# Regression

## Hung-yi Lee
## 李宏毅

# Regression: Output a scalar

Regression 問題的「函式範圍」屬於 Linear！也就是找到一個 Linear Function 來解決 Regression 問題

- Stock Market Forecast

$$f(\quad\text{<image>}\quad) = \text{Dow Jones Industrial Average at tomorrow}$$

- Self-driving Car

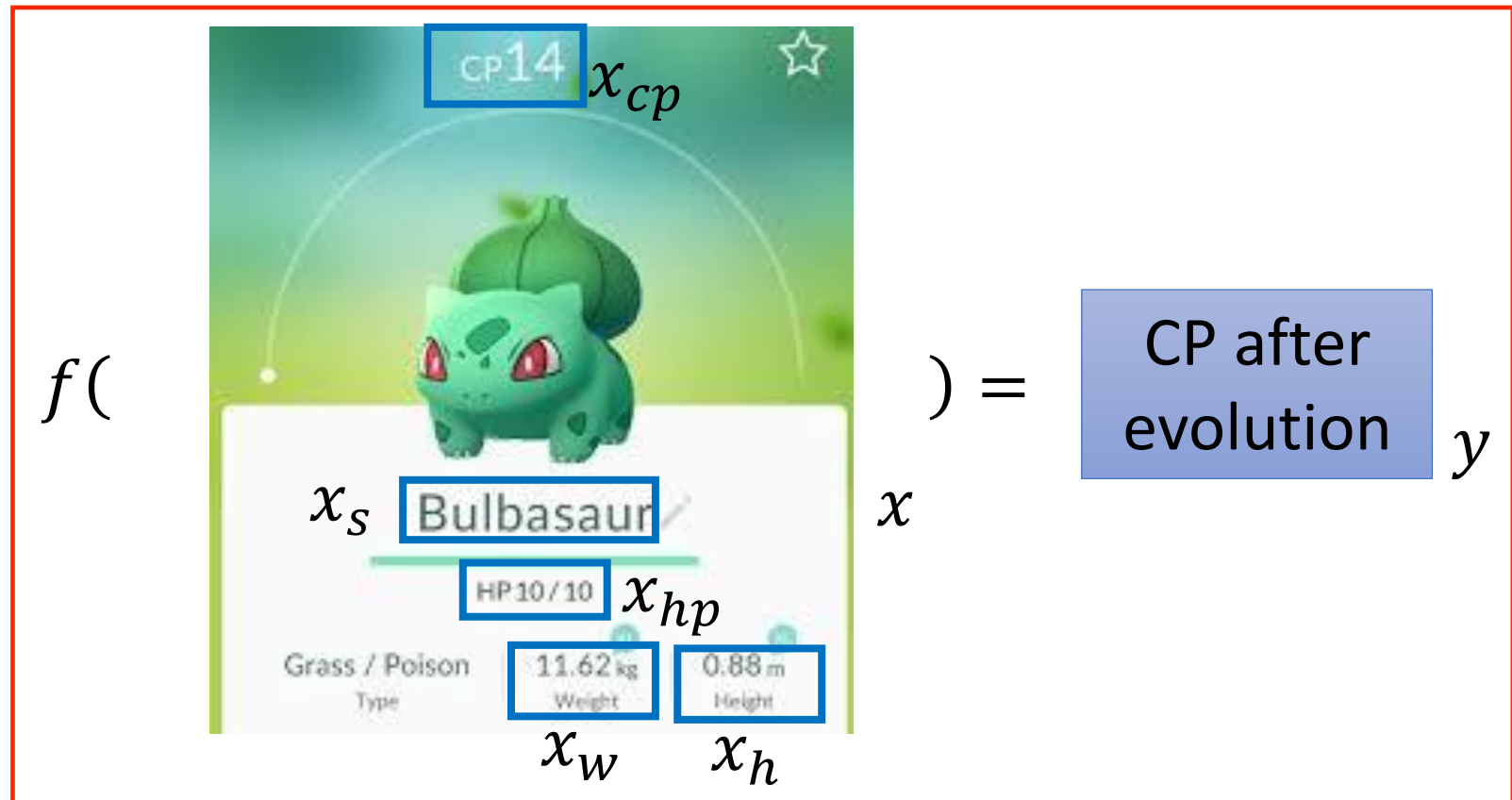$$f(\quad\text{<image>}\quad) = \text{方向盤角度}$$

- Recommendation

$$f(\quad 使用者 \text{ A} \quad 商品 \text{ B} \quad) = 購買可能性$$

# Example Application

- Estimating the Combat Power (CP) of a pokemon after evolution 輸入一隻寶可夢 => 輸出進化後的 CP

$f($  $x_{cp}$ $x_s$ $x_{hp}$ $x_w$ $x_h$ $x$ $) =$ CP after evolution $y$

# Step 1: Model

$y = b + w \cdot x_{cp}$

A set of function

Model

$f_1, f_2 \cdots$

w and b are parameters
(can be any value)

$f_1: y = 10.0 + 9.0 \cdot x_{cp}$

$f_2: y = 9.8 + 9.2 \cdot x_{cp}$

$f_3: y = -0.8 - 1.2 \cdot x_{cp}$

…… infinite

$f($ 
CP14
☆
Bulbasaur
HP 10/10
Grass / Poison    11.62 kg    0.88 m
Type              Weight      Height
$x ) =$ CP after evolution $y$

Linear model: $y = b + \sum w_i x_i$

$x_i : x_{cp}, x_{hp}, x_w, x_h \ldots$

feature

$w_i$: weight, b: bias

# Step 2: Goodness of Function

$y = b + w \cdot x_{cp}$

A set of function

Model

$f_1, f_2 \cdots$

Training Data

function input:

function Output (scalar):

$x^1$

CP612

CP979 $\hat{y}^1$

$x^2$

CP706

CP1420 $\hat{y}^2$

1Eev1ee1

HP 81/81

sparky

HP 94/94

# Step 2: Goodness of Function

Training Data:
10 pokemons

$(x^1, \hat{y}^1)$

$(x^2, \hat{y}^2)$

$\vdots$

$(x^{10}, \hat{y}^{10})$

This is real data.

$(x_{cp}^n, \hat{y}^n)$

CP after evoluation

$\hat{y}$

Original CP

$x_{cp}$

Source: https://www.openintro.org/stat/data/?data=pokemon

# Step 2: Goodness of Function

$y = b + w \cdot x_{cp}$

A set of function

Model

$f_1, f_2 \cdots$

Goodness of function f

Training Data

Loss function $L$:

Input: a function, output: how bad it is

Estimation error

$$L(f) = \sum_{n=1}^{10} \left( \hat{y}^n - f(x_{cp}^n) \right)^2$$

Sum over examples

Estimated y based on input function

$$L(w, b) = \sum_{n=1}^{10} \left( \hat{y}^n - (b + w \cdot x_{cp}^n) \right)^2$$
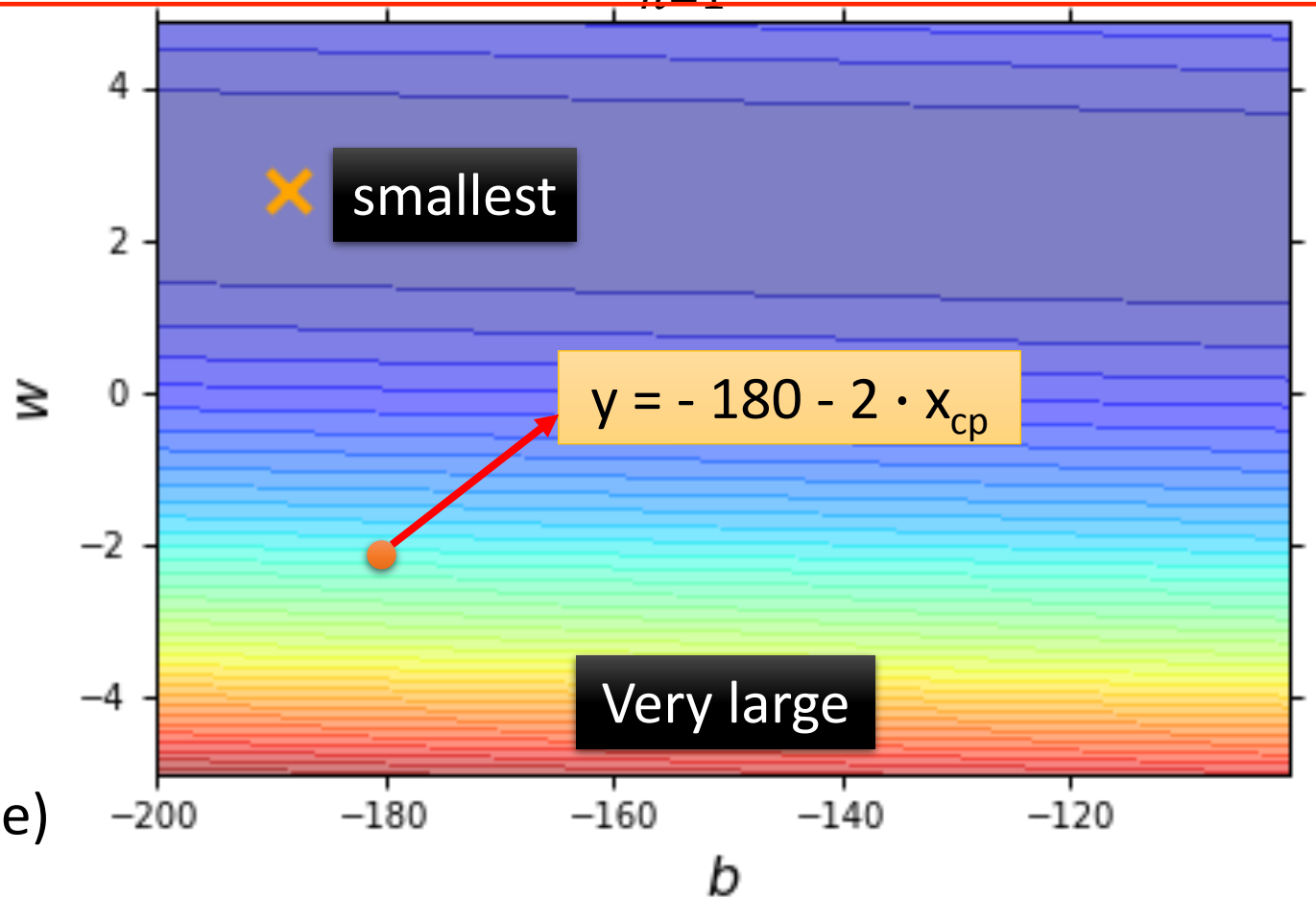
# Step 2: Goodness of Function

