# Regression

## Hung-yi Lee
## 李宏毅

# Regression: Output a scalar

- Stock Market Forecast

$$f\left( \quad \right) = \text{Dow Jones Industrial Average at tomorrow}$$

- Self-driving Car

$$f\left( \quad \right) = \text{方向盤角度}$$

- Recommendation

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

# Example Application

- Estimating the Combat Power (CP) of a pokemon after evolution

$f($  $x$ $) =$ CP after evolution $y$

$x_{cp}$

$x_s$

$x_{hp}$

$x_w$  $x_h$

？CP是这么用的？

# Step 1: Model

w and b are parameters
(can be any value)

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



A set of function

Model

$f_1, f_2 \cdots$

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

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

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

...... infinite

$f( \quad x ) = $ CP after evolution $\quad y$

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

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

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.
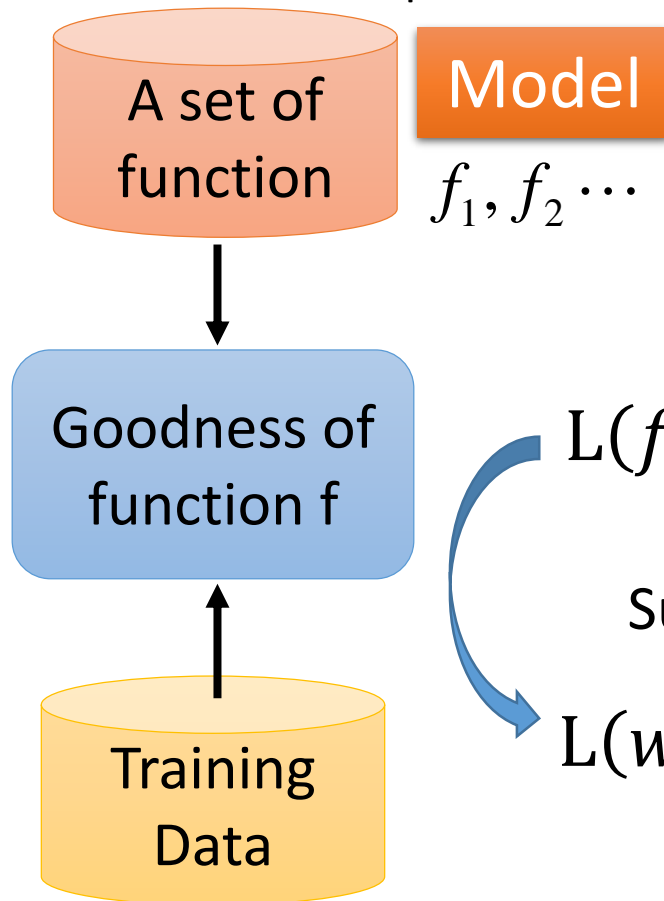


$(x_{cp}^n, \hat{y}^n)$

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

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

Sum over examples

Estimated y based on input function

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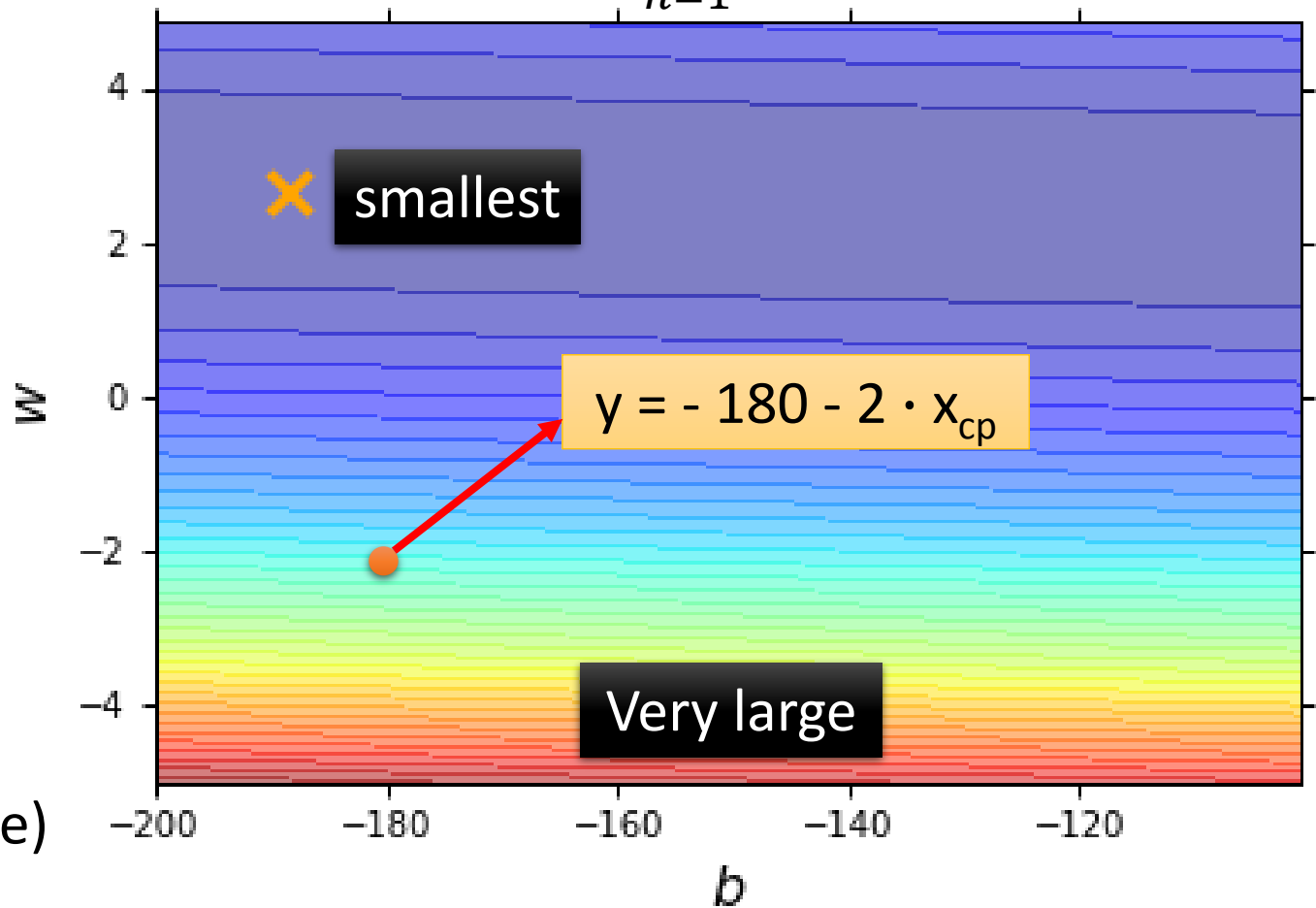
# Step 2: Goodness of Function

