

# Regression

## Hung-yi Lee

### 李宏毅

# Regression: Output a scalar

- Stock Market Forecast 股票预测

$$f(\text{过去10年... 的资料}) = \text{明天 明天道琼斯工业平均指数 Dow Jones Industrial Average at tomorrow}$$

- Self-driving Car 无人车 / 自动驾驶

$$f(\text{无人车感受到的信息}) = \text{方向盘角度} \rightarrow \text{如何行驶}$$

- Recommendation 推荐系统

$$f(\text{使用者 A 商品 B}) = \text{购买可能性}$$

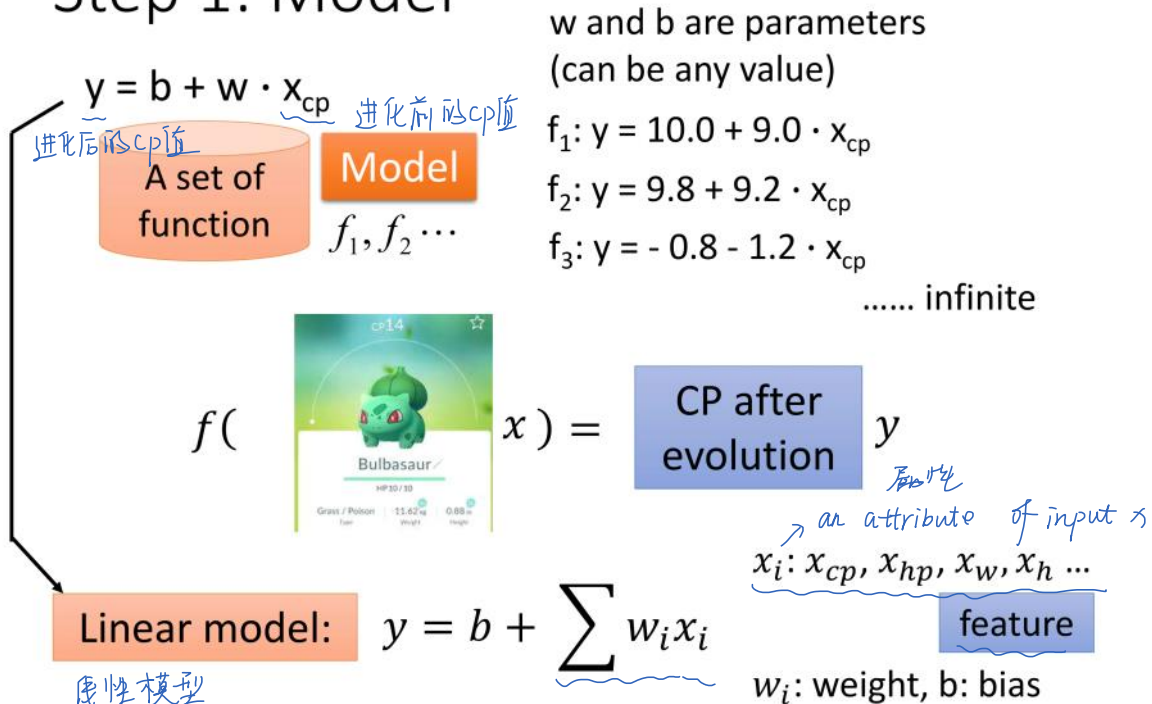
## Example Application

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

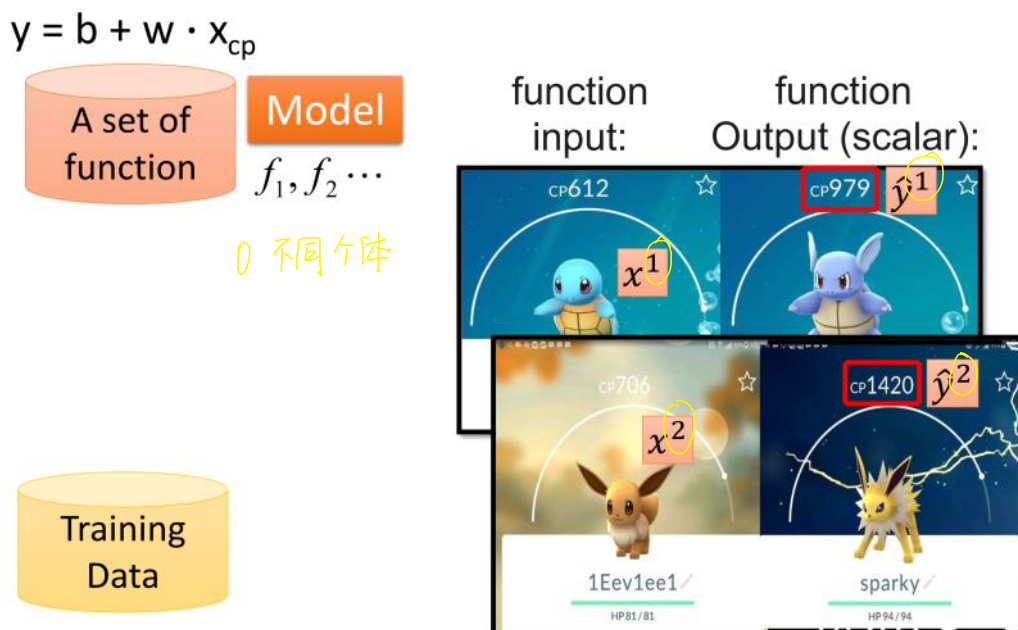
$$f(\text{某一只宝可梦 } x_s) = \text{进化后的 CP 值 } y$$

