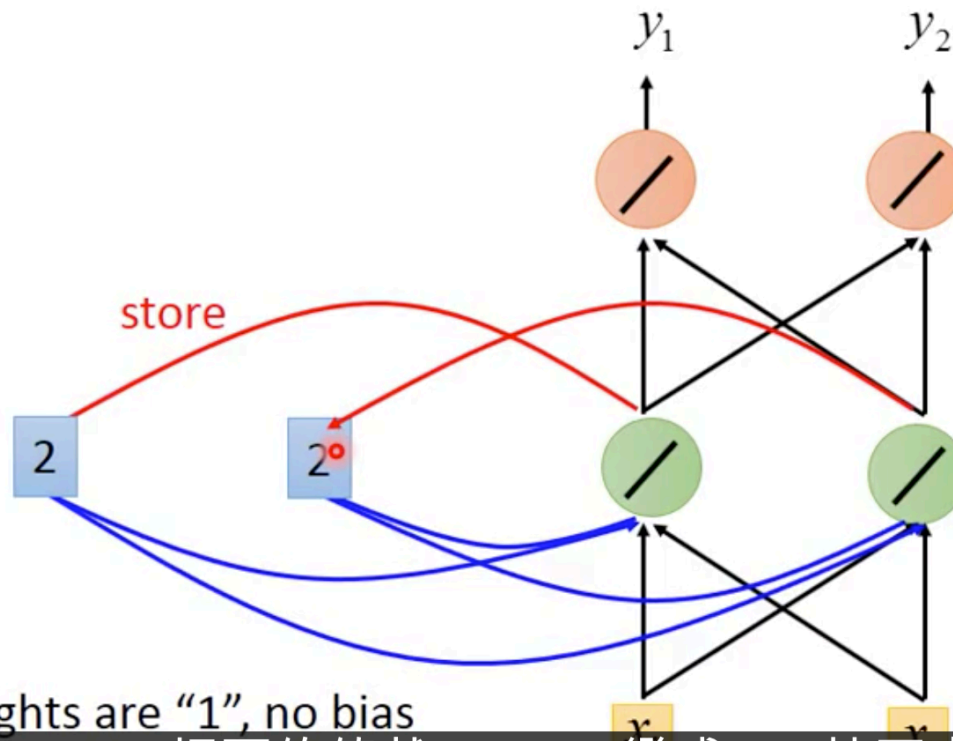
store

$a_1$  $a_2$

$y_1$  $y_2$

$x_1$  $x_2$

All the weights are "1", no bias
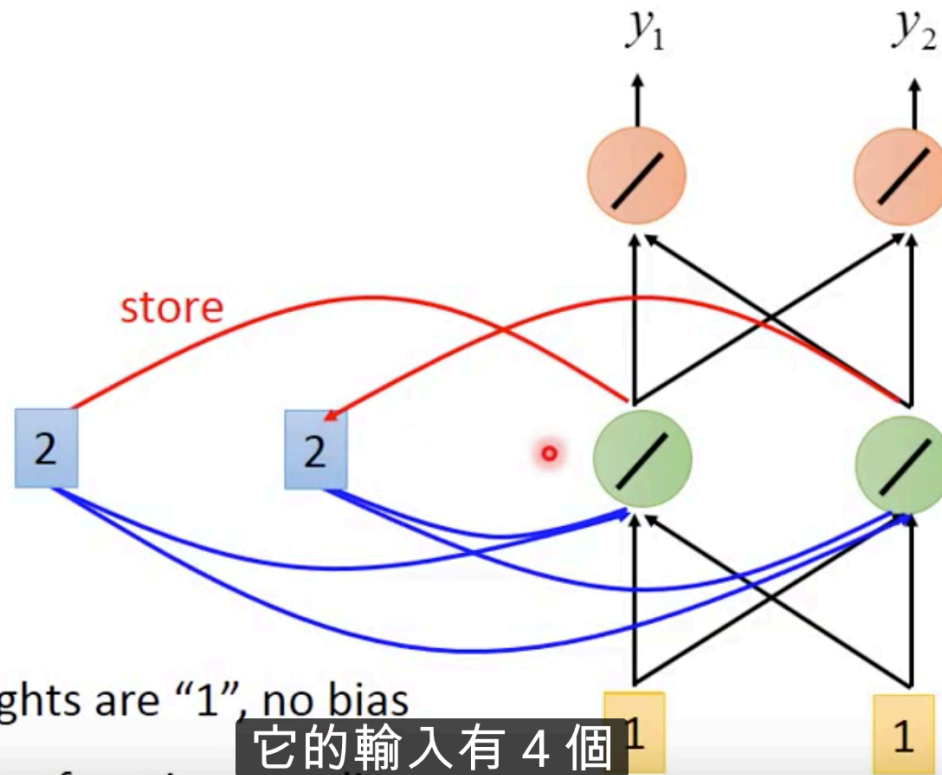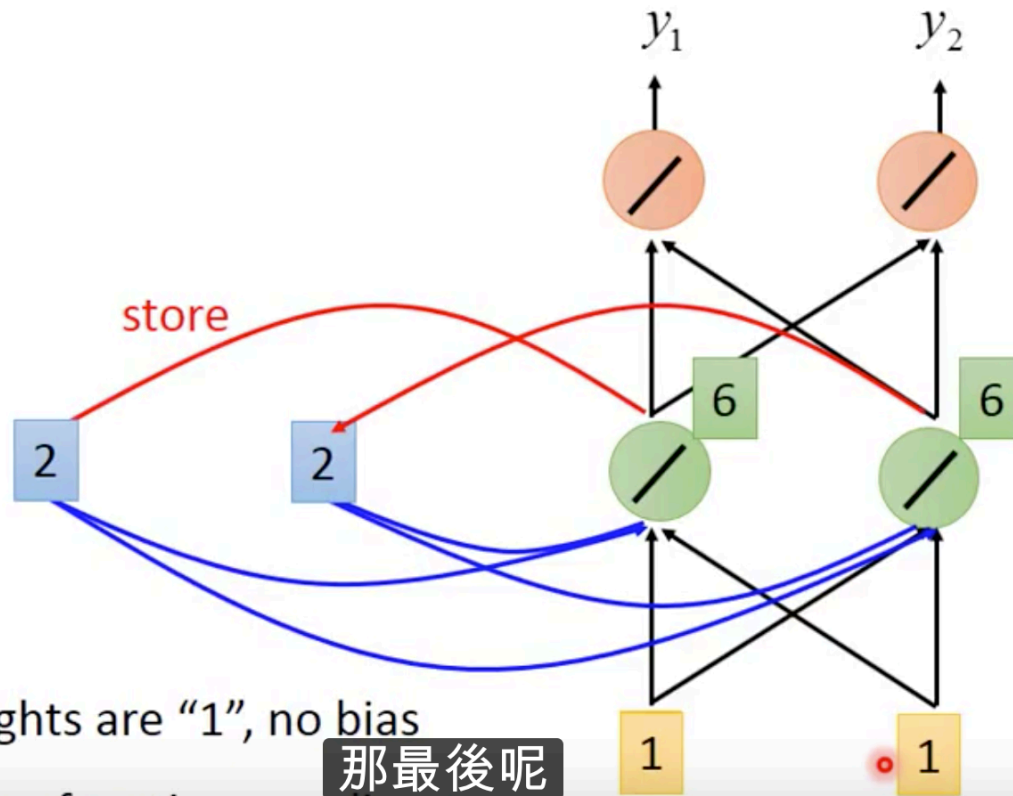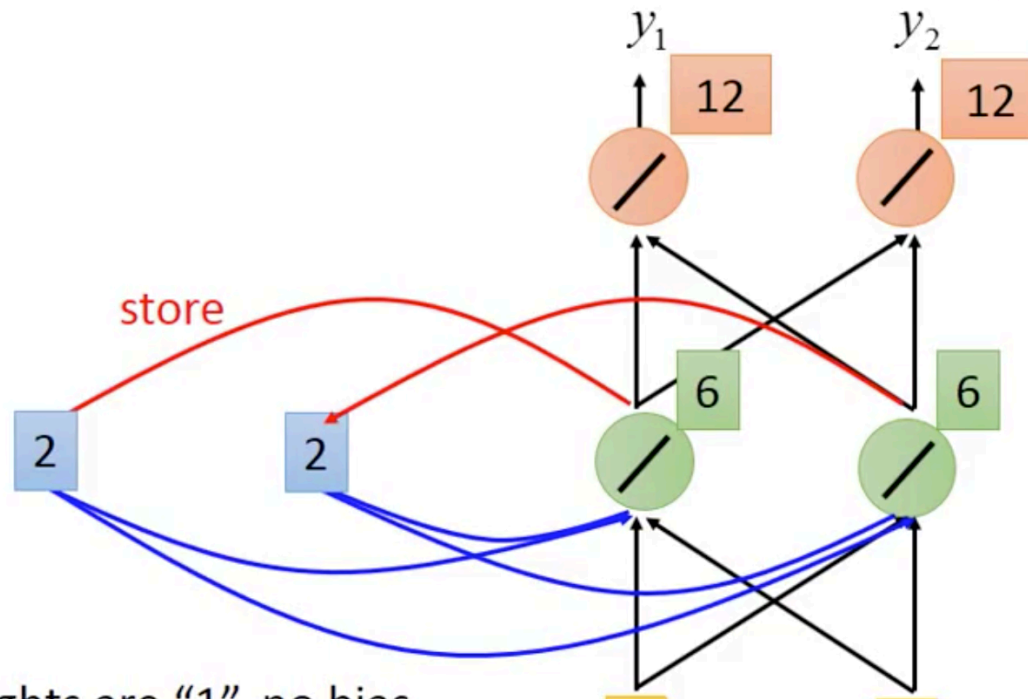
All activation functions are linear