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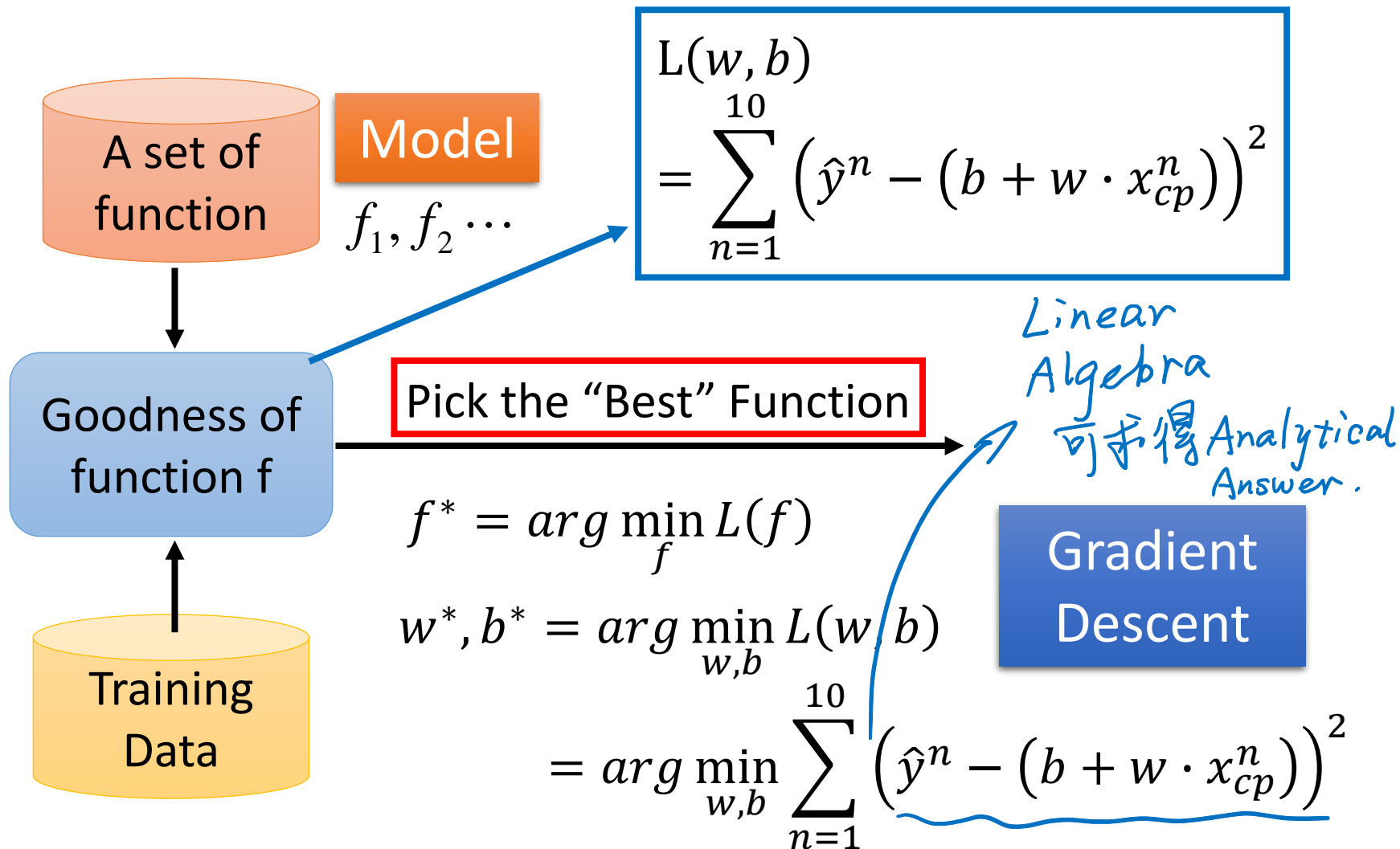
$$L(w, b) = \sum_{n=1}^{10} \left( \hat{y}^n - (b + w \cdot x_{cp}^n) \right)^2$$

Each point in the figure is a function

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

(true example)

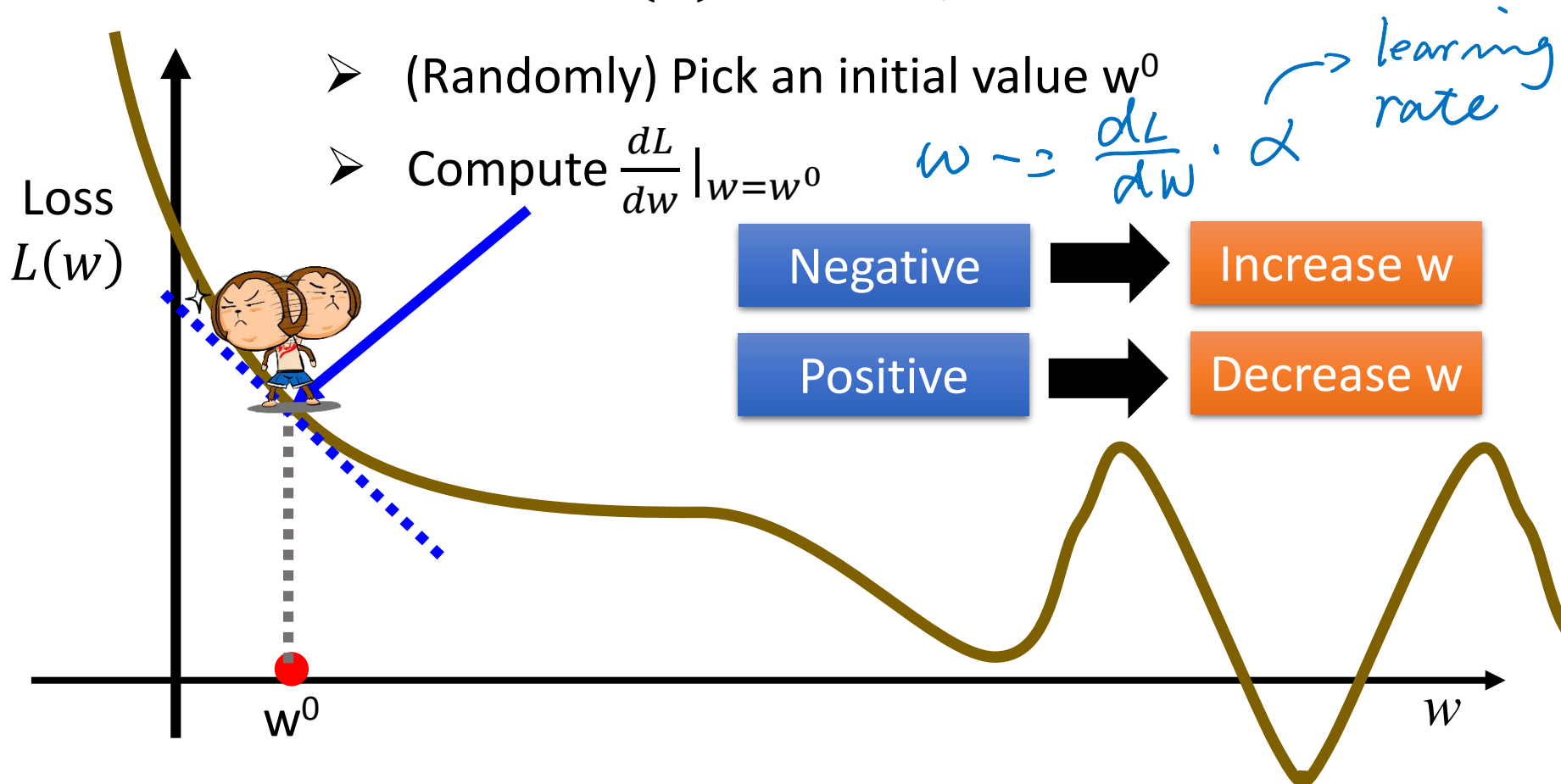# Step 3: Best Function



A set of function

**Model**

$f_1, f_2 \cdots$

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

Goodness of function f

Pick the "Best" Function

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

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

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

Training Data

Linear Algebra
可求得 Analytical Answer.
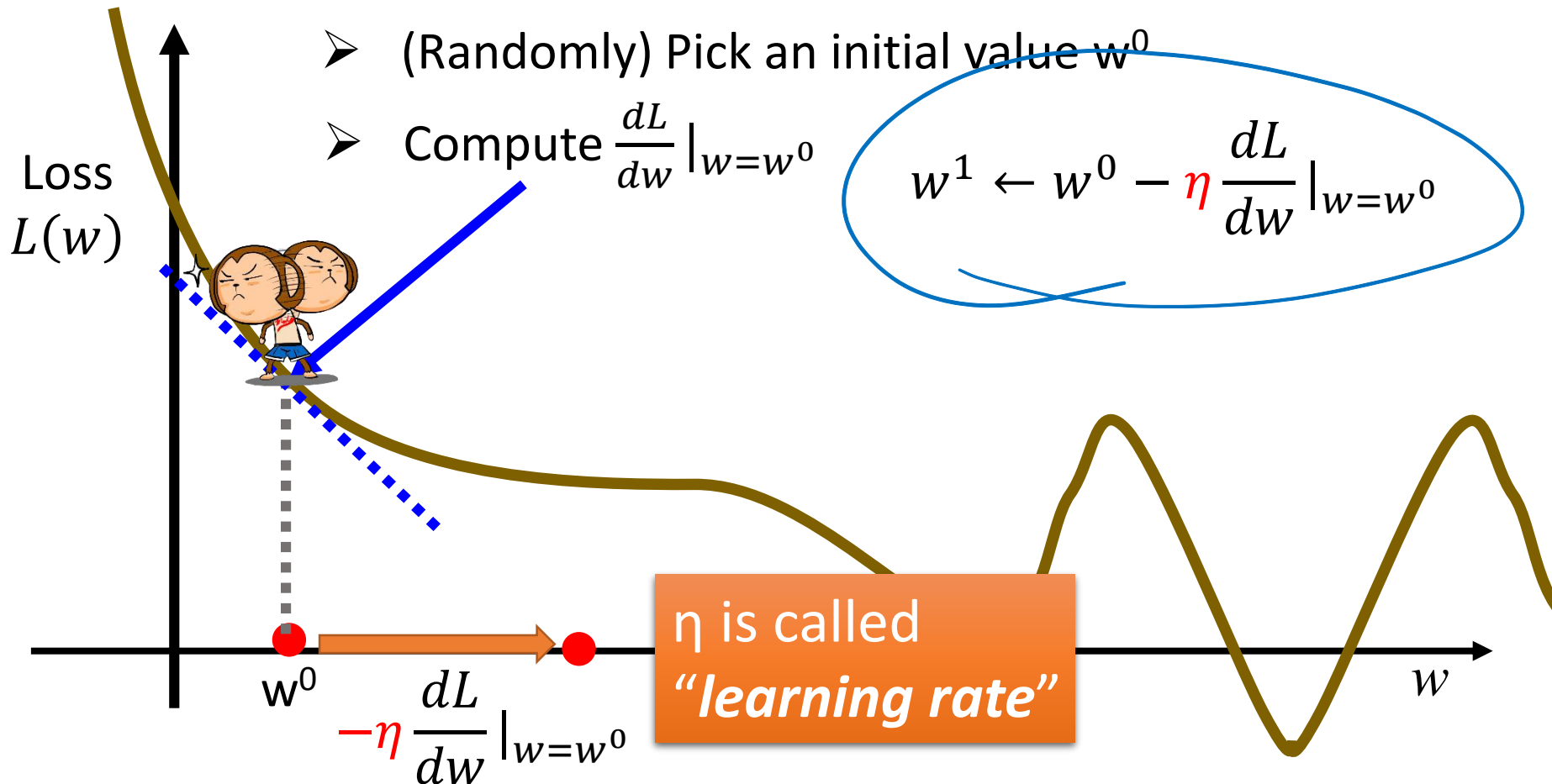
**Gradient Descent**