$x_{cp}$  某一只宝可梦之间的 CP 值  
 $x_{hp}$  HP 值  
 $x_w$  重量  
 $x_h$  高度  
 $x_s$  宝可梦种类 (Bulbasaur)  
 $x_{cp}$  进化前的 CP 值  
 $y$  进化后的 CP 值

## Step 1: Model



## Step 2: Goodness of Function



## Step 2: Goodness of Function

Training Data:  
10 pokemons

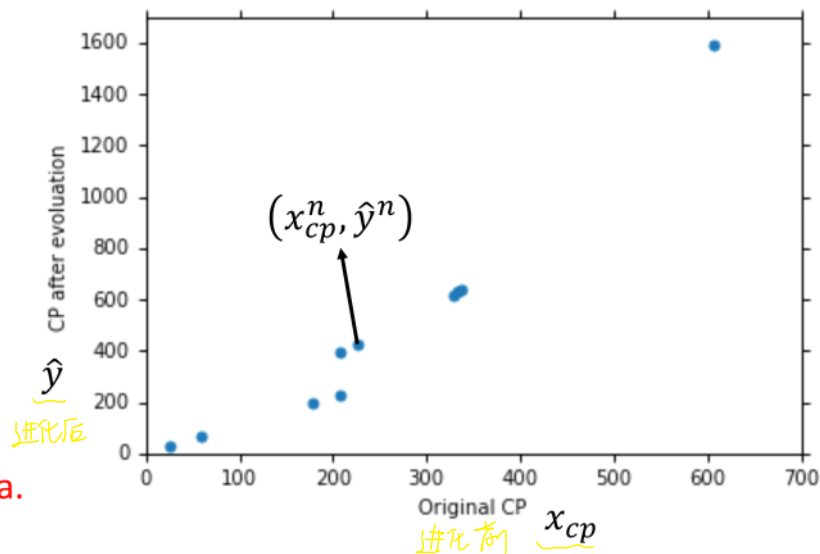
$(x^1, \hat{y}^1)$

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

$\vdots$

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

This is real data.



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

$f_1, f_2, \dots$

Loss function  $L$ :

Input: a function, output: how bad it is

$$L(f) = L(w, b)$$

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

Goodness of function  $f$

Training Data

$$L(f) = \sum_{n=1}^{10} (\hat{y}^n - \underbrace{f(x_{cp}^n)}_{\text{Estimated } y \text{ based on input function}})^2$$

Sum over examples

Estimated  $y$  based on input function

$$L(w, b) = \sum_{n=1}^{10} (\underbrace{\hat{y}^n}_{\text{真的}} - \underbrace{(b + w \cdot x_{cp}^n)}_{\text{预测的}})^2$$