# Example

Input sequence: $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$ $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$ $\begin{bmatrix} 2 \\ 2 \end{bmatrix}$ ... ...



$y_1$    $y_2$

store

given Initial values

0    0

All the weights are "1", no bias

All activation functions are linear

$x_1$    $x_2$

Input sequence: $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$ $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$ $\begin{bmatrix} 2 \\ 2 \end{bmatrix}$ ... ...

# Example

$y_1$　　$y_2$

store

given Initial values

0　　0

All the weights are "1", no bias

All activation functions are linear

1　　1

那現在輸入第一個輸入 [1 1]

Input sequence: $\begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \end{bmatrix}$ ... ...

# Example



given Initial values

store

$y_1$     $y_2$

All the weights are 1, no bias

All activation functions are linear

所以它的 output 就是 2，
那這個 neuron 的 output 也一樣是 2

# Example

Input sequence: $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$ $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$ $\begin{bmatrix} 2 \\ 2 \end{bmatrix}$ ... ...

output sequence: $\begin{bmatrix} 4 \\ 4 \end{bmatrix}$
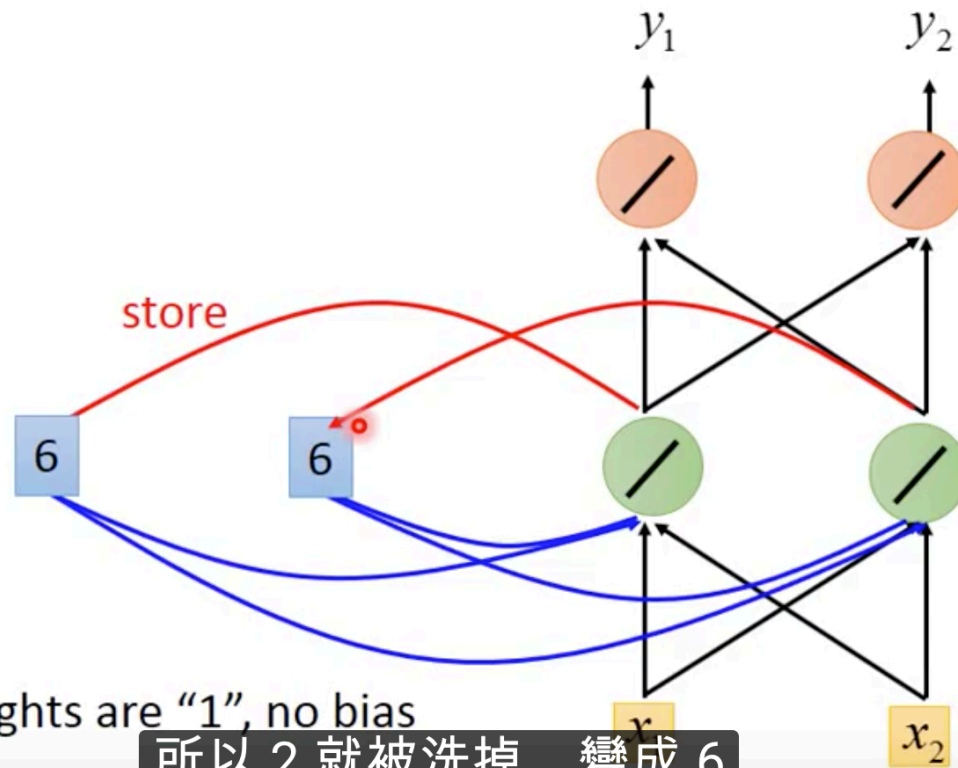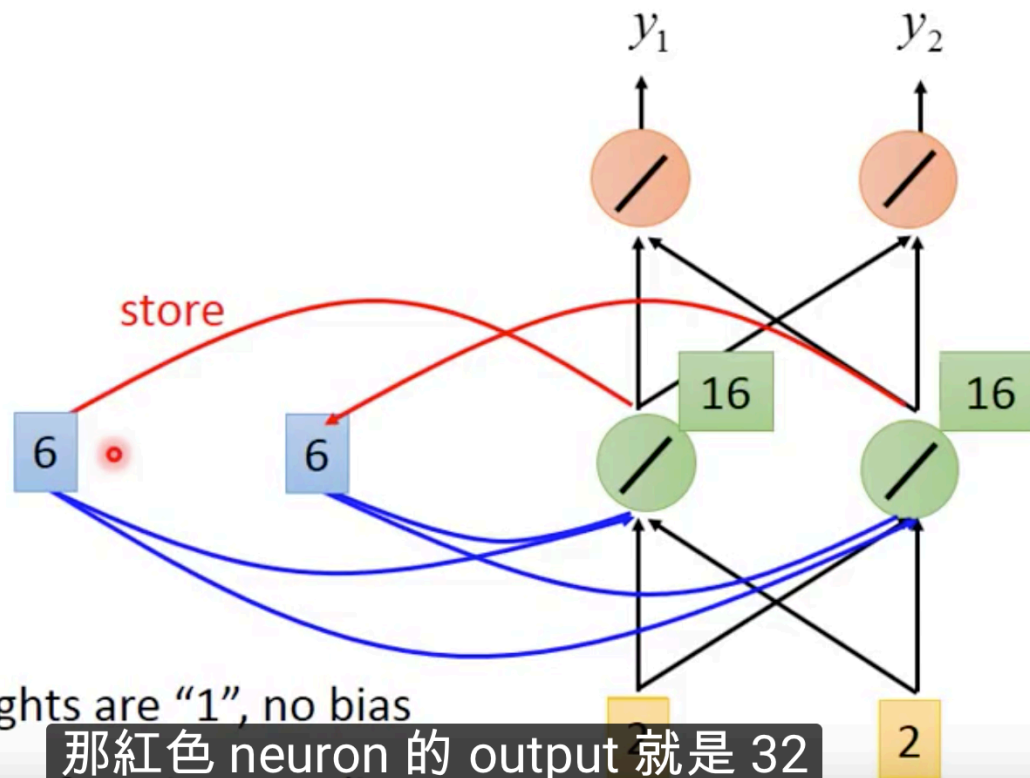