- Loss Function

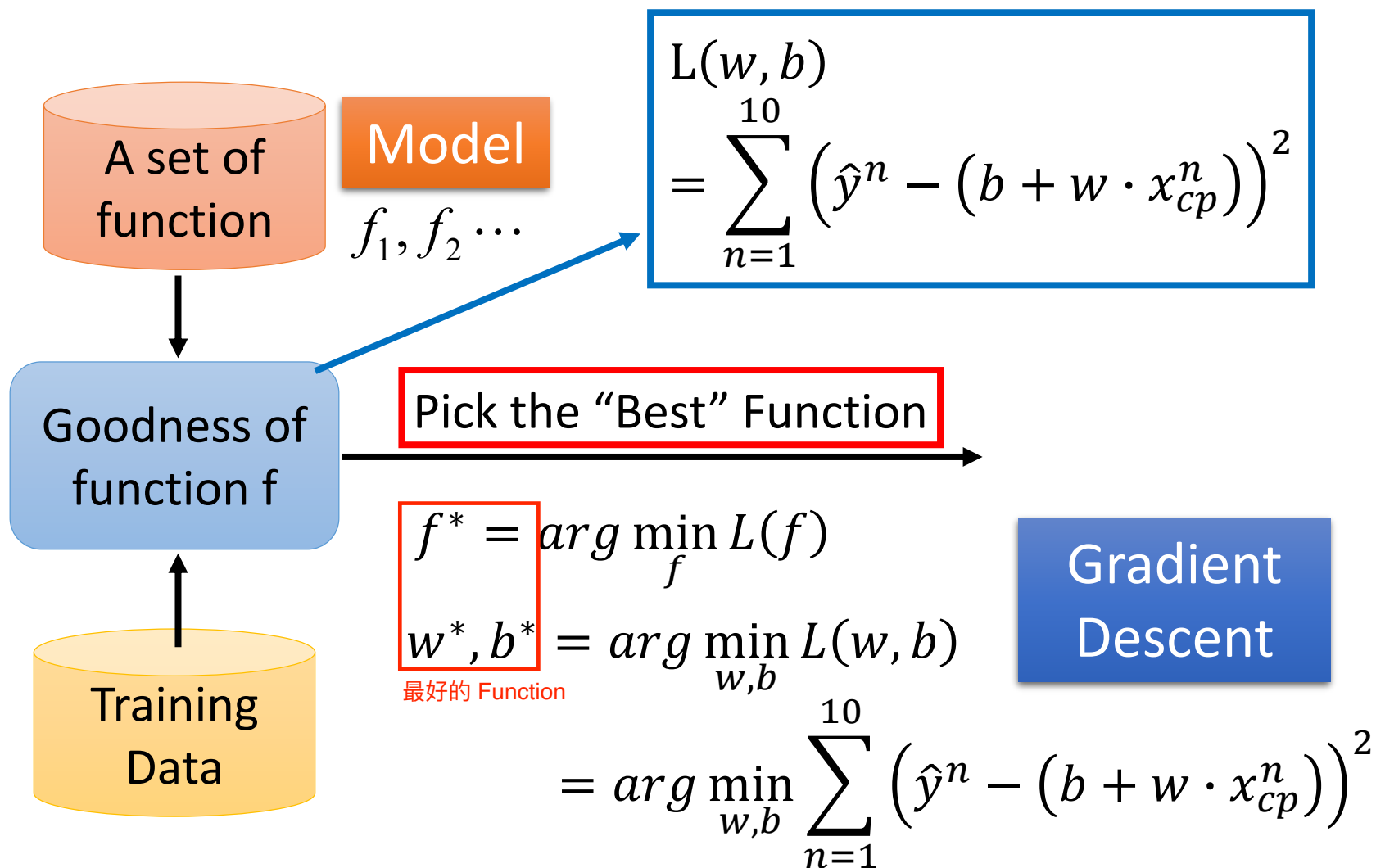$$L(w,b) = \sum_{n=1}^{10} \left( \hat{y}^n - (b + w \cdot x_{cp}^n) \right)^2$$

Each point in the figure is a function

The color represents $L(w,b)$.

(true example)



smallest

$y = -180 - 2 \cdot x_{cp}$

Very large

# Step 3: Best Function

**A set of function**

**Model**

$f_1, f_2 \cdots$

$$\text{L}(w, b) = \sum_{n=1}^{10} \left( \hat{y}^n - (b + w \cdot x_{cp}^n) \right)^2$$

**Goodness of function f**

Pick the "Best" Function

**Training Data**

$$f^* = arg \min_f L(f)$$

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

最好的 Function

$$= arg \min_{w,b} \sum_{n=1}^{10} \left( \hat{y}^n - (b + w \cdot x_{cp}^n) \right)^2$$

**Gradient Descent**

從前頁可知：要利用 Loss Function 從 Function Set 中找到最好的 Function => Loss(F) = Loss(w, b) 最小！要如何改變 w 與 b 才能得到最小的 Loss ？ => 透過 Gradient Descent

# Step 3: Gradient Descent

$$w^* = arg \min_w L(w)$$

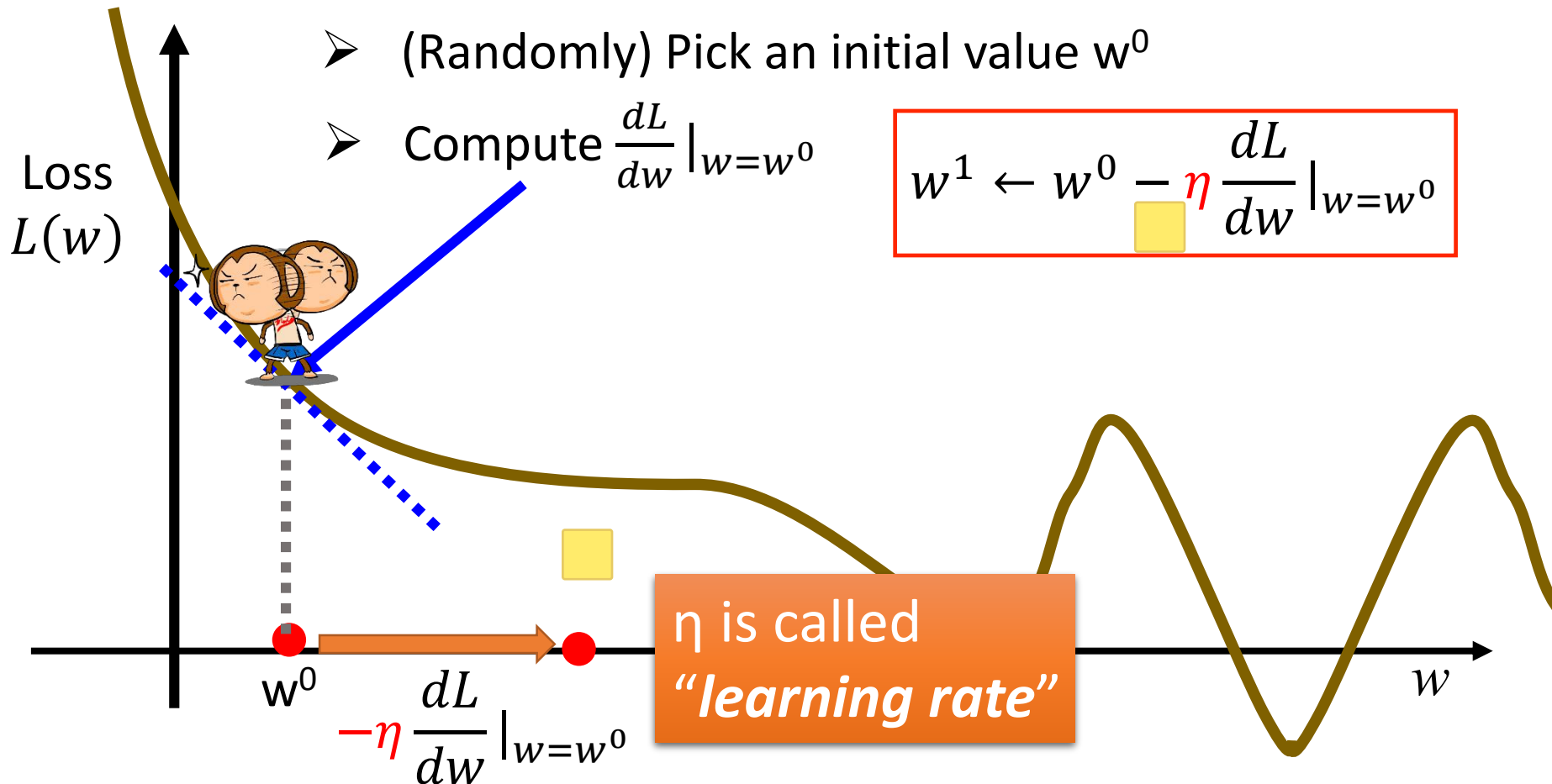- Consider loss function $L(w)$ with one parameter w:

