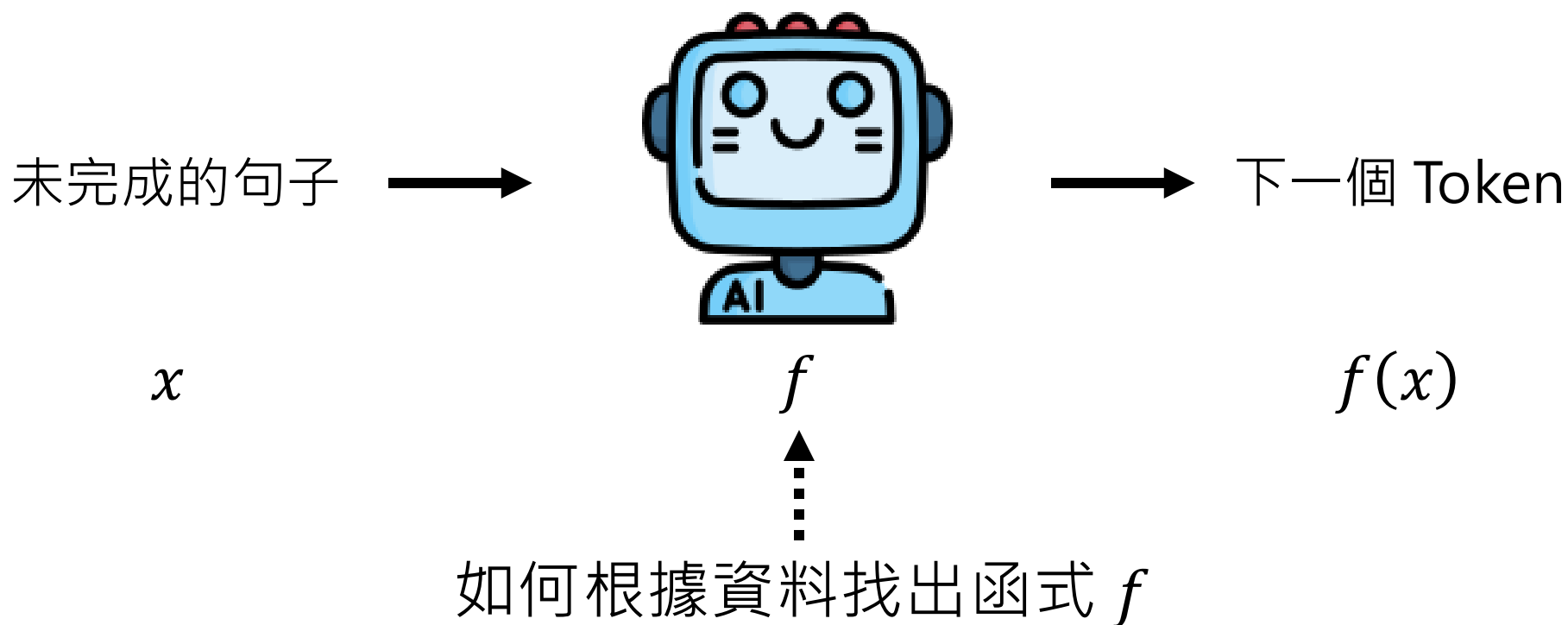


請各位同學稍待片刻
我們 14:23 開始上課

一堂課搞懂 機器學習和深度學習 的基本概念

生成式人工智慧的基本原理



機器學習 (Machine Learning)

課程規劃



原理

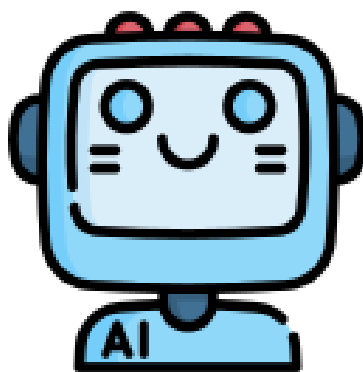


實作

可以找各式各樣的函式



x



f



如何根據資料找出函式 f



Regression

數字
(這堂課的長度)

$f(x)$

這個函式有什麼用呢？

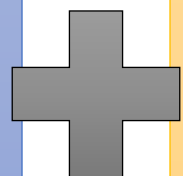
- 這個函式 f 回答一個關鍵問題

相信大家上課常常都會想的

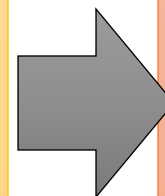


找函式步驟 3 + 1

步驟一：
我要找什麼



步驟二：
我有哪些選擇



步驟三：
選一個最好的

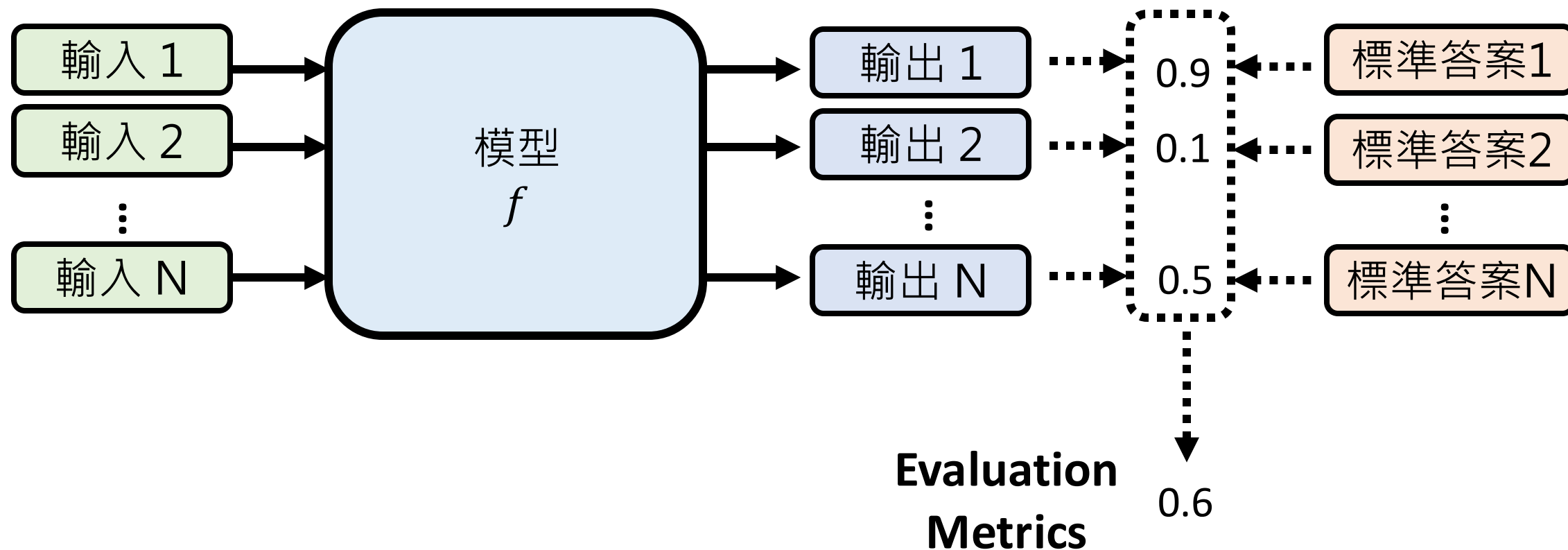
學習 (Learning)、訓練 (Training)

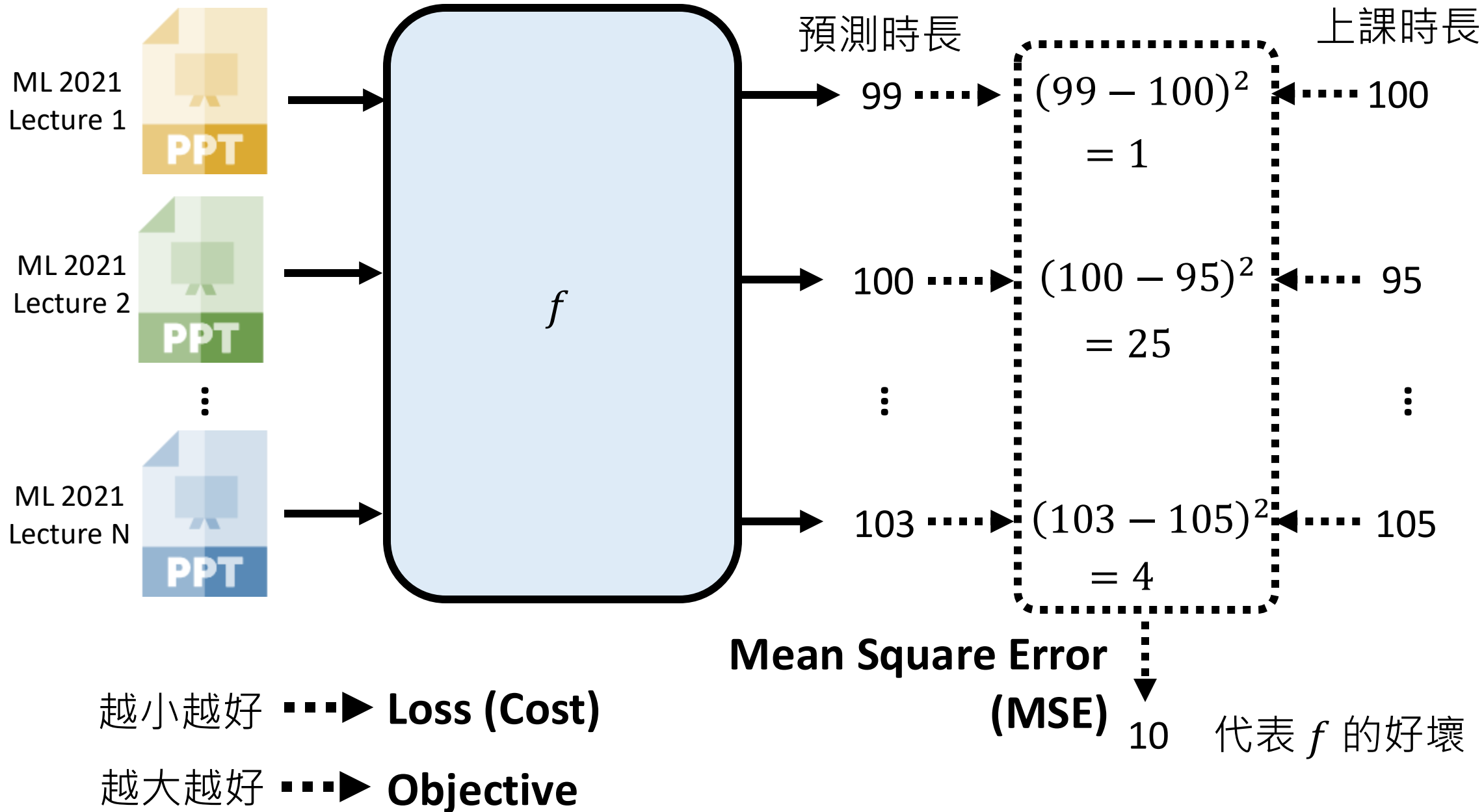
找函式步驟 3 + 1

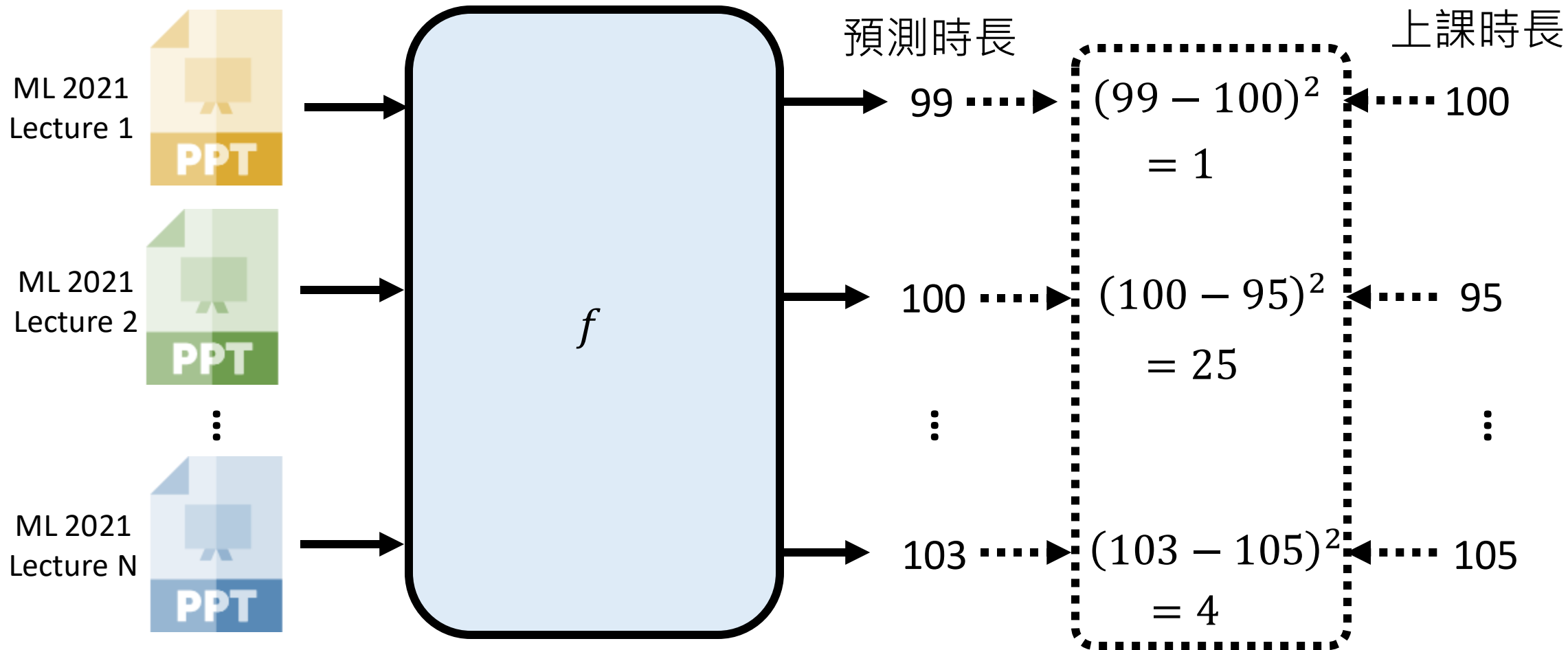


給我一個 f ，我要知道它是不是我要的

上一講：生成式人工智慧的能力檢定



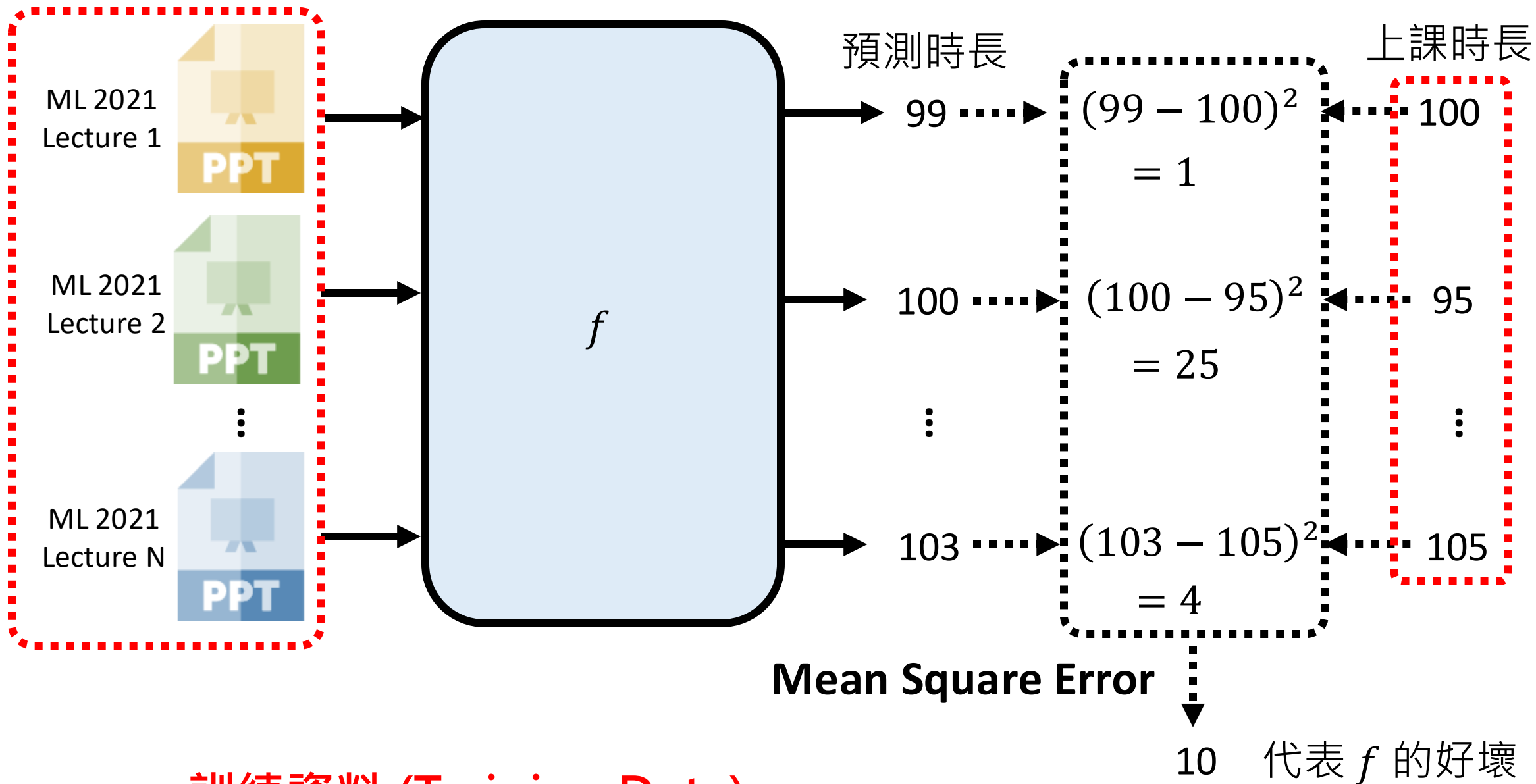




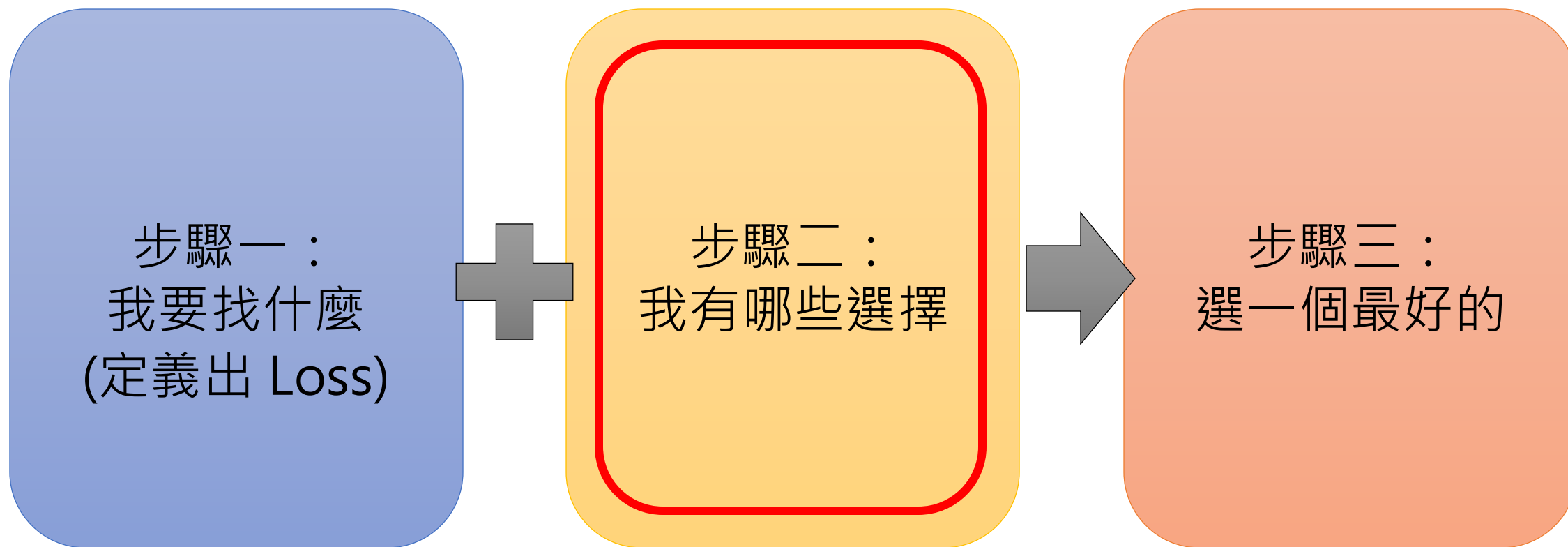
越小越好 $\cdots \blacktriangleright$ **Loss (Cost)**

越大越好 $\cdots \blacktriangleright$ **Objective**

這跟 Evaluation 的過程是一樣的
能不能把 Evaluation Metrics 當作 Loss (Objective) ?



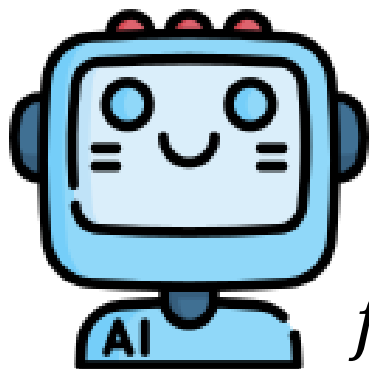
找函式步驟 3 + 1



訂出候選的函式集合



x

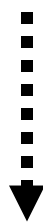


f

數字

y

函式的輸入只
能是數字



Feature

頁數 x_1

內容總字數 x_2

標題長度 x_3

有沒有 "Learning" x_4



(0 or 1)

Linear Regression

$$y = w_1 x_1 + b$$

+

b



課程長度跟投影片頁數
成某種比例關係



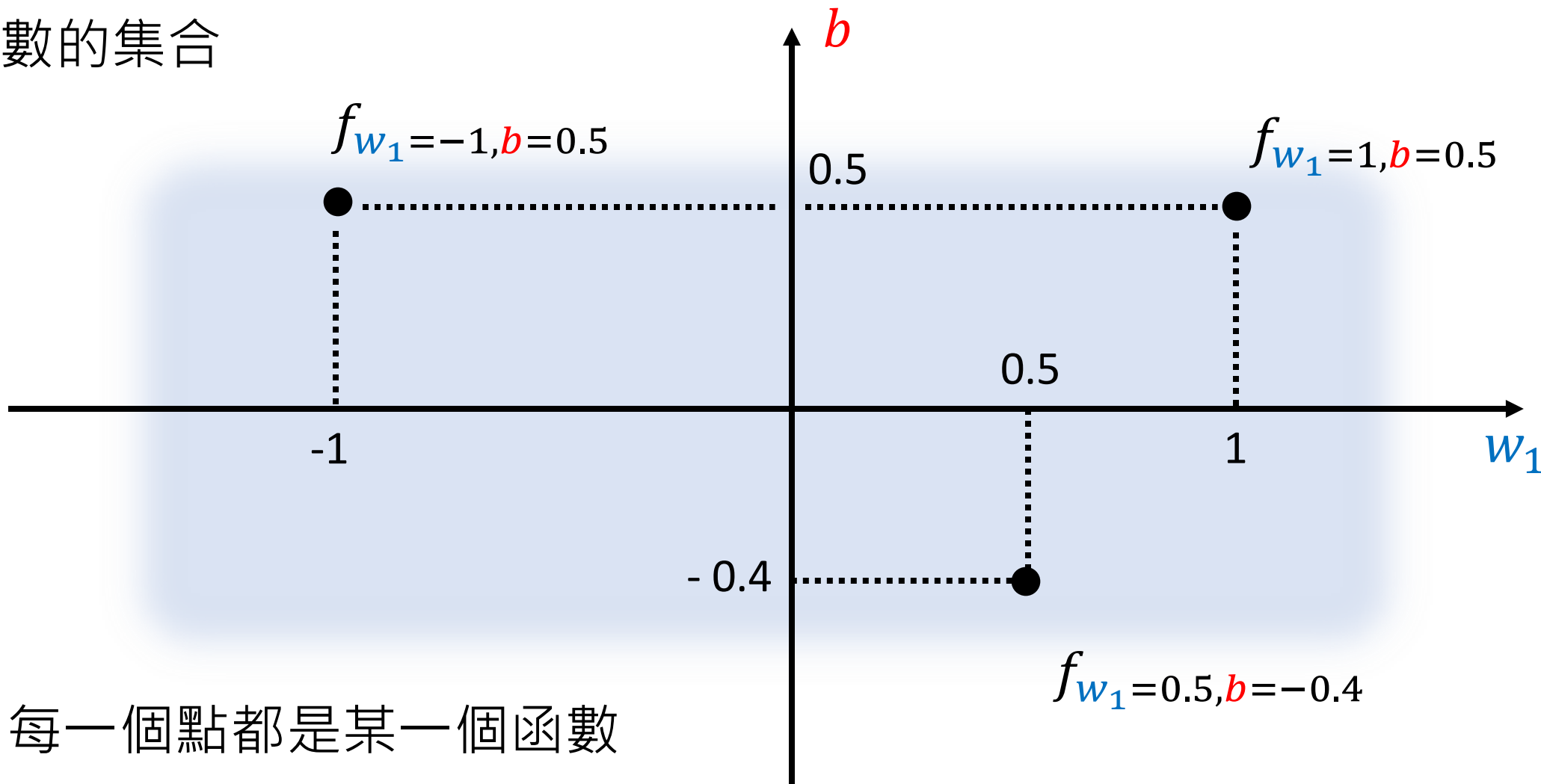
固定多出一點時間
(開場、結尾)

w_1, b 數值未知

參數 (Parameter)

$$y = w_1 x_1 + b \quad w_1, b \text{ 數值未知} \quad \text{參數 (Parameter)}$$

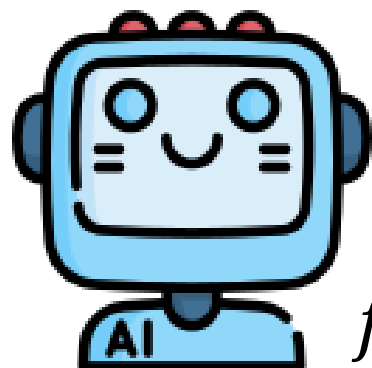
函數的集合



每一個點都是某一個函數



x

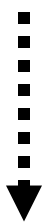


f

數字

y

函式的輸入只能是數字



頁數

x_1

內容總字數

x_2

標題長度

x_3

有沒有 "Learning"

x_4



(0 or 1)

$$y = w_1 x_1 + b$$

出自人類對於任務的理解
(domain knowledge)

$$y = w_1 x_1 + w_2 x_2 + b$$

模型 (Model)