given Initial values

store

All the weights are "1", no bias
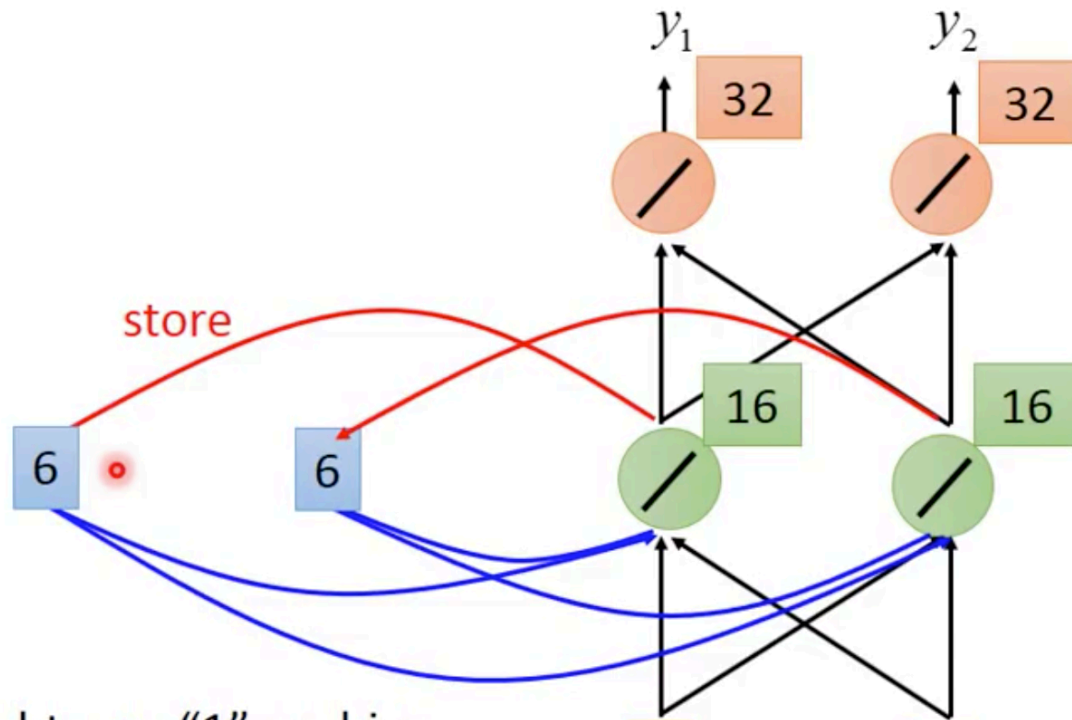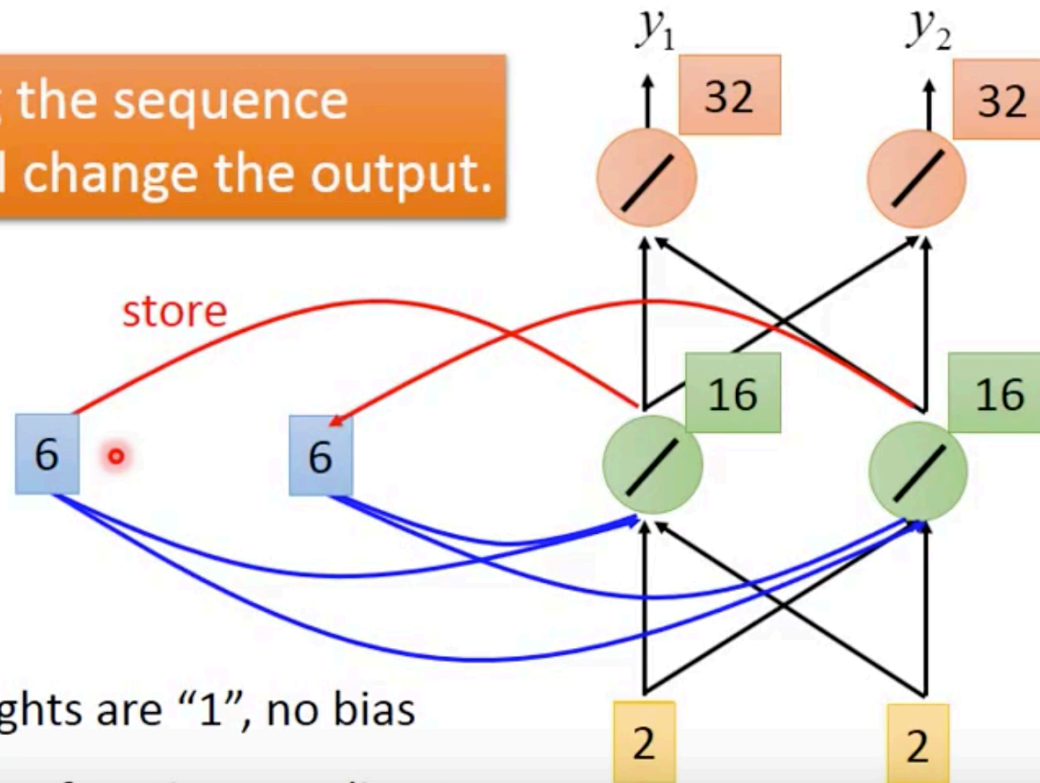
All activation functions are linear

# Example

Input sequence: $\begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \end{bmatrix}$ ... ...

output sequence: $\begin{bmatrix} 4 \\ 4 \end{bmatrix}$



store

All the weights are "1", no bias

All activation functions are linear

所以 memory 裡面的值就 update 變成 2，接下來呢

# Example

Input sequence: $\begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \end{bmatrix}$ … …

output sequence: $\begin{bmatrix} 4 \\ 4 \end{bmatrix}$



$y_1$ $y_2$

store

2    2

All the weights are "1", no bias

All activation functions are linear

它的輸入有 4 個

1    1

# Example

Input sequence: $\begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \end{bmatrix}$ ... ...