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

$w \sim \frac{dL}{dw} \cdot \alpha \rightarrow$ learning rate

Loss $L(w)$

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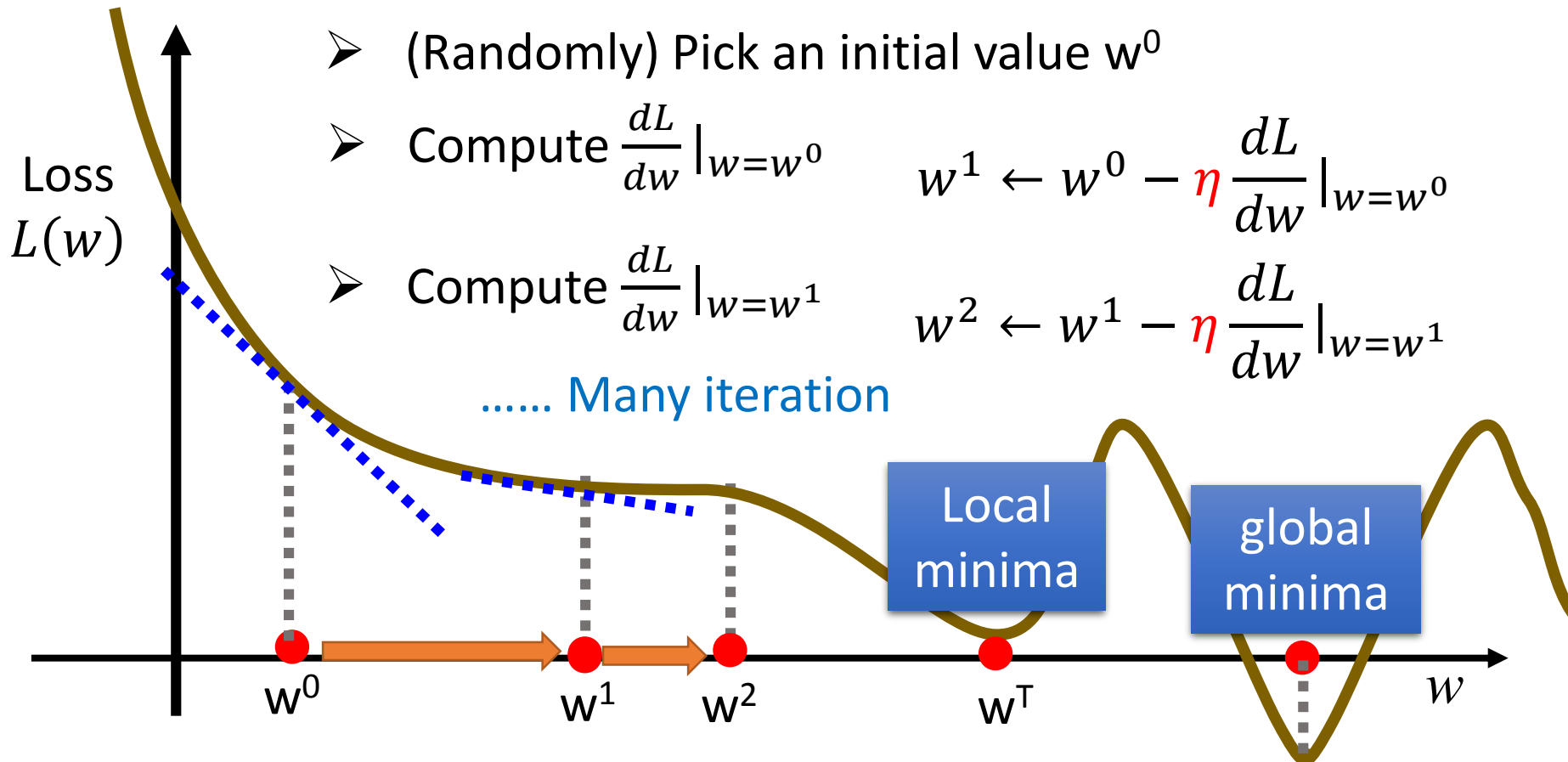
  ➢ (Randomly) Pick an initial value $w^0$

  ➢ Compute $\frac{dL}{dw}|_{w=w^0}$

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

Loss $L(w)$

$w^0$

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

η is called **"learning rate"**

$w$

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

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

…… Many iteration

Loss $L(w)$

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

Local minima

global minima

$w$

# Step 3: Gradient Descent

$$\nabla L = \begin{bmatrix} \frac{\partial L}{\partial w} \\ \frac{\partial L}{\partial b} \end{bmatrix} \text{梯度}$$
gradient

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

(P) review
戓代！

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

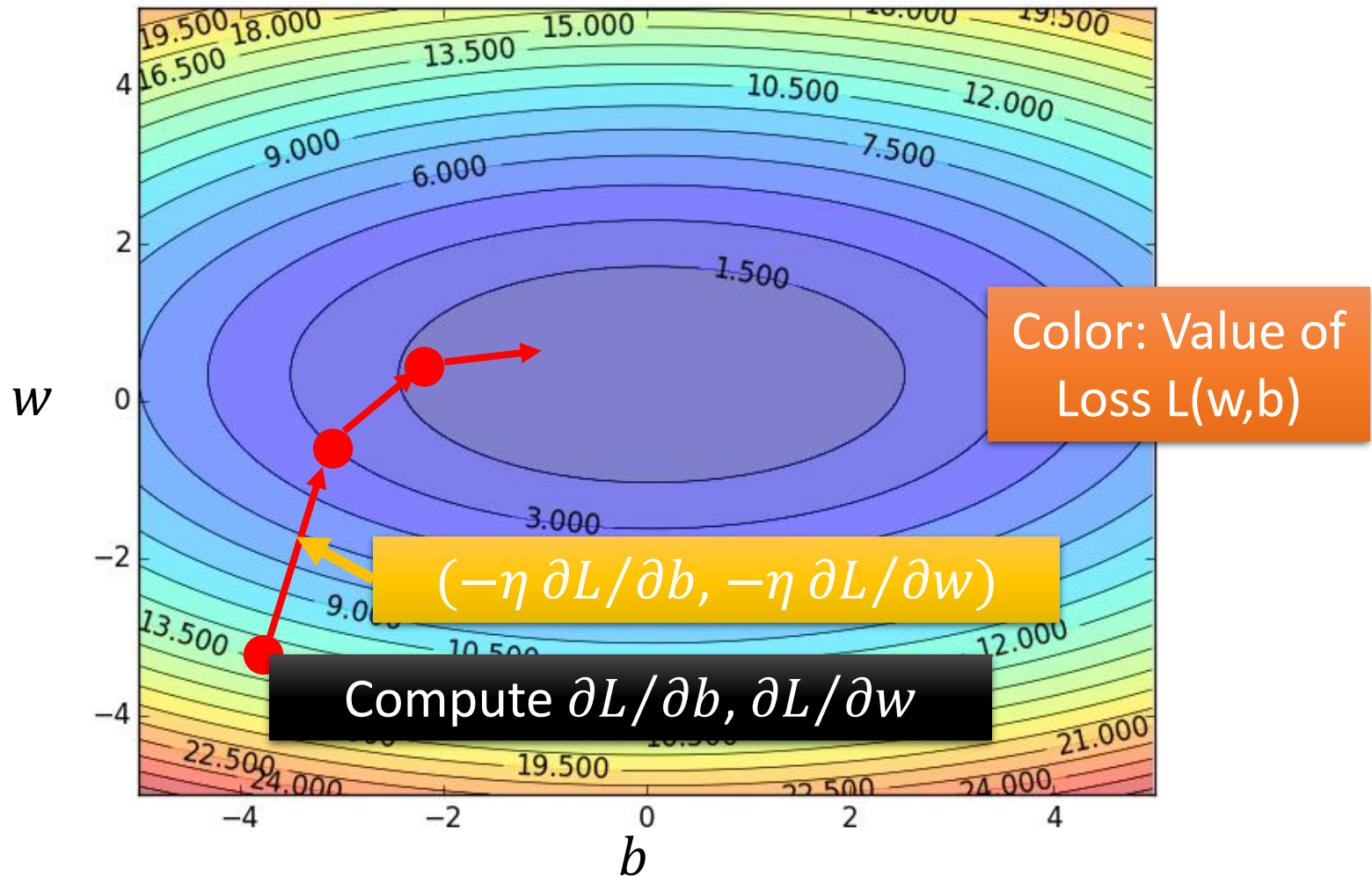  ➤ Compute $\frac{\partial L}{\partial w}|_{w=w^0,b=b^0}, \frac{\partial L}{\partial b}|_{w=w^0,b=b^0}$