第一種

窮舉所有的 w 找到最小的 L(w)

第二種

Loss
$L(w)$

➤ (Randomly) Pick an initial value $w^0$

➤ Compute $\frac{dL}{dw}\big|_{w=w^0}$

| Negative | ➡ | Increase w |
| Positive | ➡ | Decrease w |

$w^0$

$w$

# Step 3: Gradient Descent

$$w^* = arg \min_w L(w)$$

- Consider loss function $L(w)$ with one parameter w:
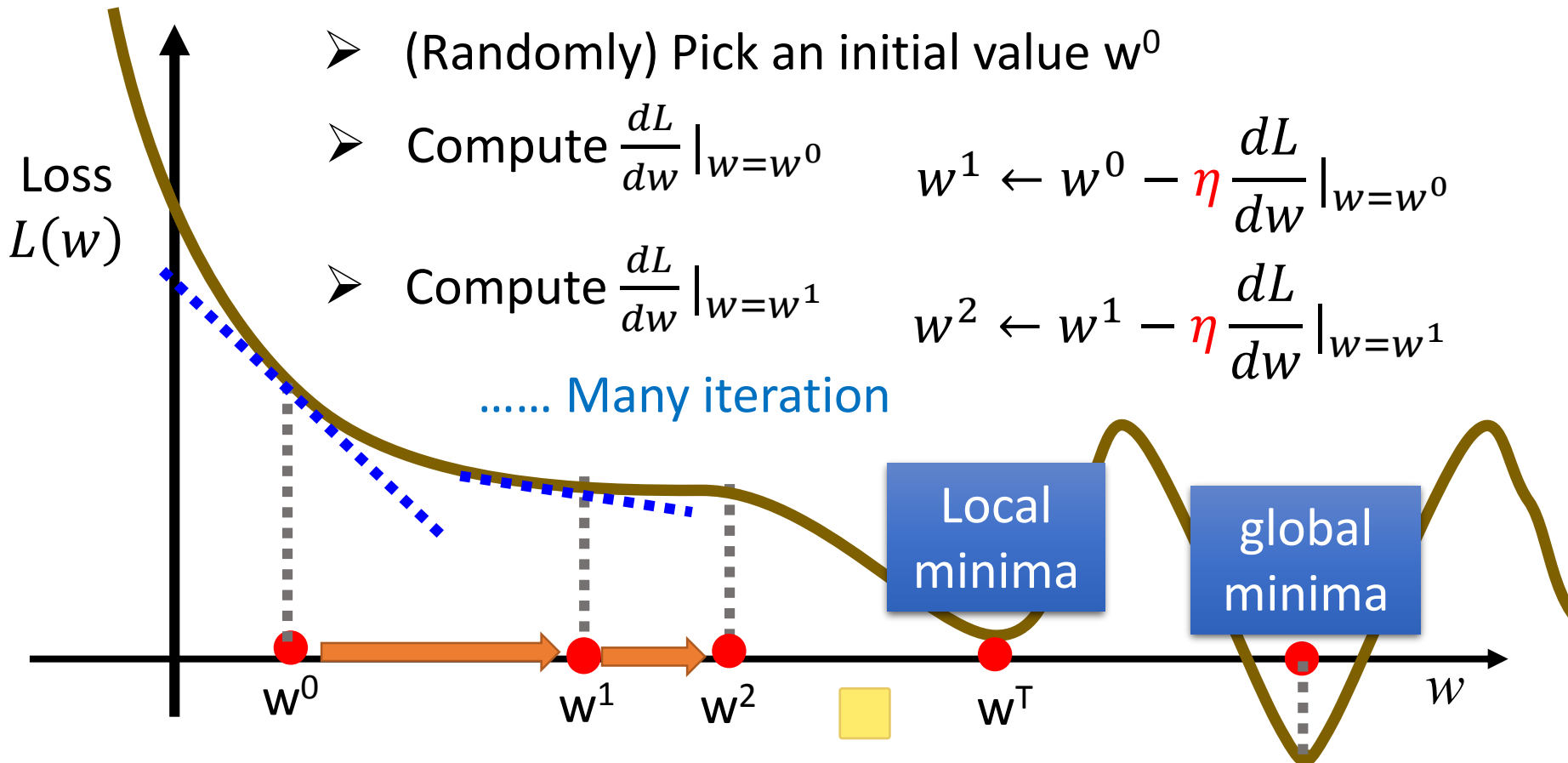
➢ (Randomly) Pick an initial value $w^0$

➢ Compute $\frac{dL}{dw}|_{w=w^0}$

$$w^1 \leftarrow w^0 - \eta \frac{dL}{dw}|_{w=w^0}$$

Loss $L(w)$

$w^0$

$-\eta \frac{dL}{dw}|_{w=w^0}$

$w$

η is called **"learning rate"**

# Step 3: Gradient Descent

$$w^* = arg \min_w L(w)$$

- Consider loss function $L(w)$ with one parameter w:

  ➢ (Randomly) Pick an initial value $w^0$

  ➢ Compute $\frac{dL}{dw}|_{w=w^0}$     $w^1 \leftarrow w^0 - \eta \frac{dL}{dw}|_{w=w^0}$

  ➢ Compute $\frac{dL}{dw}|_{w=w^1}$     $w^2 \leftarrow w^1 - \eta \frac{dL}{dw}|_{w=w^1}$

  …… Many iteration

Loss $L(w)$

Local minima

global minima

$w^0$   $w^1$   $w^2$   $w^T$   $w$

# Step 3: Gradient Descent

$$\begin{bmatrix} \dfrac{\partial L}{\partial w} \\ \dfrac{\partial L}{\partial b} \end{bmatrix} \text{gradient}$$

- How about two parameters? $w^*, b^* = arg \min_{w,b} L(w,b)$

  ➢ (Randomly) Pick an initial value $w^0, b^0$

  ➢ Compute $\dfrac{\partial L}{\partial w}\big|_{w=w^0,b=b^0}, \dfrac{\partial L}{\partial b}\big|_{w=w^0,b=b^0}$

$$w^1 \leftarrow w^0 - \eta \dfrac{\partial L}{\partial w}\big|_{w=w^0,b=b^0} \qquad b^1 \leftarrow b^0 - \eta \dfrac{\partial L}{\partial b}\big|_{w=w^0,b=b^0}$$

  ➢ Compute $\dfrac{\partial L}{\partial w}\big|_{w=w^1,b=b^1}, \dfrac{\partial L}{\partial b}\big|_{w=w^1,b=b^1}$

$$w^2 \leftarrow w^1 - \eta \dfrac{\partial L}{\partial w}\big|_{w=w^1,b=b^1} \qquad b^2 \leftarrow b^1 - \eta \dfrac{\partial L}{\partial b}\big|_{w=w^1,b=b^1}$$

# Step 3: Gradient Descent

將「兩個參數」的 Loss Function 的 Gradient Descent 視覺化！



Color: Value of Loss L(w,b)

$(-\eta \ \partial L / \partial b, \ -\eta \ \partial L / \partial w)$

Compute $\partial L / \partial b, \partial L / \partial w$

參數更新的大小

計算切線斜率

# Step 3: Gradient Descent

- Formulation of $\partial L/\partial w$ and $\partial L/\partial b$

$$L(w,b) = \sum_{n=1}^{10} \left(\hat{y}^n - \left(b + w \cdot x_{cp}^n\right)\right)^2$$

$$\frac{\partial L}{\partial w} = ? \sum_{n=1}^{10} 2\left(\hat{y}^n - \left(b + w \cdot x_{cp}^n\right)\right)\left(-x_{cp}^n\right)$$

$$\frac{\partial L}{\partial b} = ? \sum_{n=1}^{10} 2\left(\hat{y}^n - \left(b + w \cdot x_{cp}^n\right)\right)$$

# How's the results?

$$y = b + w \cdot x_{cp}$$

b = -188.4

w = 2.7

Average Error on Training Data

$$= \frac{1}{10} \sum_{n=1}^{10} e^n = 31.9$$

Training Data

$e^1$

$e^2$

CP after evoluation

Original CP

# How's the results? - Generalization

What we really care about is the error on new data (testing data)

$y = b + w \cdot x_{cp}$

b = -188.4

w = 2.7

Average Error on Testing Data

$$= \frac{1}{10} \sum_{n=1}^{10} e^n = 35.0$$

> Average Error on Training Data (31.9)

Another 10 pokemons as testing data

How can we do better?