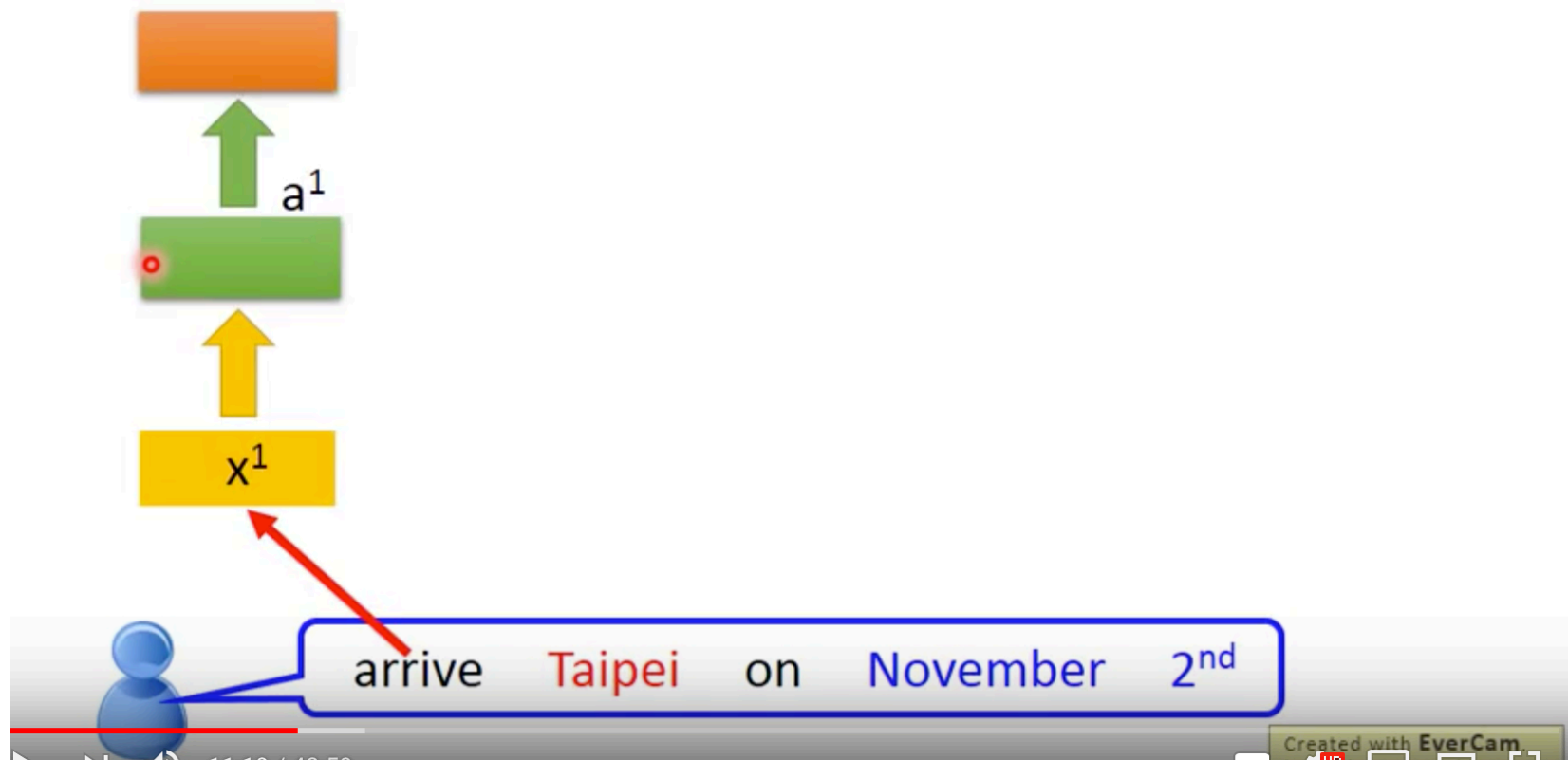output sequence: $\begin{bmatrix} 4 \\ 4 \end{bmatrix}$



$y_1$ $y_2$

store

2    2    6    6

All the weights are "1", no bias

1    1

那最後呢

All activation functions are linear

# Example

Input sequence: $\begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \end{bmatrix}$ ... ...

output sequence: $\begin{bmatrix} 4 \\ 4 \end{bmatrix} \begin{bmatrix} 12 \\ 12 \end{bmatrix}$



All the weights are "1", no bias

All activation functions are linear

第二次再輸 [1 1] 的時候，輸出就是 [12 12]

# Example

Input sequence: $\begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \end{bmatrix}$ ......

output sequence: $\begin{bmatrix} 4 \\ 4 \end{bmatrix} \begin{bmatrix} 12 \\ 12 \end{bmatrix}$



All the weights are "1", no bias

All activation functions are linear

所以 2 就被洗掉，變成 6

# Example

Input sequence: $\begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \end{bmatrix} \cdots \cdots$

output sequence: $\begin{bmatrix} 4 \\ 4 \end{bmatrix} \begin{bmatrix} 12 \\ 12 \end{bmatrix}$



store

$y_1$    $y_2$

16    16

6    6

All the weights are "1", no bias

All activation functions are linear

那紅色 neuron 的 output 就是 32

2

# Example

Input sequence: $\begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \end{bmatrix}$ ... ...

output sequence: $\begin{bmatrix} 4 \\ 4 \end{bmatrix} \begin{bmatrix} 12 \\ 12 \end{bmatrix} \begin{bmatrix} 32 \\ 32 \end{bmatrix}$



All the weights are "1", no bias

All activation functions are linear

所以 input 2 跟 2 的時候呢，output 是 32

# RNN



arrive  Taipei  on  November  2nd

# RNN

Probability of
"arrive" in each slot



store

$y^1$

$a^1$

$a^1$

$x^1$

arrive   Taipei   on   November   2nd

# RNN

Probability of "arrive" in each slot

# RNN

Probability of "arrive" in each slot

Probability of "Taipei" in each slot



然後再根據 a2 產生 y2
y2 是 Taipei 屬於哪一個 slot 的機率

# RNN

Different

Prob of "leave" in each slot

Prob of "Taipei" in each slot

Prob of "arrive" in each slot

Prob of "Taipei" in each slot

$y^1$

$y^2$ ......

$y^1$

$y^2$ ......

store

store

$a^1$

$a^1$

$a^2$

$a^1$

$a^1$

$a^2$

$x^1$

$x^2$ ......

$x^1$

$x^2$ ......

leave

Taipei

arrive

Taipei

The values stored in the memory is different.

Of course it can be deep ...