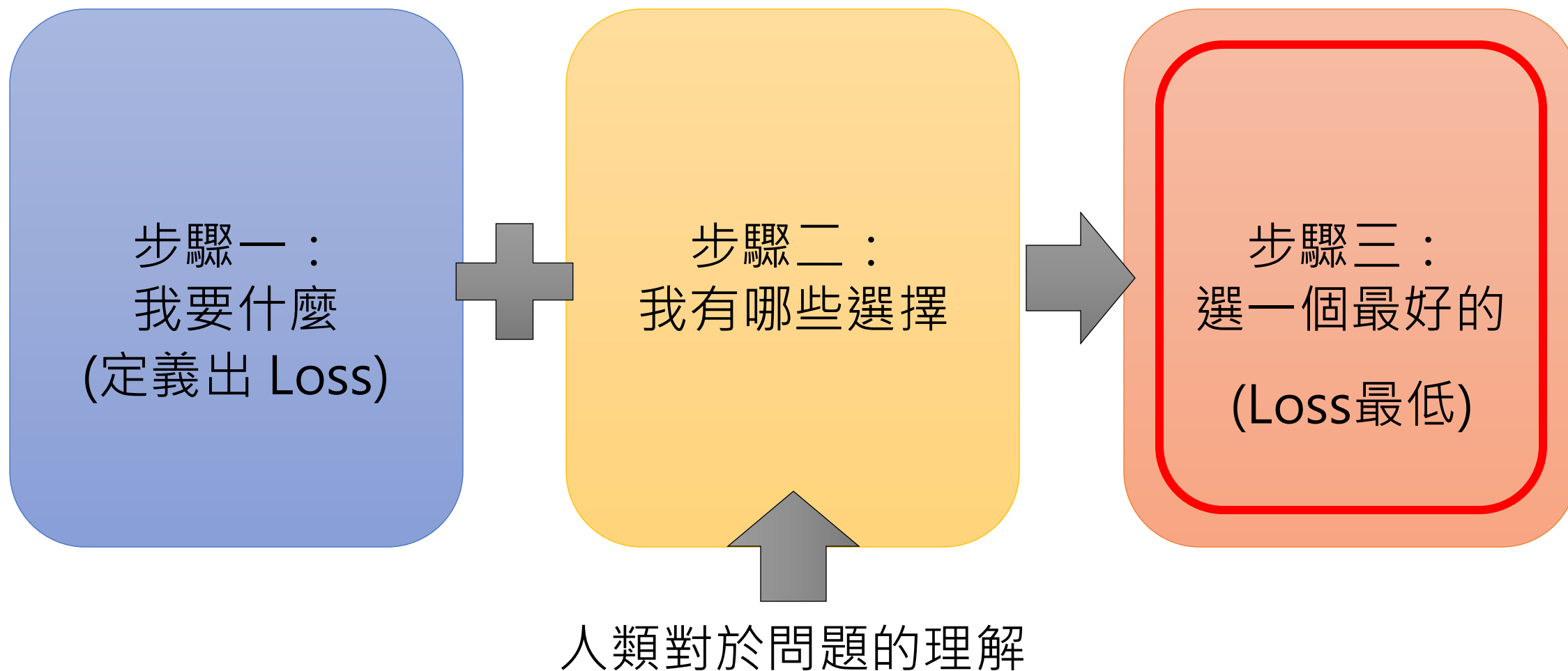
$$y = w_1 x_1 + w_2 x_2 + w_3 x_3 + b$$

$$y = w_1 x_1 + w_2 x_2 x_4 + w_3 x_3^2 + b$$



有無盡的可能

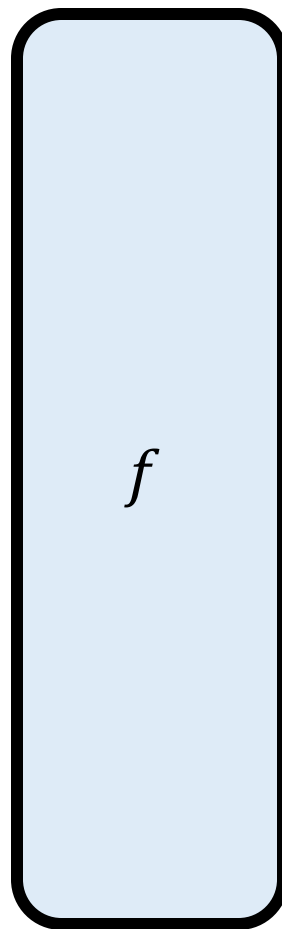
找函式步驟 3 + 1



ML 2021
Lecture 1



x_1^1



預測時長

y^1



$(y^1 - \hat{y}^1)^2$

上課時長

\hat{y}^1

ML 2021
Lecture 2



x_1^2



y^2



$(y^2 - \hat{y}^2)^2$

\hat{y}^2

\vdots

\vdots

\vdots

ML 2021
Lecture N



x_1^N



y^N



$(y^N - \hat{y}^N)^2$

\hat{y}^N

先把 Loss 的數學
式寫出來

$$L = \frac{1}{N} \sum_{i=1}^N (y^i - \hat{y}^i)^2 = \frac{1}{N} \sum_{i=1}^N (\underset{\substack{\uparrow \\ y = w_1 x_1 + b}}{w_1 x_1^i + b} - \hat{y}^i)^2$$

選一個最好 (Loss最低) 的函式

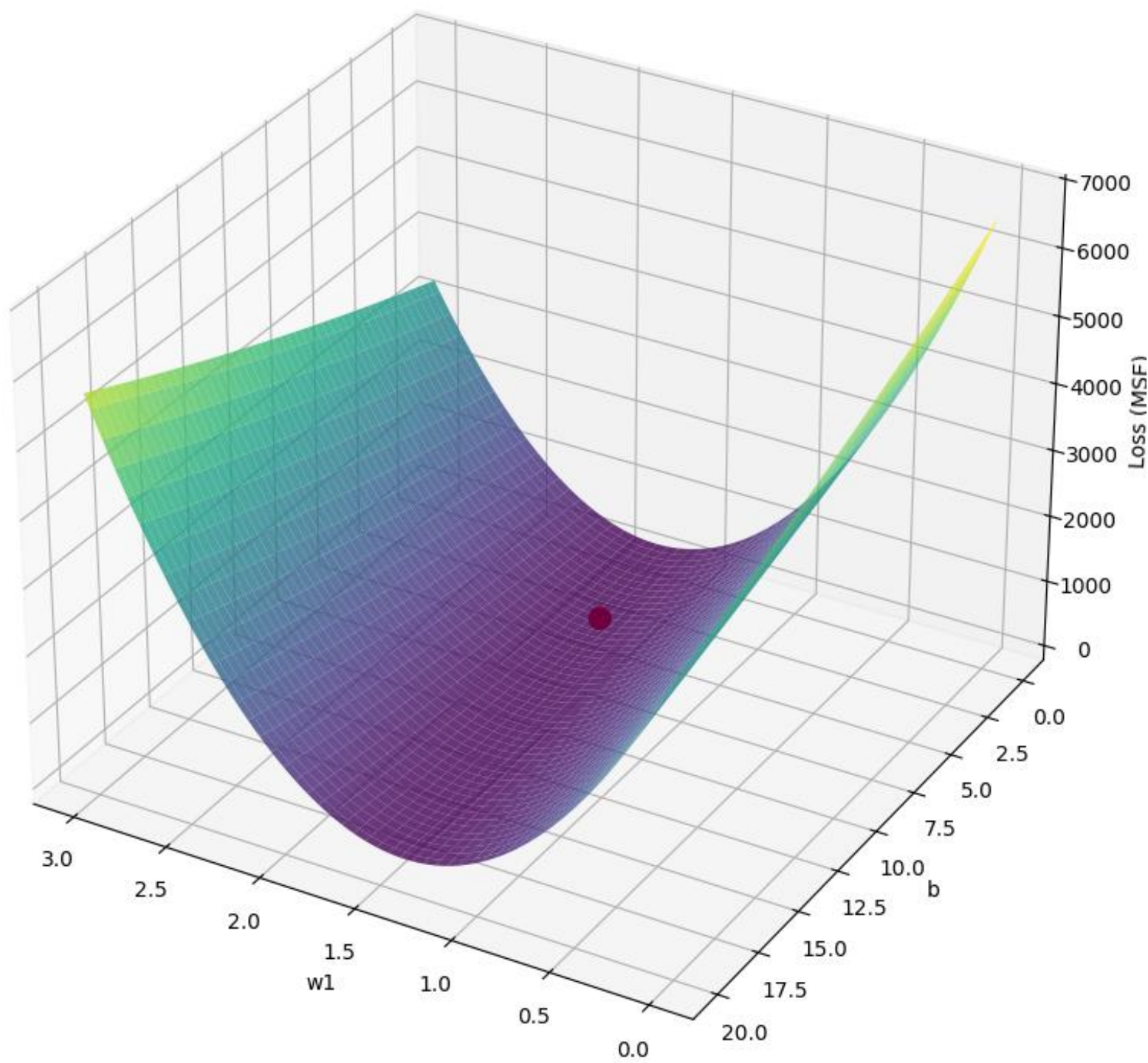
$$L(\underline{w_1}, b) = \frac{1}{N} \sum_{i=1}^N (w_1 x_1^i + b - \hat{y}^i)^2$$

$$w_1^*, b^* = \arg \min_{w_1, b} L(w_1, b)$$

這是我們真正要解的問題
Optimization

暴力算出所有候選
函式的 Loss (MSE)

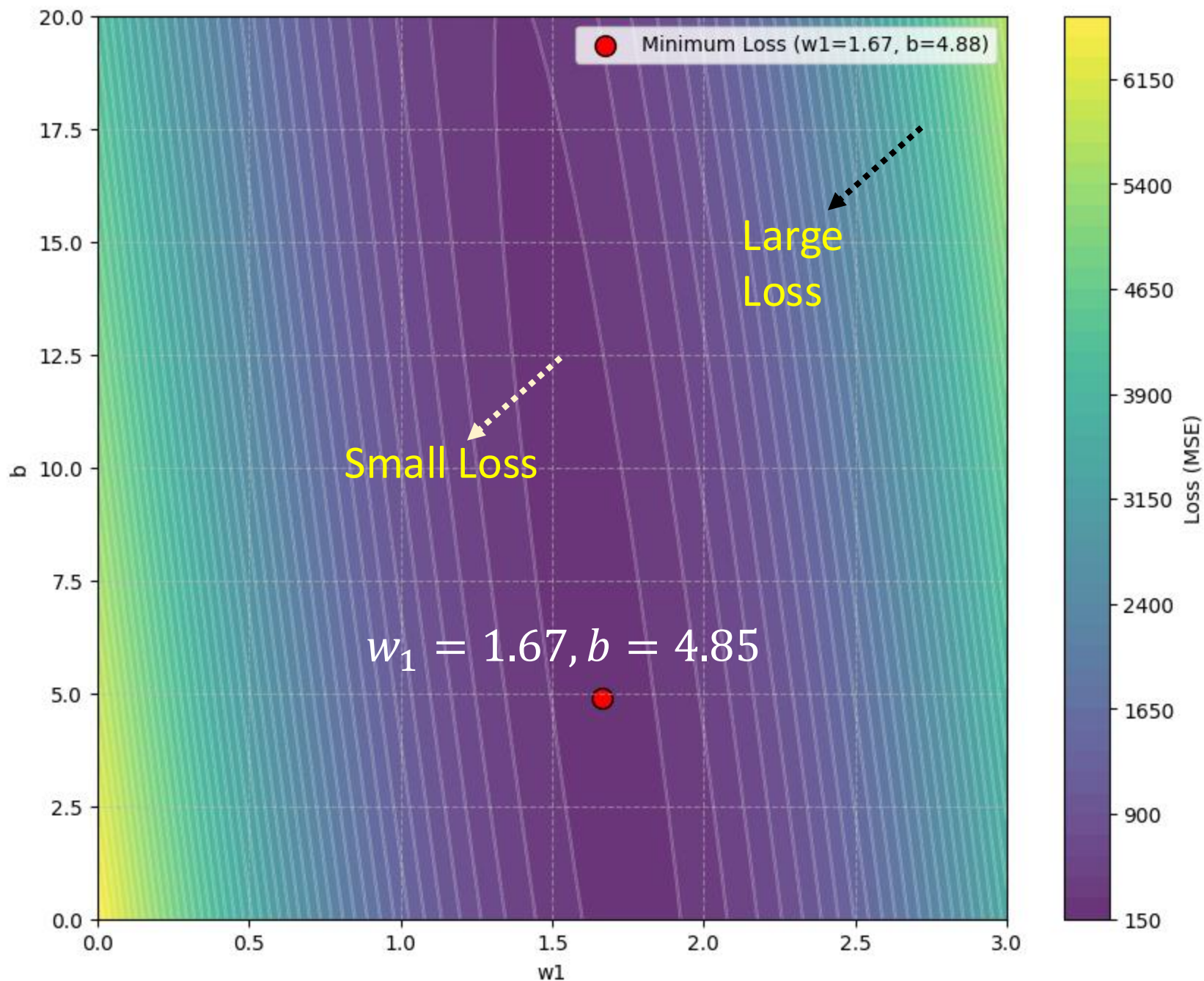
$$w_1^*, b^* = \arg \min_{w_1, b} L(w_1, b)$$



暴力算出所有候選
函式的 Loss (MSE)

$$y = w_1 x_1 + b$$

Loss Surface
(Loss 等高線圖)



選一個最好 (Loss最低) 的函式

當 Loss 是 MSE

$$L(w_1, b) = \frac{1}{N} \sum_{i=1}^N (w_1 x_1^i + b - \hat{y}^i)^2$$

當函式集合寫成這樣

$$y = w_1 x_1 + b$$

Linear Regression

$$w_1^*, b^* = \arg \min_{w_1, b} L(w_1, b)$$

線性代數告訴我們有這個問題有
Closed-form Solution
(有公式解)

我們需要更通用的做法

Gradient Descent

梯度下降法

$$w_1^* = \arg \min_{w_1} L(w_1)$$

Loss
 L

左右各踏一小步
往右一小步可以讓 L 變小

Local
Minimum

往左往右都不會
讓 L 變小

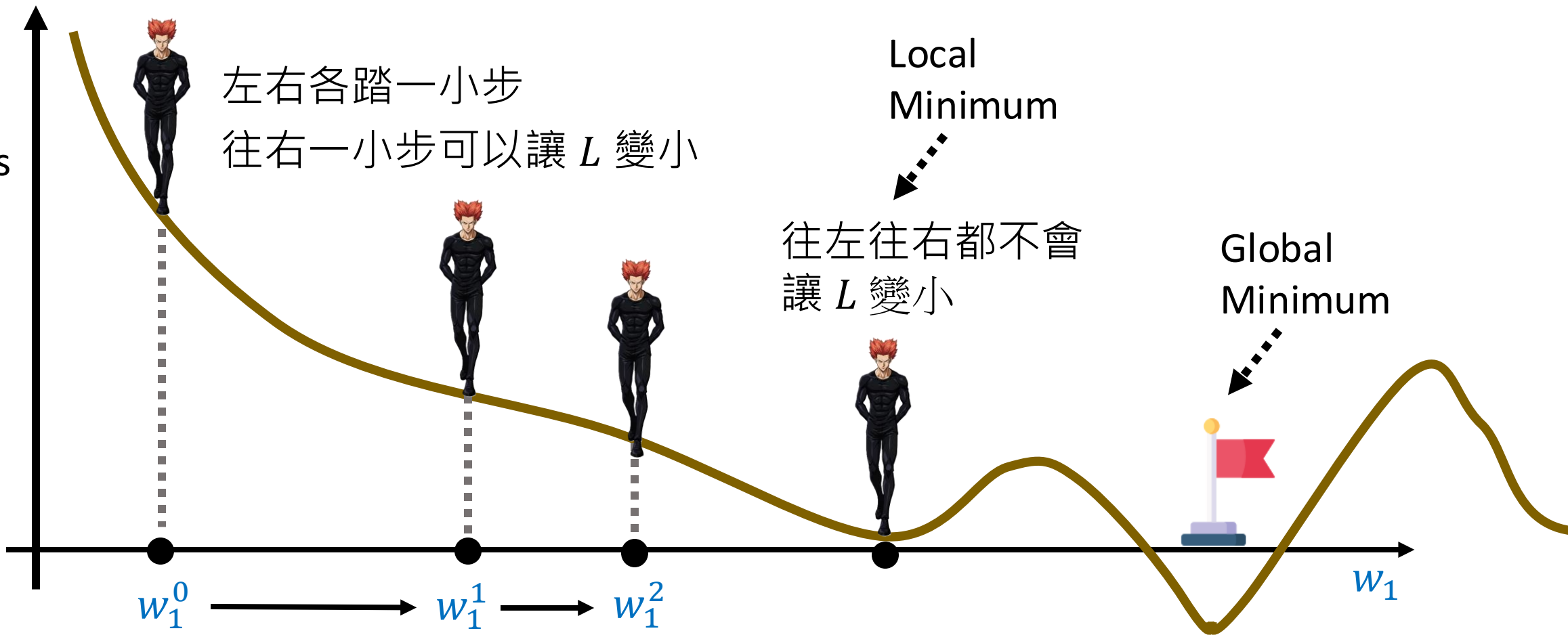
Global
Minimum

w_1^0

w_1^1

w_1^2

w_1



Gradient Descent

梯度下降法

$$w_1^* = \arg \min_{w_1} L(w_1)$$

Loss
 L

計算切線斜率

(斜率是負的) → 向右走

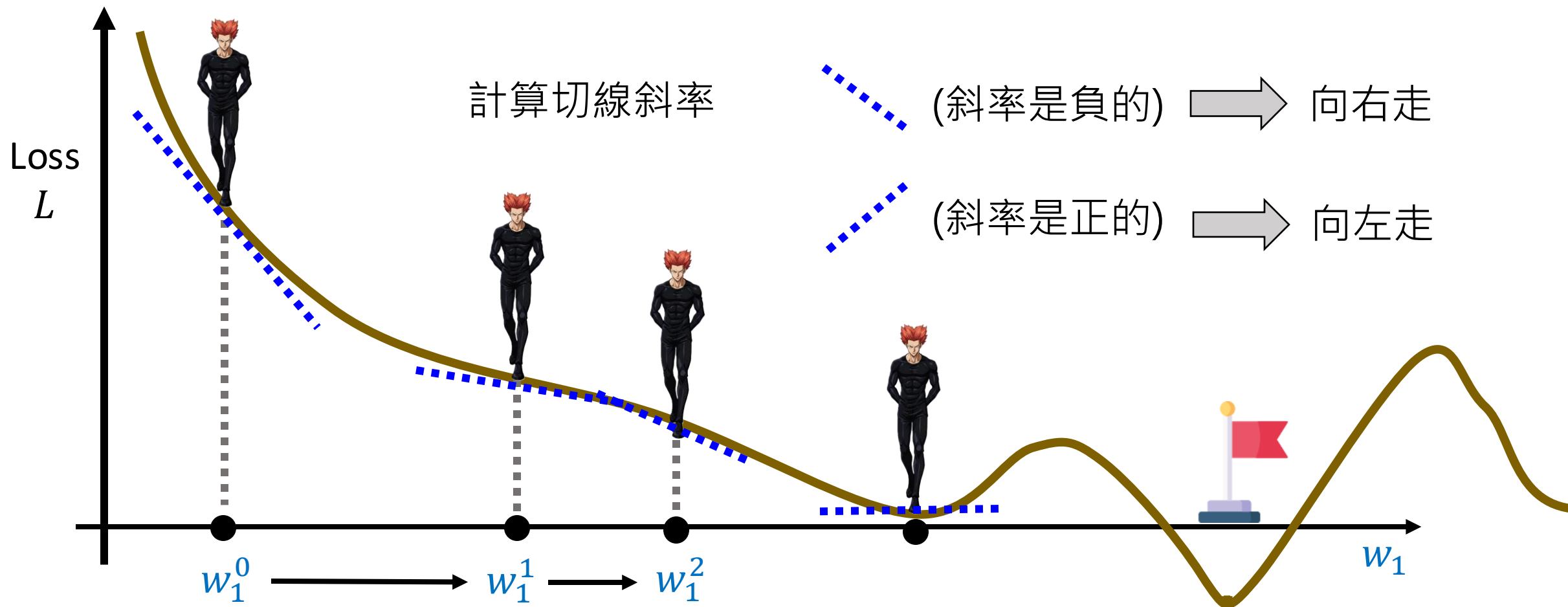
(斜率是正的) → 向左走

w_1^0

w_1^1

w_1^2

w_1



Gradient Descent

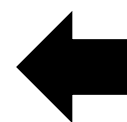
梯度下降法

$$w_1^* = \arg \min_{w_1} L(w_1)$$

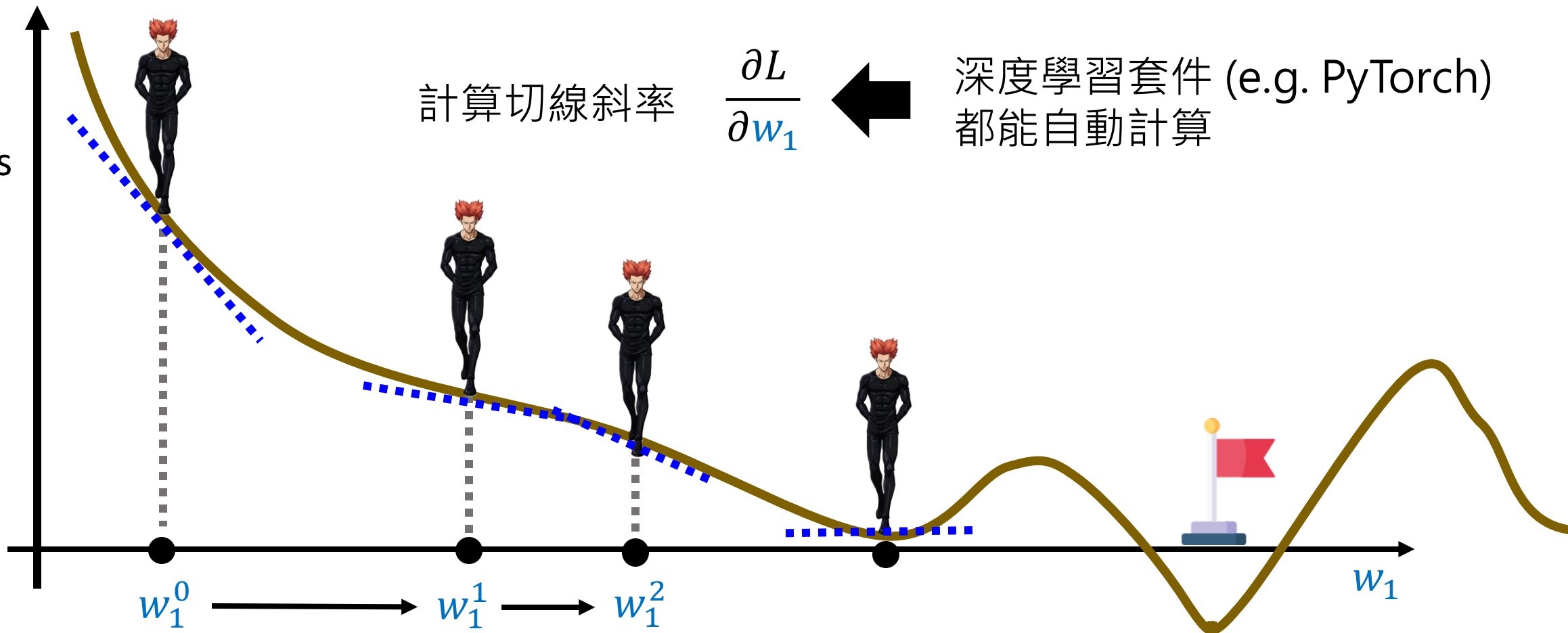
Loss
 L

計算切線斜率

$$\frac{\partial L}{\partial w_1}$$



深度學習套件 (e.g. PyTorch)
都能自動計算

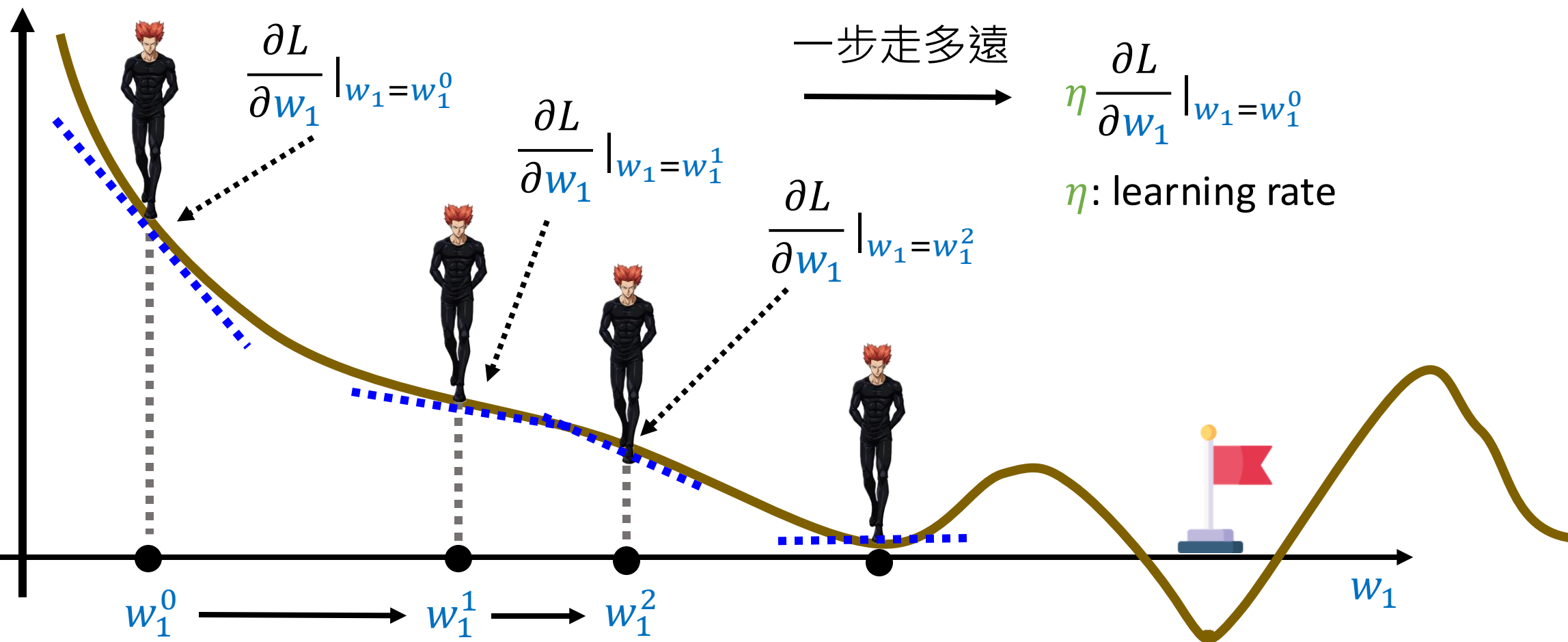


Gradient Descent

梯度下降法

$$w_1^* = \arg \min_{w_1} L(w_1)$$

Loss
 L



Gradient Descent

$$w_1^*, b^* = \arg \min_{w_1, b} L(w_1, b)$$

- (Randomly) Pick initial values w_1^0, b^0
- Compute

Gradient

$$\frac{\partial L}{\partial w_1} \Big|_{w_1=w_1^0, b=b^0}$$

$$\frac{\partial L}{\partial b} \Big|_{w_1=w_1^0, b=b^0}$$

$$w_1^1 \leftarrow w_1^0 - \eta \frac{\partial L}{\partial w_1} \Big|_{w_1=w_1^0, b=b^0}$$

$$b^1 \leftarrow b^0 - \eta \frac{\partial L}{\partial b} \Big|_{w_1=w_1^0, b=b^0}$$