## Selecting another Model

$$y = b + w_1 \cdot x_{cp} + w_2 \cdot (x_{cp})^2$$

重新定義一個 Function Set

## Best Function

b = -10.3

$w_1 = 1.0$, $w_2 = 2.7 \times 10^{-3}$

Average Error = 15.4

## Testing:

Average Error = 18.4

Better! Could it be even better?

## Selecting another Model

$$y = b + w_1 \cdot x_{cp} + w_2 \cdot (x_{cp})^2 + w_3 \cdot (x_{cp})^3$$

重新定義一個 Function Set => 項數更多

## Best Function

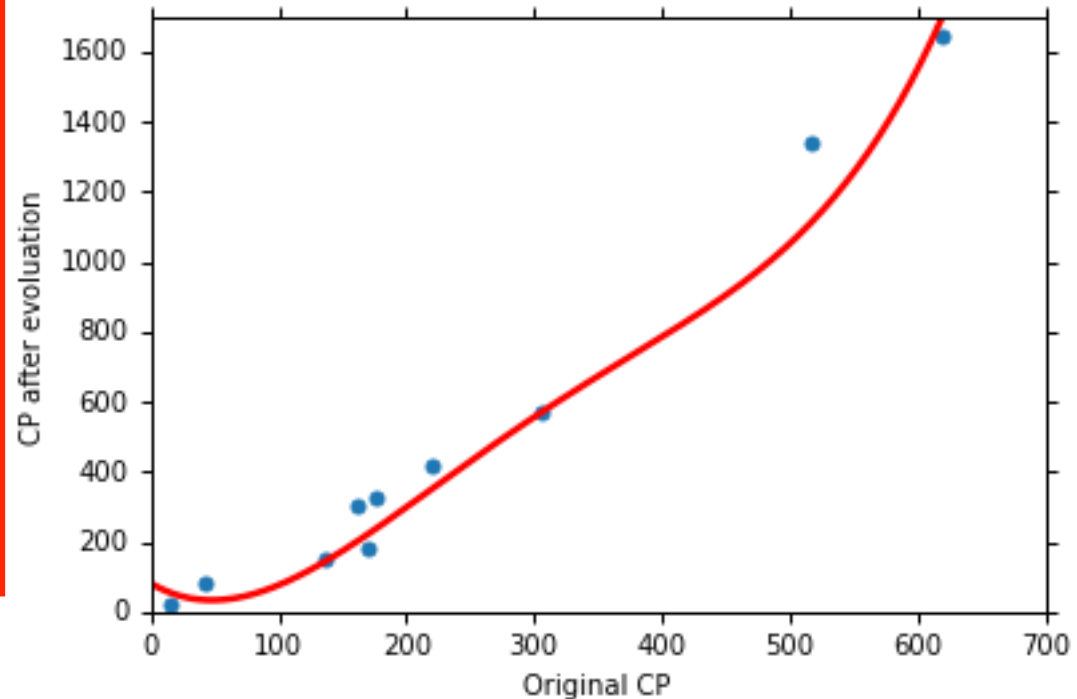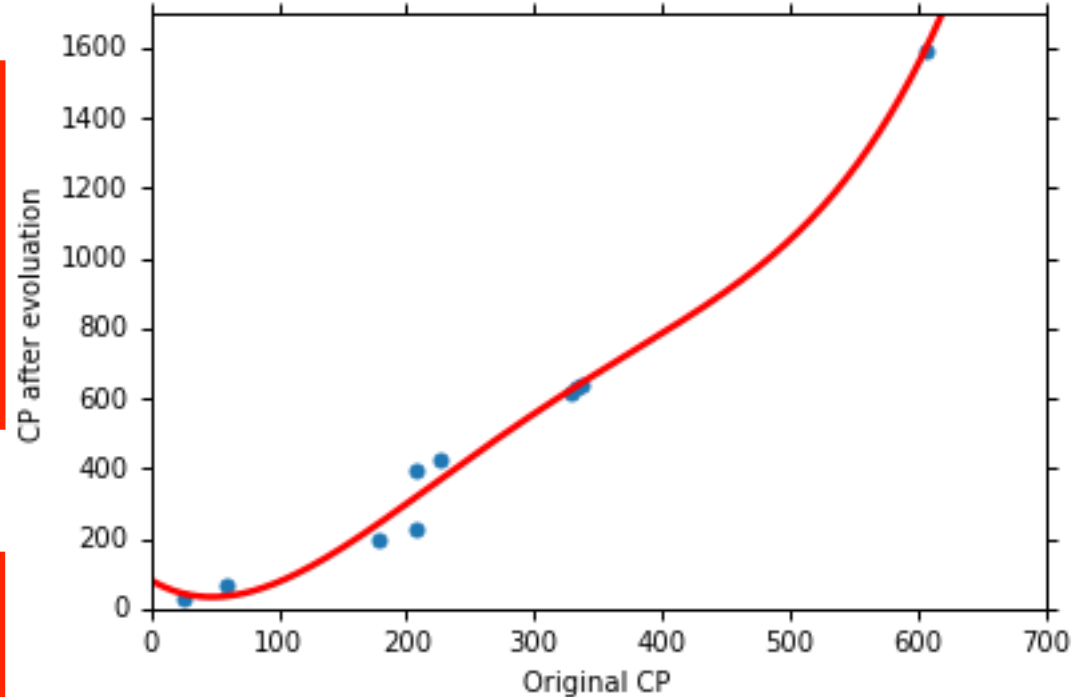b = 6.4, $w_1$ = 0.66

$w_2$ = 4.3 x $10^{-3}$

$w_3$ = -1.8 x $10^{-6}$

Average Error = 15.3

## Testing:

Average Error = 18.1

Slightly better. How about more complex model?

## Selecting another Model

$$y = b + w_1 \cdot x_{cp} + w_2 \cdot (x_{cp})^2 + w_3 \cdot (x_{cp})^3 + w_4 \cdot (x_{cp})^4$$

重新定義一個 Function Set => 項數更多

## Best Function

Average Error = 14.9

## Testing:

Average Error = 28.8

The results become worse ...

## *Selecting another Model*

$$y = b + w_1 \cdot x_{cp} + w_2 \cdot (x_{cp})^2 + w_3 \cdot (x_{cp})^3 + w_4 \cdot (x_{cp})^4 + w_5 \cdot (x_{cp})^5$$
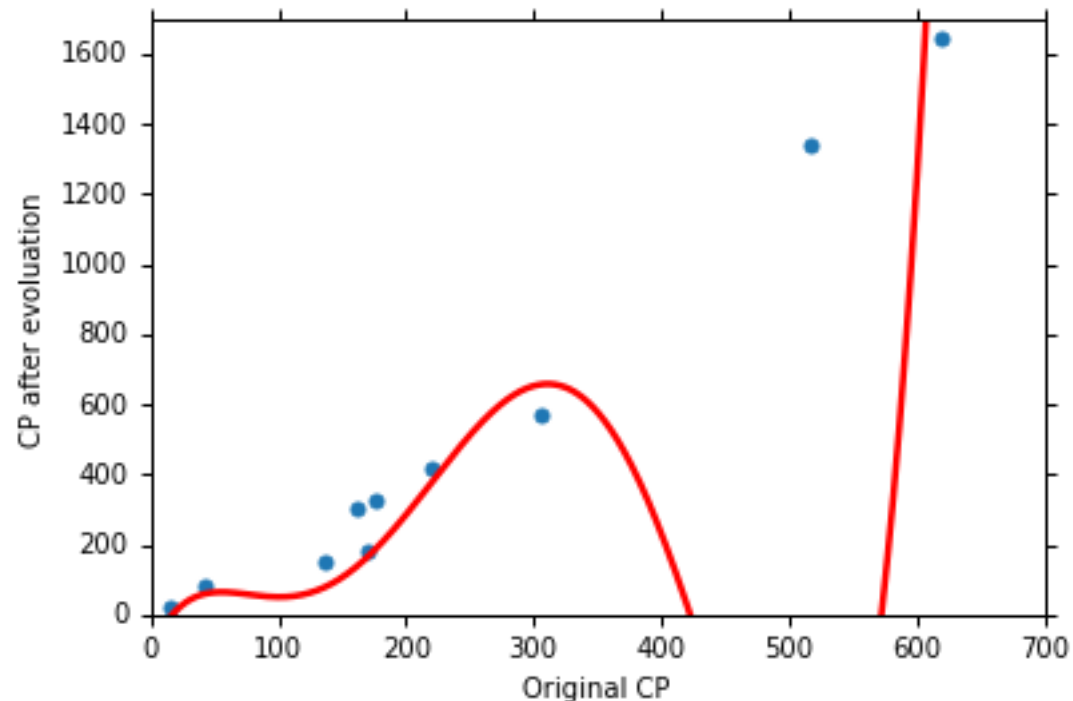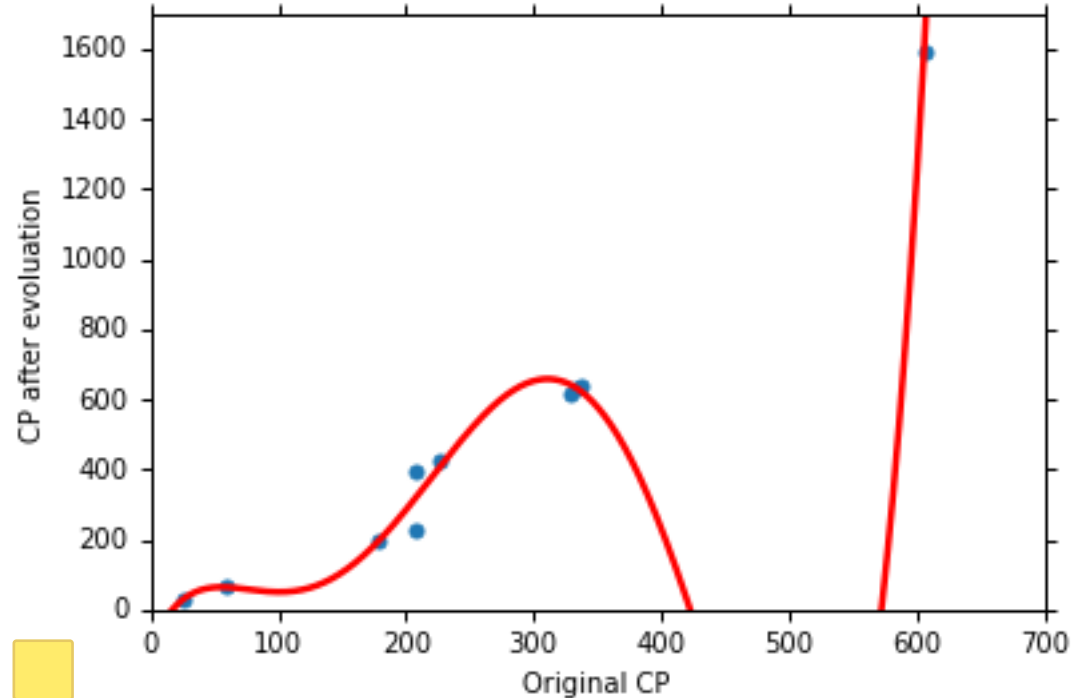
重新定義一個 Function Set => 項數更多