# Of course it can be deep ...

# Of course it can be deep ...



在下一個時間點的時候呢，每一個 hidden layer

# Of course it can be deep ...



最後的 output，這個 process 就一直持續下去

# Of course it can be deep ...



這個 deep 你要疊幾層呢

# Elman Network & Jordan Network



**_Elman Network_**

$y^t$            $y^{t+1}$

$W^o$        $W^o$

$W^h$

$W^i$        $W^i$

$x^t$         $x^{t+1}$

# Elman Network & Jordan Network

# Bidirectional RNN

# Bidirectional RNN

# Bidirectional RNN

# Bidirectional RNN

# Bidirectional RNN

Benefit: every part of output considers the whole input sequence

The above is actually just a simple version of RNN (called SimpleRNN).

Issues with the SimpleRNN: training is difficult, due to issues including "exploding gradient" or "vanishing gradient" in the gradient descent method.

More advanced types of RNN:
LSTM, and GRU (a simpler version than LSTM).

When people use RNN, they mostly use LSTM or GRU.

Keras let you create SimpleRNN, LSTM or GRU using just one line of code.

# Long Short-term Memory (LSTM)

# Long Short-term Memory (LSTM)



這個 Long Short-term 的 memory 呢
它有 3 個 gate

# Long Short-term Memory (LSTM)



Memory Cell

Input Gate

它必須先通過一個閘門，通過一個 input gate

Other part of the network

# Long Short-term Memory (LSTM)

Memory Cell

Signal control the input gate
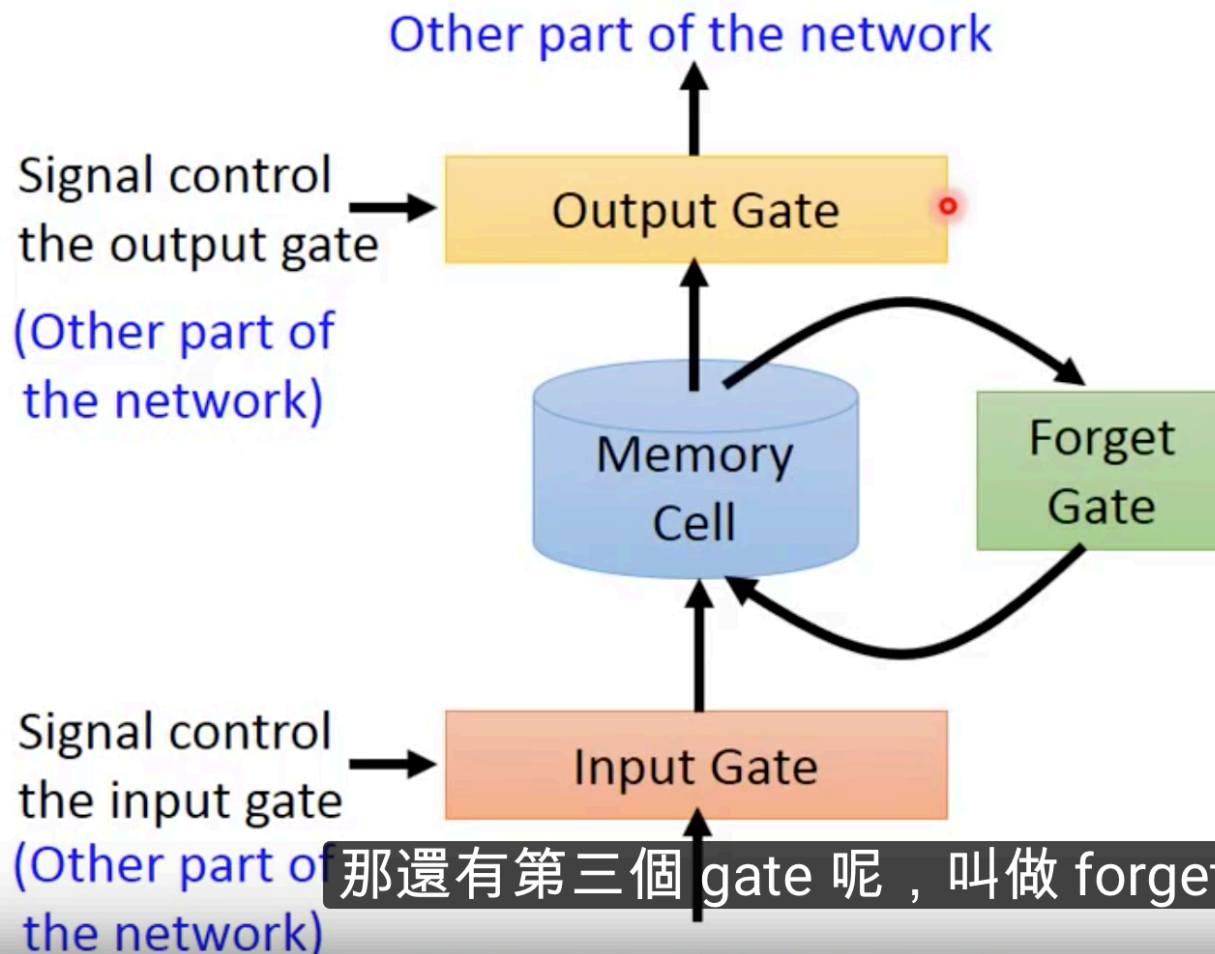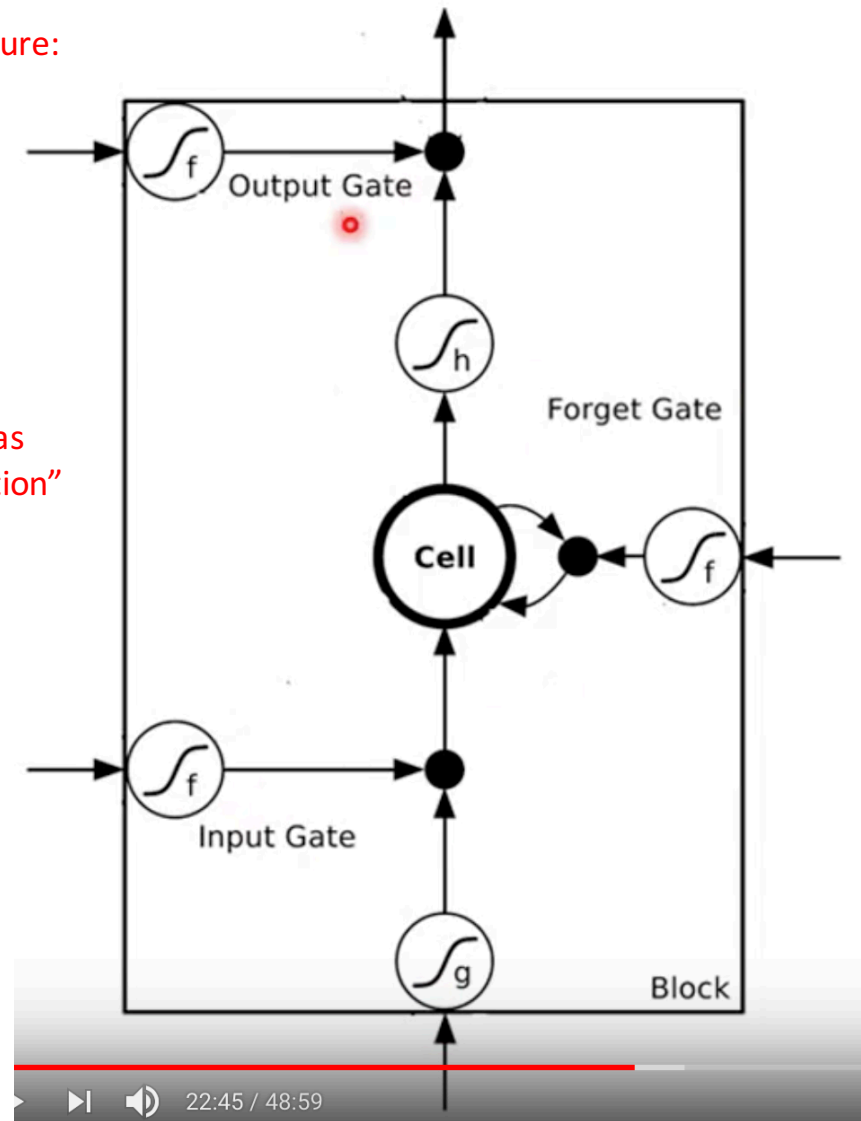(Other part of the network)