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

  ➤ Compute $\frac{\partial L}{\partial w}|_{w=w^1,b=b^1}, \frac{\partial L}{\partial b}|_{w=w^1,b=b^1}$

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

# Step 3: Gradient Descent



Color: Value of Loss L(w,b)

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

Compute $\partial L/\partial b, \partial L/\partial w$
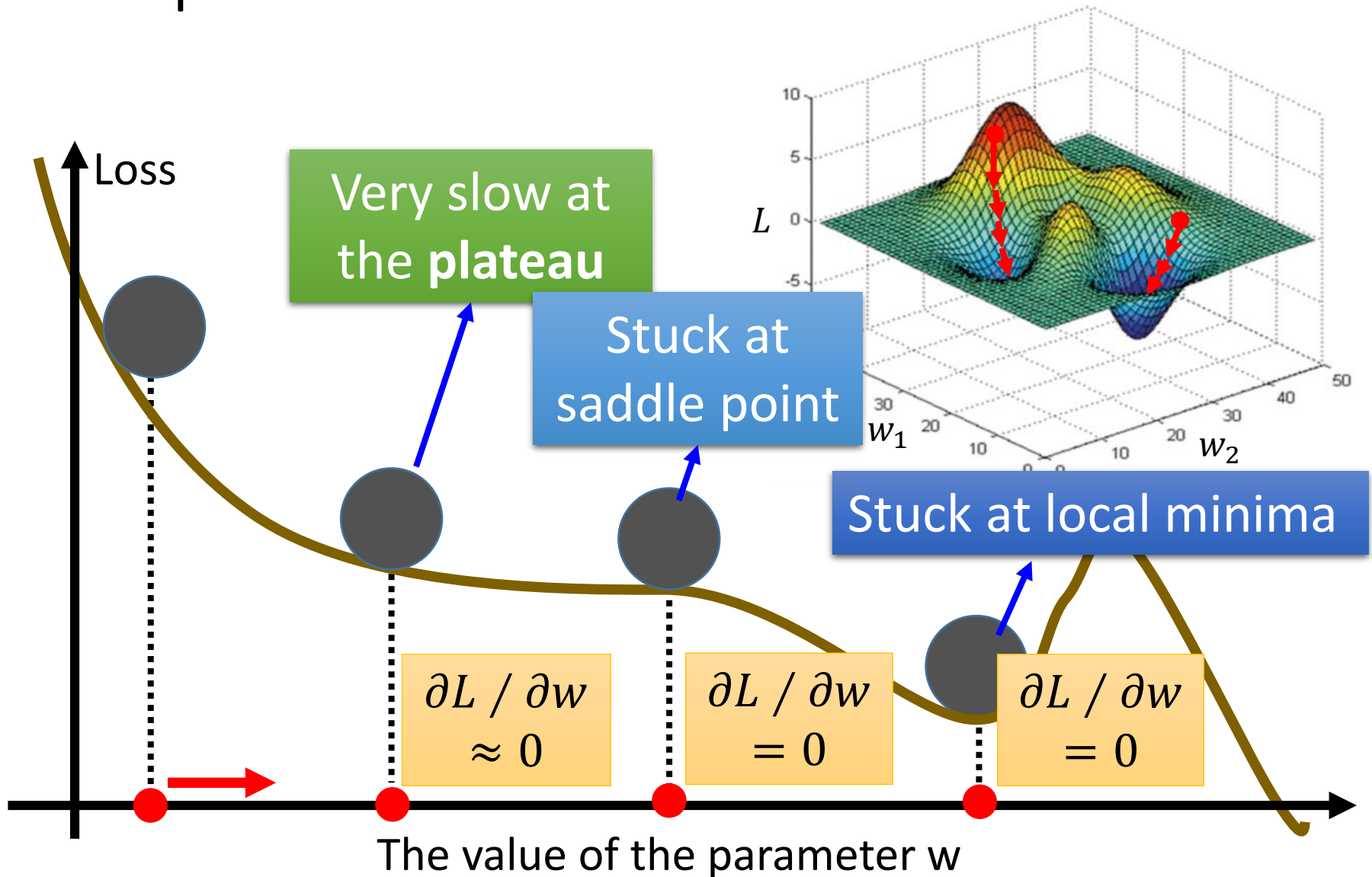
# Step 3: Gradient Descent

- When solving:

$$\theta^* = \arg \max_\theta L(\theta) \quad \text{by gradient descent}$$

- Each time we update the parameters, we obtain $\theta$ that makes $L(\theta)$ smaller.
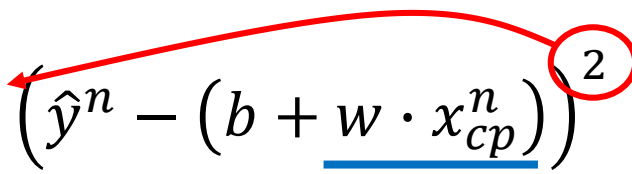
$$L(\theta^0) > L(\theta^1) > L(\theta^2) > \cdots$$

Is this statement correct?

# Step 3: Gradient Descent



Loss

Very slow at the **plateau**

Stuck at saddle point

Stuck at local minima

$\partial L \,/\, \partial w \approx 0$

$\partial L \,/\, \partial w = 0$

$\partial L \,/\, \partial w = 0$

The value of the parameter w

# Step 3: Gradient Descent

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

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

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

$$\frac{\partial L}{\partial b} = ?$$

# Step 3: Gradient Descent

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

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

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

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

# Step 3: Gradient Descent

# How's the results?

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

b = -188.4

w = 2.7

Average Error on
Training Data

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

Loss是
MSE
≠Loss!
too小?



Training Data

# How's the results? - Generalization

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

→ 这么好抓的吗？

Another 10 pokemons as 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)

How can we do better?

## Selecting another Model

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



## 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

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

## 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

*Overfitting*

Slightly better. How about more complex model?

# *Selecting another Model*

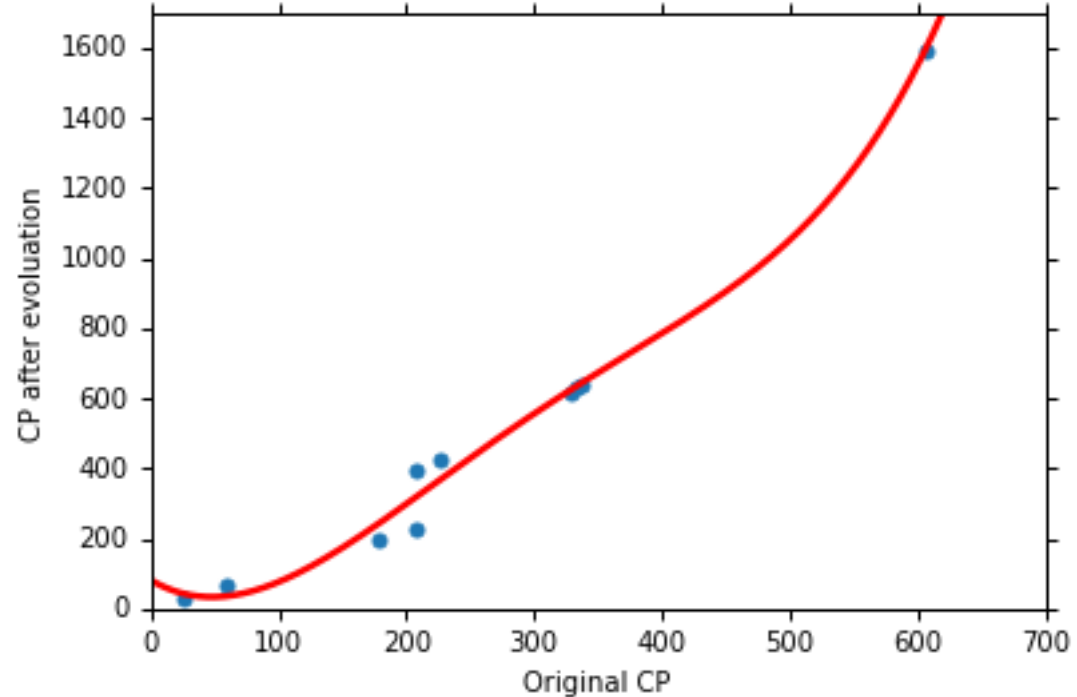$$y = b + w_1 \cdot x_{cp} + w_2 \cdot (x_{cp})^2 + w_3 \cdot (x_{cp})^3 + w_4 \cdot (x_{cp})^4$$