Can be done in one line in most deep learning frameworks

- Update w_1 and b iteratively

Gradient Descent

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) \quad \boldsymbol{\theta} = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \\ \vdots \end{bmatrix}$$

➤ (Randomly) Pick initial values $\boldsymbol{\theta}^0$

$$\text{gradient } \mathbf{g}^0 = \begin{bmatrix} \frac{\partial L}{\partial \theta_1} |_{\boldsymbol{\theta}=\boldsymbol{\theta}^0} \\ \frac{\partial L}{\partial \theta_2} |_{\boldsymbol{\theta}=\boldsymbol{\theta}^0} \\ \vdots \end{bmatrix} \quad \begin{bmatrix} \theta_1^1 \\ \theta_2^1 \\ \vdots \end{bmatrix} \leftarrow \begin{bmatrix} \theta_1^0 \\ \theta_2^0 \\ \vdots \end{bmatrix} - \begin{bmatrix} \eta \frac{\partial L}{\partial \theta_1} |_{\boldsymbol{\theta}=\boldsymbol{\theta}^0} \\ \eta \frac{\partial L}{\partial \theta_2} |_{\boldsymbol{\theta}=\boldsymbol{\theta}^0} \\ \vdots \end{bmatrix}$$

$$\mathbf{g}^0 = \nabla L(\boldsymbol{\theta}^0)$$

$$\boldsymbol{\theta}^1 \leftarrow \boldsymbol{\theta}^0 - \eta \mathbf{g}^0$$

Gradient Descent

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} L(\boldsymbol{\theta})$$

➤ (Randomly) Pick initial values $\boldsymbol{\theta}^0$

➤ Compute gradient $\boldsymbol{g}^0 = \nabla L(\boldsymbol{\theta}^0)$

$$\boldsymbol{\theta}^1 \leftarrow \boldsymbol{\theta}^0 - \eta \boldsymbol{g}^0 \quad \leftarrow \quad \text{1 iteration (1 update)}$$

➤ Compute gradient $\boldsymbol{g}^1 = \nabla L(\boldsymbol{\theta}^1)$

$$\boldsymbol{\theta}^2 \leftarrow \boldsymbol{\theta}^1 - \eta \boldsymbol{g}^1$$

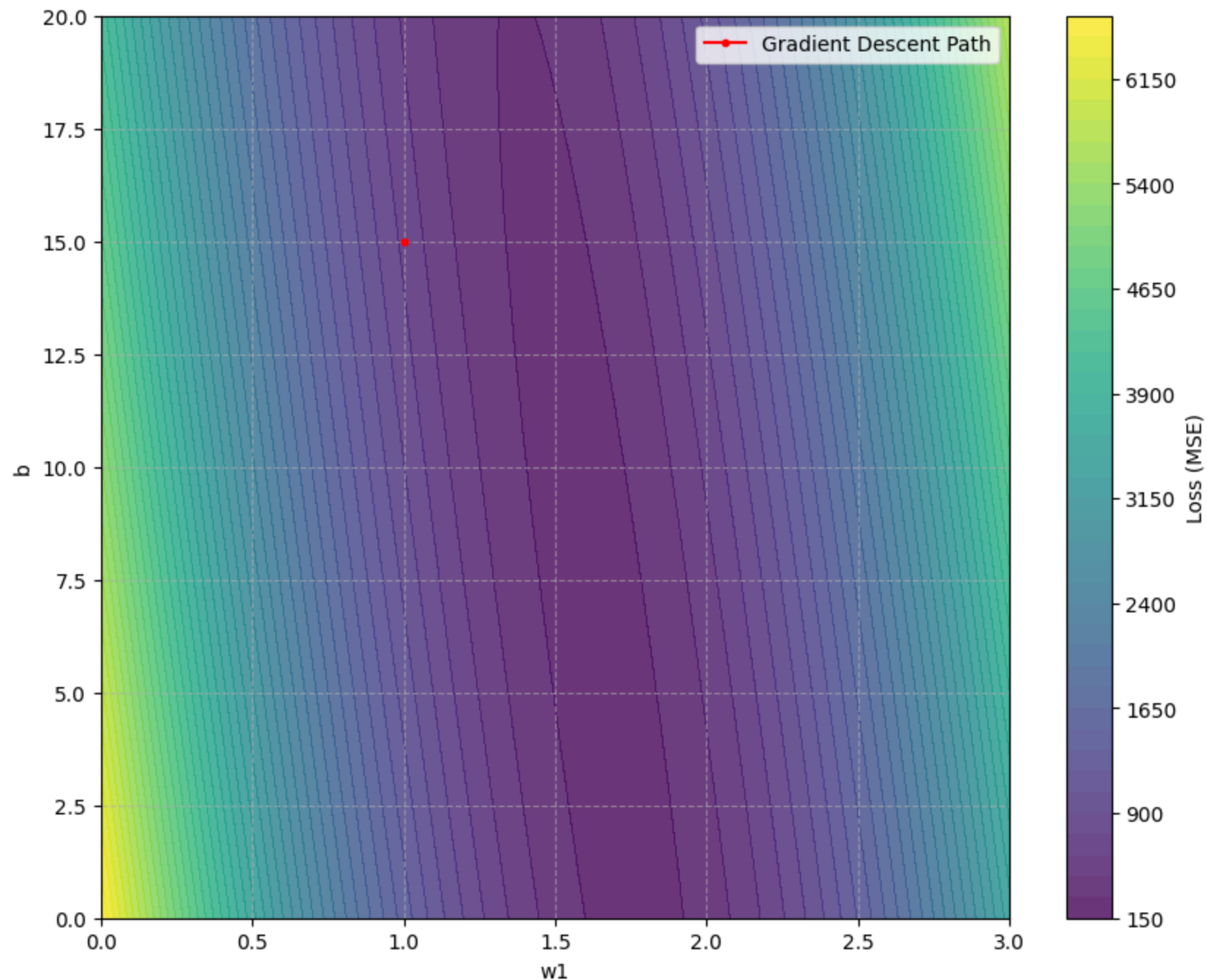
➤ Compute gradient $\boldsymbol{g}^2 = \nabla L(\boldsymbol{\theta}^2)$

$$\boldsymbol{\theta}^3 \leftarrow \boldsymbol{\theta}^2 - \eta \boldsymbol{g}^2$$