Input Gate

Neural network learns when to open/close the gate

Other part of the network

那這個 input gate 呢，它要被打開的時候

# Long Short-term Memory (LSTM)

Other part of the network

Memory Cell

Signal control the input gate (Other part of the network)

Input Gate

Other part of the network

# Long Short-term Memory (LSTM)

Other part of the network

Output Gate

Memory Cell

Signal control
the input gate
(Other part of
the network)

Input Gate

Other part of the network

輸出的地方也有一個 output gate

# Long Short-term Memory (LSTM)

Other part of the network

Signal control the output gate

(Other part of the network)

Output Gate

Memory Cell

Signal control the input gate

(Other part of the network)

Input Gate

Other part of the network

# Long Short-term Memory (LSTM)

Other part of the network

Signal control
the output gate

(Other part of
the network)

Output Gate

Memory
Cell

Forget
Gate

Signal control
the input gate

(Other part of
the network)

Input Gate

Other part of the network

那還有第三個 gate 呢，叫做 forget gate

# Long Short-term Memory (LSTM)

Other part of the network

Signal control
the output gate

(Other part of
the network)

Signal control
the input gate
(Other part of
the network)

Output Gate

Memory
Cell

Forget
Gate

Input Gate

Signal control
the forget gate

(Other part of
the network)

Other part of the network

把它 format 掉

# Long Short-term Memory (LSTM)

A more detailed look at its structure:

Every input and output here
is a real number.

The whole thing can be seen as
replacing the "activation function"
of an ordinary neuron.

好，假設我們現在 cell 裡面呢

23:27 / 48:59

Activation function f is usually a sigmoid function

Between 0 and 1

Mimic open and close gate

Activation function f is usually a sigmoid function

Between 0 and 1

Mimic open and close gate

$f(zi)$ 得到 $g(z)*f(zi)$

Activation function f is usually a sigmoid function

Between 0 and 1

Mimic open and close gate

得到 f(zf)，接下來呢

Activation function f is usually a sigmoid function

Between 0 and 1

Mimic open and close gate

$z_o$ — Output Gate

$a$

$f$ (Sh)

Forget Gate

$c$   $f(z_f)$

$c$   $f$   $z_f$

$cf(z_f)$

$f(z_i)$   $g(z)f(z_i)$

$z_i$ — Input Gate   multiply

$g(z)$

$g$

c 乘上 f(zf)

$z$

25:05 / 48:59

Activation function f is usually a sigmoid function

Between 0 and 1

Mimic open and close gate

$$c' = g(z)f(z_i) + cf(z_f)$$

Activation function f is usually a sigmoid function

Between 0 and 1

Mimic open and close gate

$$c' = g(z)f(z_i) + cf(z_f)$$

新的存在 memory 裡面的值就是 c'，所以

Activation function f is usually a sigmoid function

Between 0 and 1

Mimic open and close gate

$$c' = g(z)f(z_i) + cf(z_f)$$

好那把這個 c' 通過 h，得到 h(c')

Activation function f is usually a sigmoid function

Between 0 and 1

Mimic open and close gate

$$c' = g(z)f(z_i) + cf(z_f)$$

zo 通過 f 得到 f(zo)

$$a = h(c')f(z_o)$$

multiply

Output Gate

$$f(z_o)$$

$$h(c')$$

Activation function f is
usually a sigmoid function

Between 0 and 1

Mimic open and close gate

Forget Gate

$$c \quad f(z_f)$$

$$c'$$

$$cf(z_f)$$

$$f(z_i) \quad g(z)f(z_i)$$

multiply

Input Gate

$$g(z)$$

$$c' = g(z)f(z_i) + cf(z_f)$$

$$z$$

跟這個 h(c') 乘起來，如果 f(zo)是 1，就等於是

# LSTM - Example

| $x_1$ | 1 | 3 | 2 | 4 | 2 | 1 | 3 | 6 | 1 |
|-------|---|---|---|---|---|---|---|---|---|
| $x_2$ | 0 | 1 | 0 | 1 | 0 | 0 | -1 | 1 | 0 |
| $x_3$ | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 |

| $y$ | 0 | 0 | 0 | 0 | 0 | 7 | 0 | 0 | 6 |
|-----|---|---|---|---|---|---|---|---|---|

$x_1$

$x_2$ → LSTM cell → $y$

$x_3$

# LSTM - Example

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| $x_1$ | 1 | 3 | 2 | 4 | 2 | 1 | 3 | 6 | 1 |
| $x_2$ | 0 | 1 | 0 | 1 | 0 | 0 | -1 | 1 | 0 |
| $x_3$ | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 |

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| $y$ | 0 | 0 | 0 | 0 | 0 | 7 | 0 | 0 | 6 |

When $x_2 = 1$, add the numbers of $x_1$ into the memory

When $x_2 = -1$, reset the memory

When $x_3 = 1$, output the number in the memory.