# *Best Function*

Average Error = 14.9

# *Testing:*

Average Error = 28.8

The results become worse …

# Selecting another Model

$$y = b + w_1 \cdot x_{cp} + w_2 \cdot (x_{cp})^2 + w_3 \cdot (x_{cp})^3 + w_4 \cdot (x_{cp})^4 + w_5 \cdot (x_{cp})^5$$
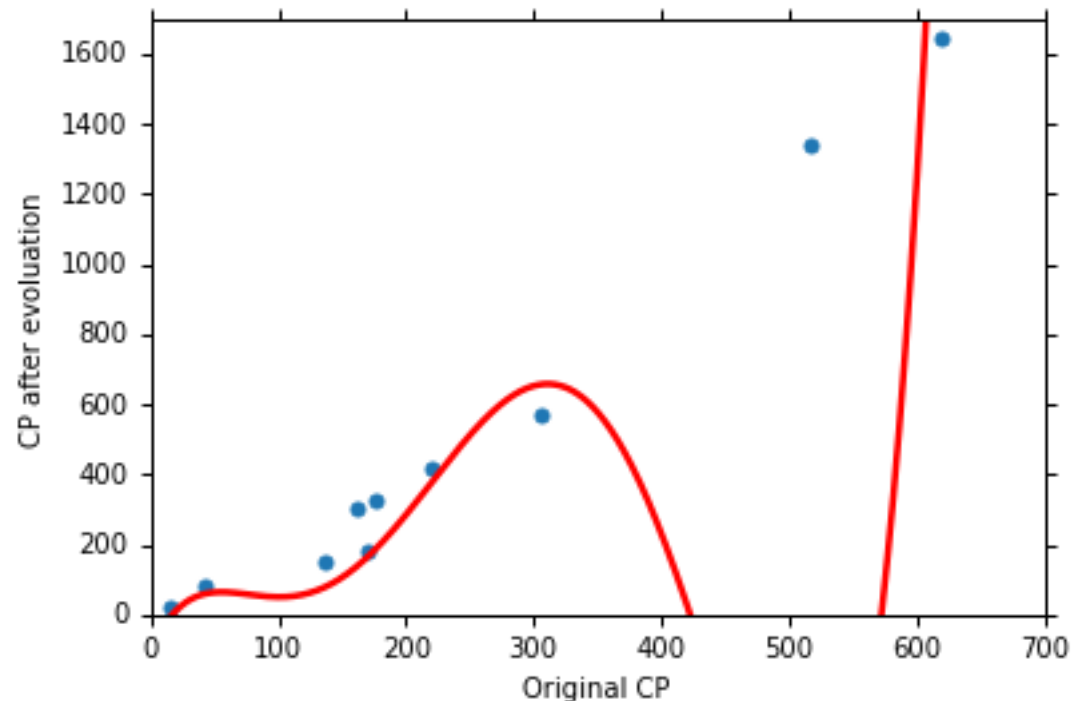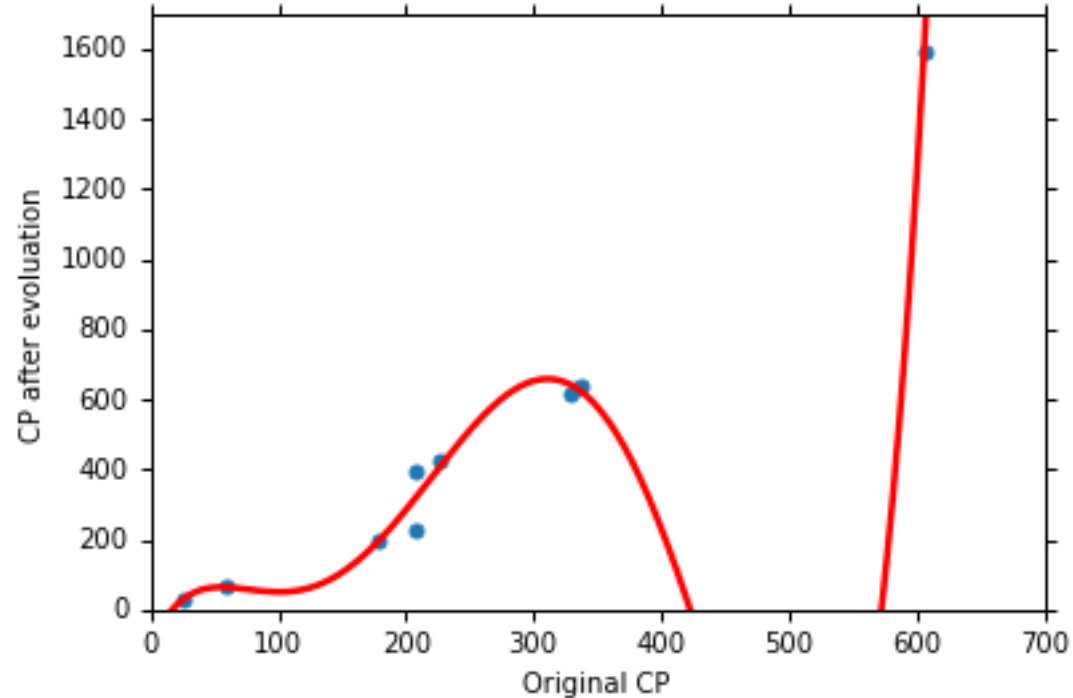
# Best Function

Average Error = 12.8

# Testing:

Average Error = 232.1

The results are so bad.
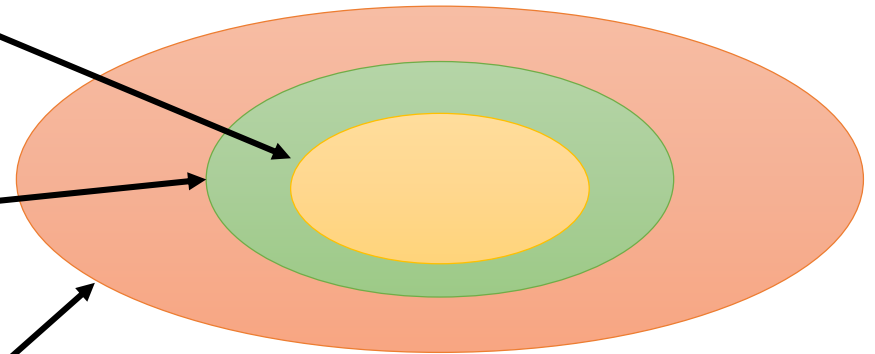
# Model Selection