概念很簡單，
做起來不容易

$$w_1^0 = 1.0, b^0 = 15.0$$

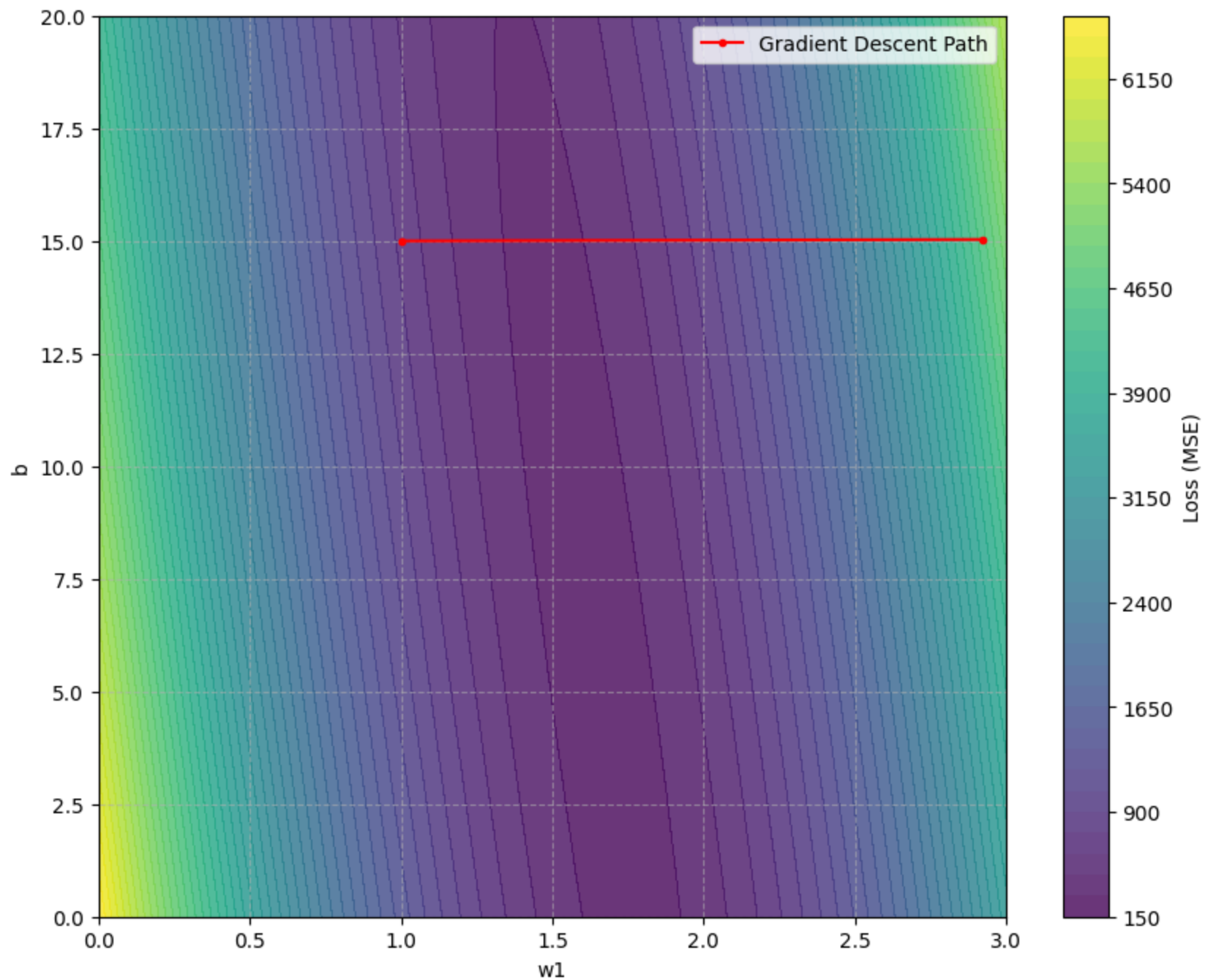
$$\eta = 0.001$$



概念很簡單，
做起來不容易

$$w_1^0 = 1.0, b^0 = 15.0$$

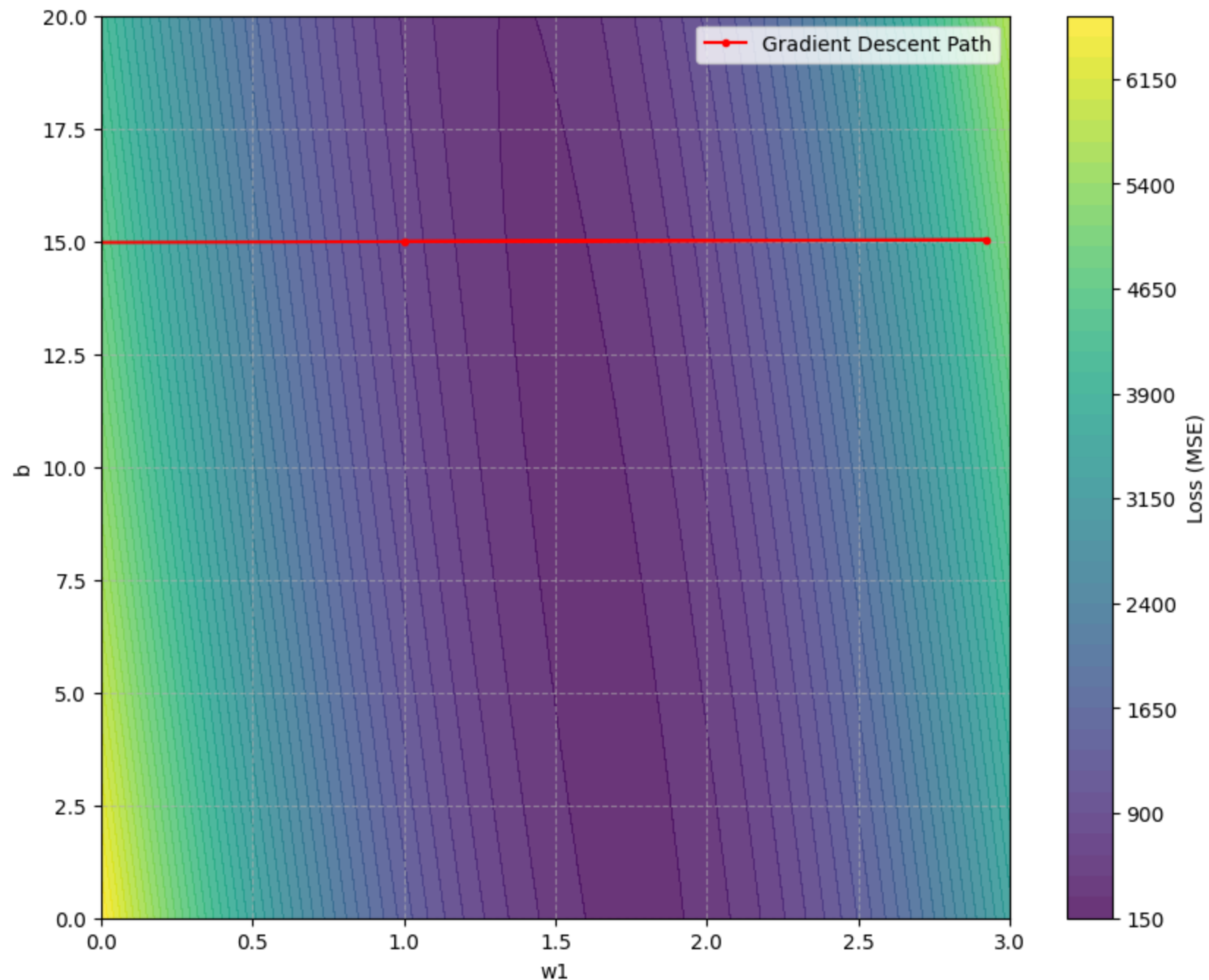
$$\eta = 0.001$$



概念很簡單，
做起來不容易

$$w_1^0 = 1.0, b^0 = 15.0$$

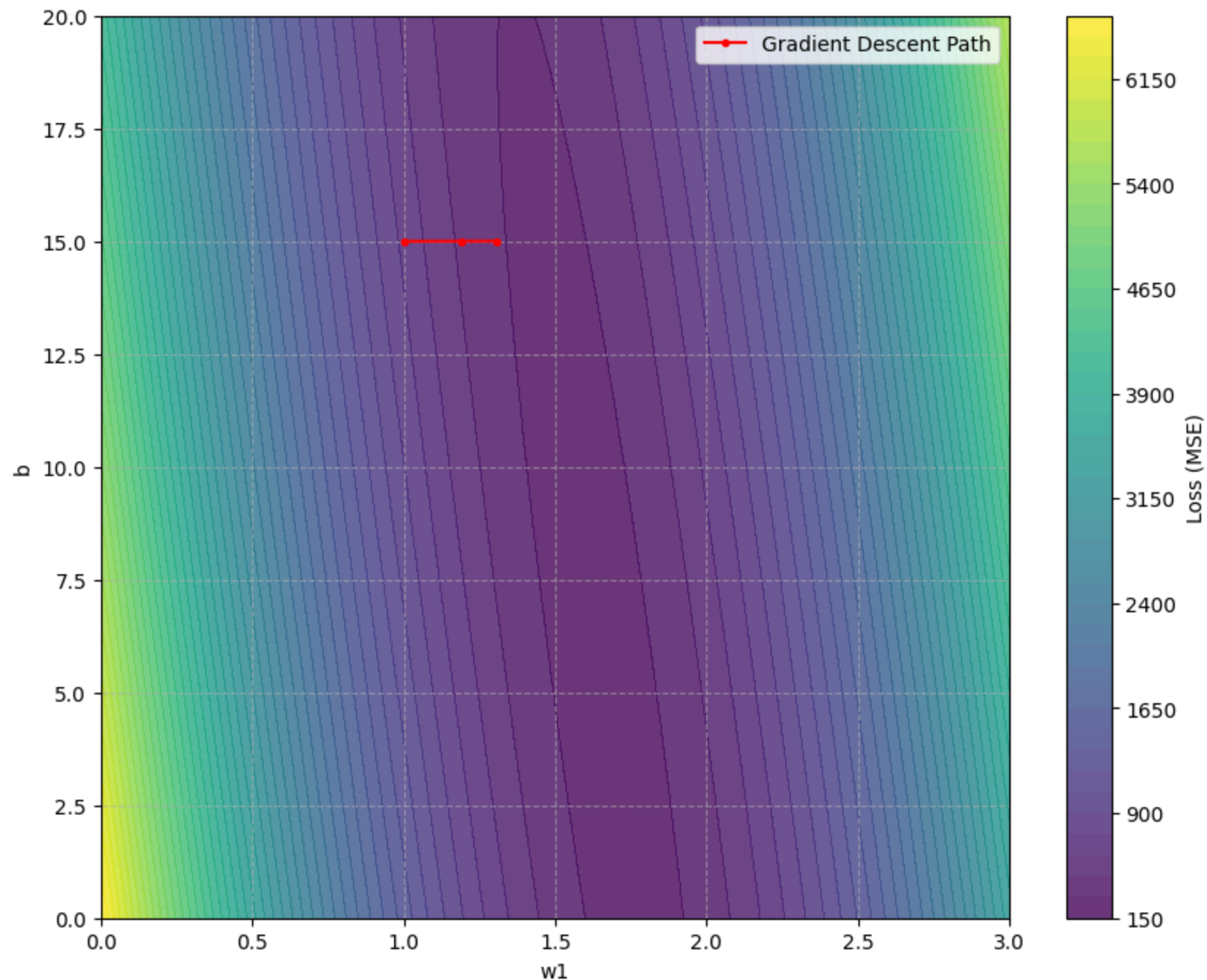
$$\eta = 0.001$$



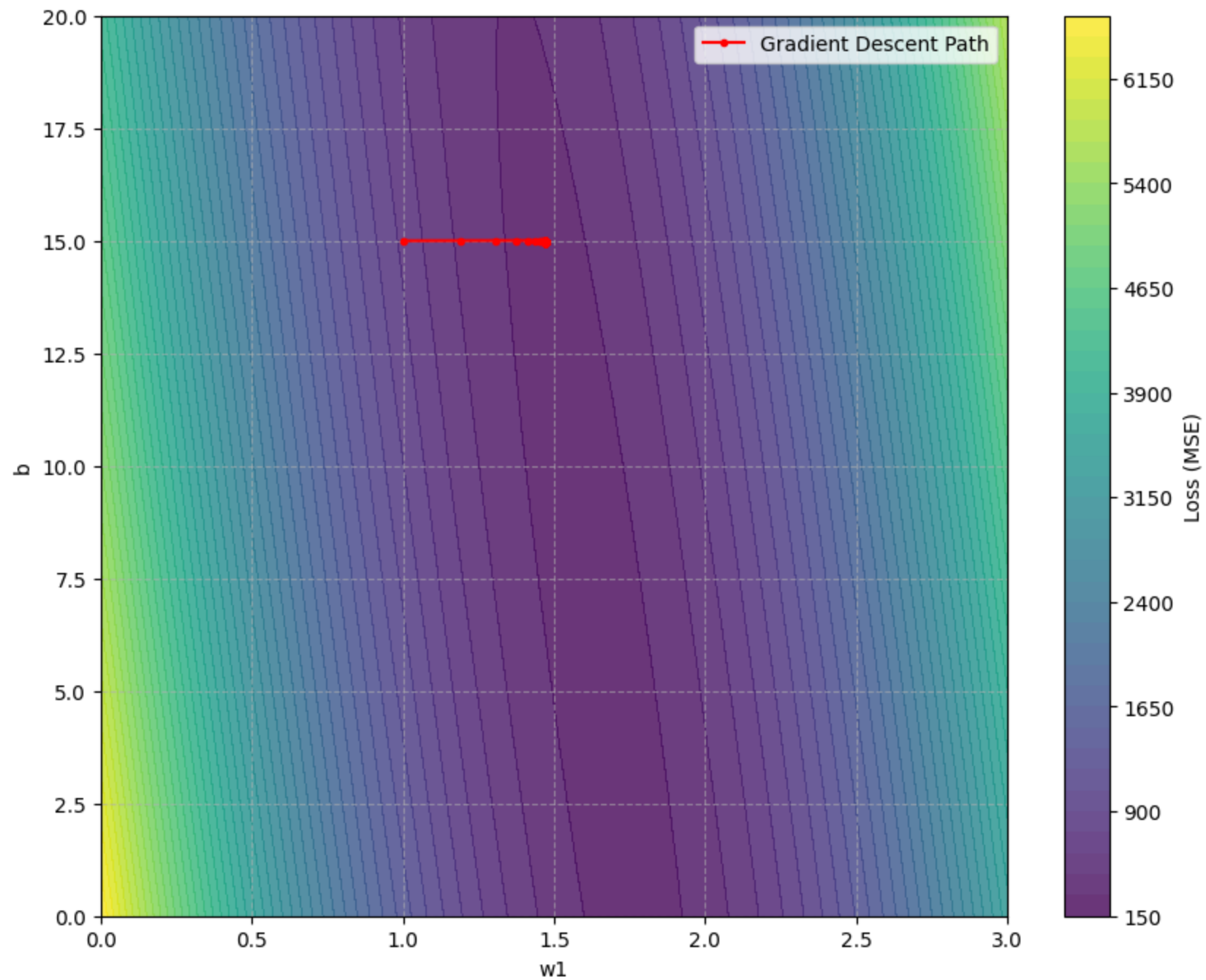
概念很簡單，
做起來不容易

$$w_1^0 = 1.0, b^0 = 15.0$$

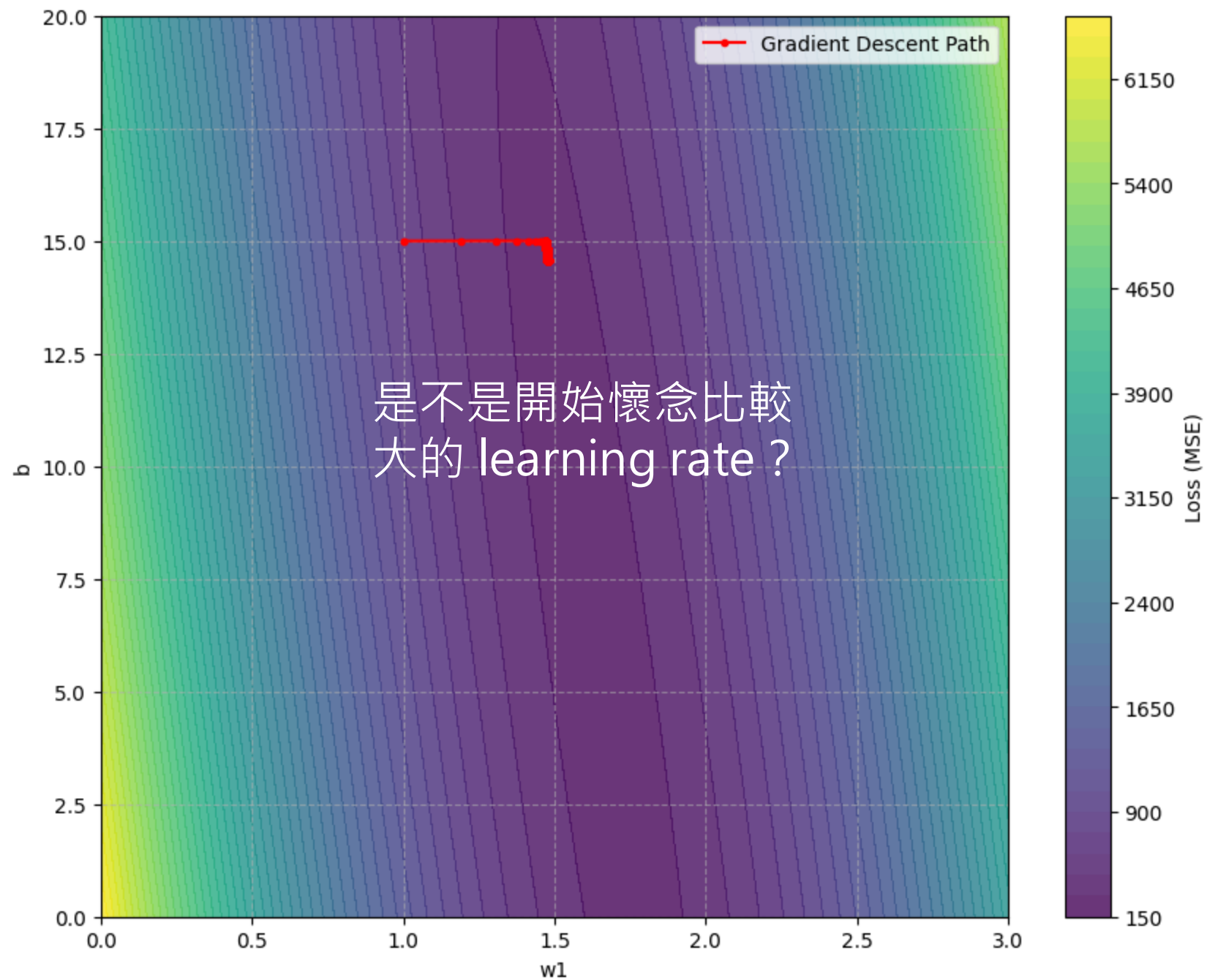
$$\eta = 0.0001$$



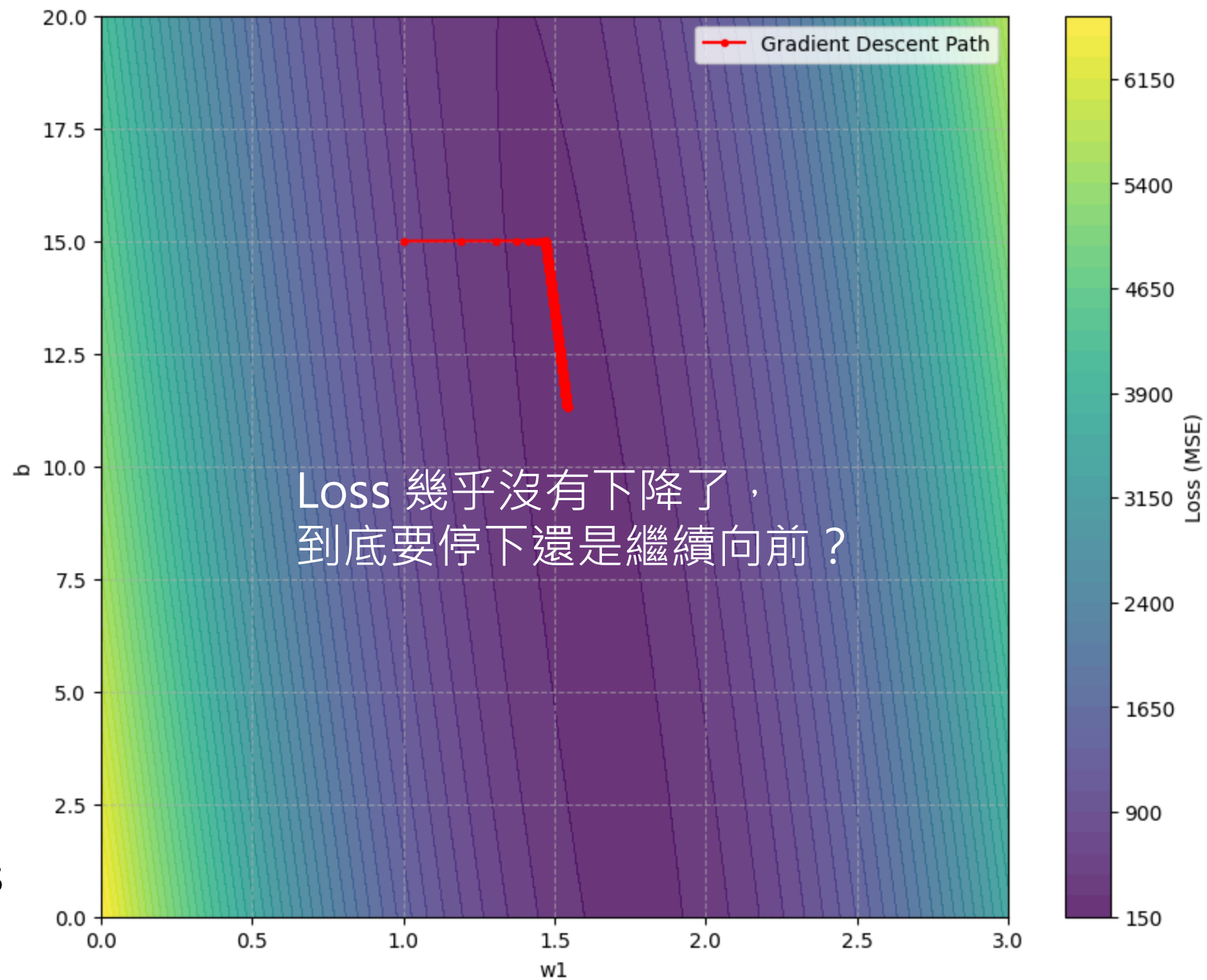
Update 100 times



Update 1,000 times



Update 10,000 times



參數更新太慢？

$$\theta^* = \arg \min_{\theta} L(\theta)$$

➤ (Randomly) Pick initial values θ^0

➤ Compute gradient $g^0 = \nabla L(\theta^0)$

$$\theta^1 \leftarrow \theta^0 - \eta g^0$$

➤ Compute gradient $g^1 = \nabla L(\theta^1)$

$$\theta^2 \leftarrow \theta^1 - \eta g^1$$

➤ Compute gradient $g^2 = \nabla L(\theta^2)$

$$\theta^3 \leftarrow \theta^2 - \eta g^2$$

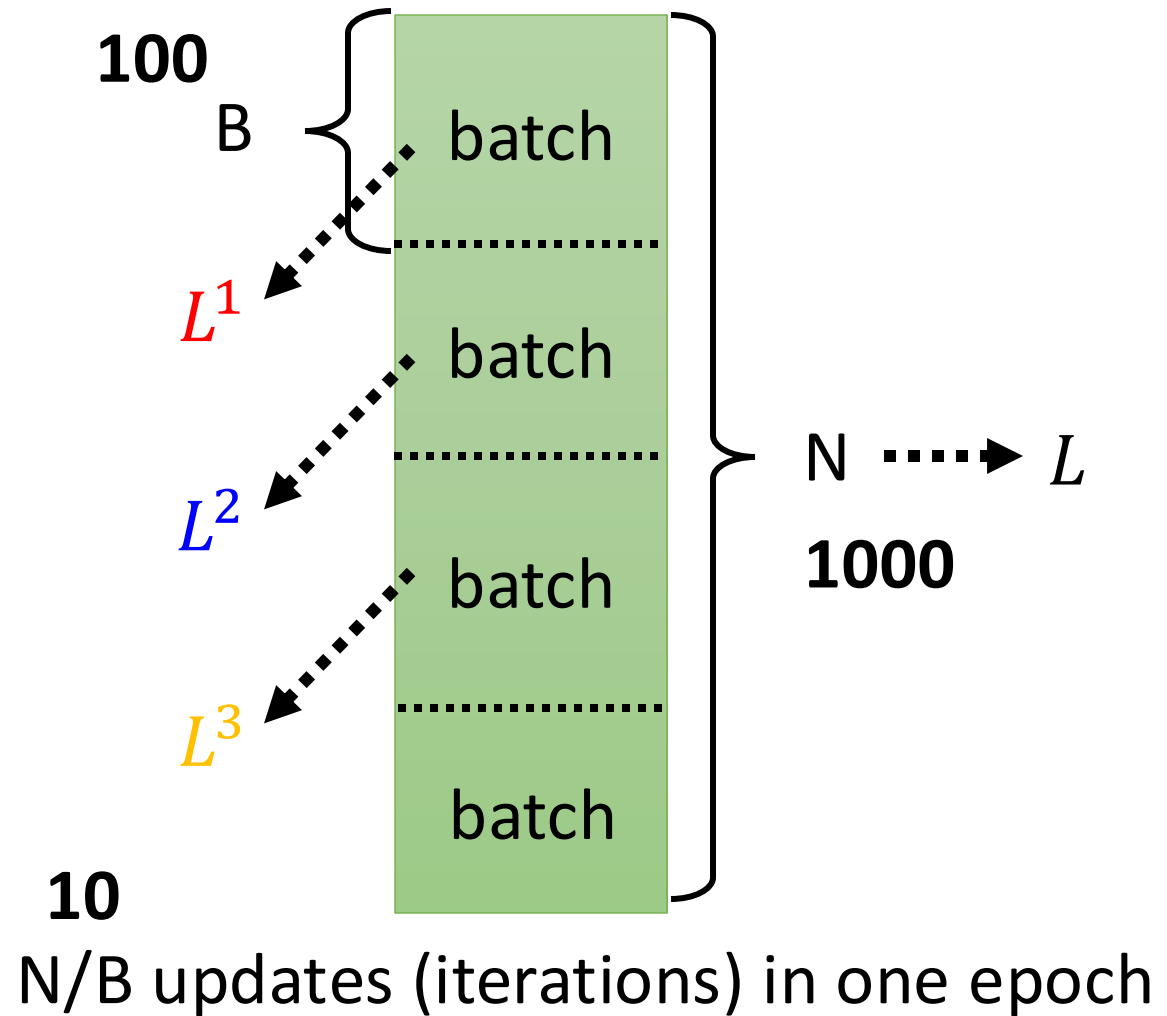
$$L = \frac{1}{N} \sum_{i=1}^N \dots \dots$$

如果訓練資料很多，要等很久才能更新一次參數

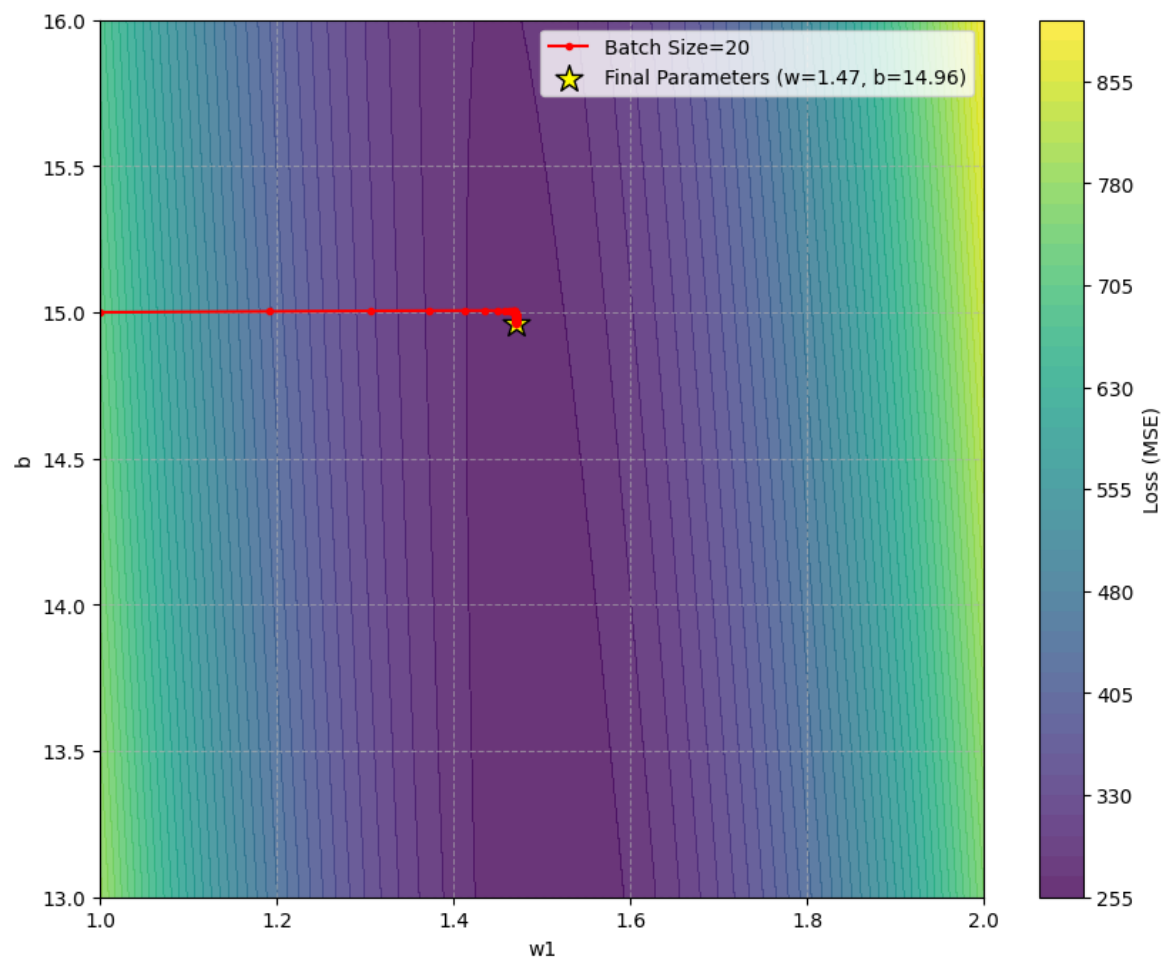
迫不及待更新參數

- (Randomly) Pick initial values θ^0
- Compute gradient $g^0 = \nabla L^1(\theta^0)$
 $\theta^1 \leftarrow \theta^0 - \eta g^0$
- Compute gradient $g^1 = \nabla L^2(\theta^1)$
 $\theta^2 \leftarrow \theta^1 - \eta g^1$
- Compute gradient $g^2 = \nabla L^3(\theta^2)$
 $\theta^3 \leftarrow \theta^2 - \eta g^2$

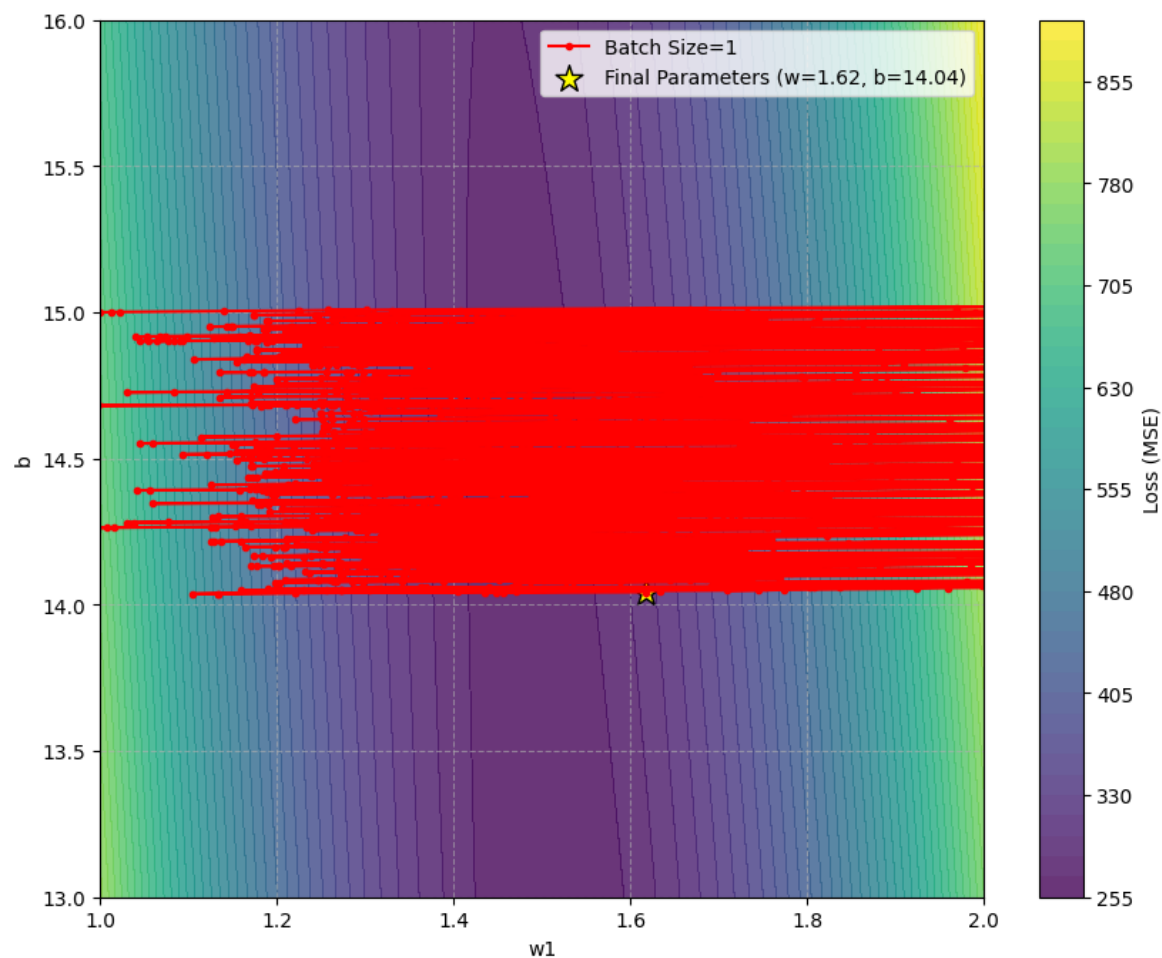
1 **epoch** = see all the batches once



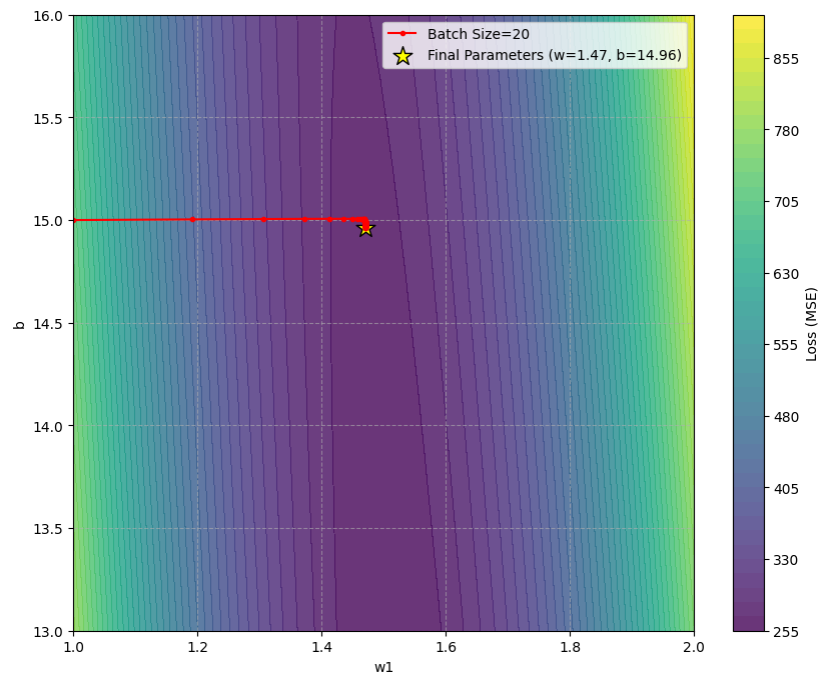
epochs = 100



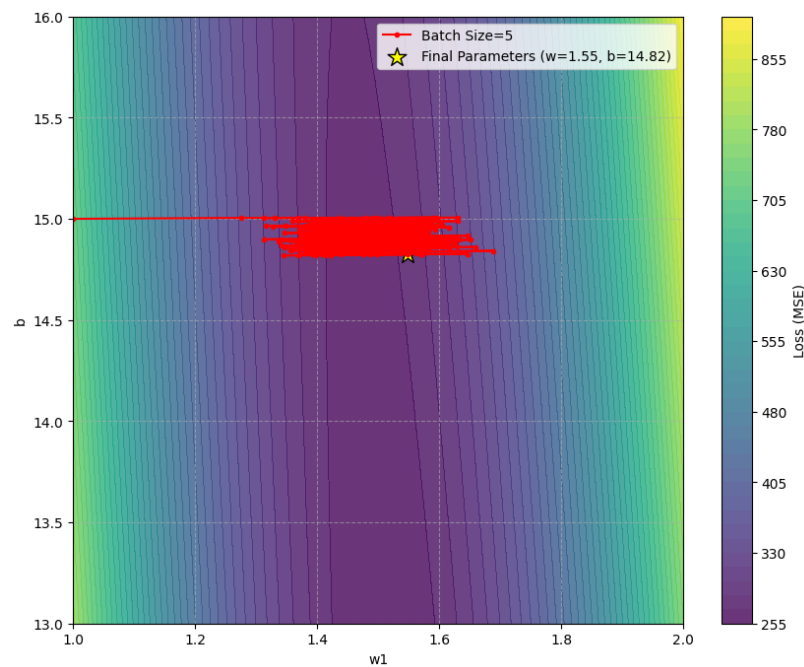
Batch size = all training data
(Full Batch)



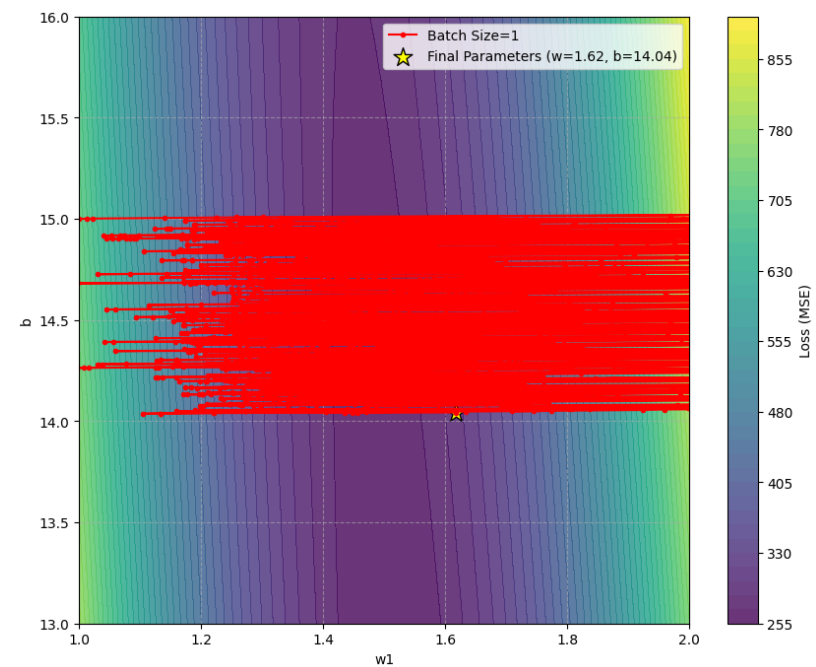
Batch size = 1
(Stochastic Gradient Descent, SGD)



Batch size
= all training data



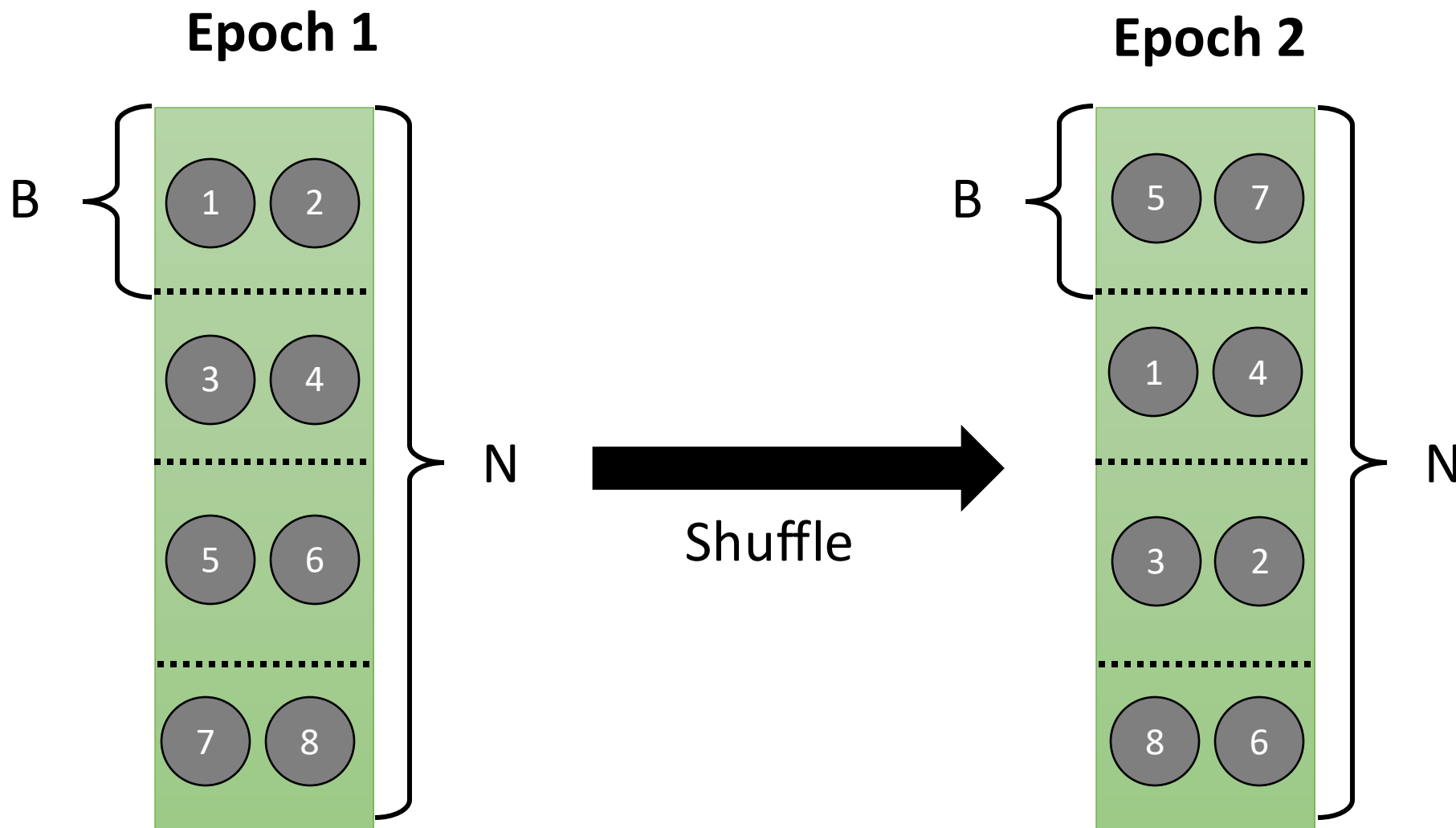
Batch size
= 5



Batch size
= 1

又多了一個可以調的 hyperparameter

Shuffle



步驟一：
我要什麼

+

步驟二：
我有哪些選擇



步驟三：
選一個最好的

$$L = \frac{1}{N} \sum_{i=1}^N (y^i - \hat{y}^i)^2 \text{ MSE}$$

$$L(w_1^*, b^*) = 240$$

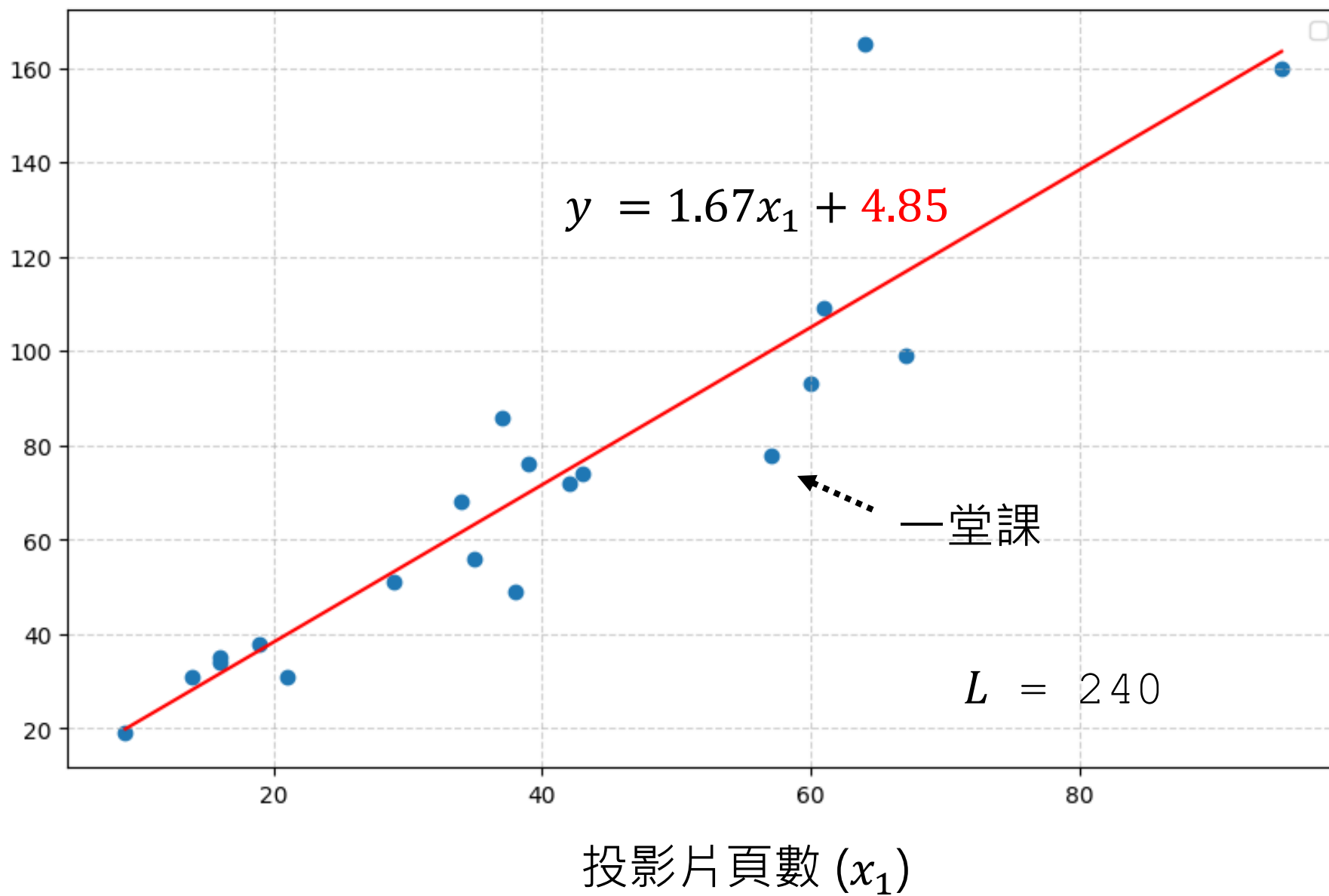
$$y = w_1 x_1 + b$$

$$y = 1.67x_1 + 4.85$$

$$w_1^*, b^* = \arg \min_{w_1, b} L(w_1, b)$$

$$w_1^* = 1.67, b^* = 4.85$$

課程
時長
(y)