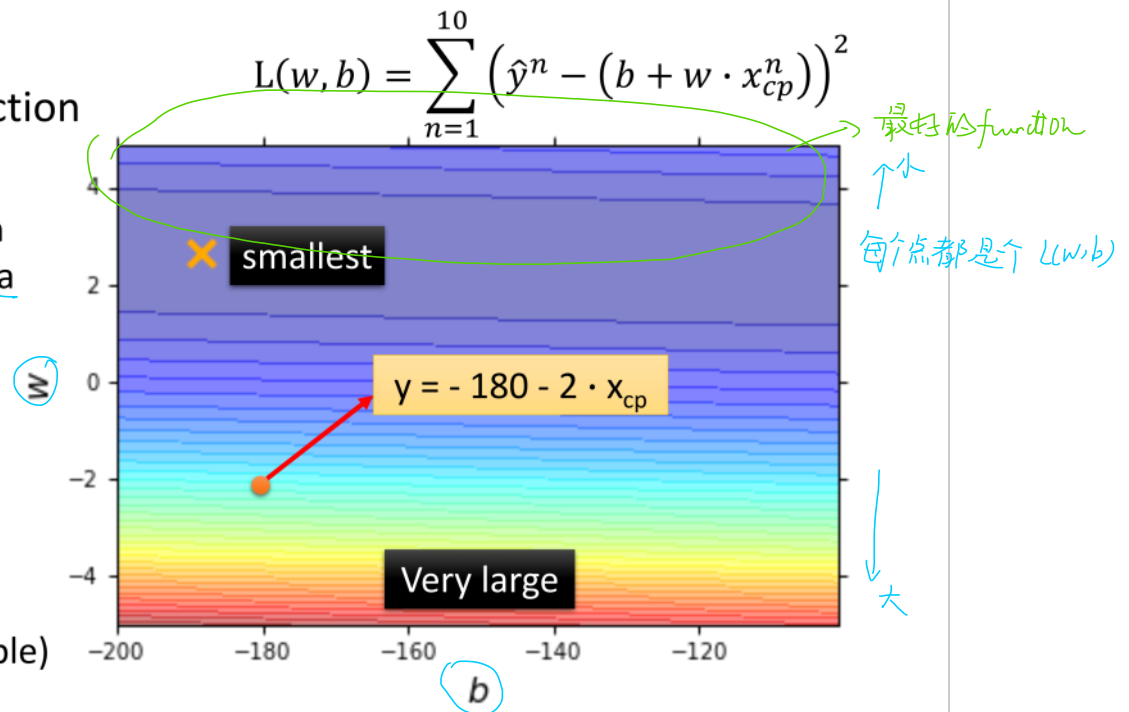
## Step 2: Goodness of Function

### • Loss Function

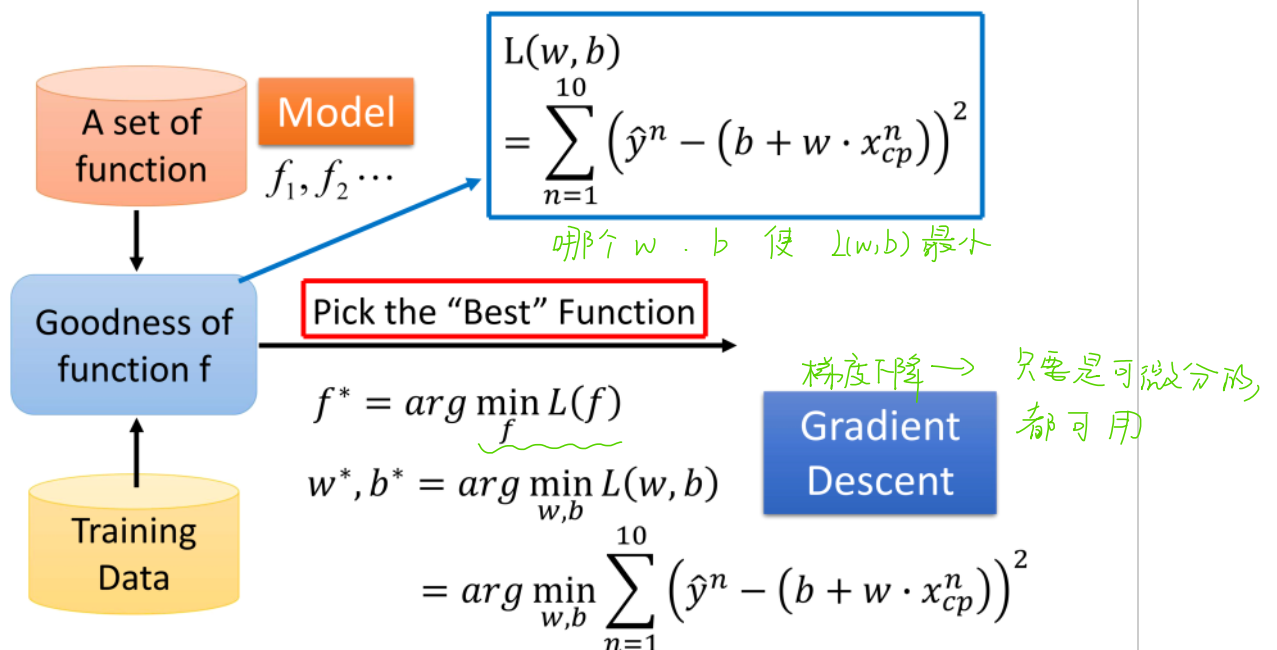
Each point in the figure is a function

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

(true example)