## *Best Function*

Average Error = 12.8

## *Testing:*

Average Error = 232.1

The results are so bad.

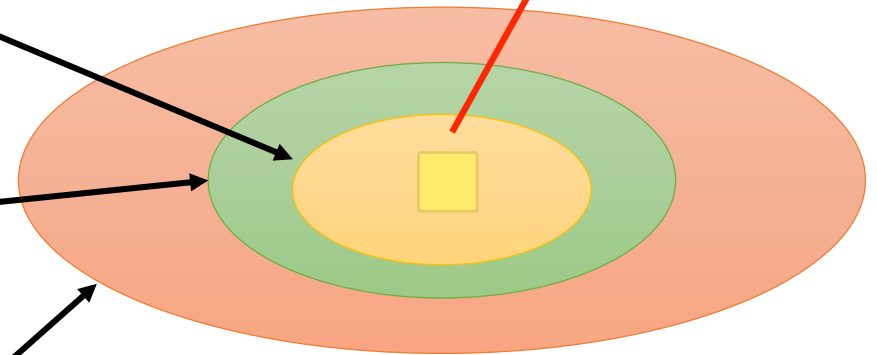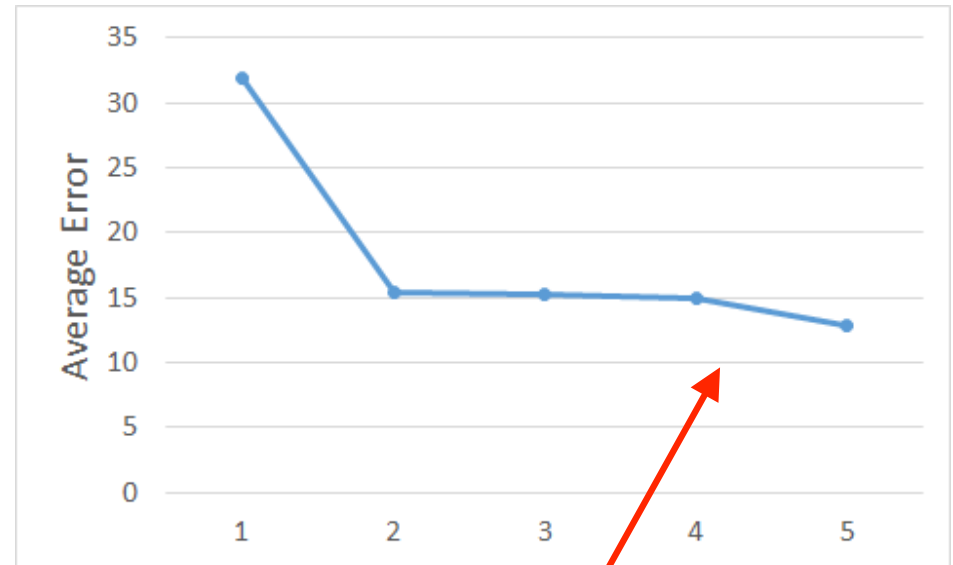# Model Selection

1. $y = b + w \cdot x_{cp}$

2. $y = b + w_1 \cdot x_{cp} + w_2 \cdot (x_{cp})^2$

3. $y = b + w_1 \cdot x_{cp} + w_2 \cdot (x_{cp})^2 + w_3 \cdot (x_{cp})^3$

4. $y = b + w_1 \cdot x_{cp} + w_2 \cdot (x_{cp})^2 + w_3 \cdot (x_{cp})^3 + w_4 \cdot (x_{cp})^4$

5. $y = b + w_1 \cdot x_{cp} + w_2 \cdot (x_{cp})^2 + w_3 \cdot (x_{cp})^3 + w_4 \cdot (x_{cp})^4 + w_5 \cdot (x_{cp})^5$
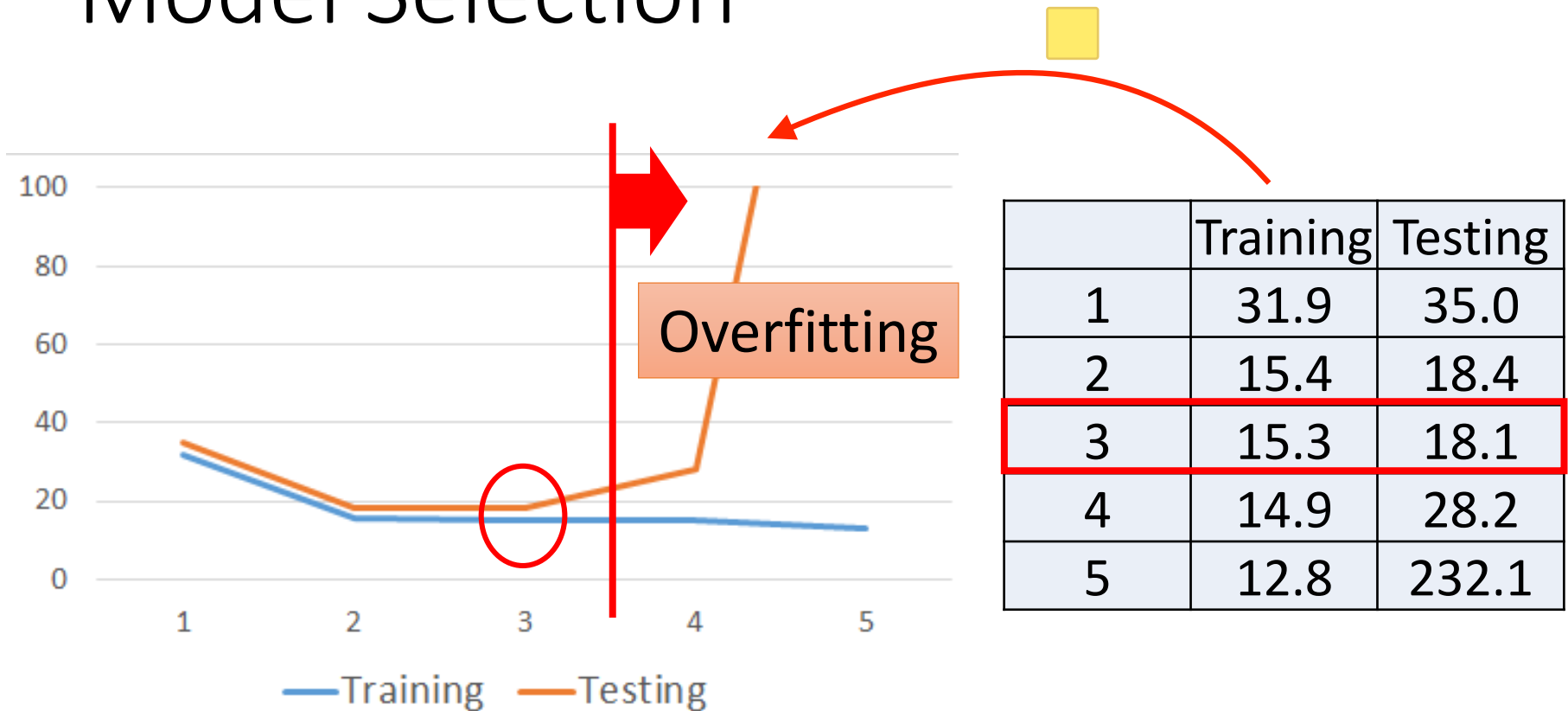
A more complex model yields lower error on training data.