步驟一：
我要什麼

+

步驟二：
我有哪些選擇



步驟三：
選一個最好的

$$L = \frac{1}{N} \sum_{i=1}^N (y^i - \hat{y}^i)^2 \text{ MSE}$$

$$L(w_1^*, b^*) = 240$$

$$y = w_1 x_1 + b$$

$$y = 1.67x_1 + 4.85$$

$$w_1^*, b^* = \arg \min_{w_1, b} L(w_1, b)$$

$$w_1^* = 1.67, b^* = 4.85$$

測試今天這堂課

測試 (Testing)



步驟一：
我要什麼

+

步驟二：
我有哪些選擇



步驟三：
選一個最好的

$$L = \frac{1}{N} \sum_{i=1}^N (y^i - \hat{y}^i)^2 \text{ MSE}$$

$$L(w_1^*, b^*) = 240$$

$$y = w_1 x_1 + b$$

$$y = 1.67x_1 + 4.85$$

$$w_1^*, b^* = \arg \min_{w_1, b} L(w_1, b)$$

$$w_1^* = 1.67, b^* = 4.85$$



測試今天這堂課

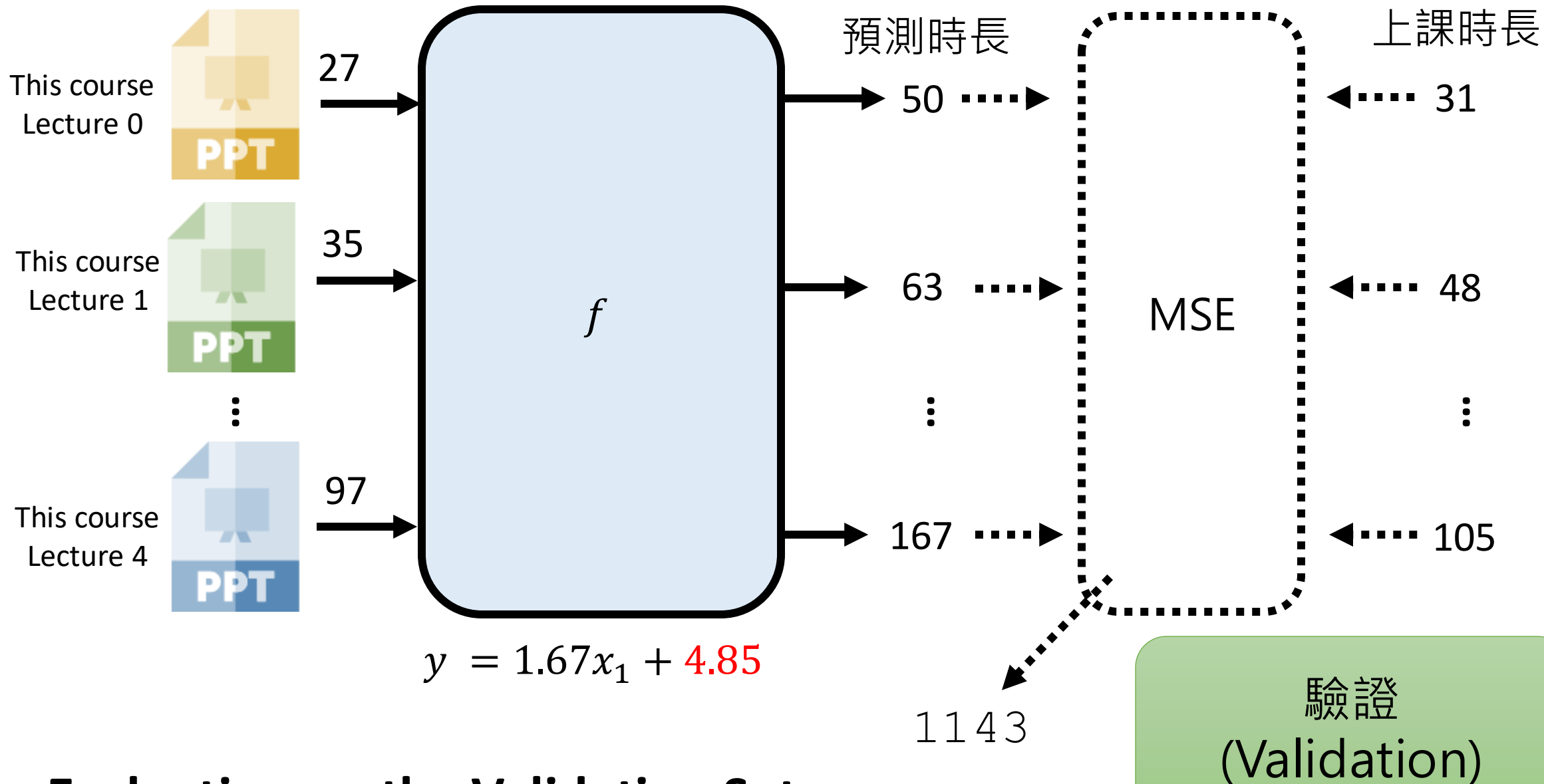
測試 (Testing)

真的大考

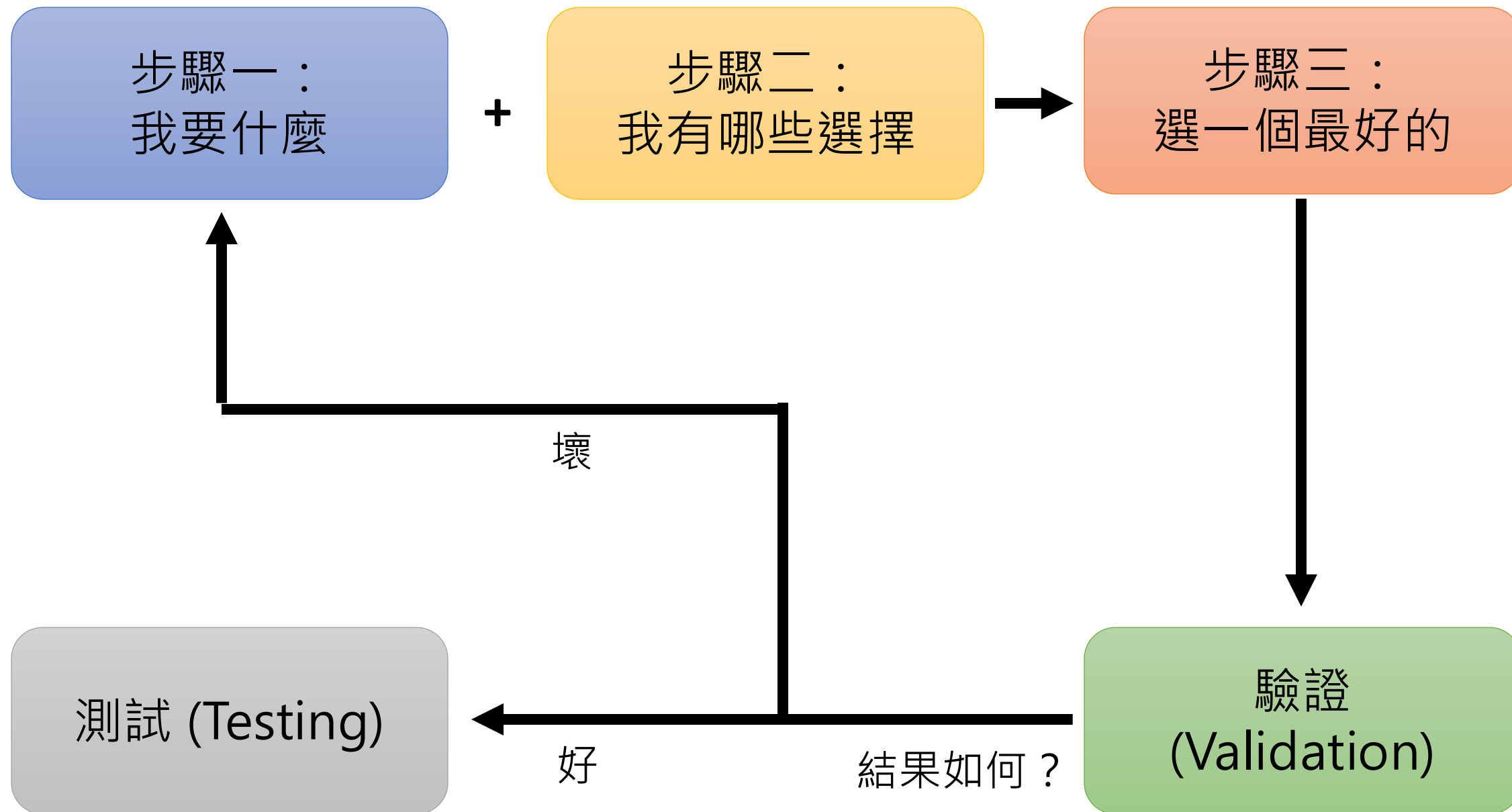
驗證
(Validation)

模擬考





模擬考



步驟一：
我要什麼

+

步驟二：
我有哪些選擇



步驟三：
選一個最好的

在《機器學習》2021 上計算 MSE

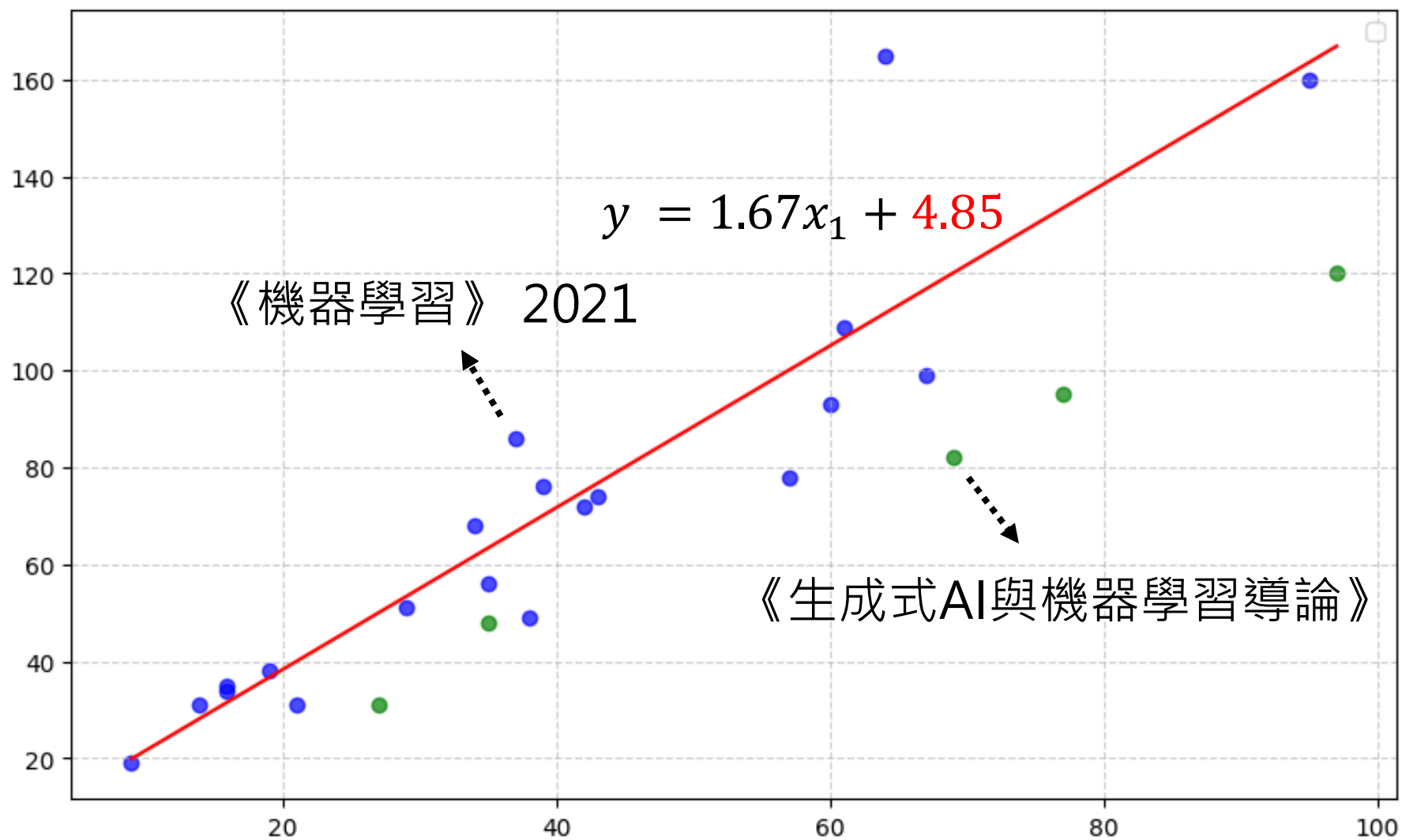
會不會有巨大差異？

在《生成式AI與機器學習
導論》2025 上計算 MSE

你以為你要的目標，
跟實際上的目標不一致

驗證
(Validation)

課程
時長
(y)



投影片頁數 (x_1)

更換 訓練資料

【機器學習2021】(中文版)

Hung-yi Lee - 1/40

🔄 🔗

▶  **49:59** **【機器學習2021】預測本頻道觀看人數(上) - 機器學習基本...**
Hung-yi Lee

2  **58:35** **【機器學習2021】預測本頻道觀看人數(下) - 深度學習基本...**
Hung-yi Lee

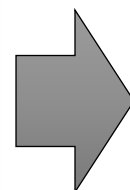
3  **51:23** **【機器學習2021】機器學習任務攻略**
Hung-yi Lee

4  **33:45** **【機器學習2021】類神經網路訓練不起來怎麼辦(一)：局...**
Hung-yi Lee

5  **30:59** **【機器學習2021】類神經網路訓練不起來怎麼辦(二)：批...**
Hung-yi Lee

6  **37:42** **【機器學習2021】類神經網路訓練不起來怎麼辦(三)：自動...**
Hung-yi Lee

7  **19:27** **【機器學習2021】類神經網路訓練不起來怎麼辦(四)：損失...**
Hung-yi Lee



【生成式AI導論 2024】

Hung-yi Lee - 1/20

🔄 🔗

▶  **25:40** **【生成式AI導論 2024】第0講：課程說明 (17:15 有茉莉...**
Hung-yi Lee

2  **29:29** **【生成式AI導論 2024】第1講：生成式AI是什麼？**
Hung-yi Lee

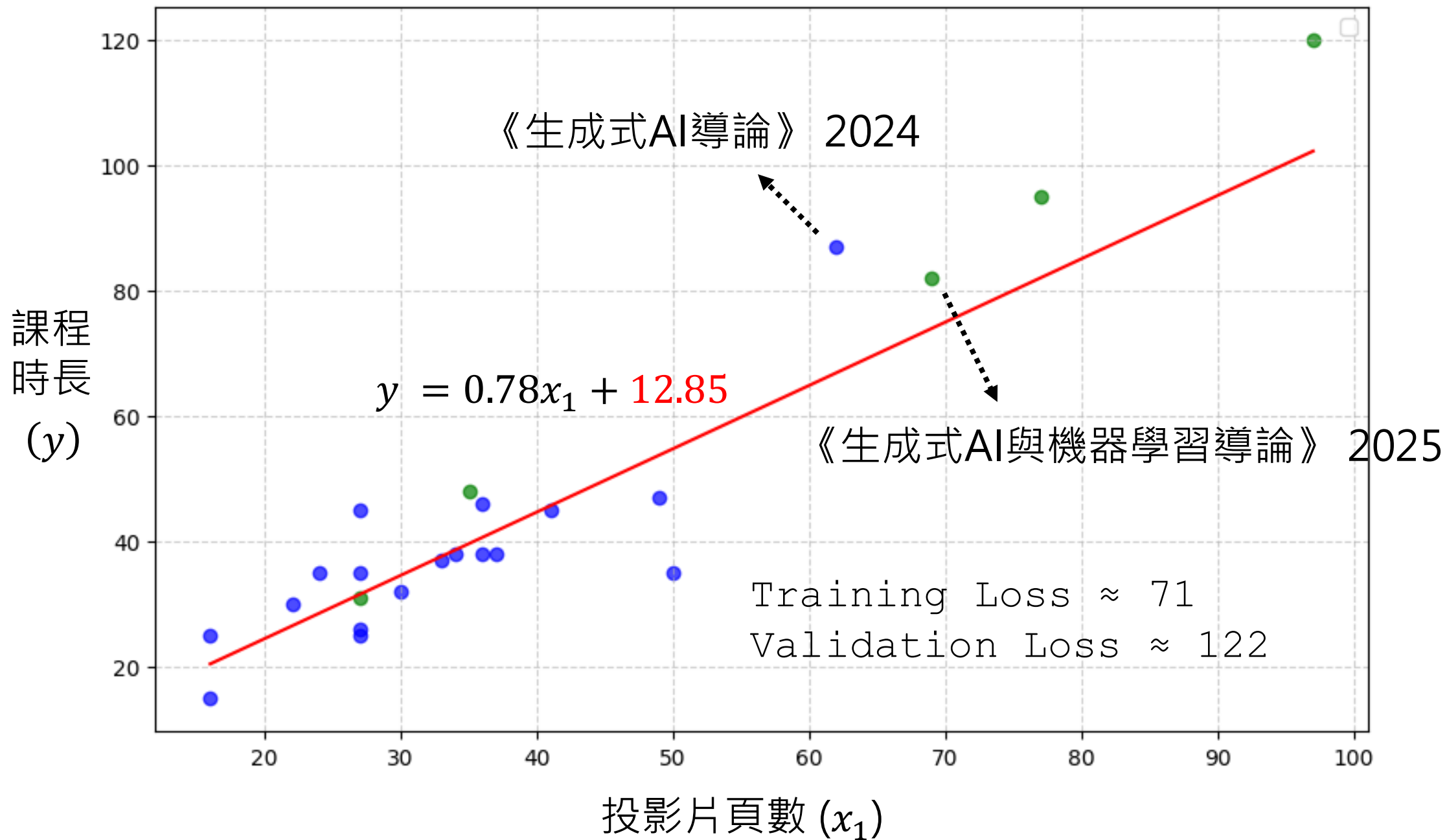
3  **26:06** **【生成式AI導論 2024】第2講：今日的生成式人工智慧...**
Hung-yi Lee

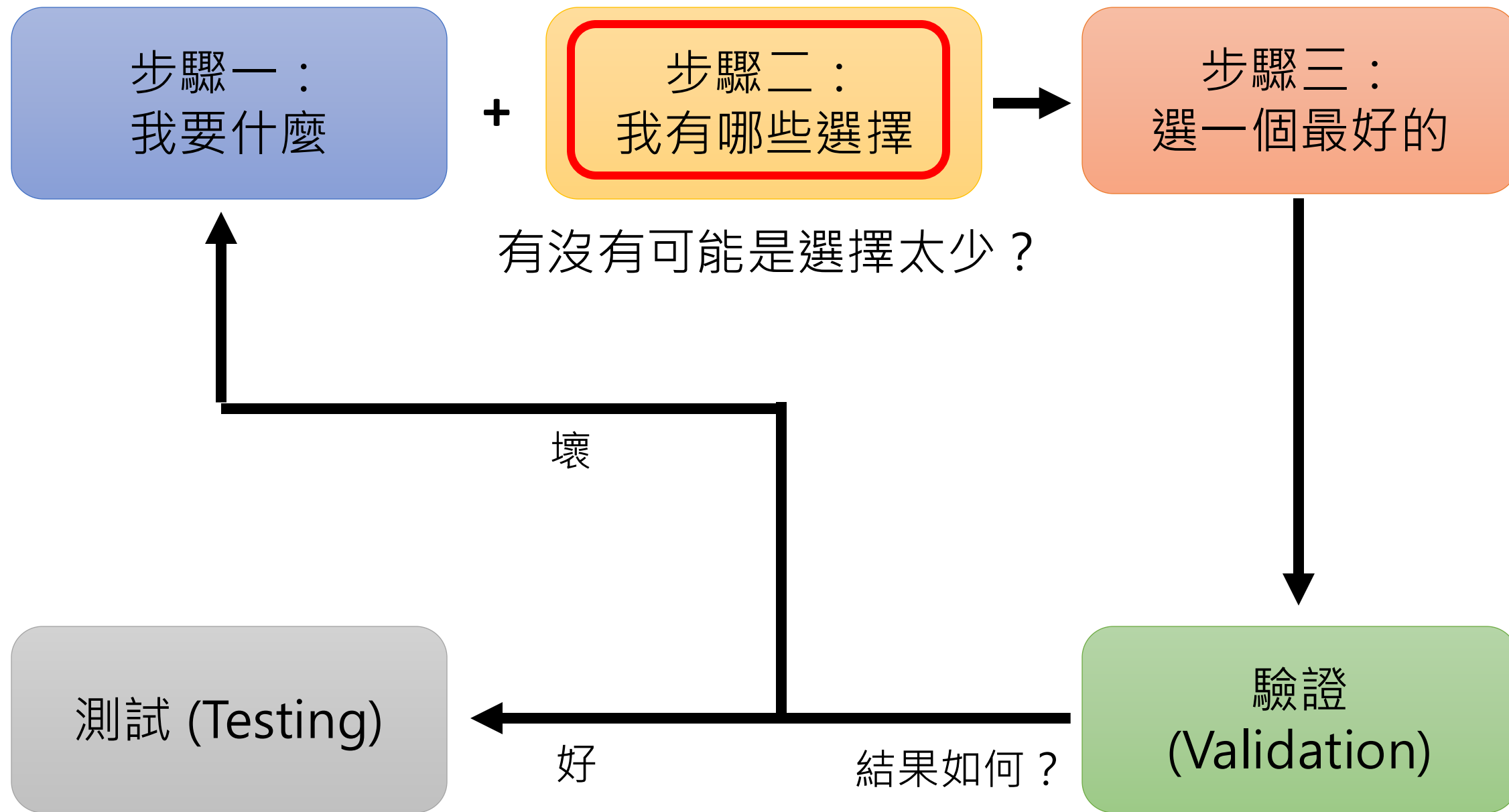
4  **34:35** **【生成式AI導論 2024】第3講：訓練不了人工智慧？你...**
Hung-yi Lee

5  **47:22** **【生成式AI導論 2024】第4講：訓練不了人工智慧？你...**
Hung-yi Lee

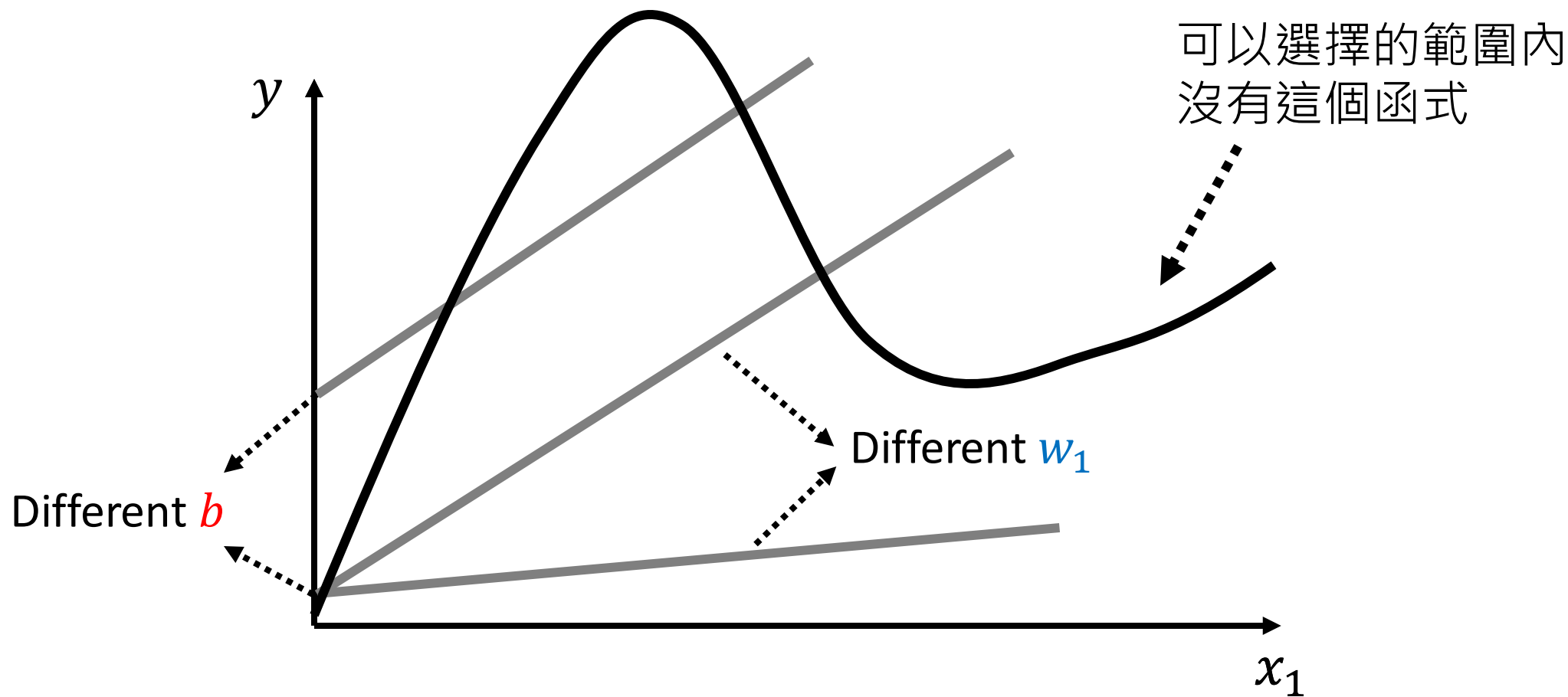
6  **25:20** **【生成式AI導論 2024】第5講：訓練不了人工智慧？你...**
Hung-yi Lee