对 training
Data
算是一直下降

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

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

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

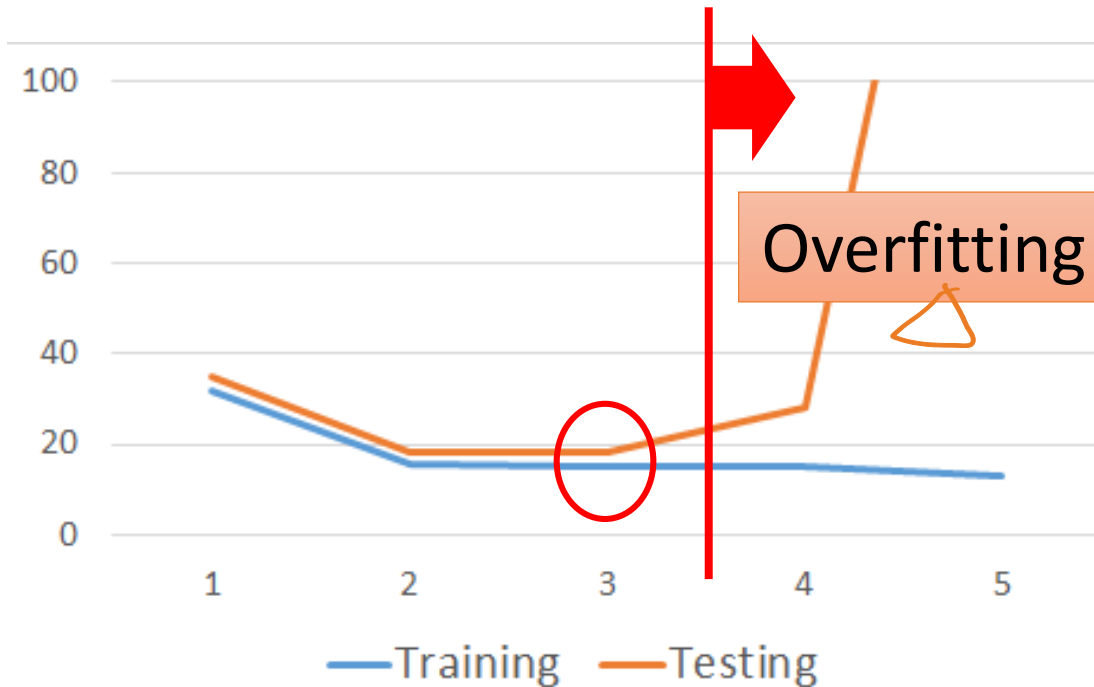4. $y = b + w_1 \cdot x_{cp} + w_2 \cdot (x_{cp})^2 + w_3 \cdot (x_{cp})^3 + w_4 \cdot (x_{cp})^4$

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.

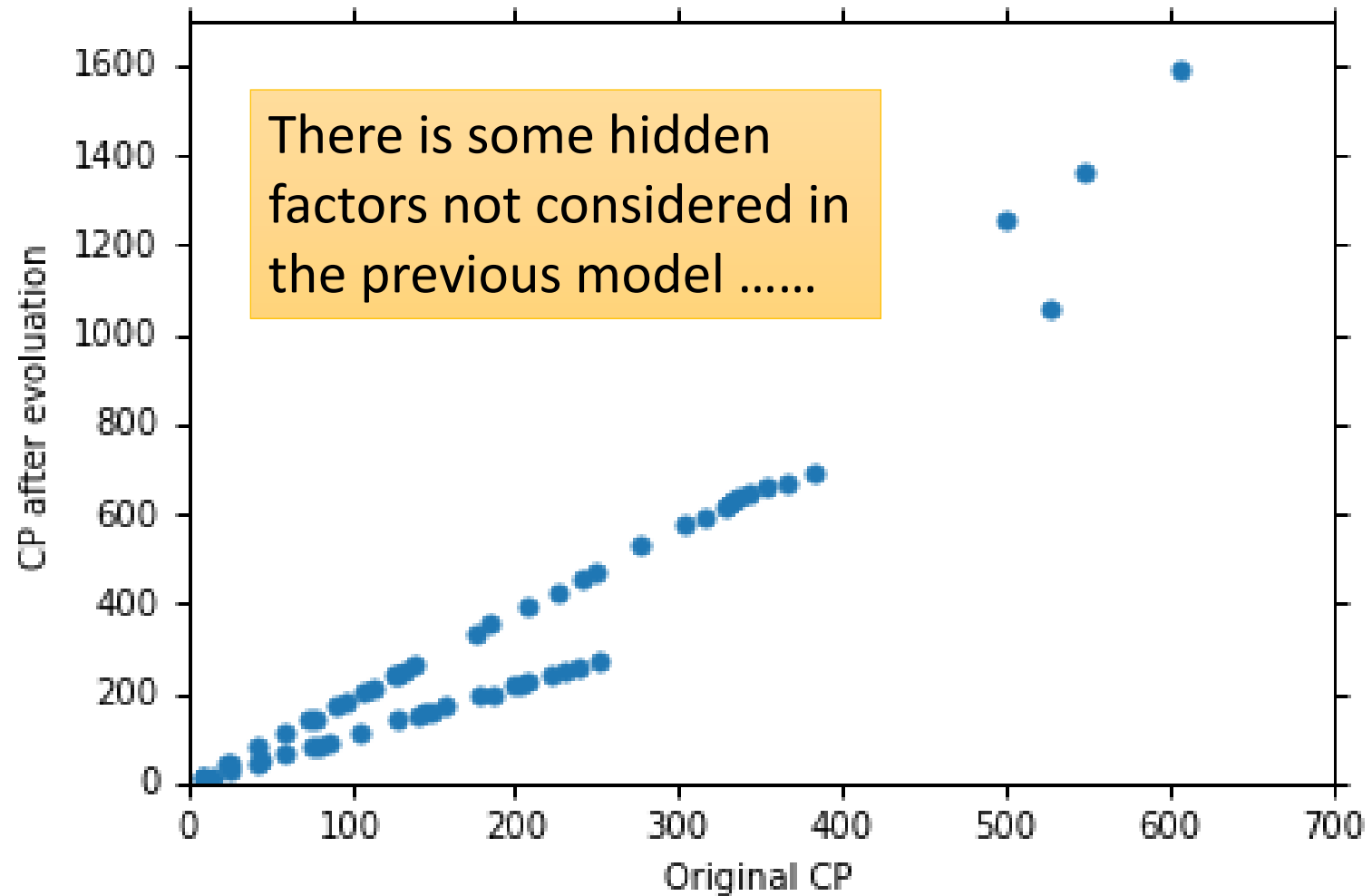If we can truly find the best function

# Model Selection



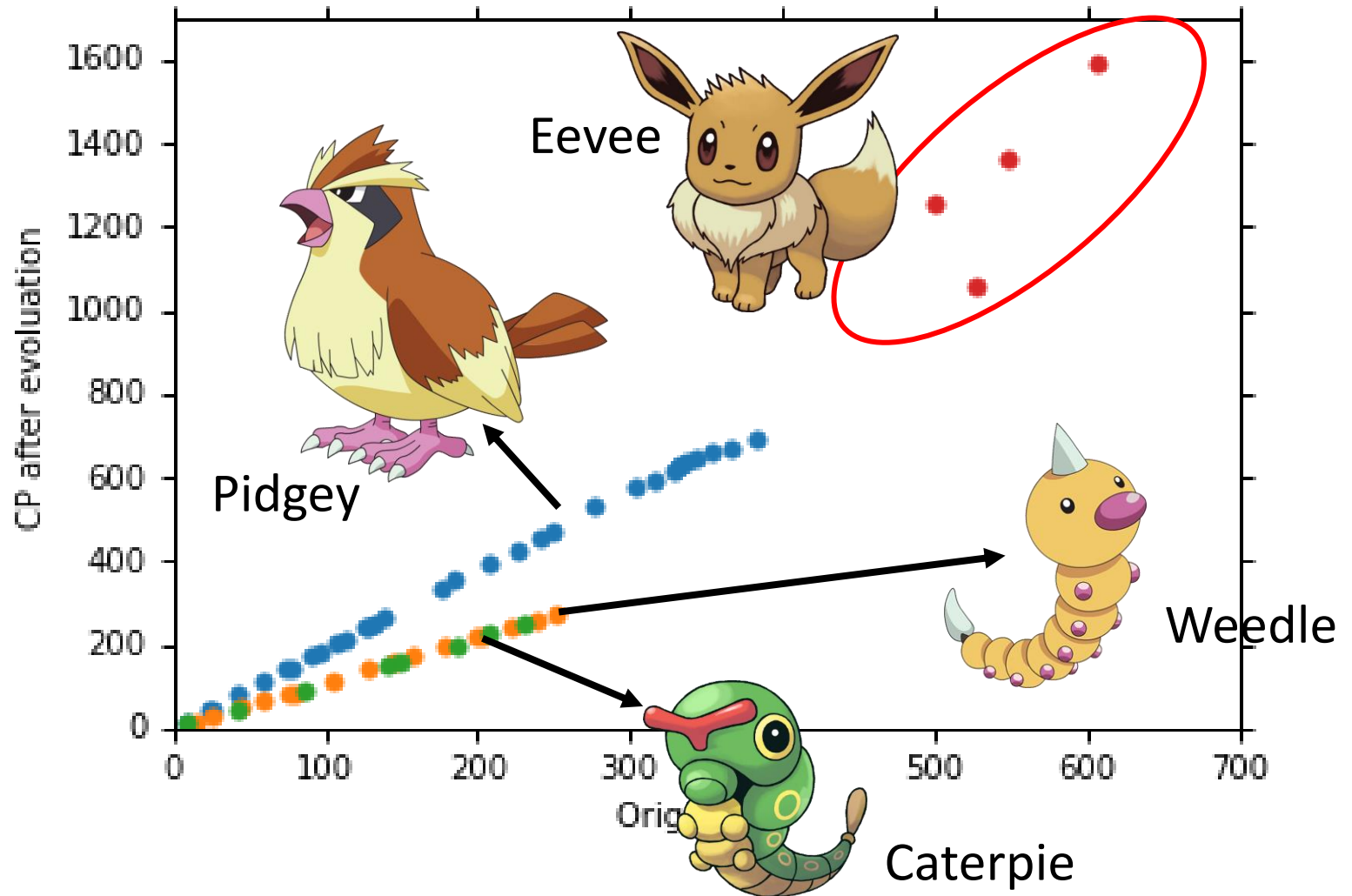| | Training | Testing |
|---|---|---|
| 1 | 31.9 | 35.0 |
| 2 | 15.4 | 18.4 |
| 3 | 15.3 | 18.1 |
| 4 | 14.9 | 28.2 |
| 5 | 12.8 | 232.1 |