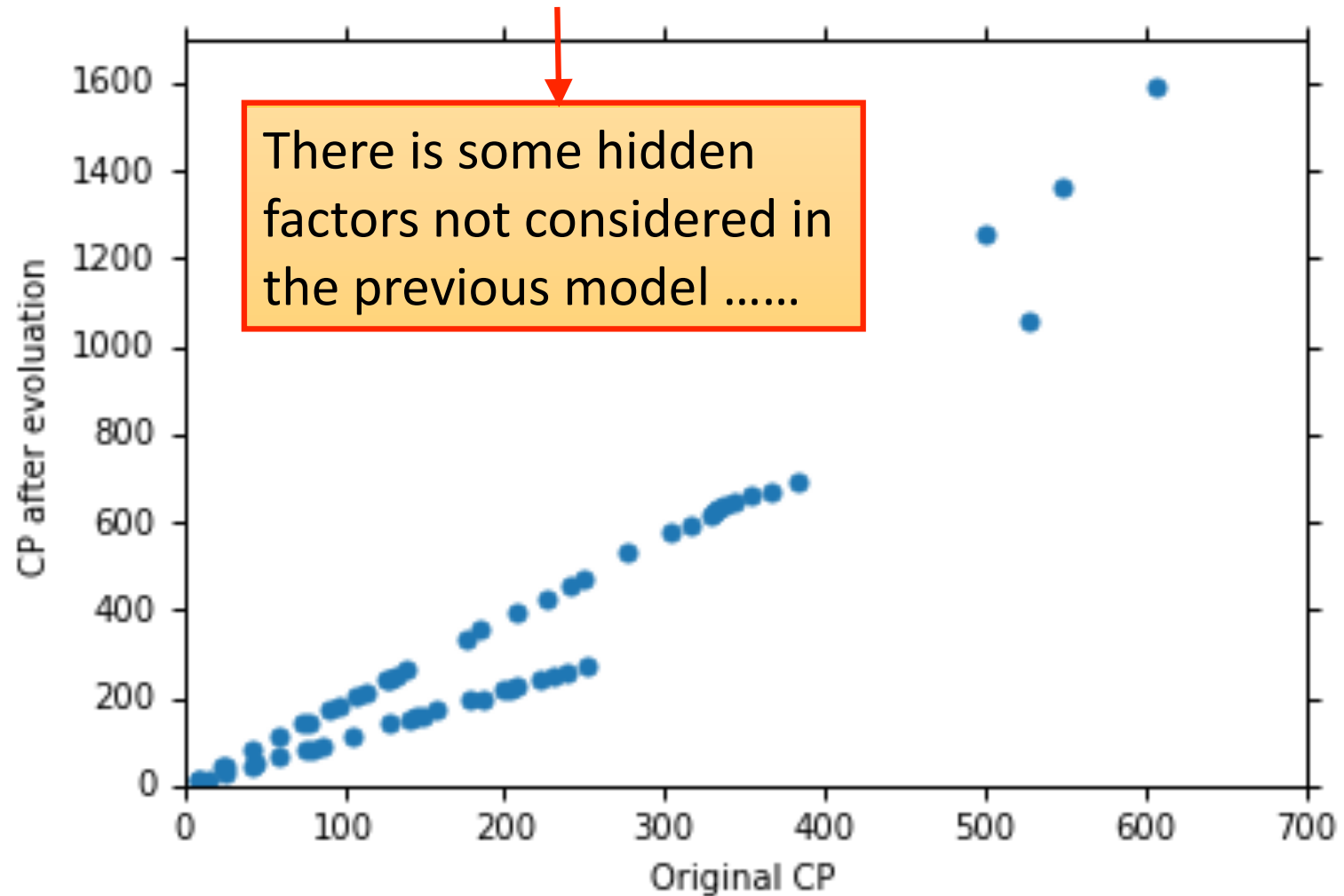If we can truly find the best function

# Model Selection



| | Training | Testing |
|---|---|---|
| 1 | 31.9 | 35.0 |
| 2 | 15.4 | 18.4 |
| 3 | 15.3 | 18.1 |
| 4 | 14.9 | 28.2 |
| 5 | 12.8 | 232.1 |

Overfitting

Training  Testing

A more complex model does not always lead to better performance on ***testing data***.

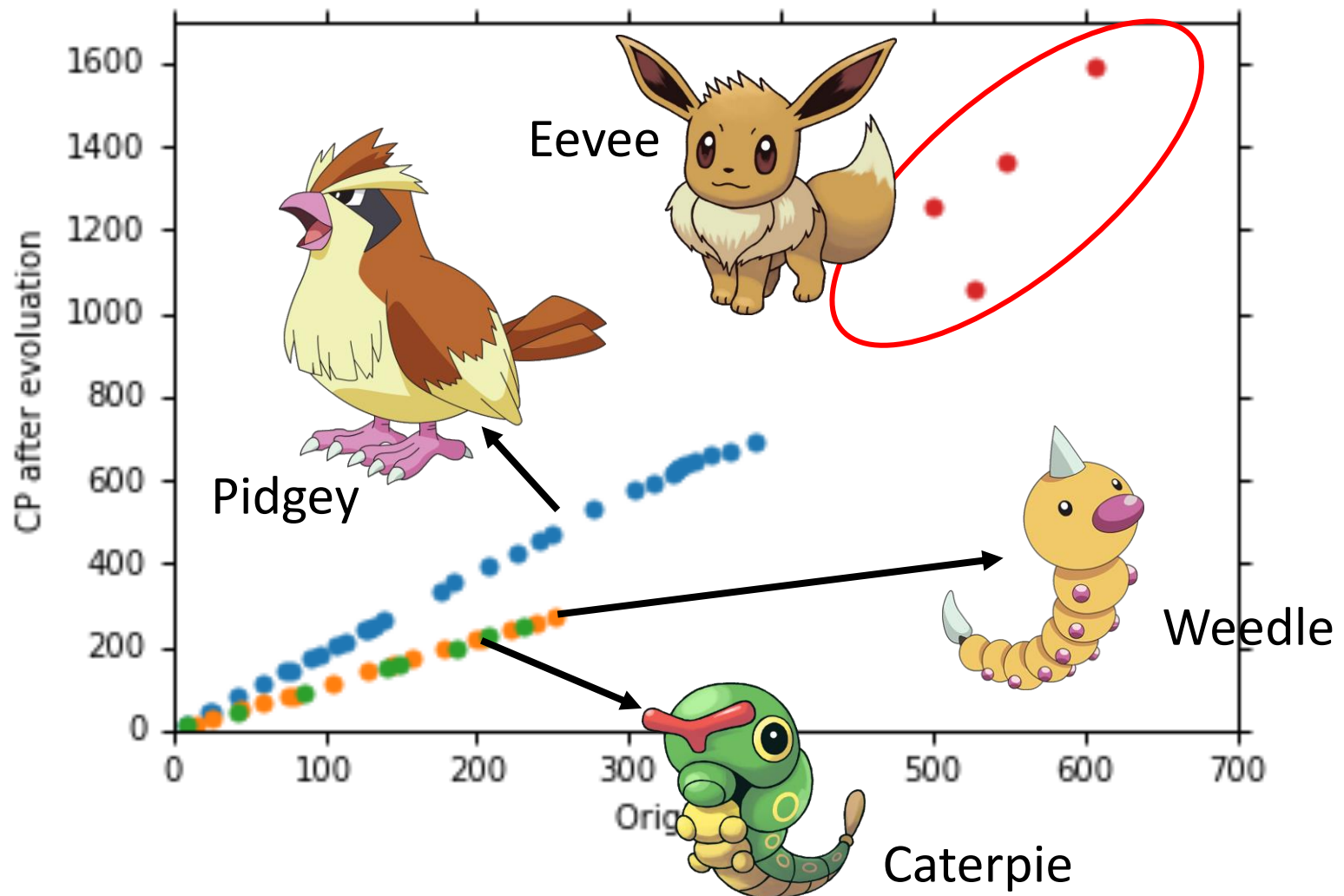This is ***Overfitting***.  ➡  Select suitable model

# Let's collect more data

收集更多 data 後發現：原來的 CP 與進化後的 CP 不是單單一條 Regression 的關係



There is some hidden factors not considered in the previous model ……

# What are the hidden factors?

# Back to step 1: Redesign the Model

$$y = b + \sum w_i x_i$$

Linear model?