7  **34:26** **【生成式AI導論 2024】第6講：大型語言模型修練史 - ...**
Hung-yi Lee



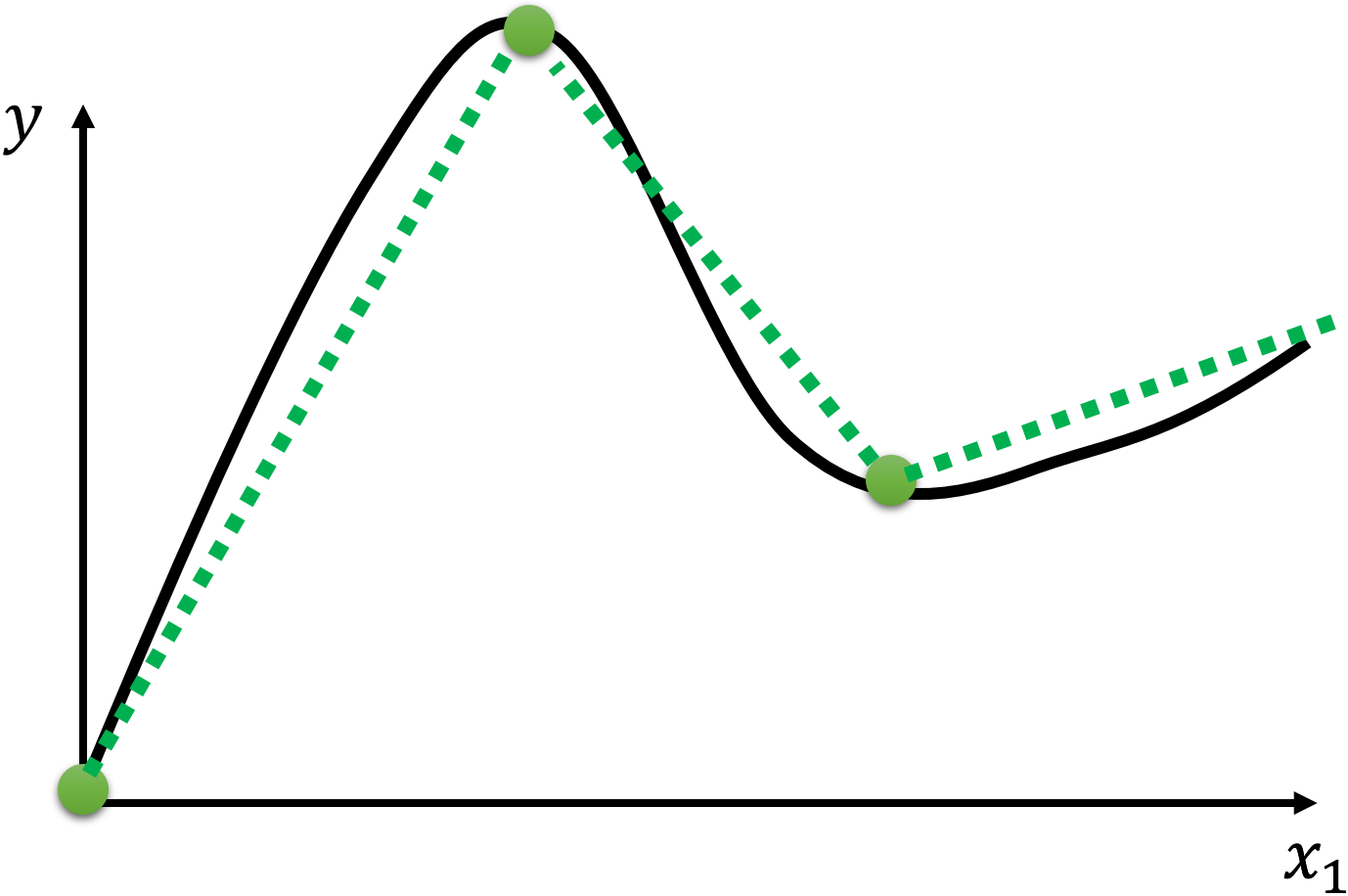


$$y = w_1 x_1 + b$$

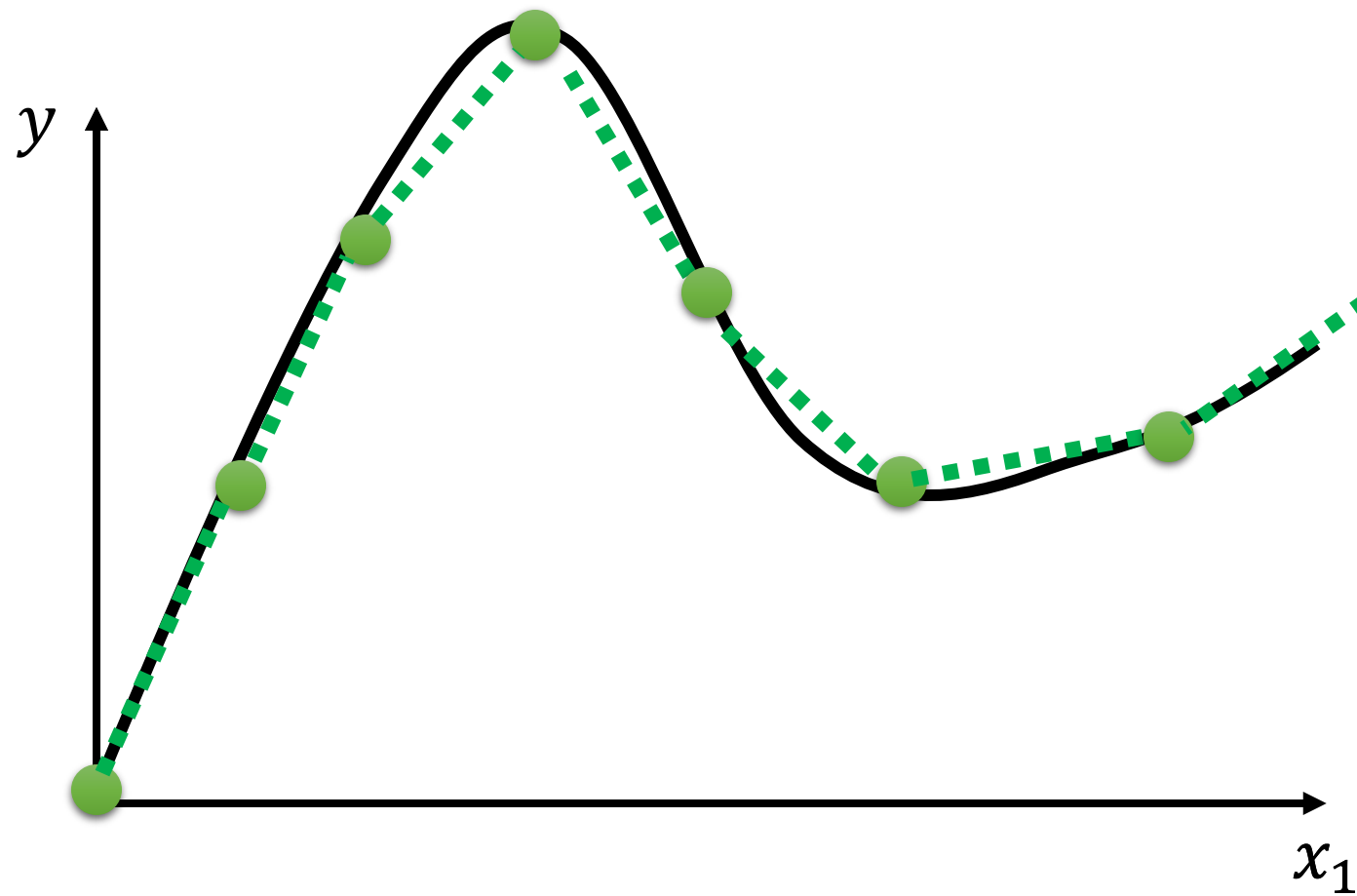


畫一個有機會包含所有函數的範圍

Piecewise
Linear

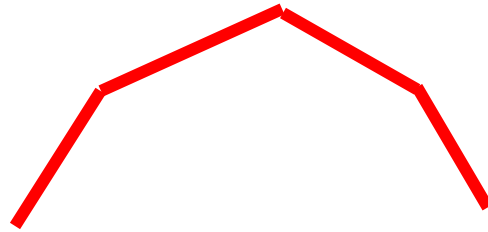


**Piecewise
Linear**



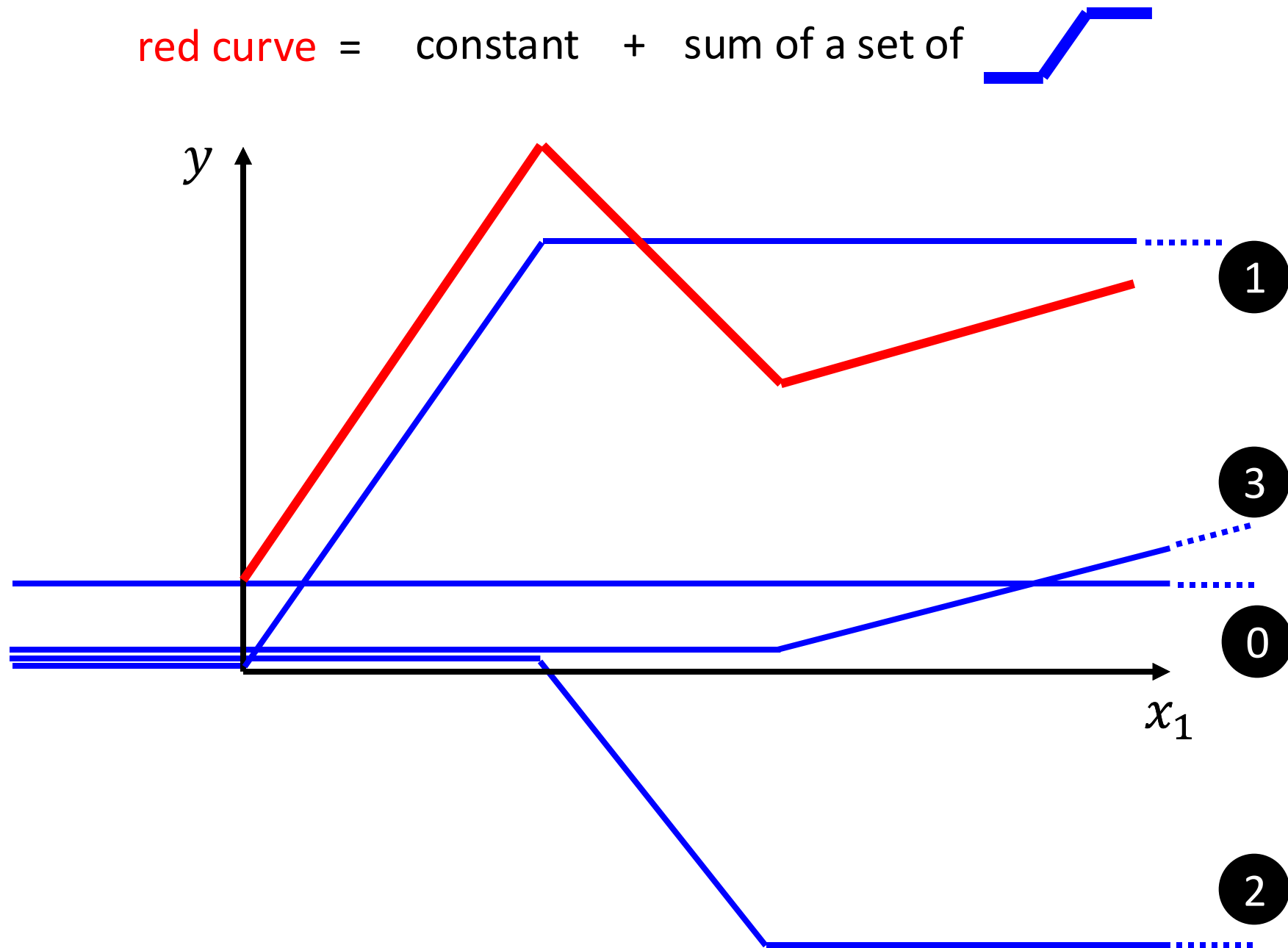
All Piecewise Linear Curves

= constant + sum of a set of 

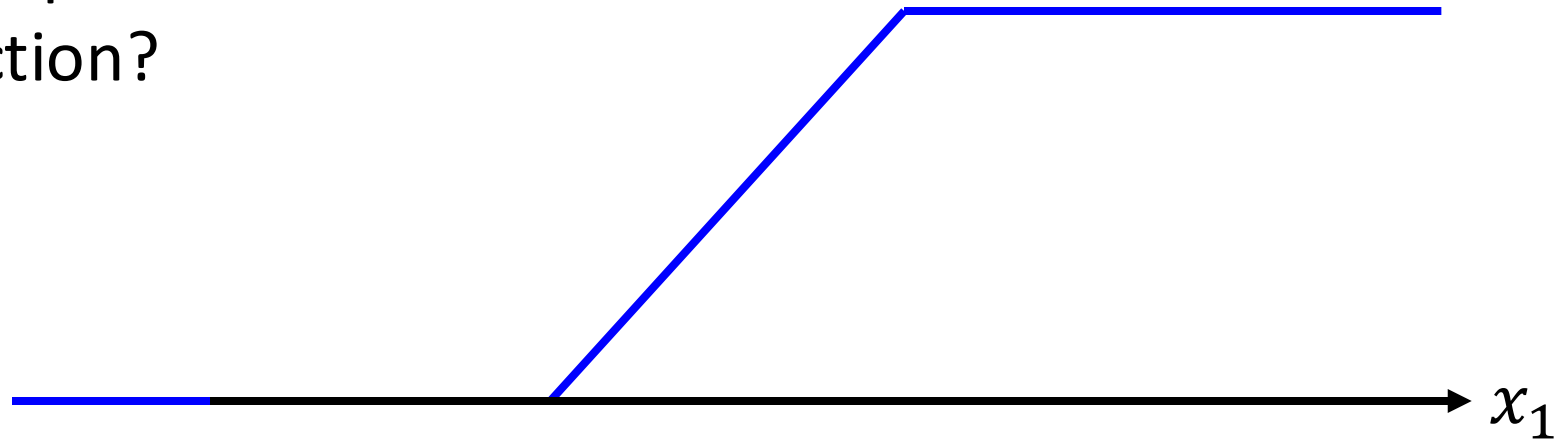


More pieces require more 

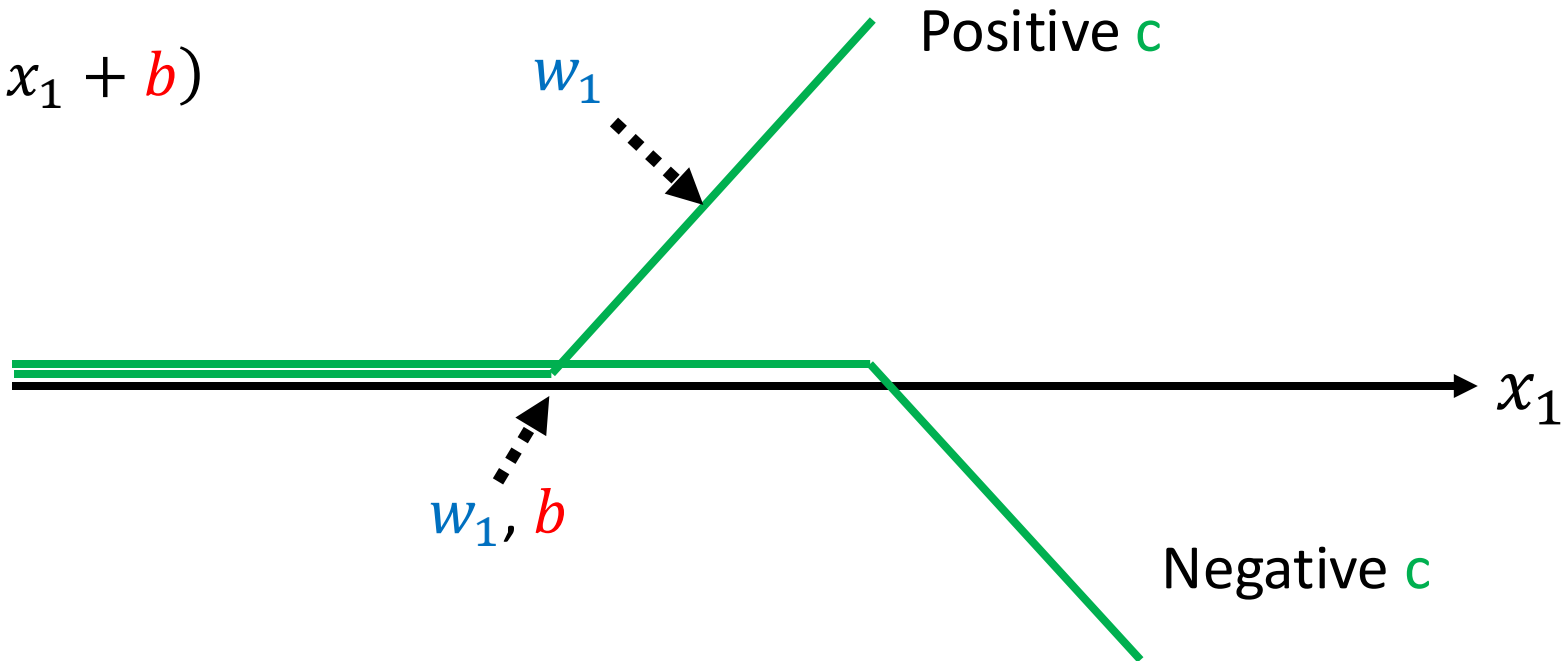
red curve = constant + sum of a set of 



How to represent
this function?




$$c \max(0, w_1 x_1 + b)$$



Any Curves \approx Piecewise Linear Curves

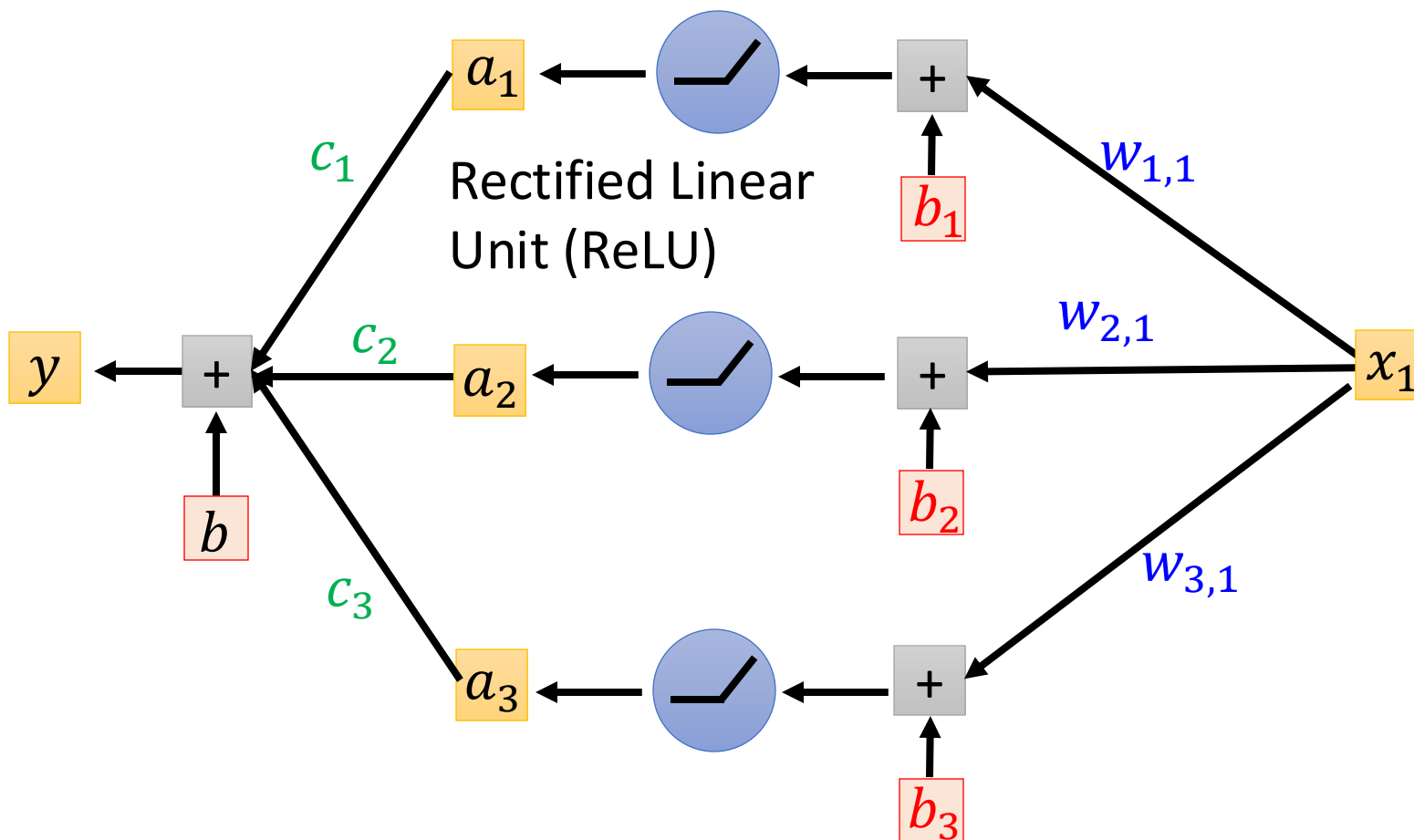
= constant + sum of a set of 

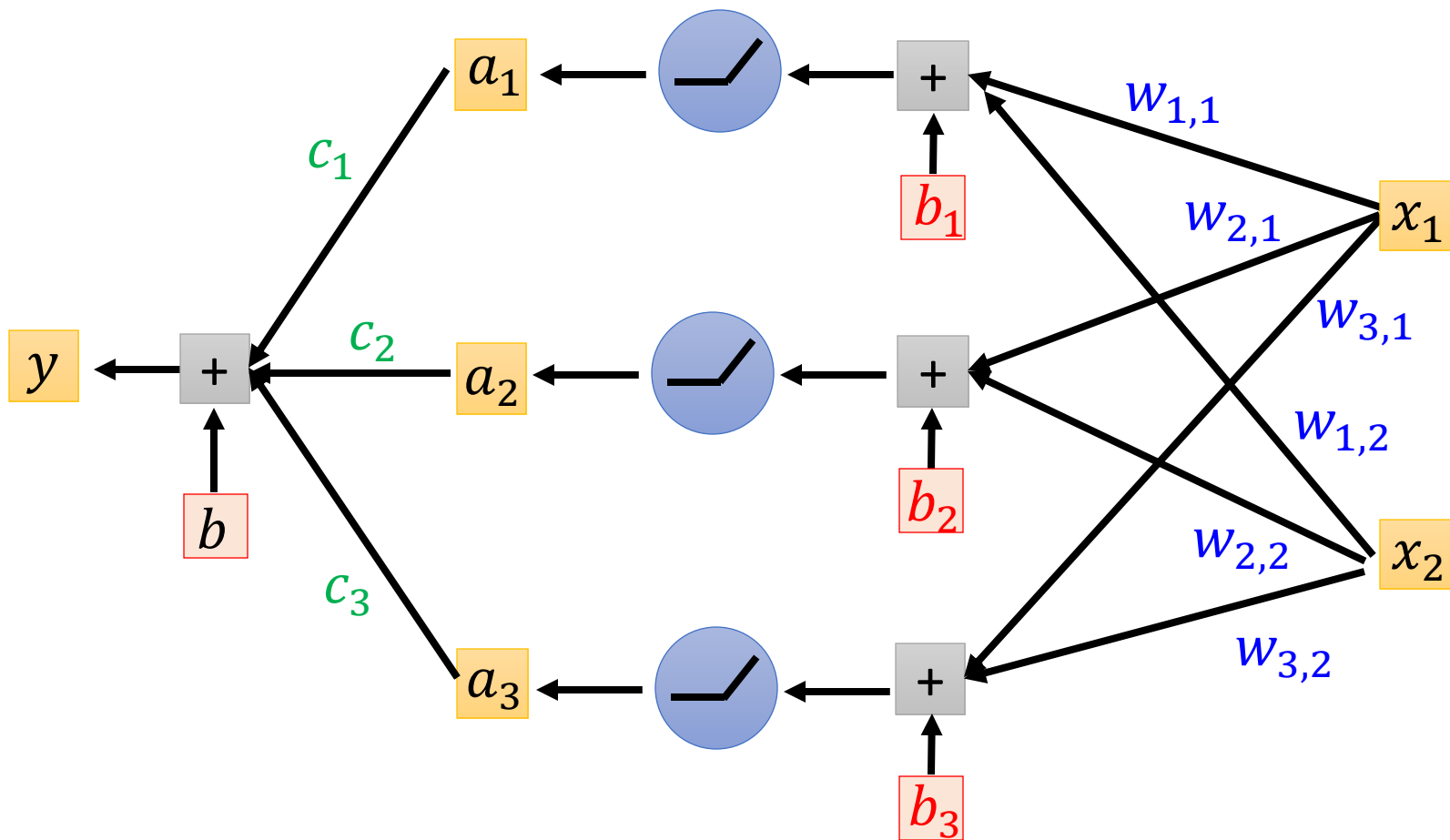
= constant + sum of a larger set of  $c \max(0, w_1 x_1 + b)$