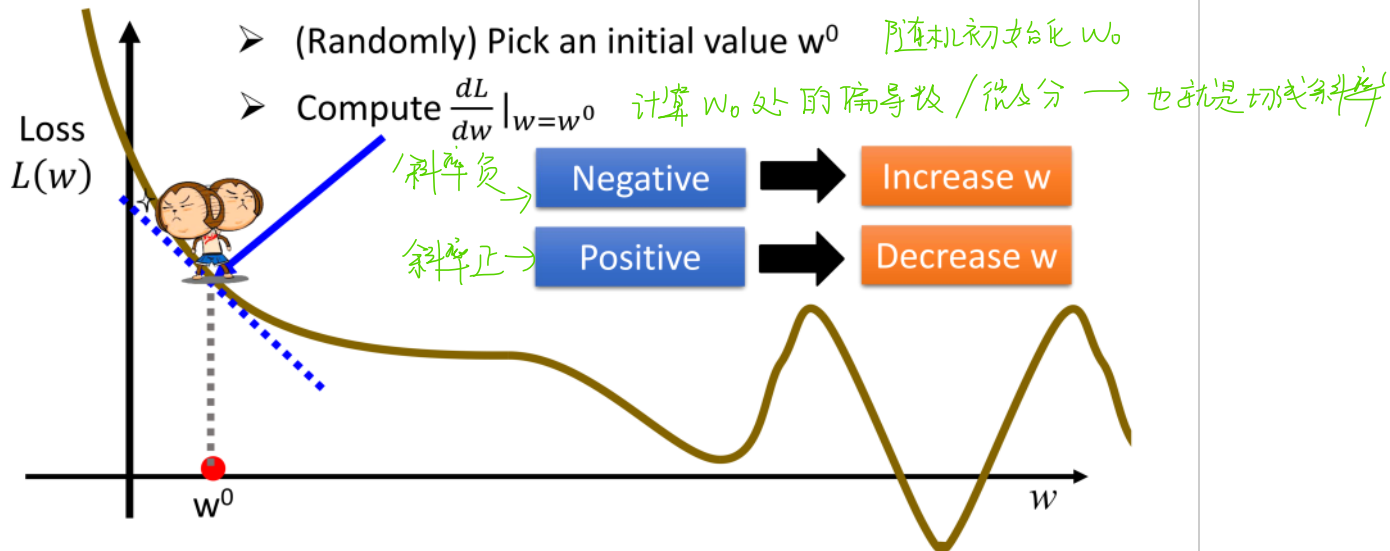
## Step 3: Best Function



## Step 3: Gradient Descent

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

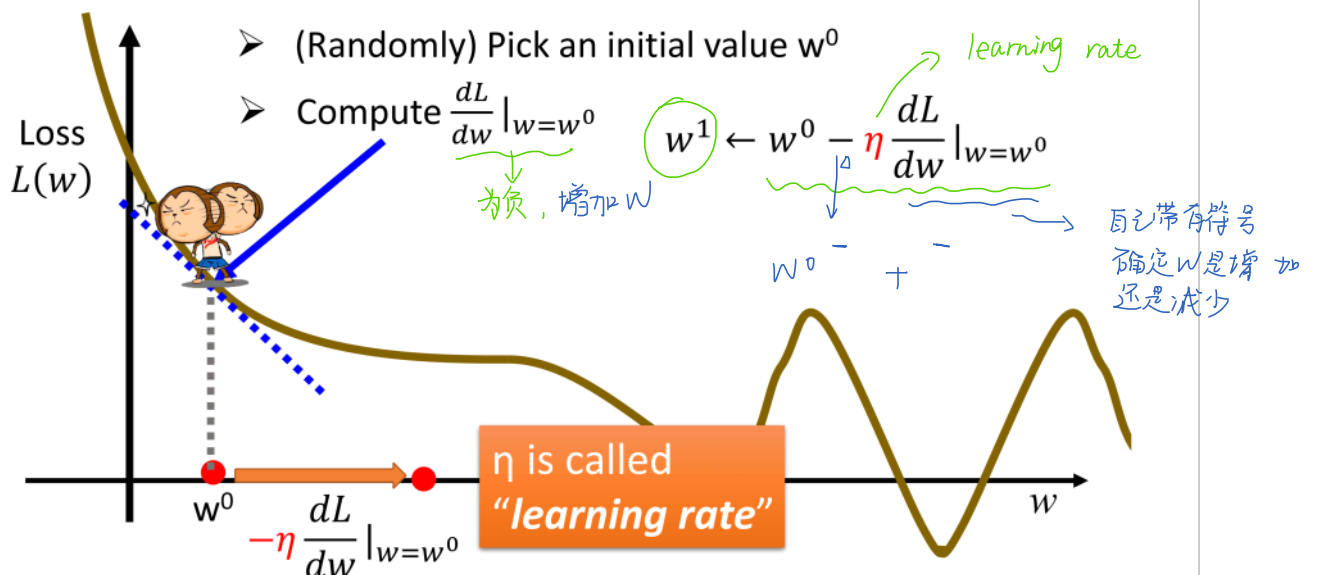
- Consider loss function  $L(w)$  with one parameter  $w$ :



## Step 3: Gradient Descent

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

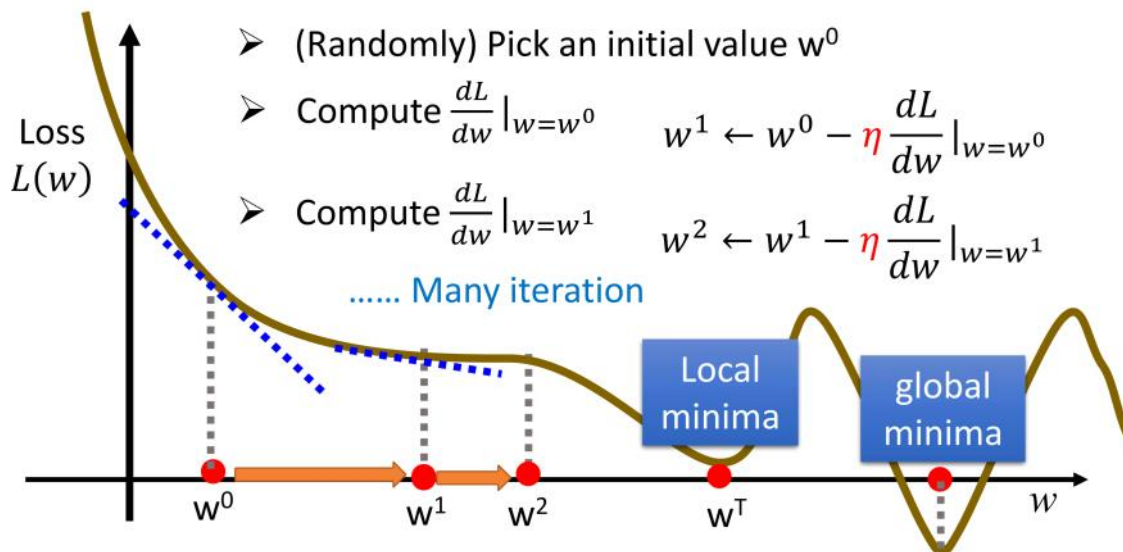
- Consider loss function  $L(w)$  with one parameter  $w$ :