$x_s$ = species of x

x

| If $x_s$ = Pidgey: | $y = b_1 + w_1 \cdot x_{cp}$ |
|---|---|
| If $x_s$ = Weedle: | $y = b_2 + w_2 \cdot x_{cp}$ |
| If $x_s$ = Caterpie: | $y = b_3 + w_3 \cdot x_{cp}$ |
| If $x_s$ = Eevee: | $y = b_4 + w_4 \cdot x_{cp}$ |

不同的物種，就會有不同的 weight 與 bias

y

# Back to step 1: Redesign the Model

$$y = b + \sum w_i x_i$$

Linear model?

$$y = b_1 \cdot \delta(x_s = \text{Pidgey})$$
$$+ w_1 \cdot \delta(x_s = \text{Pidgey})x_{cp}$$
$$+ b_2 \cdot \delta(x_s = \text{Weedle})$$
$$+ w_2 \cdot \delta(x_s = \text{Weedle})x_{cp}$$
$$+ b_3 \cdot \delta(x_s = \text{Caterpie})$$
$$+ w_3 \cdot \delta(x_s = \text{Caterpie})x_{cp}$$
$$+ b_4 \cdot \delta(x_s = \text{Eevee})$$
$$+ w_4 \cdot \delta(x_s = \text{Eevee})x_{cp}$$

$$\delta(x_s = \text{Pidgey})$$

$$\begin{cases} =1 & \text{If } x_s = \text{Pidgey} \\ =0 & \text{otherwise} \end{cases}$$

將上面的四條式子結合成一條

所以，新的 Model 仍屬於 Linear Model：「進化後 CP」同時考慮「原 CP」與「物種」

# Back to step 1: Redesign the Model

$$y = b + \sum w_i x_i$$

Linear model?

$$y = b_1 \cdot \boxed{1}$$

$$+ w_1 \cdot \boxed{1} \quad x_{cp}$$

$$+ b_2 \cdot \boxed{0}$$

$$+ w_2 \cdot \boxed{0}$$

$$+ b_3 \cdot \boxed{0}$$

$$+ w_3 \cdot \boxed{0}$$

$$+ b_4 \cdot \boxed{0}$$

$$+ w_4 \cdot \boxed{0}$$

$$\delta(x_s = \text{Pidgey})$$

$$\begin{cases} =1 & \text{If } x_s = \text{Pidgey} \\ \\ =0 & \text{otherwise} \end{cases}$$
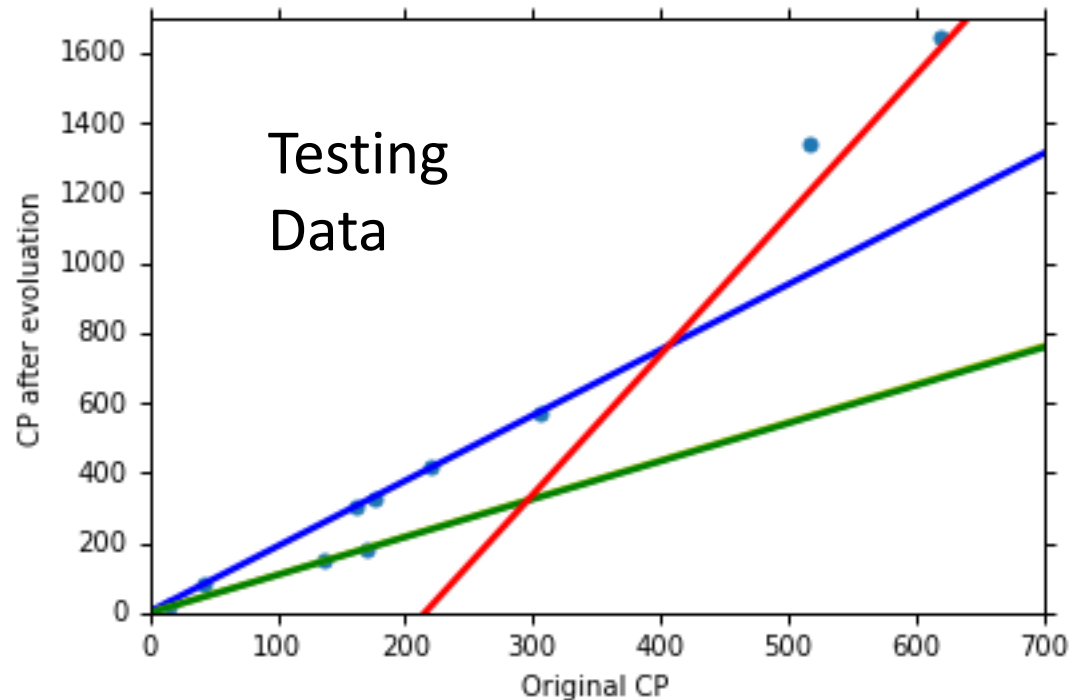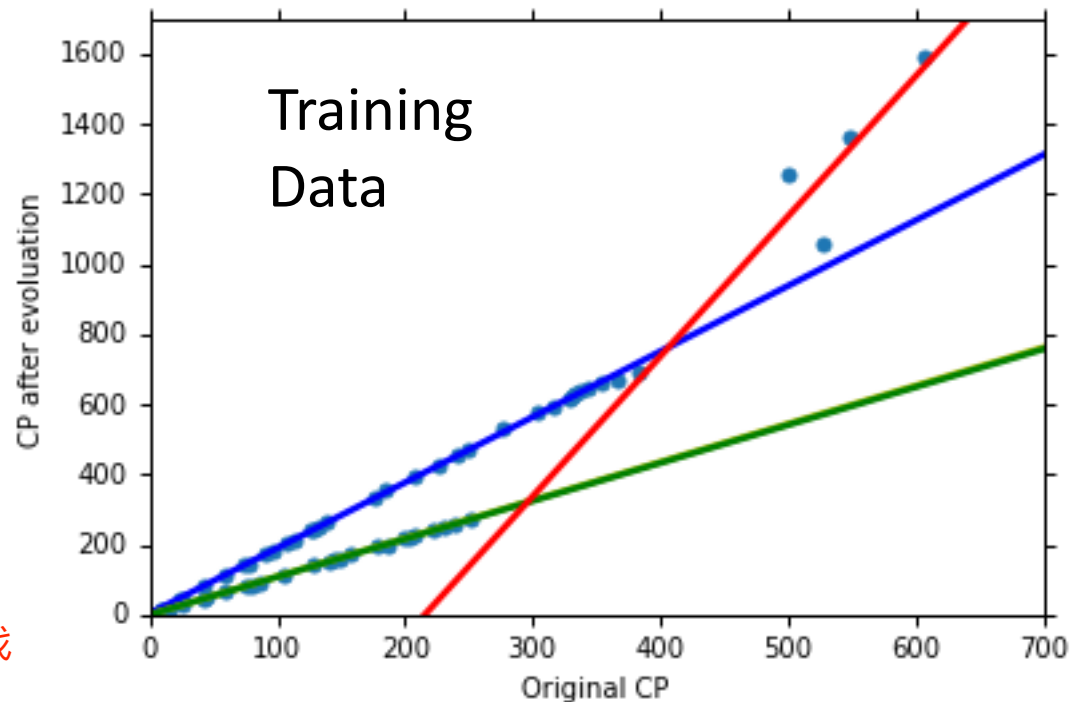
If $x_s = \text{Pidgey}$