└──┘
H

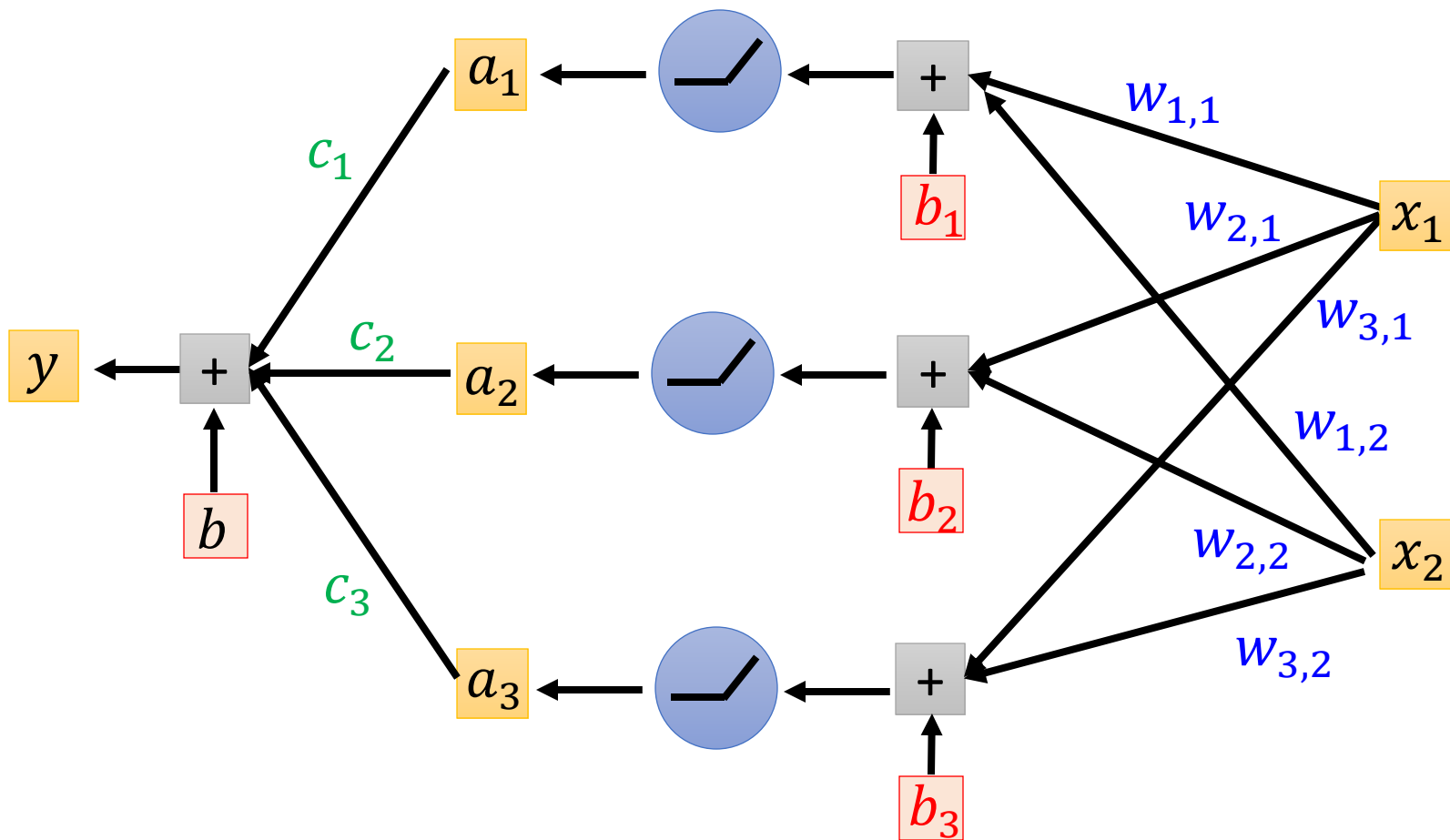
$$y = b + \sum_{i=1}^H c_i \max(0, w_{i,1} x_1 + b_i)$$

$$y = b + \sum_{i=1}^H \underbrace{c_i \max(0, w_{i,1}x_1 + b_i)}_{a_i}$$



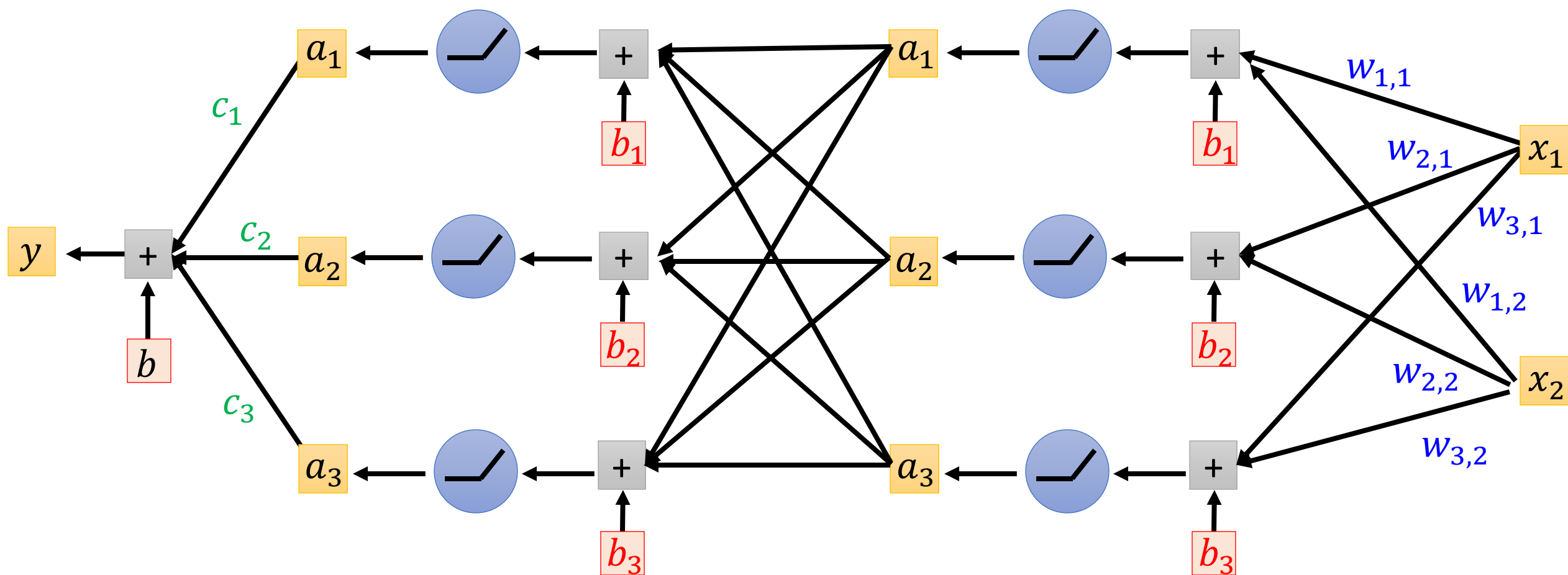


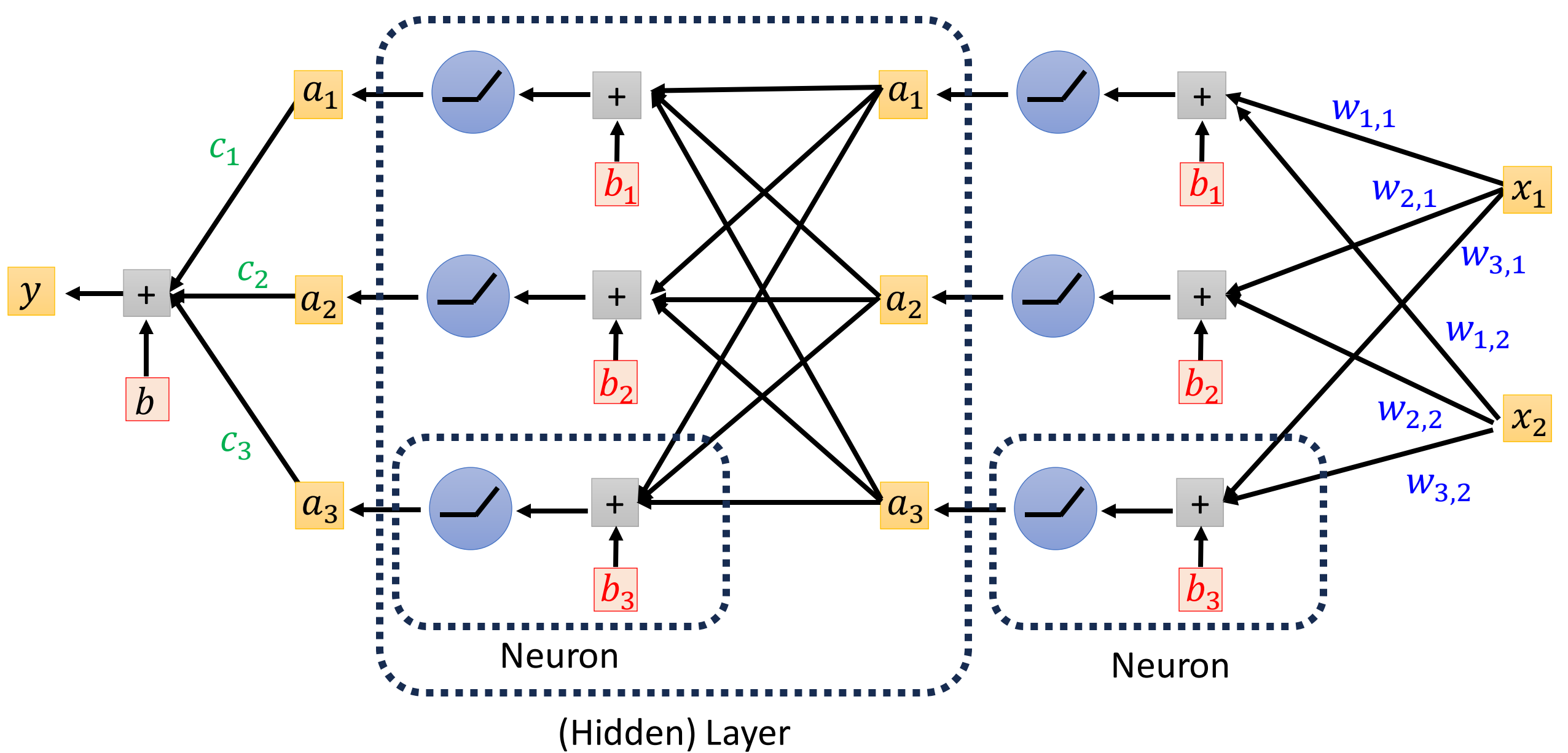
$$y = \textcolor{red}{b} + [\textcolor{green}{c}_1 \quad \textcolor{green}{c}_2 \quad \textcolor{green}{c}_3] \begin{bmatrix} \textcolor{brown}{a}_1 \\ \textcolor{brown}{a}_2 \\ \textcolor{brown}{a}_3 \end{bmatrix} \quad \begin{bmatrix} \textcolor{brown}{a}_1 \\ \textcolor{brown}{a}_2 \\ \textcolor{brown}{a}_3 \end{bmatrix} = \sigma \left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right)$$



$$y = b + c^T a \quad a = \sigma(b + W x)$$

$$\mathbf{a}' = \sigma(\mathbf{b}' + \mathbf{W}' \mathbf{a}) \quad \mathbf{a} = \sigma(\mathbf{b} + \mathbf{W} \mathbf{x})$$



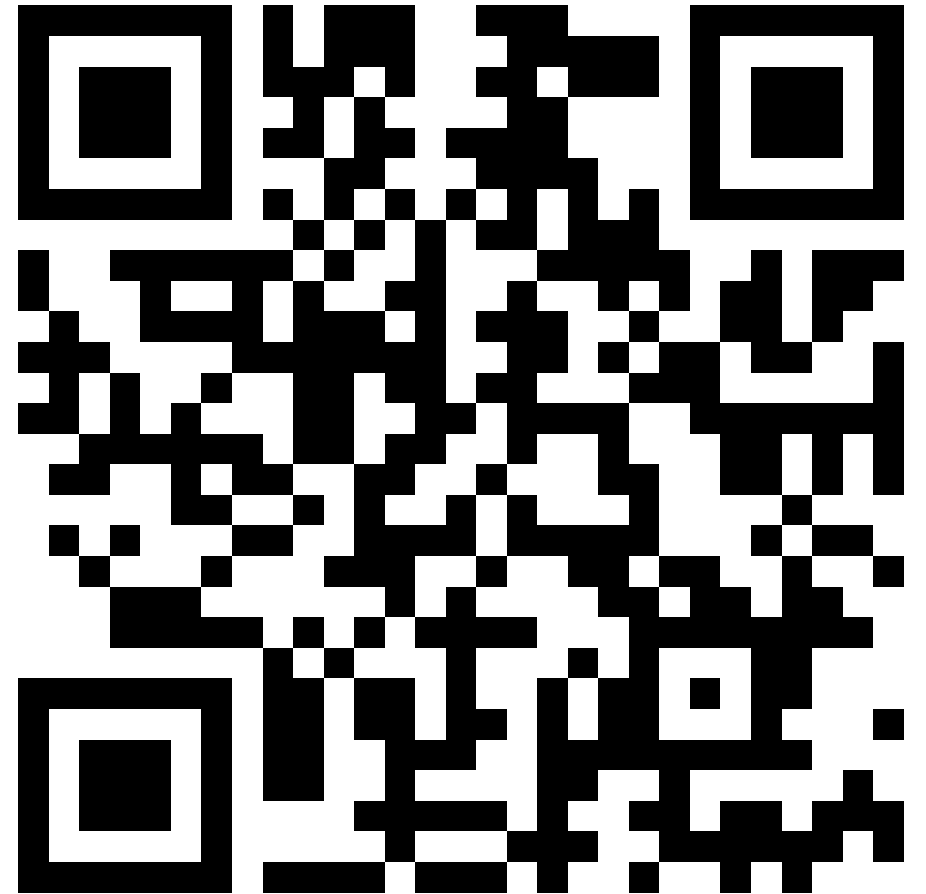


Neural Network

Many Layers \longrightarrow 深度學習 (Deep Learning)

Backpropagation

Computing gradients in an efficient way

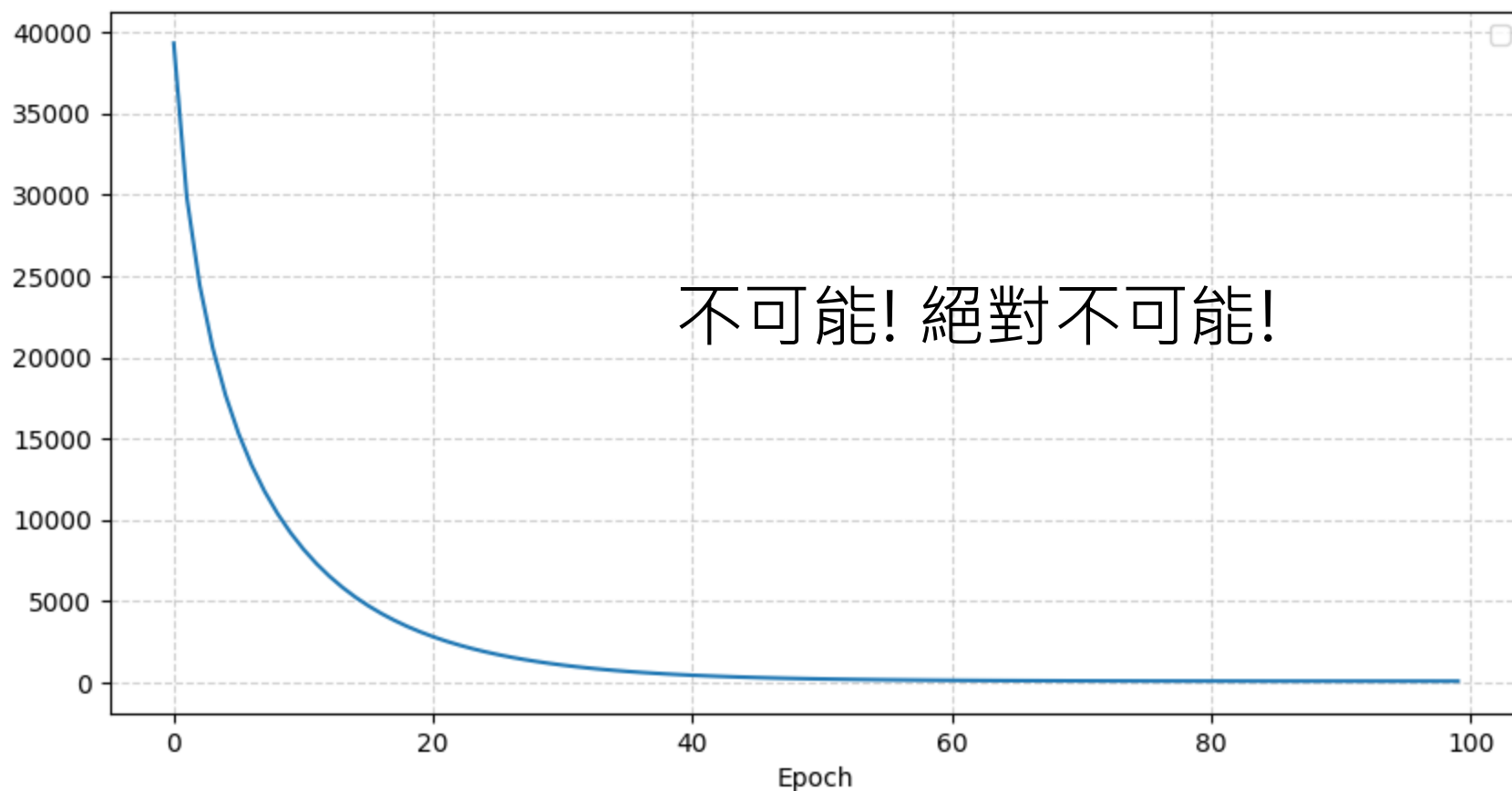


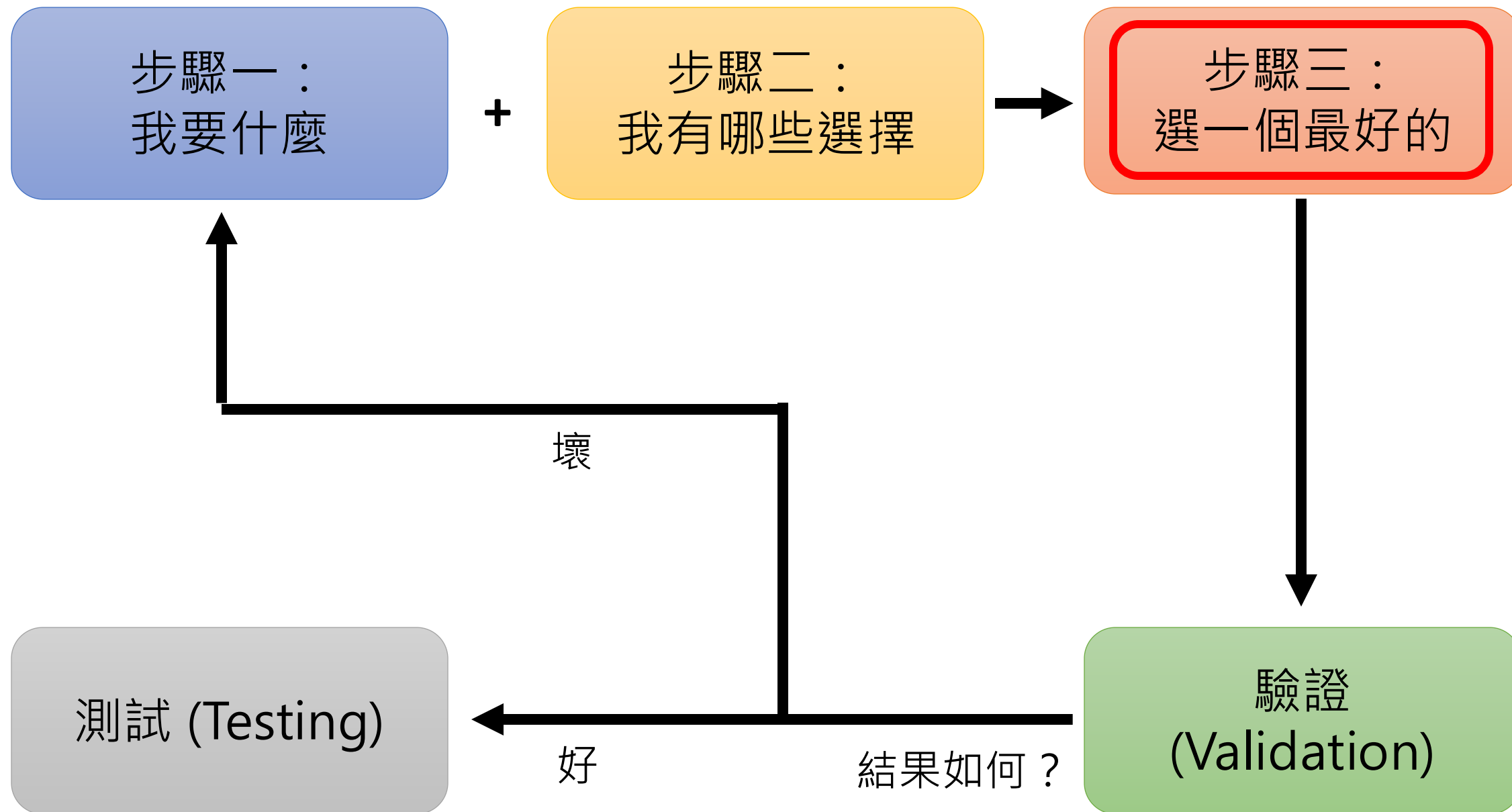
<https://youtu.be/ibJpTrp5mcE>

$$y = w_1 x_1 + b \longrightarrow y = b + \sum_{i=1}^H c_i \max(0, w_{i,1} x_1 + b_i)$$

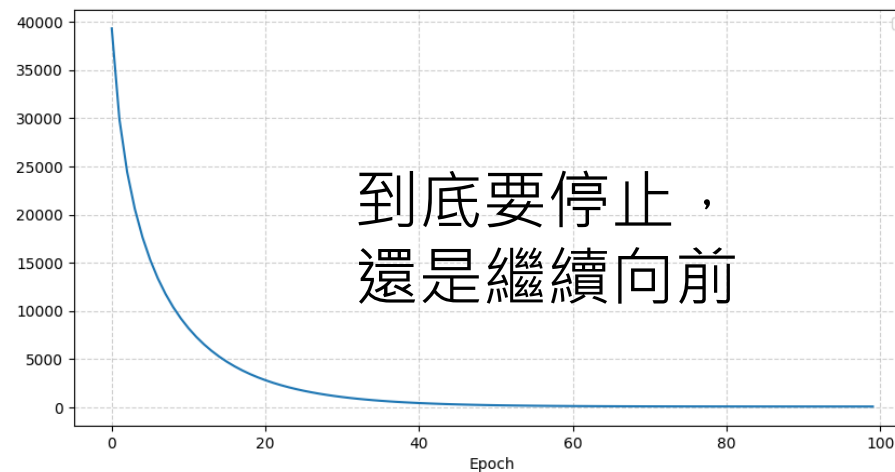
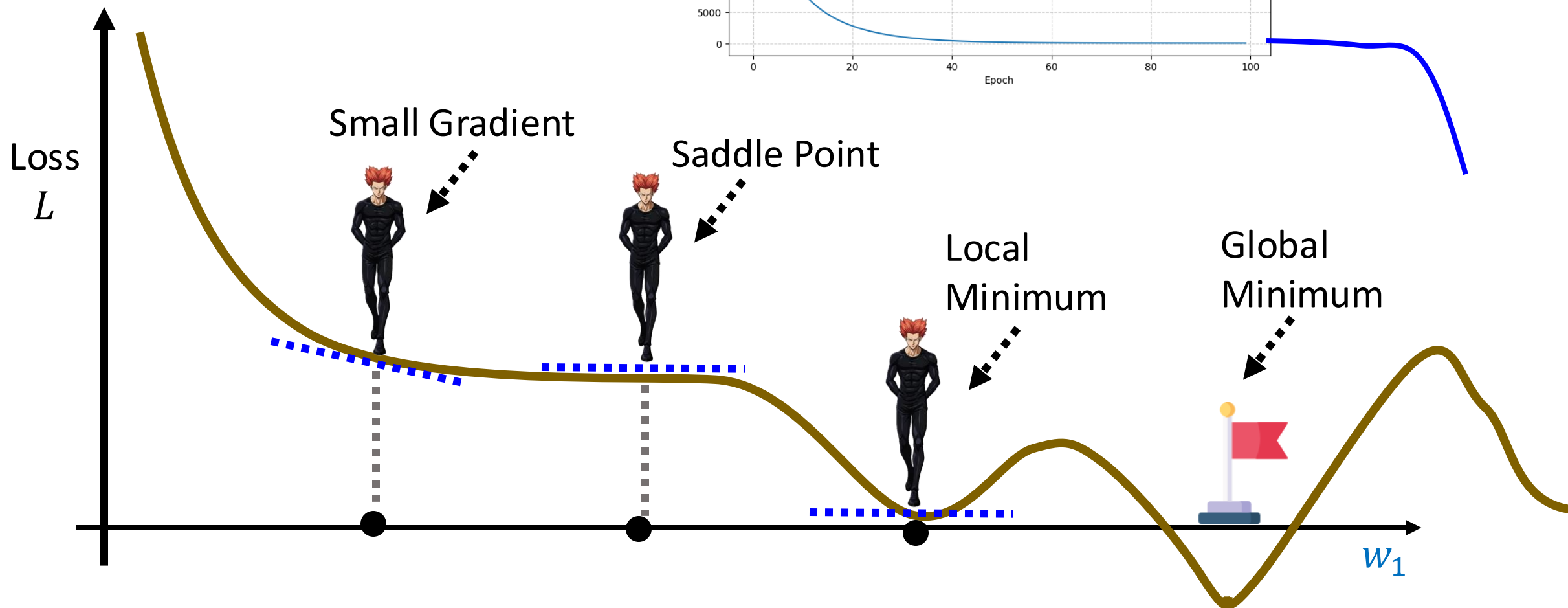
Training Loss ≈ 71
 Training Loss ≈ 80
 $H = 100$