## Step 3: Gradient Descent

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

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



## Step 3: Gradient Descent

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

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

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

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

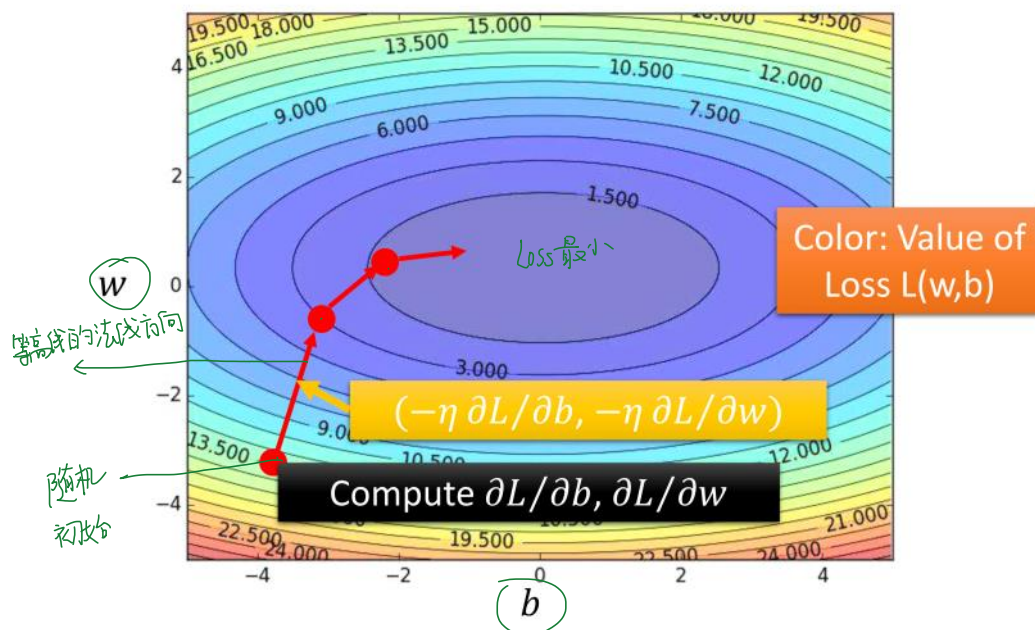
$$\underline{w^1} \leftarrow w^0 - \eta \frac{\partial L}{\partial w} \big|_{\underline{w=w^0}, b=b^0} \quad \underline{b^1} \leftarrow b^0 - \eta \frac{\partial L}{\partial b} \big|_{w=w^0, \underline{b=b^0}} \quad \text{更新 } w^1, b^1$$

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

$$\underline{w^2} \leftarrow w^1 - \eta \frac{\partial L}{\partial w} \big|_{\underline{w=w^1}, b=b^1} \quad \underline{b^2} \leftarrow b^1 - \eta \frac{\partial L}{\partial b} \big|_{w=w^1, \underline{b=b^1}} \quad \text{更新 } w^2, b^2$$



## Step 3: Gradient Descent



## Step 3: Gradient Descent

- When solving:

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

- Each time we update the parameters, we obtain  $\theta$  that makes  $L(\theta)$  smaller.

$$L(\theta^0) > L(\theta^1) > L(\theta^2) > \dots$$

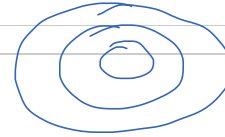
Is this statement correct?

In Linear regression the loss function is Convex  $\rightarrow$  无 local optimal  
线性回归 凸函数 无局部最优解