那我們就來實際做一下運算

我們假設 g 跟 h 都是 linear 的，這樣計算比較方便

好那我們現在 input 第一個 vector [3 1 0]

那 forget gate 呢，input 是 [3 1 0]，forget gate 呢

也沒有甚麼影響，0*1+3
所以存在 memory 裡面的值變成 3

output gate 還是被關起來的，3無法通過，所以輸出就是 0

forget gate 被打開的關係，所以 3 * 1 + 4
所以 memory 裡面存的值會變成 7

所以 7 呢，仍然無法被輸出，
所以整個 memory 的輸出仍然是 0

所以 forget gate 還是打開的，所以 7 * 1 + 0

那這個 7 它沒有辦法被輸出，因為 output gate
仍然是關閉的，所以 output 仍然是 0

forget gate 這個時候仍然跟原來一樣，它是被打開的

那 1 * 7 = 7 這樣子

所以呢，memory 裡面存的值會被洗掉

那它讀出來的值也是 0

Original Network:

# Original Network:



$a_1$     $a_2$

$z_1$     $z_2$

$x_1$     $x_2$    Input

但是如果是 LSTM 的話呢

Original Network:

> Simply replace the neurons with LSTM

Each edge has a different weight

$a_1$ $a_2$

Output Gate Output Gate

Forget Gate Forget Gate

Cell Cell

Input Gate Input Gate

Block Block

**4 times of parameters** $x_1$ $x_2$ Input

# How to understand LSTM as RNN

## LSTM

# LSTM

$$c^{t-1}$$

The vector of memory values at time t-1 (where each memory's value is a real number)

# LSTM

$c^{t-1}$

vector

Z

Input vector at time t

$x^t$

Linear transformation
(multiple input vector by a matrix to get a vector z)