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

This is ***Overfitting***. ➡ Select suitable model
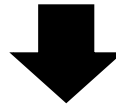
# Let's collect more data



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: | $y = b_4 + w_4 \cdot x_{cp}$ |

↓

y

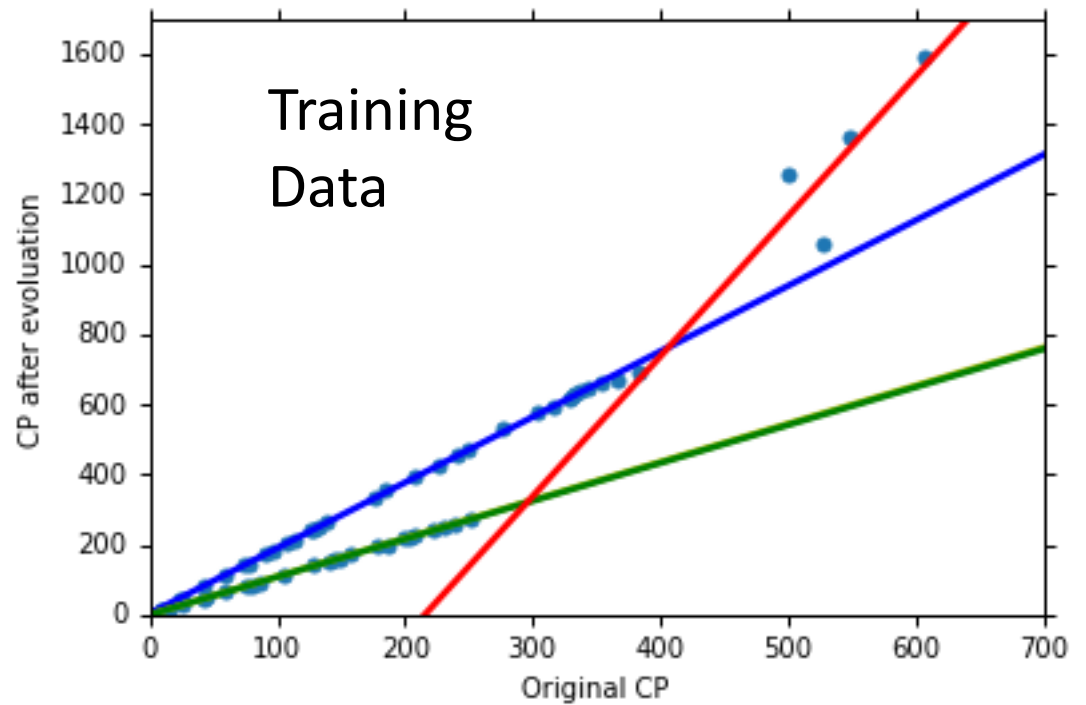# Back to step 1: Redesign the Model

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

$y = b_1 \cdot$ | 1 | $\delta(x_s = \text{Pidgey})$

$+w_1 \cdot$ | 1 | $x_{cp}$ $\delta(x_s = \text{pidgey}) \cdot x_{cp}$

$+b_2 \cdot$ | 0 | $\delta(x_s = \cdots)$    $\delta(x_s = \text{Pidgey})$

$+w_2 \cdot$ | 0 |

$+b_3 \cdot$ | 0 |

$+w_3 \cdot$ | 0 |

$+b_4 \cdot$ | 0 |

$+w_4 \cdot$ | 0 |

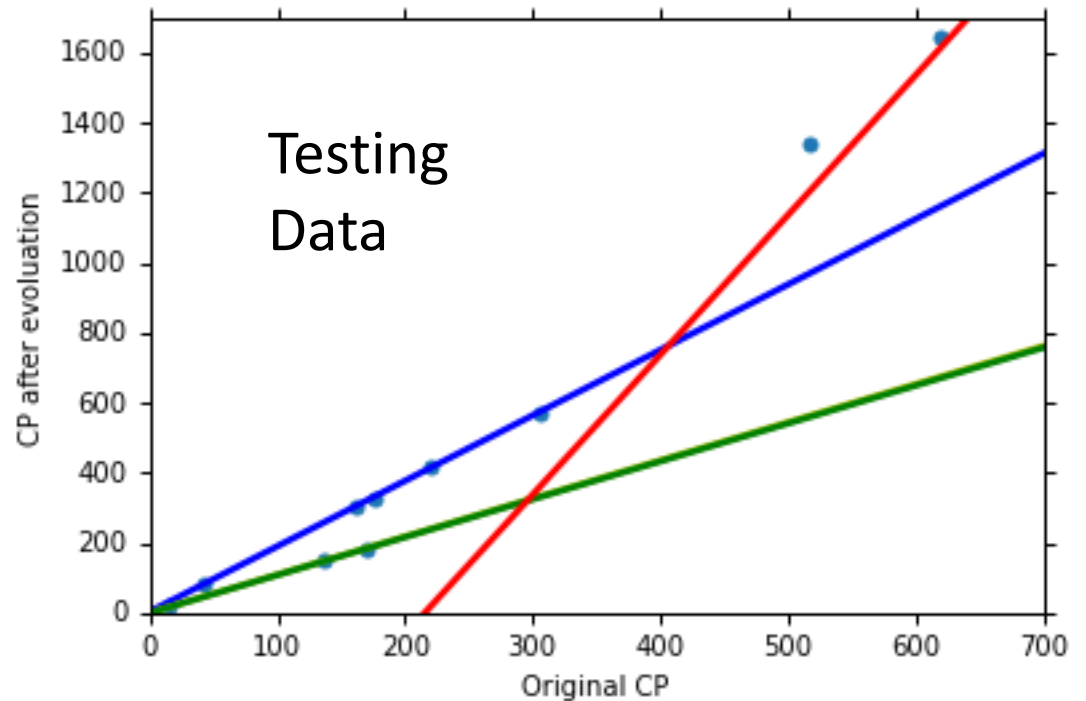$\begin{cases} =1 & \text{If } x_s = \text{Pidgey} \\ \\ =0 & \text{otherwise} \end{cases}$
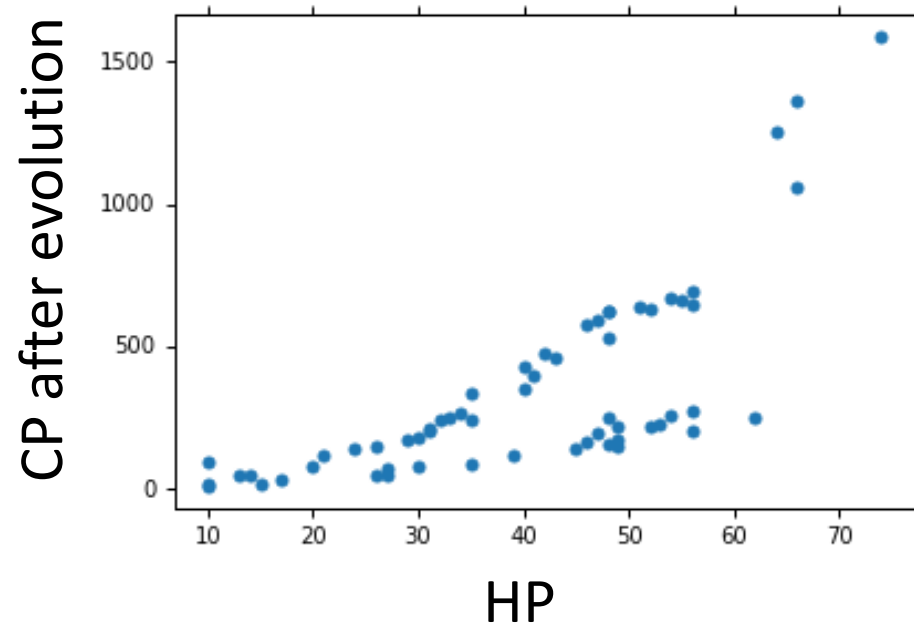
If $x_s = \text{Pidgey}$
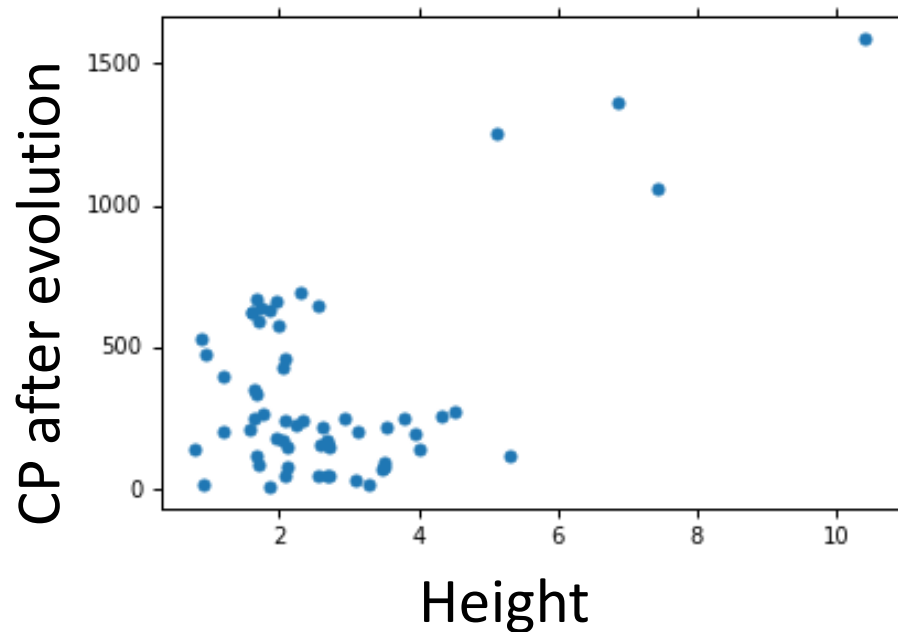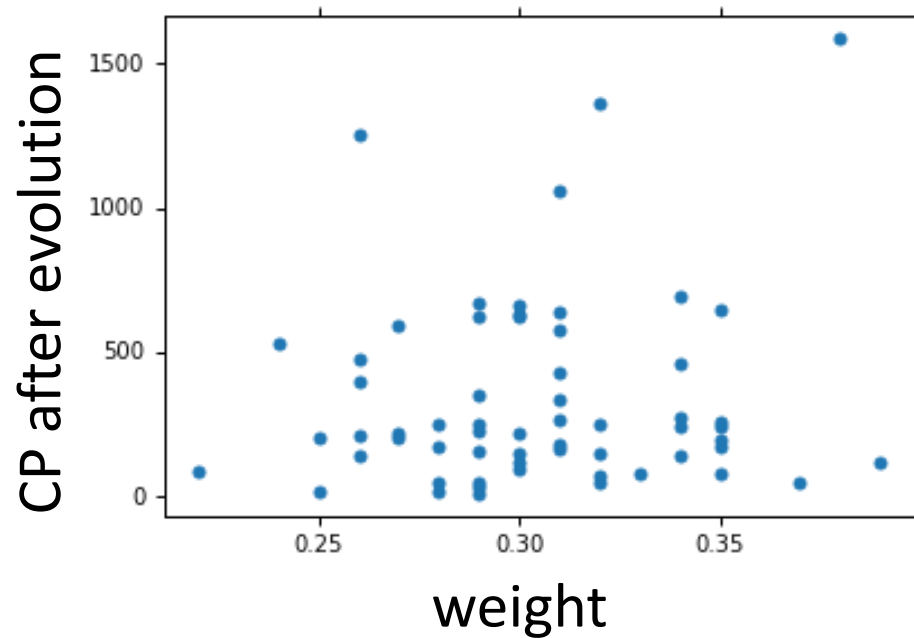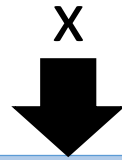
$$y = b_1 + w_1 \cdot x_{cp}$$

Average error = 3.8

Average error = 14.3