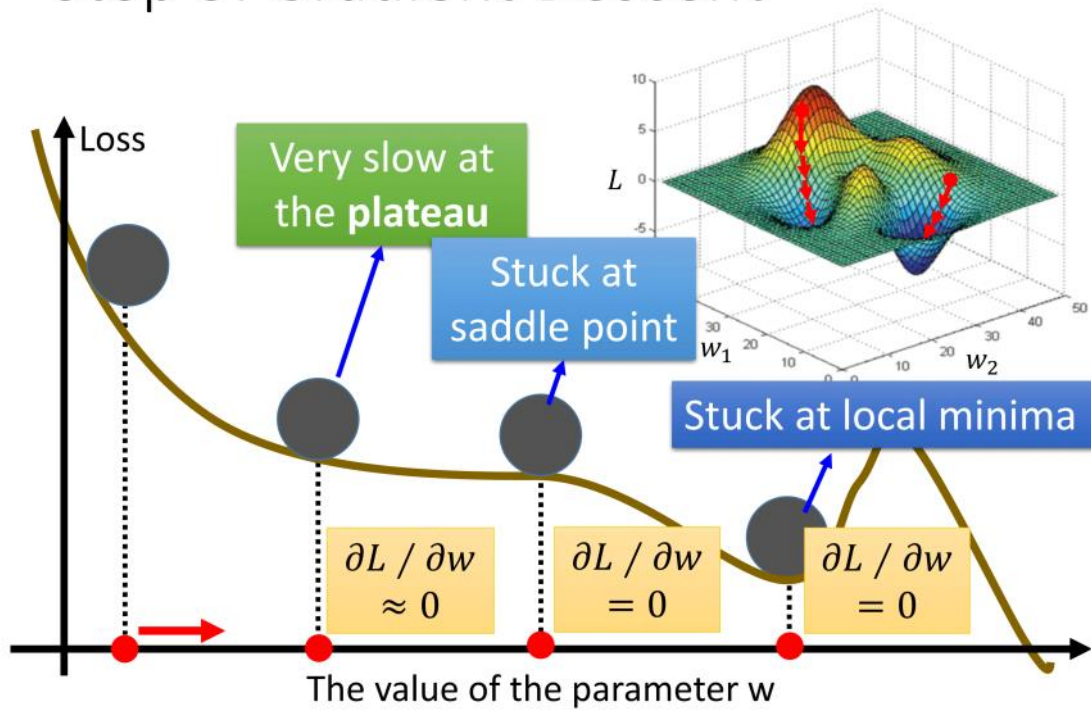
等高线是一圈一圈的





等高线是一圈  
一圈的

## Step 3: Gradient Descent



## Step 3: Gradient Descent

- Formulation of  $\frac{\partial L}{\partial w}$  and  $\frac{\partial L}{\partial b}$

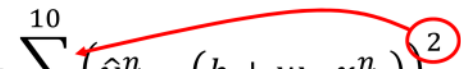
$$L(w, b) = \sum_{n=1}^{10} \left( \hat{y}^n - (b + \underline{w \cdot x_{cp}^n}) \right)^2$$

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

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

## Step 3: Gradient Descent

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

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


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

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

## Step 3: Gradient Descent

## How's the results?

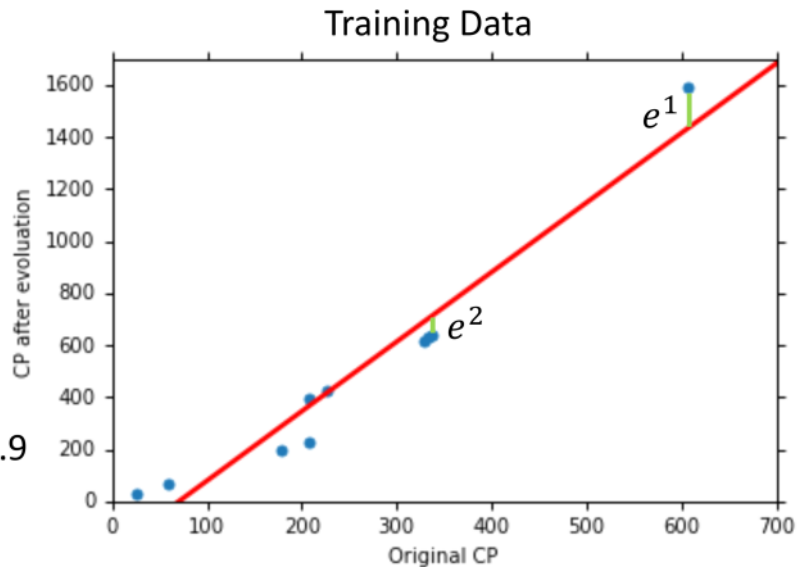
$$y = b + w \cdot x_{cp}$$

$$b = -188.4$$

$$w = 2.7$$

Average Error on Training Data

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



## How's the results? - Generalization

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

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

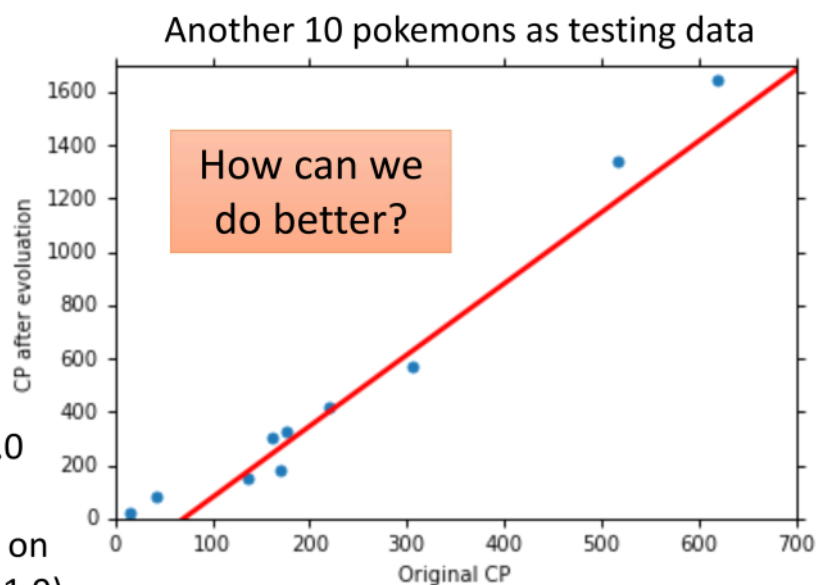
$$b = -188.4$$

$$w = 2.7$$

Average Error on Testing Data

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

> Average Error on Training Data (31.9)