# LSTM

$c^{t-1}$

vector

Z

$x^t$

# LSTM



$c^{t-1}$

vector

$z^i$   $z$

Linear
transformation

$x^t$

得到 z^i

# LSTM



$c^{t-1}$

vector

$z^f$   $z^i$   $z$   $z^o$   4 vectors

4 different linear transformations

$x^t$

# LSTM



這 4 個 z 其實都是 vector

# LSTM

Element-wise vector operation for the layer

# LSTM


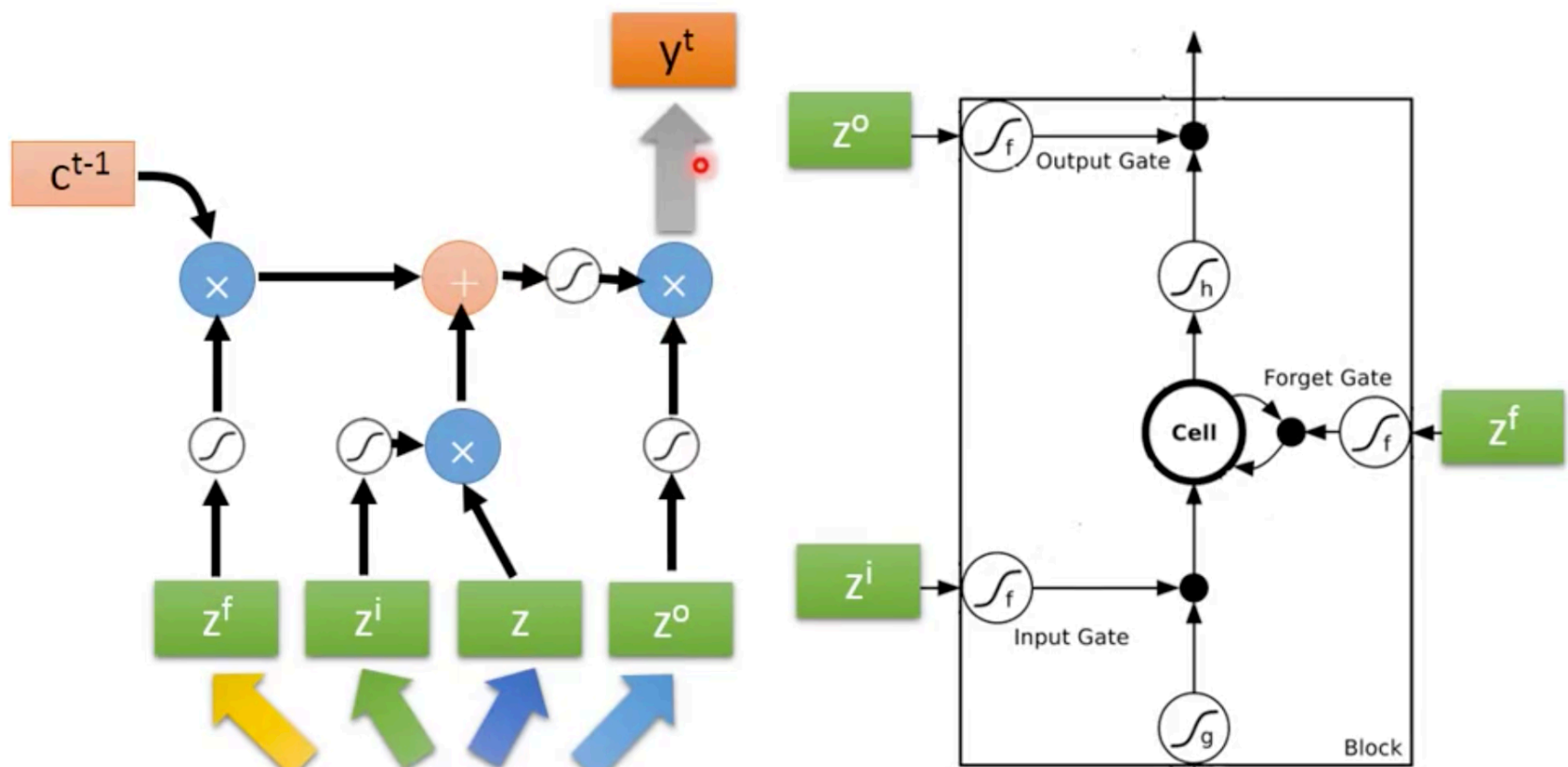
再相乘，最後就得到最後的 output 的 y

# LSTM



這個時候相加以後的結果，就是 memory 裡面存的值

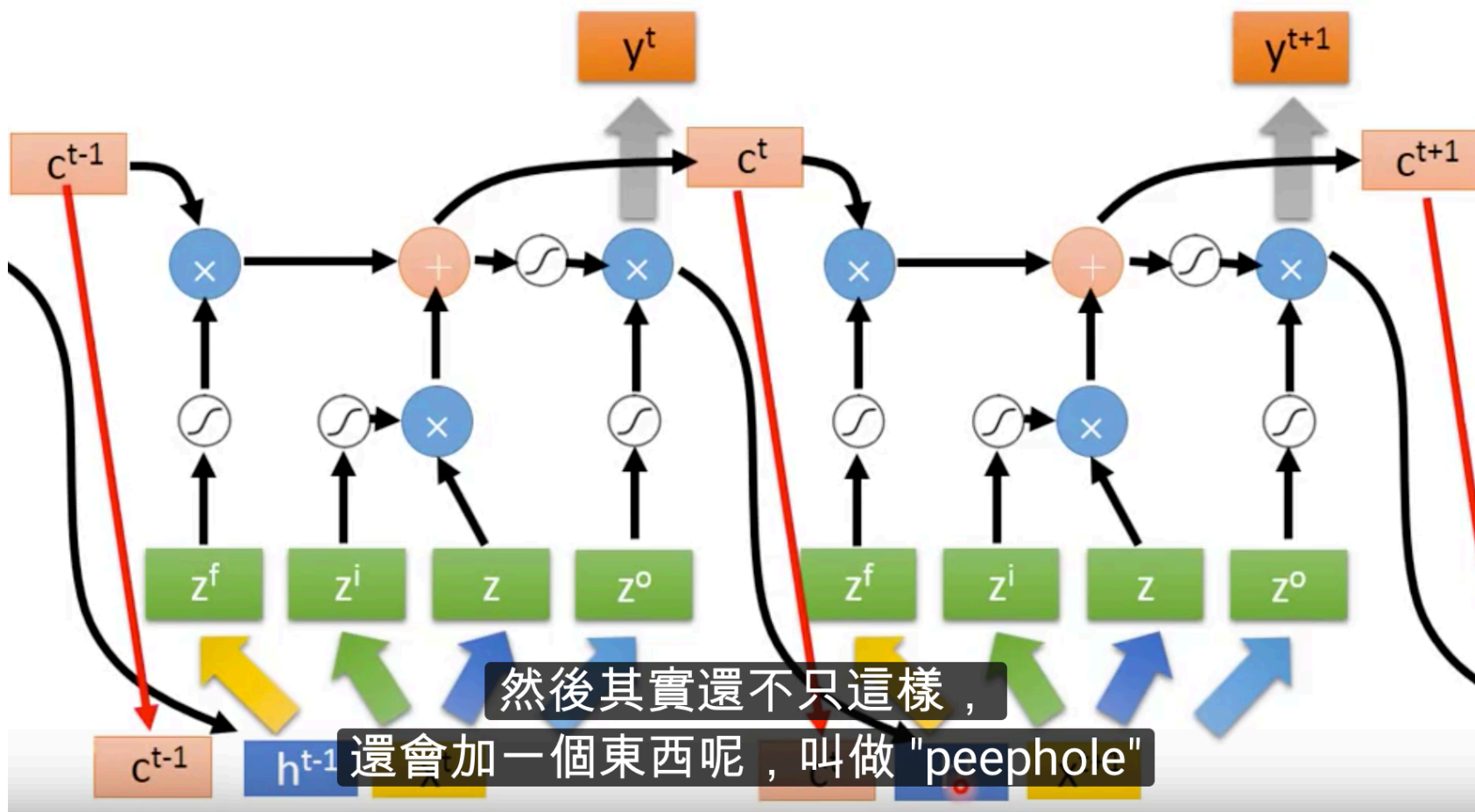# LSTM

LSTM

Real LSTM

LSTM

Real LSTM

Extension: "peephole"

然後其實還不只這樣，
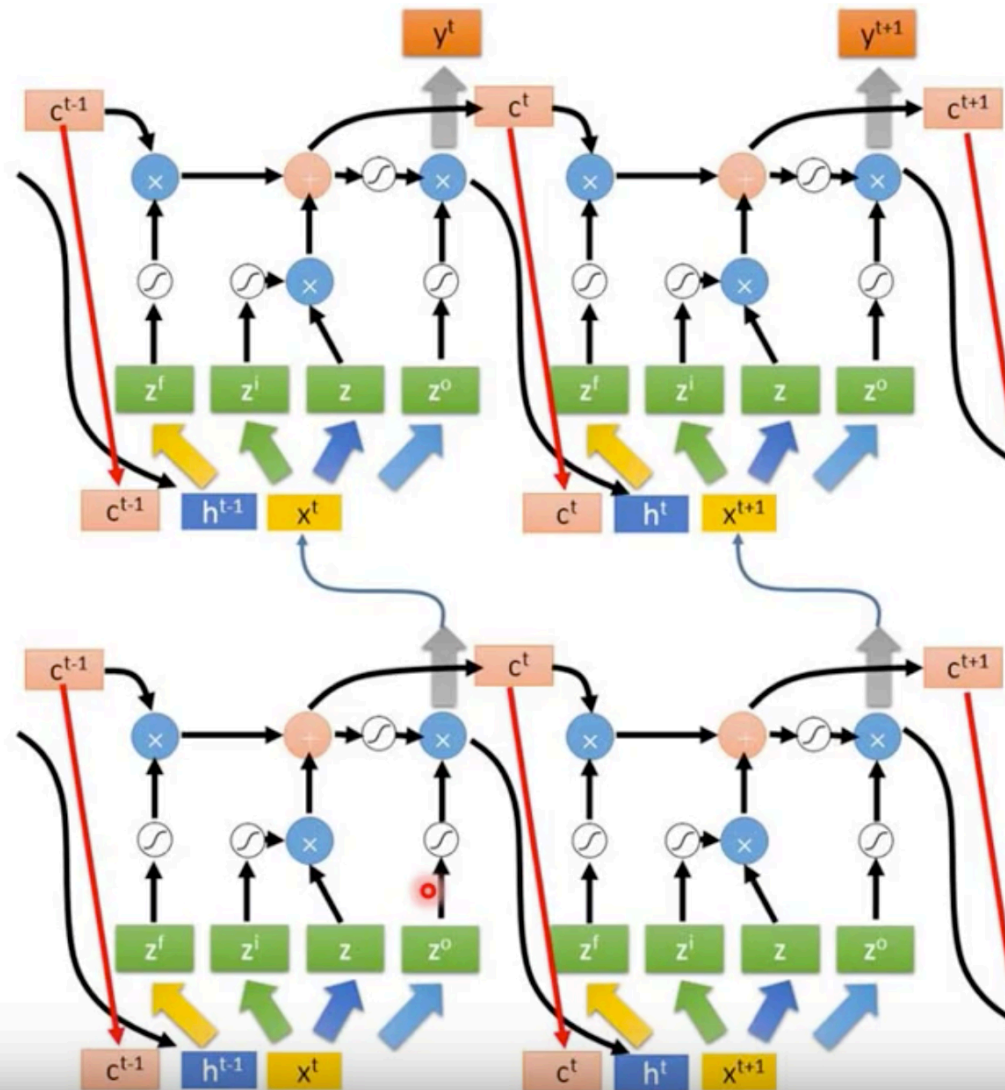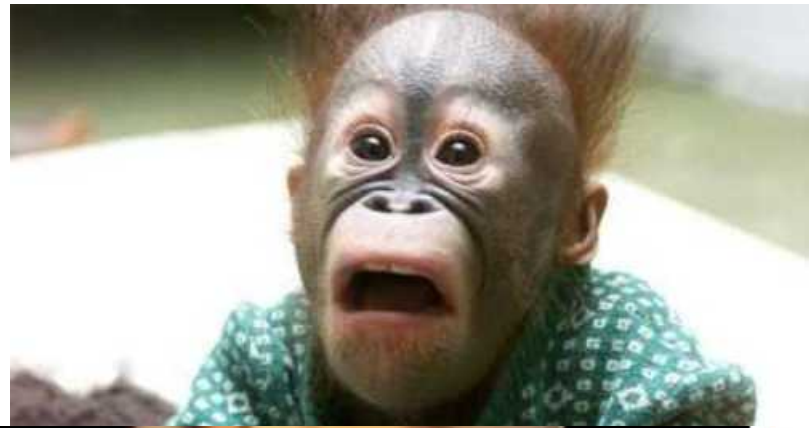還會加一個東西呢，叫做 "peephole"

# Multiple-layer LSTM

The first time a person sees

LSTM

The first time a person sees

LSTM



Don't worry if you cannot understand this.
Keras can handle it.

Keras supports
"LSTM", "GRU", "SimpleRNN" layers