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

$y = y' + w_9 \cdot x_{hp} + w_{10} \cdot (x_{hp})^2$
$+ w_{11} \cdot x_h + w_{12} \cdot (x_h)^2 + w_{13} \cdot x_w + w_{14} \cdot (x_w)^2$

Training Error = 1.9

Testing Error = 102.3

Overfitting!

↓

y

# Back to step 2: Regularization

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

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

→ 因为 expect $L \to 0$

The functions with smaller $w_i$ are better

⇒ function 更平滑.

hyper-parameter

➤ Smaller $w_i$ means … **smoother**

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

对干扰更 微虑

$$y + \sum w_i \Delta x_i = b + \sum w_i (x_i + \Delta x_i)$$
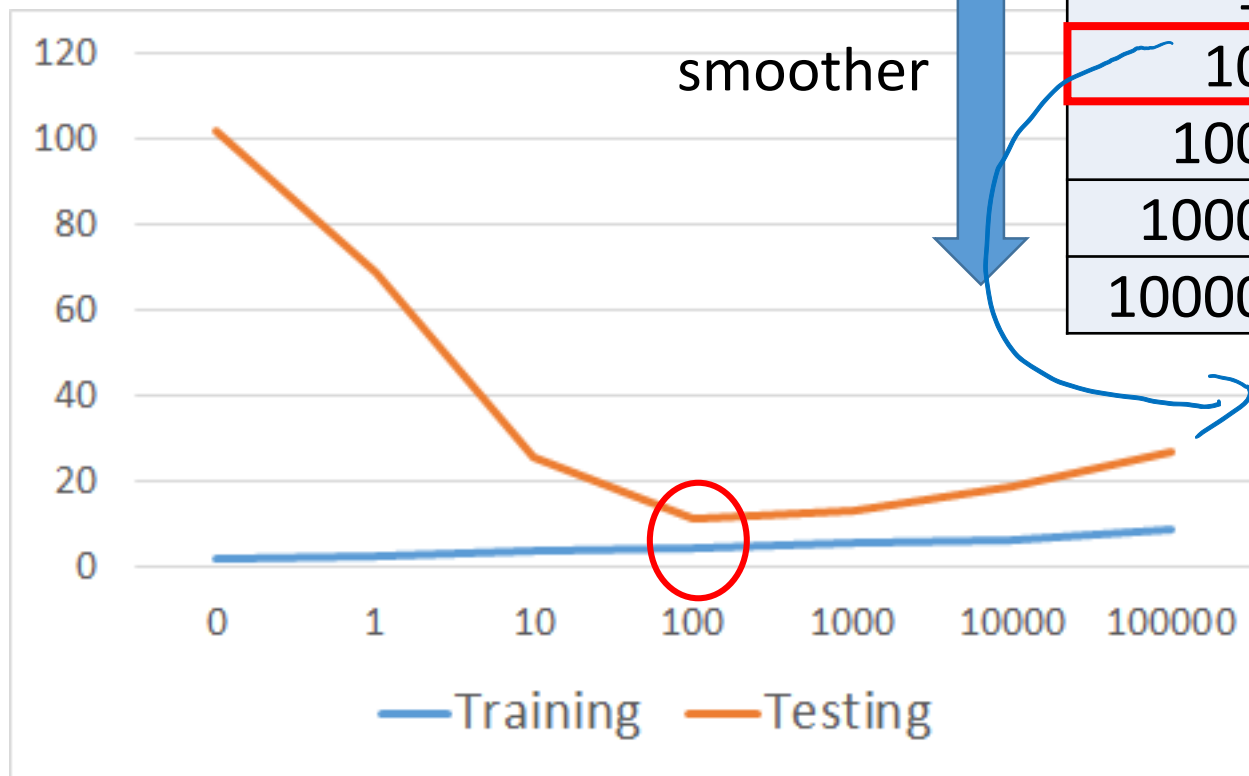
➤ We believe smoother function is more likely to be correct

Do you have to apply regularization on bias?

对函数平滑也无影响. 只是上下平移罢?

# 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



只有plot咯.

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

# Reference

- Bishop: Chapter 1.1

# Acknowledgment

- 感謝 鄭凱文 同學發現投影片上的符號錯誤
- 感謝 童寬 同學發現投影片上的符號錯誤
- 感謝 黃振綸 同學發現課程網頁上影片 連結錯誤 的符號錯誤