## Selecting another Model

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

### Best Function

$$b = -10.3$$

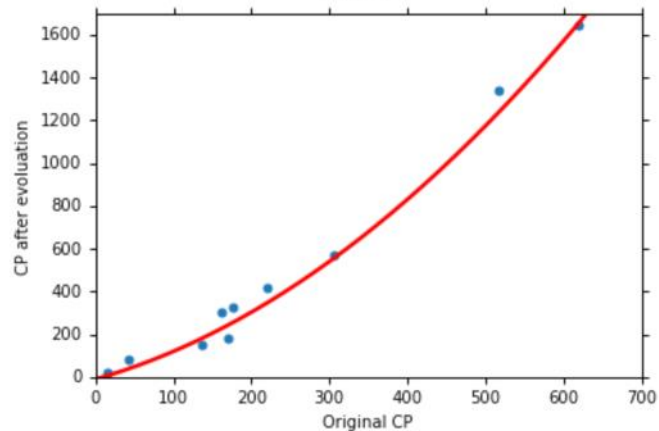
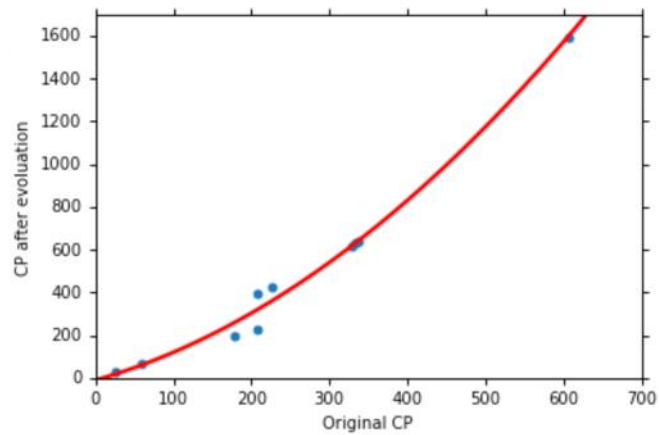
$$w_1 = 1.0, w_2 = 2.7 \times 10^{-3}$$

$$\text{Average Error} = 15.4$$

### Testing:

$$\text{Average Error} = 18.4$$

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

$$b = 6.4, w_1 = 0.66$$

$$w_2 = 4.3 \times 10^{-3}$$

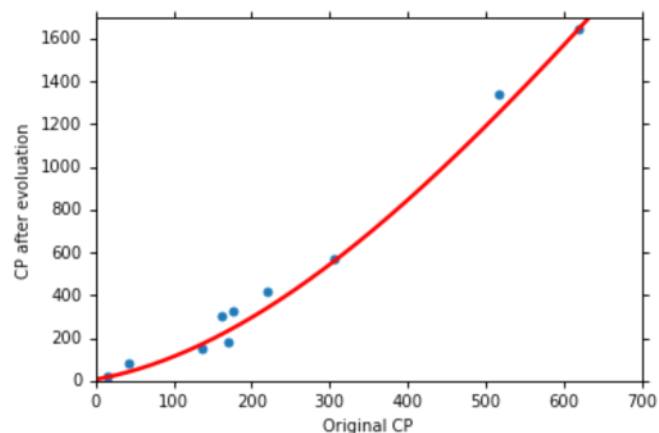
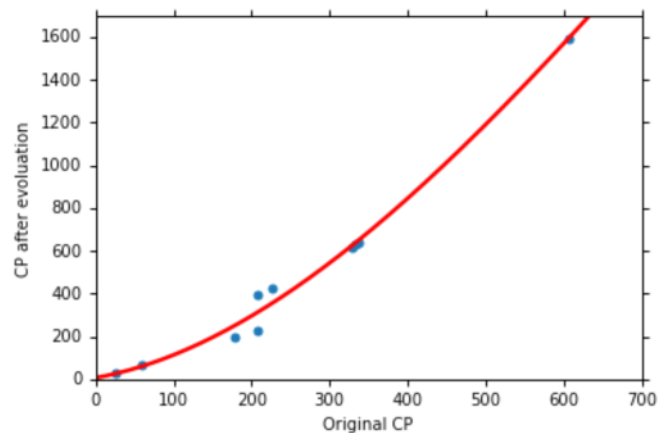
$$w_3 = -1.8 \times 10^{-6}$$

$$\text{Average Error} = 15.3$$

### Testing:

$$\text{Average Error} = 18.1$$

Slightly better.  
How about more complex model?