參數太多了，
只能看
Loss Curve



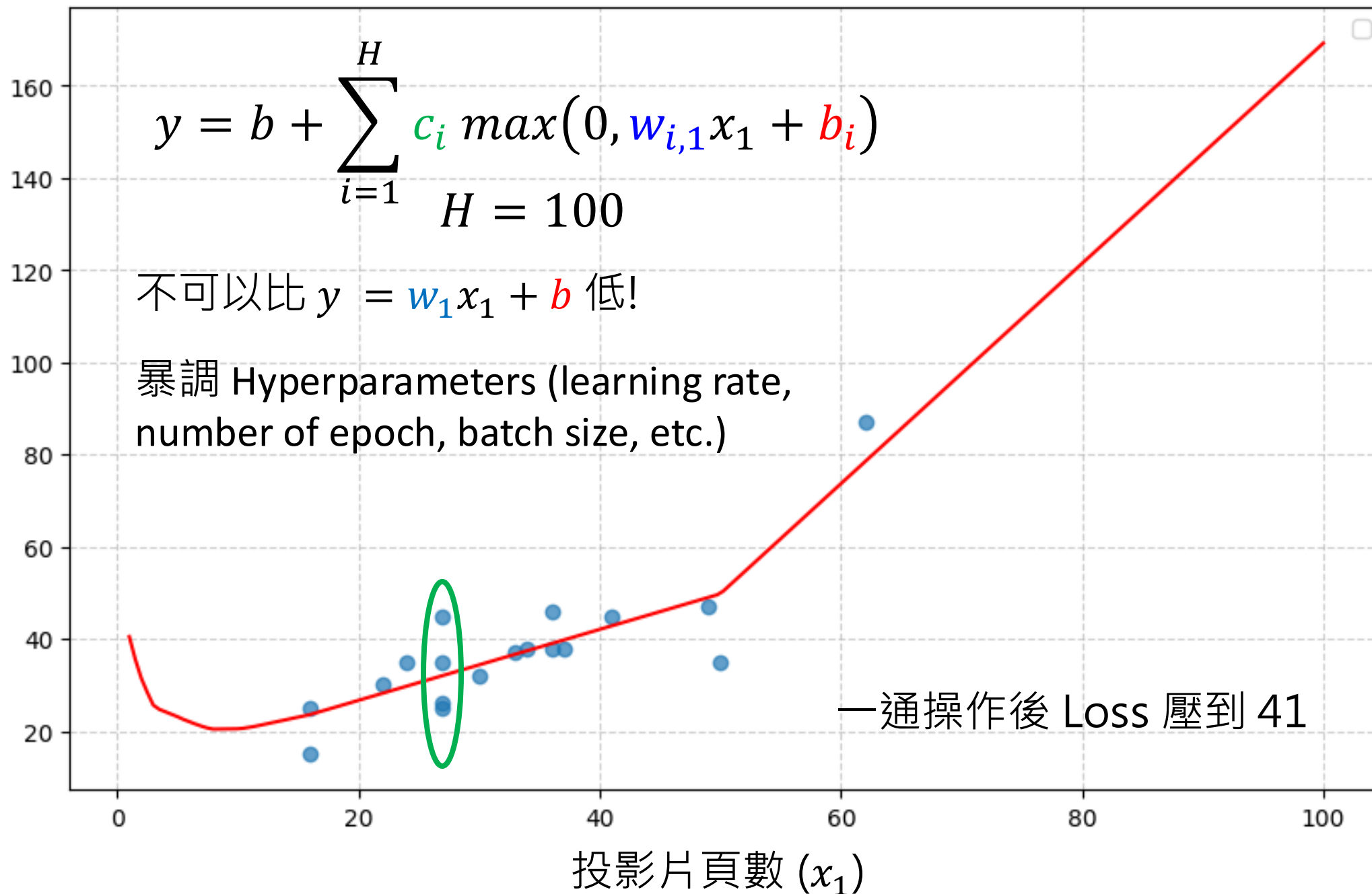


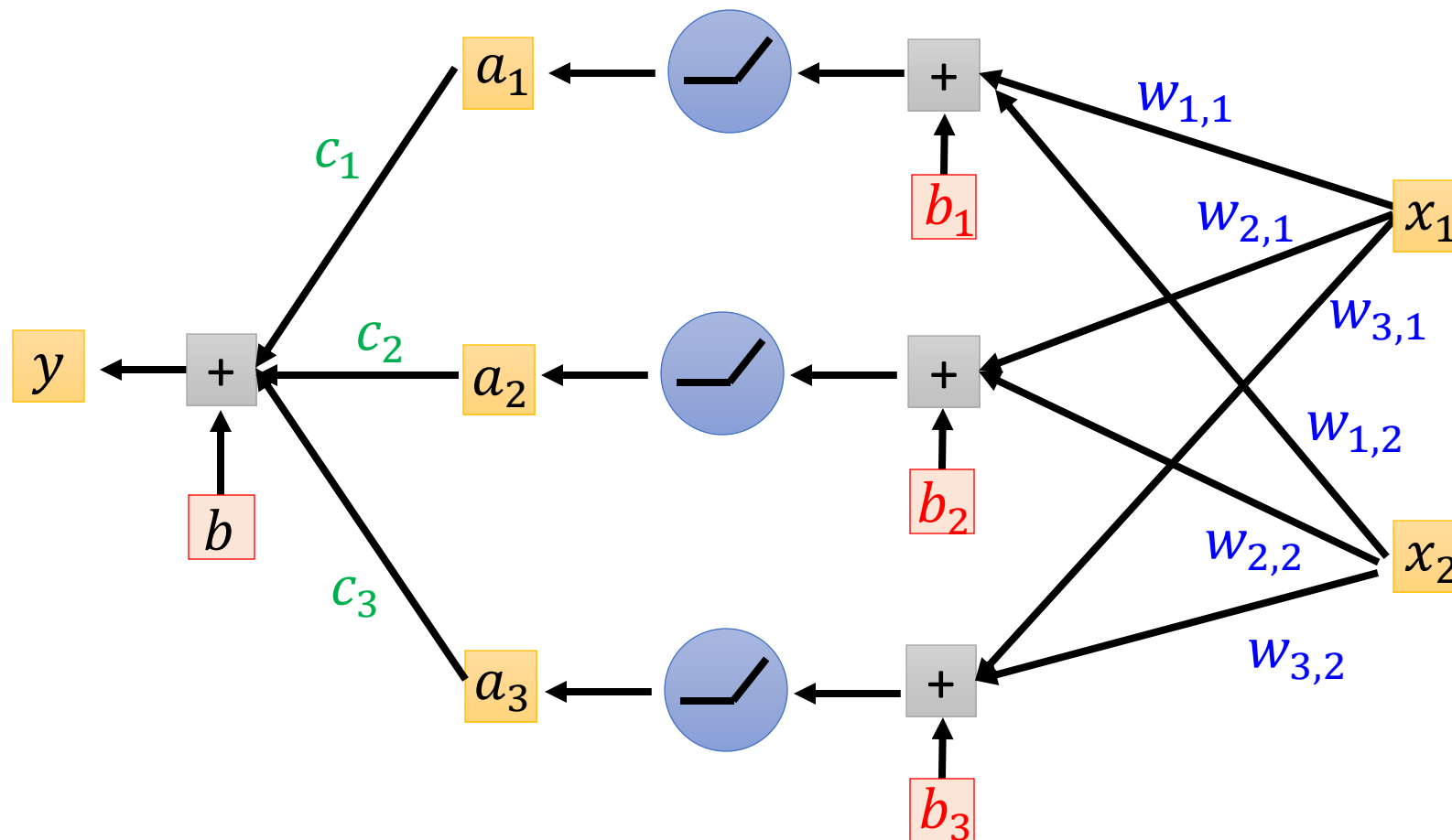
Optimization Fail!



誰知道呢

課程
時長
(y)





有多少頁投影片

~~投影片中總共有
多少字~~

每頁投影片平均
多少字

	Linear	Deep	+ No. of Word	+ Avg of Word
Training:	Loss \approx 71	41	40	22
Validation:	Loss \approx 122			1307

步驟一：
我要什麼

+

步驟二：
我有哪些選擇



步驟三：
選一個最好的



驗證
(Validation)

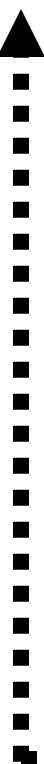
Loss \approx 21

劃定的範圍越大，越容易
Overfitting

Overfitting

差距巨大

Loss \approx 1307



如果世上一切函數都可以選，會怎麼樣？

訓練資料



Duration: 10



Duration: 20



Duration: 30

$$f_{\text{lazy}}(\text{PPT}) = 10$$

$$f_{\text{lazy}}(\text{PPT}) = 20$$

$$f_{\text{lazy}}(\text{PPT}) = 30$$

$$f_{\text{lazy}}(\text{Other}) = 0$$

在訓練資料上的
Loss 為 0

你說這是不是在訓練資料上 Loss 最低的函式？

$$f_{lazy}(\text{PPT}_{\text{green}}) = 10$$

$$f_{lazy}(\text{PPT}_{\text{light blue}}) = 20$$

$$f_{lazy}(\text{PPT}_{\text{blue}}) = 30$$

$$f_{lazy}(\text{Other}) = 0$$

在訓練資料上的
Loss 為 0

驗證資料



Duration: 15



Duration: 32



Duration: 33

$$f_{lazy}(\text{PPT}_{\text{grey}}) = 0$$

$$f_{lazy}(\text{PPT}_{\text{yellow}}) = 0$$

$$f_{lazy}(\text{PPT}_{\text{orange}}) = 0$$

在驗證資料上的
Loss 為炸裂

Function with Unknown Parameters

$$f(\text{Digimon}) = \begin{cases} \text{Digimon} & \text{If } e(\text{Digimon}) \geq h \\ \text{Pokémon} & \text{If } e(\text{Digimon}) < h \end{cases}$$

f_h : function with threshold h

$\mathcal{H} = \{1, 2, \dots, 10,000\}$ $|\mathcal{H}|$: number of candidate functions (model “complexity”)

【機器學習 2022】再探寶可夢、數碼寶貝分類器 – 淺談機器學習原理



https://youtu.be/_j9MVVcvyZI?si=cKWY8QmyS3-wX4l9

Overfitting

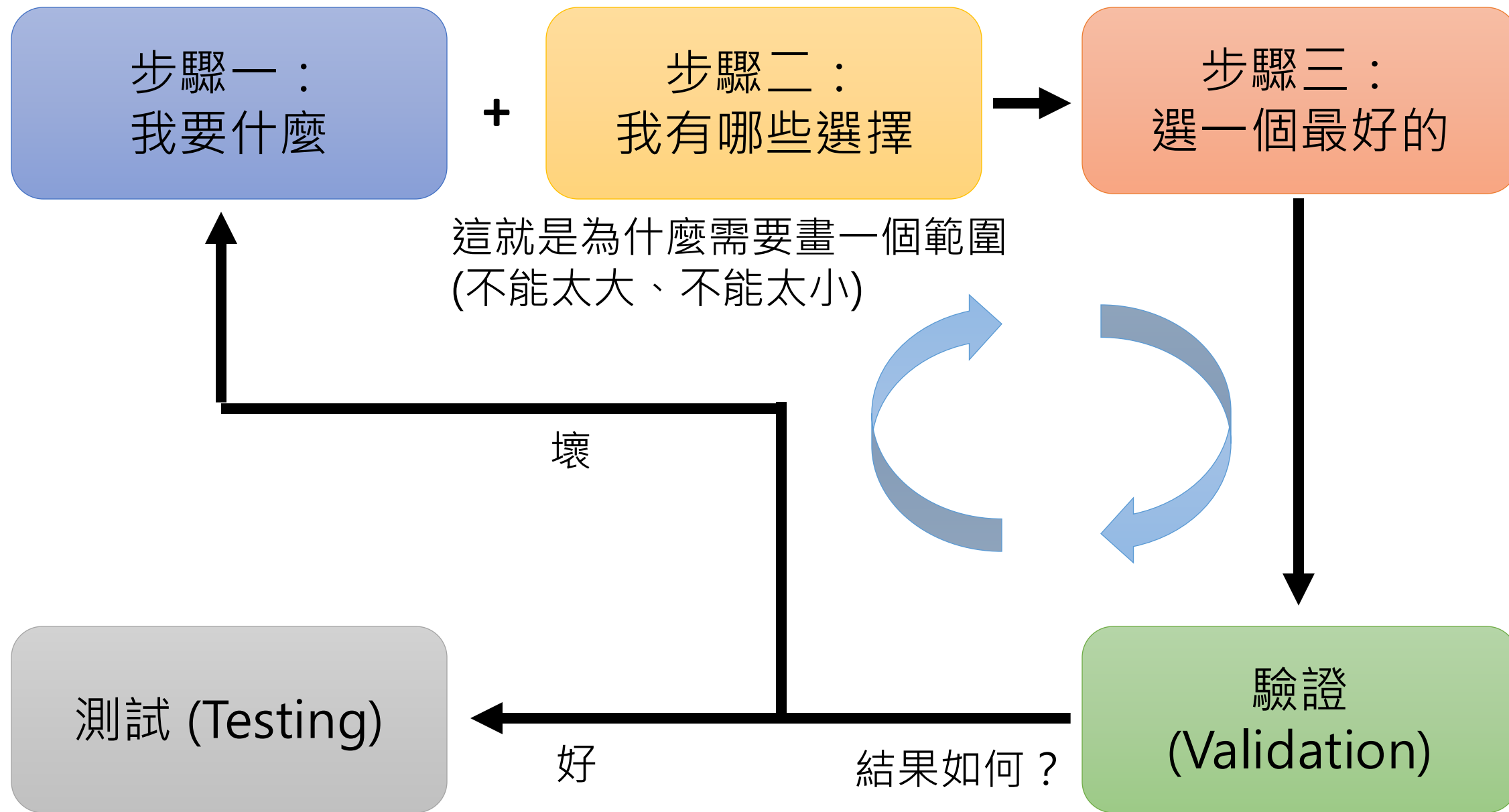
- 選擇越多，訓練和驗證的差距越大

看著路開車

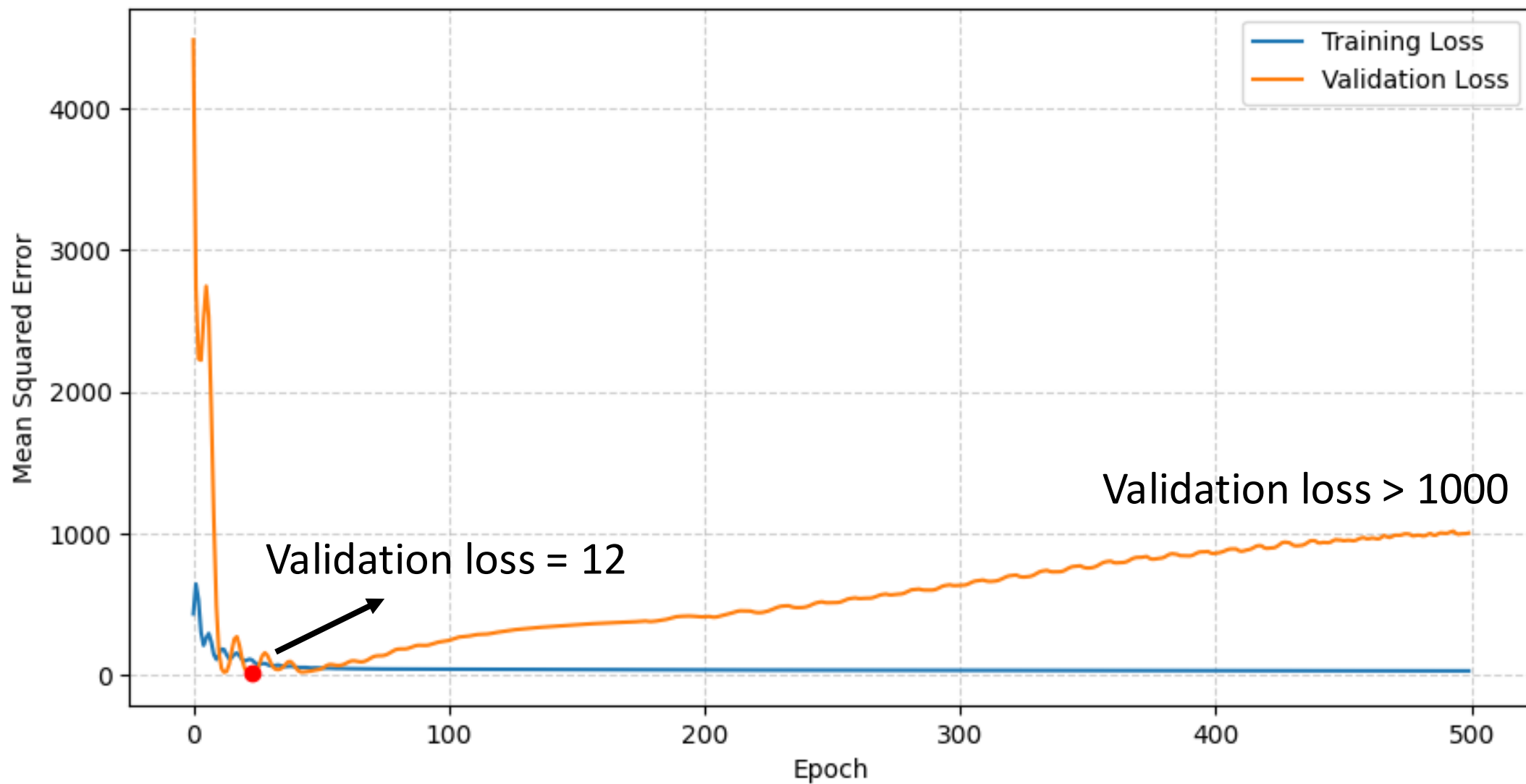
看著貼紙開車

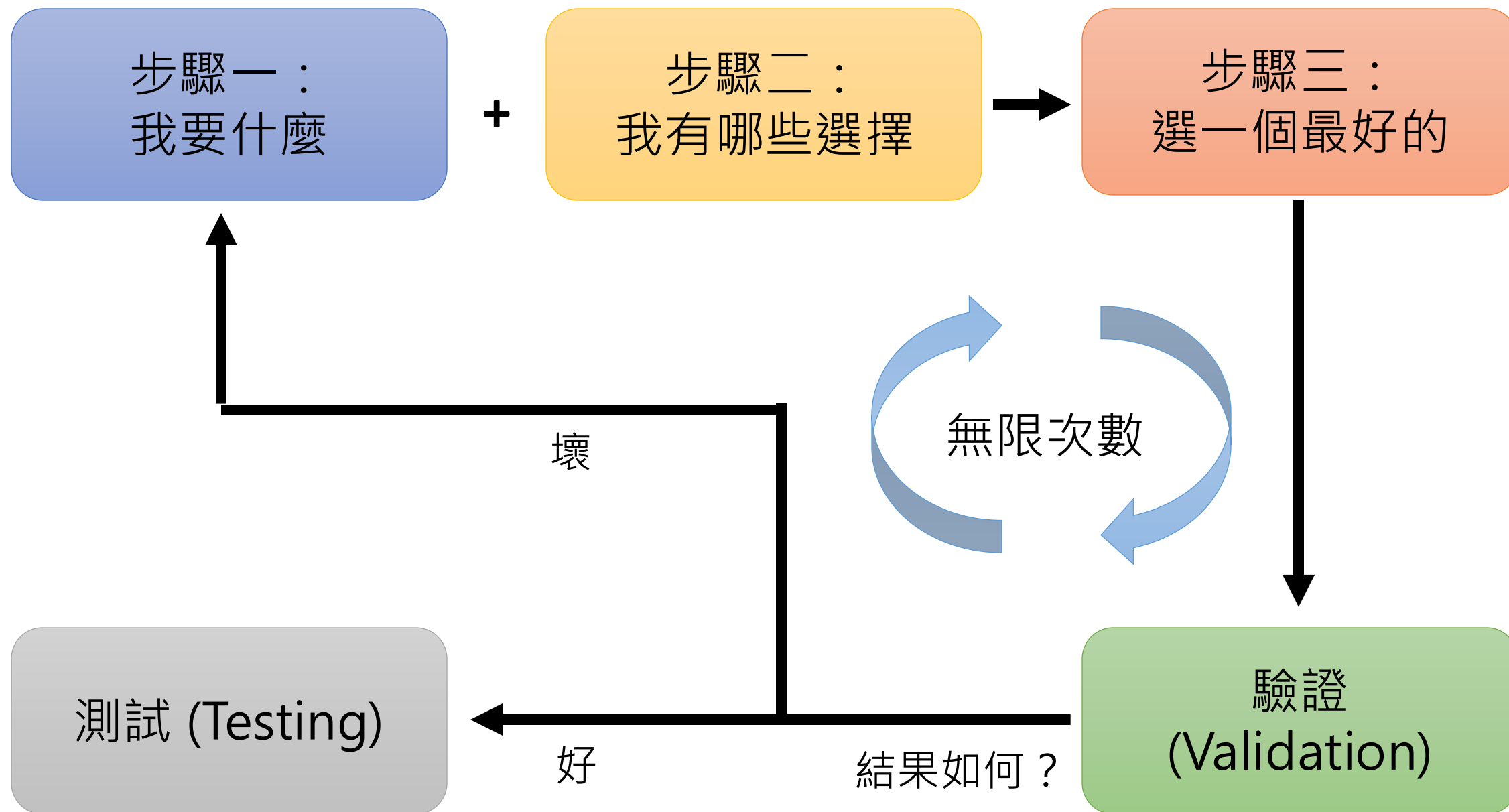
只在駕訓班才能開車





每一個 Epoch 結束都去量 Validation Loss





如果可以無限的使用驗證資料

訓練資料

$$f_{lazy2}(\text{PPT}) = 10$$

$$f_{lazy2}(\text{PPT}) = 20$$

$$f_{lazy2}(\text{PPT}) = 30$$

跟訓練資料
一樣

驗證資料

$$f_{lazy2}(\text{PPT}) = 15$$

$$f_{lazy2}(\text{PPT}) = 32$$

$$f_{lazy2}(\text{PPT}) = 33$$

跟驗證資料
一樣

測試資料

$$f_{lazy2}(\text{PPT}) = 0$$

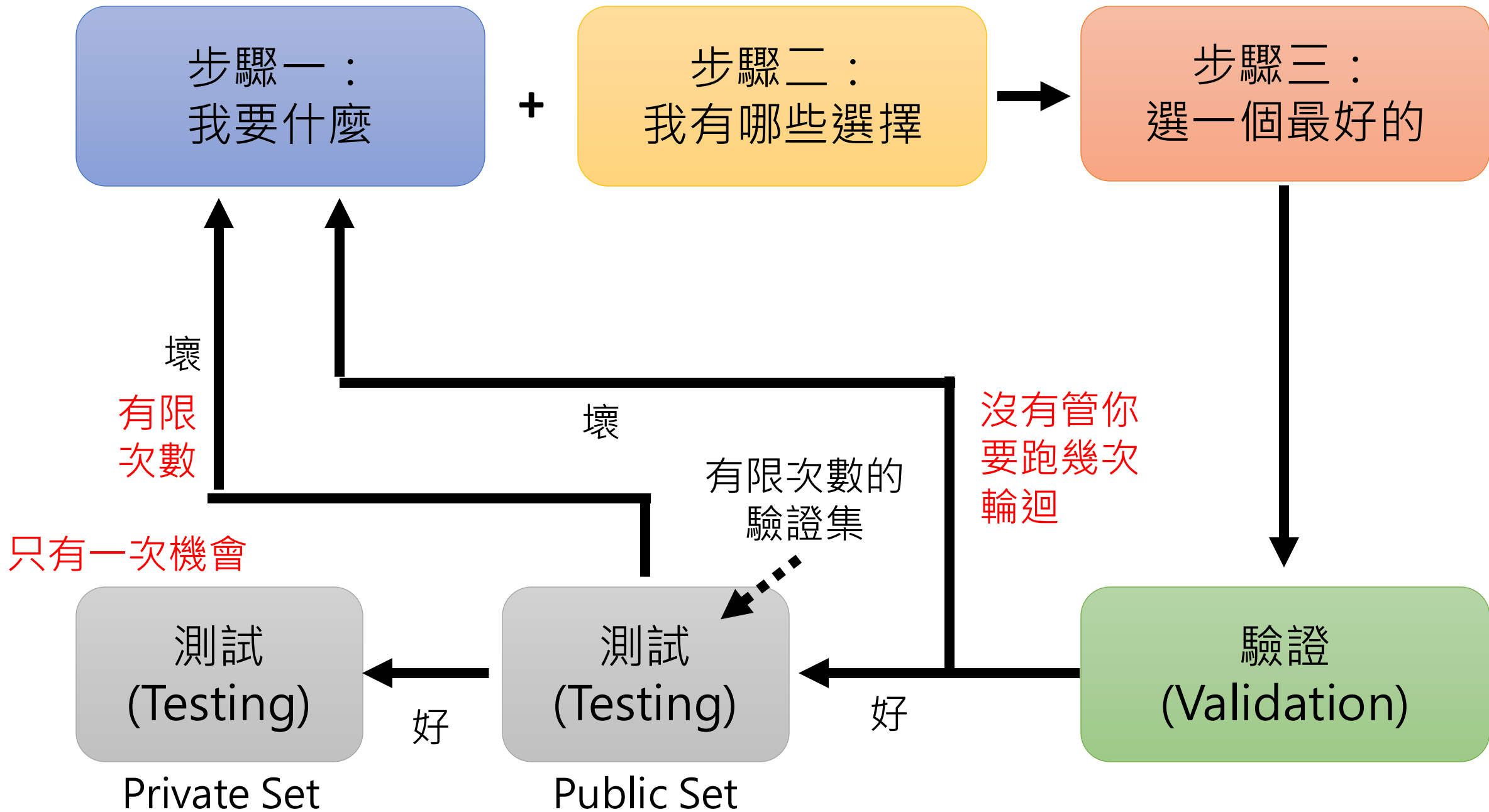
$$f_{lazy2}(\text{PPT}) = 0$$

$$f_{lazy2}(\text{PPT}) = 0$$

亂給答案

這就是為什麼人工智慧
常常在 Benchmark 上
打敗人類

Rank	Model	EM	F1
	Human Performance <i>Stanford University</i> (Rajpurkar & Jia et al. '18)	86.831	89.452
1 Mar 20, 2019	BERT + DAE + AoA (ensemble) <i>Joint Laboratory of HIT and iFLYTEK Research</i>	87.147	89.474
2 Mar 15, 2019	BERT + ConvLSTM + MTL + Verifier (ensemble) <i>Layer 6 AI</i>	86.730	89.286
3 Mar 05, 2019	BERT + N-Gram Masking + Synthetic Self- Training (ensemble) <i>Google AI Language</i> https://github.com/google-research/bert	86.673	89.147
4 May 21, 2019	XLNet (single model) <i>XLNet Team</i>	86.346	89.133
5 Apr 13, 2019	SemBERT(ensemble) <i>Shanghai Jiao Tong University</i>	86.166	88.886



課程規劃

原理

實作

範例程式

連結：

<https://colab.research.google.com/drive/1SFtkeDL9jp5LtaVsj-2JApSGpOltLi9?usp=sharing>