$$y = b_1 + w_1 \cdot x_{cp}$$

Training
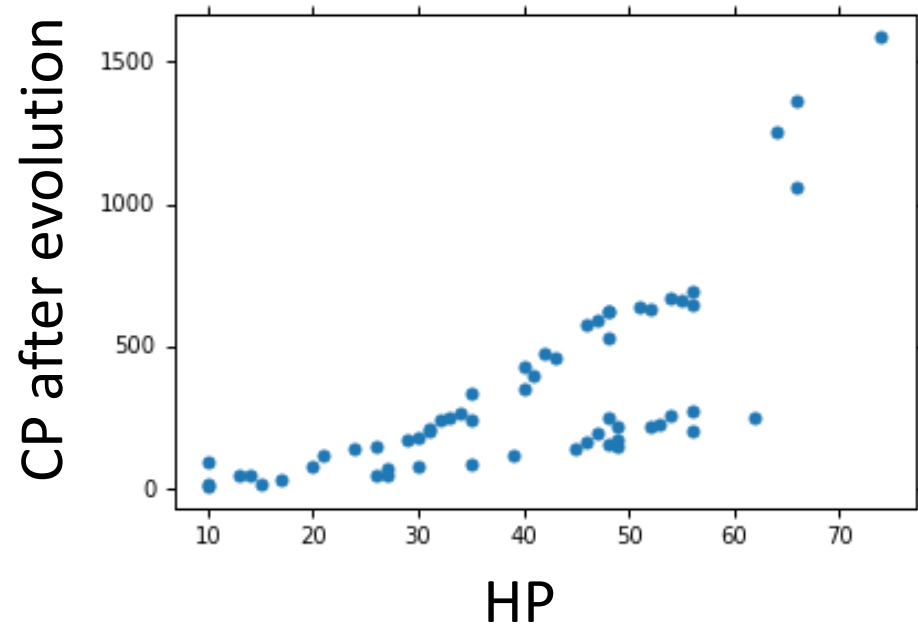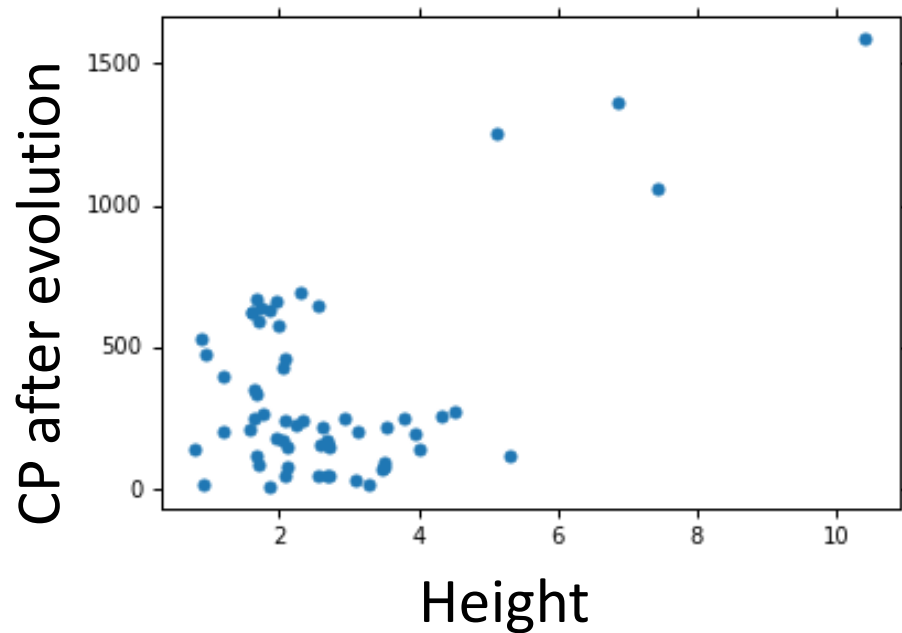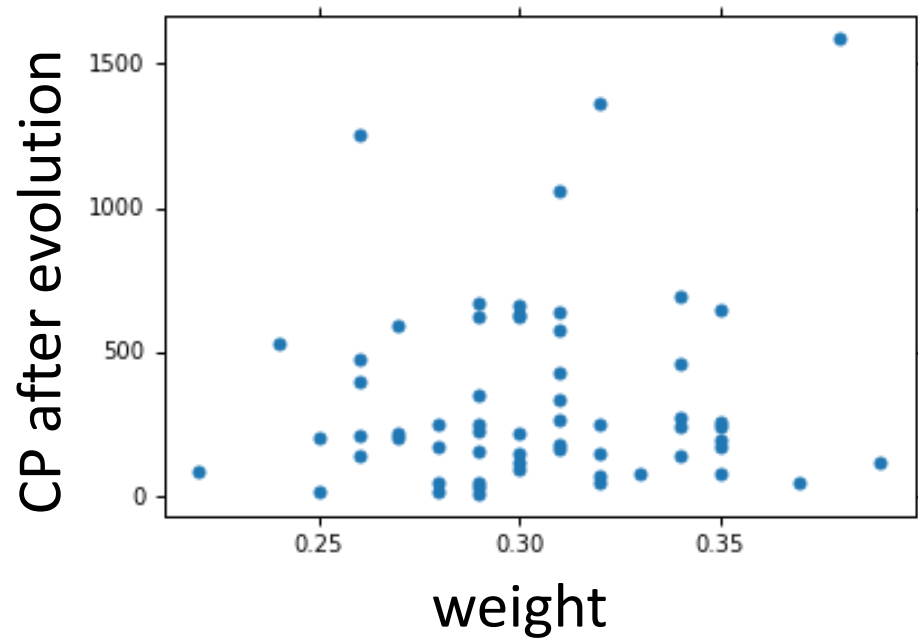Data

Average error
= 3.8

利用新的 Model (Function Set)，從中找到一個 Dest Function 後，確實在 Test Data 上得到更小的 Error。

但是，仍有 Error 代表可能有些「因素」還沒考慮進 Model 中！

Average error
= 14.3

Testing
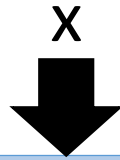Data

Are there any other hidden factors?

# Back to step 1: Redesign the Model Again

$x$
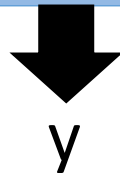
If $x_s$ = Pidgey:   $y' = b_1 + w_1 \cdot x_{cp} + w_5 \cdot (x_{cp})^2$

If $x_s$ = Weedle:   $y' = b_2 + w_2 \cdot x_{cp} + w_6 \cdot (x_{cp})^2$

If $x_s$ = Caterpie:   $y' = b_3 + w_3 \cdot x_{cp} + w_7 \cdot (x_{cp})^2$

If $x_s$ = Eevee:   $y' = b_4 + w_4 \cdot x_{cp} + w_8 \cdot (x_{cp})^2$

$y = y' + w_9 \cdot x_{hp} + w_{10} \cdot (x_{hp})^2$

$+ w_{11} \cdot x_h + w_{12} \cdot (x_h)^2 + w_{13} \cdot x_w + w_{14} \cdot (x_w)^2$

$y$

Training Error = 1.9

Testing Error = 102.3

Overfitting!

# Back to step 2: Regularization

$$y = b + \sum w_i x_i$$

$$L = \sum_n \left( \hat{y}^n - \left( b + \sum w_i x_i \right) \right)^2 \quad +\lambda \sum (w_i)^2$$

The functions with smaller $w_i$ are better
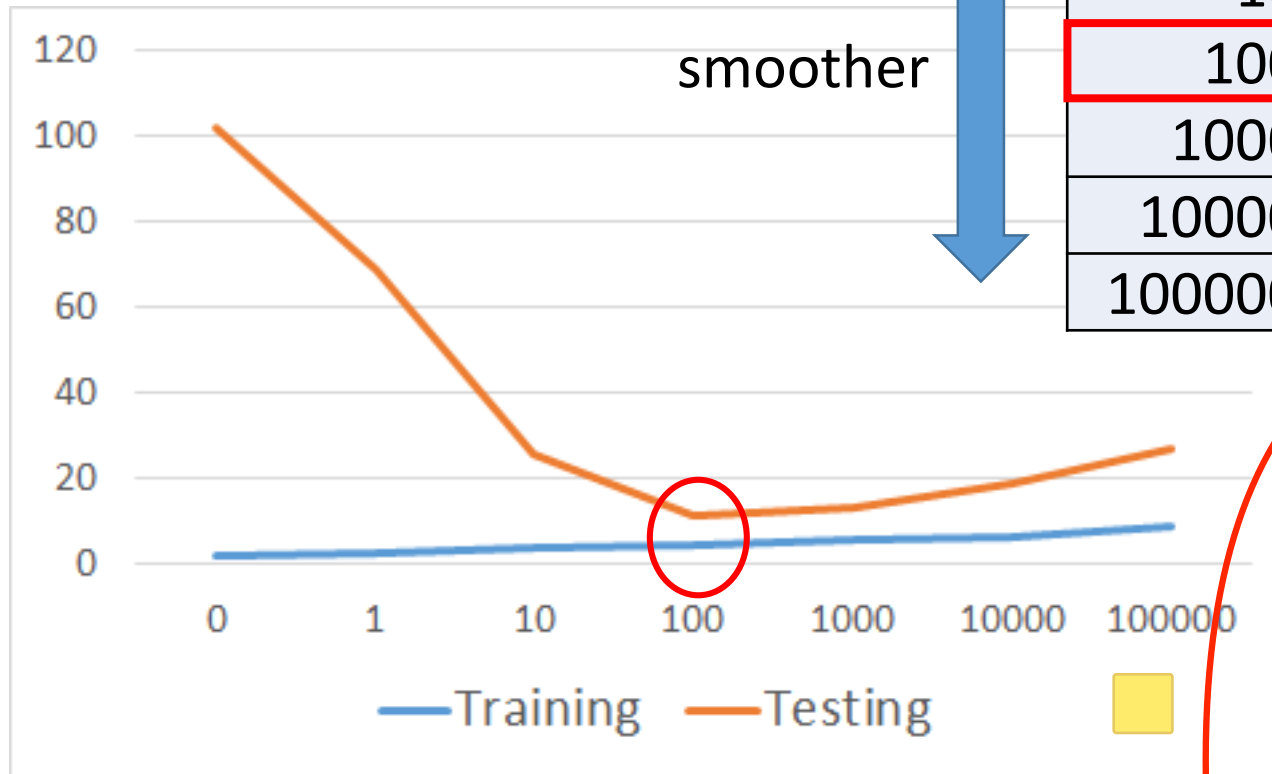
➢ Smaller $w_i$ means … smoother

$$y = b + \sum w_i x_i$$

$$y + \sum w_i \Delta x_i = b + \sum w_i (x_i + \Delta x_i)$$

➢ We believe smoother function is more likely to be correct

Do you have to apply regularization on bias?

# Regularization

| $\lambda$ | Training | Testing |
|---|---|---|
| 0 | 1.9 | 102.3 |
| 1 | 2.3 | 68.7 |
| 10 | 3.5 | 25.7 |
| 100 | 4.1 | 11.1 |
| 1000 | 5.6 | 12.8 |
| 10000 | 6.3 | 18.7 |
| 100000 | 8.5 | 26.8 |

smoother



Training — Testing

How smooth?

Select $\lambda$ obtaining the best model

➢ Training error: larger $\lambda$, considering the training error less

➢ We prefer smooth function, but don't be too smooth.

# Conclusion

- Pokémon: Original CP and species almost decide the CP after evolution
  - There are probably other hidden factors
- Gradient descent
  - More theory and tips in the following lectures
- We finally get average error = 11.1 on the testing data
  - How about new data? Larger error? Lower error?
- Next lecture: Where does the error come from?
  - More theory about overfitting and regularization
  - The concept of validation