## Selecting another Model

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

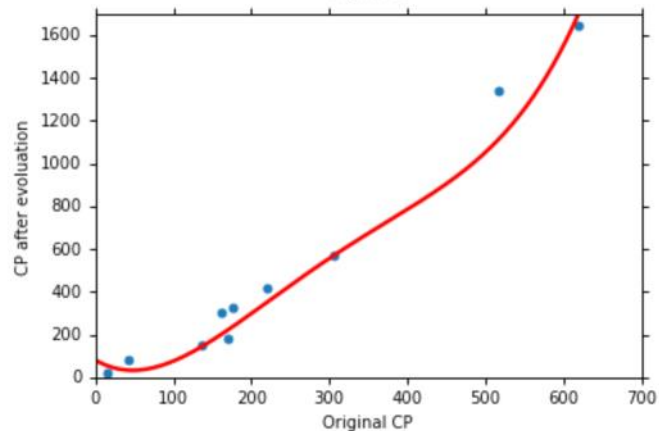
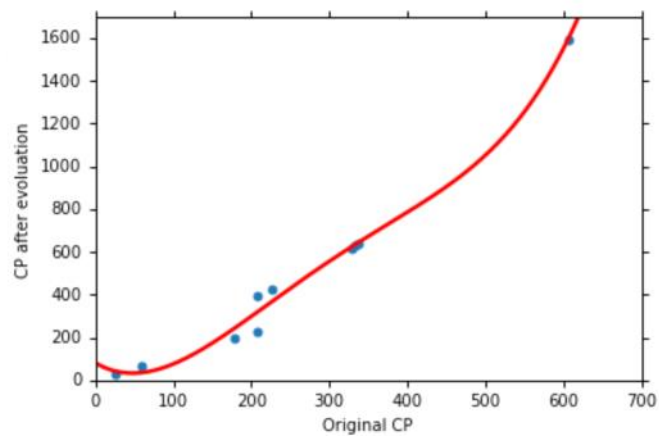
## Best Function

Average Error = 14.9

## Testing:

Average Error = 28.8

The results become worse ...



## Selecting another Model

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

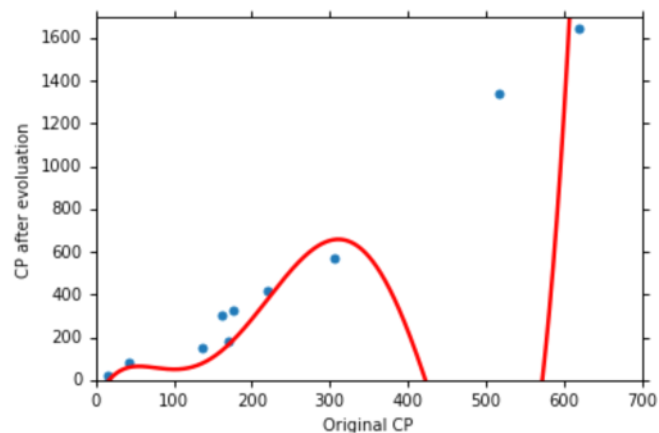
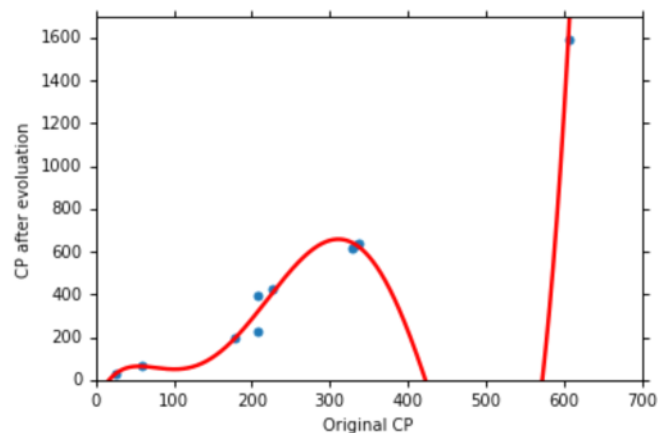
## Best Function

Average Error = 12.8

## Testing:

Average Error = 232.1

The results are so bad.



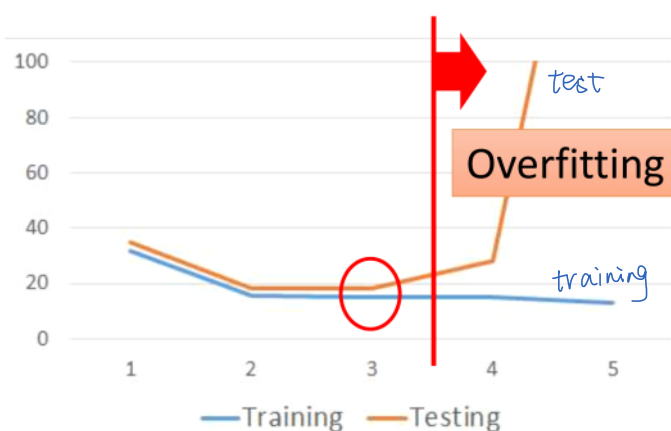
# Model Selection

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. *越复杂的 model*  
 If we can truly find the best function. *function 越复杂*

# 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

*model*

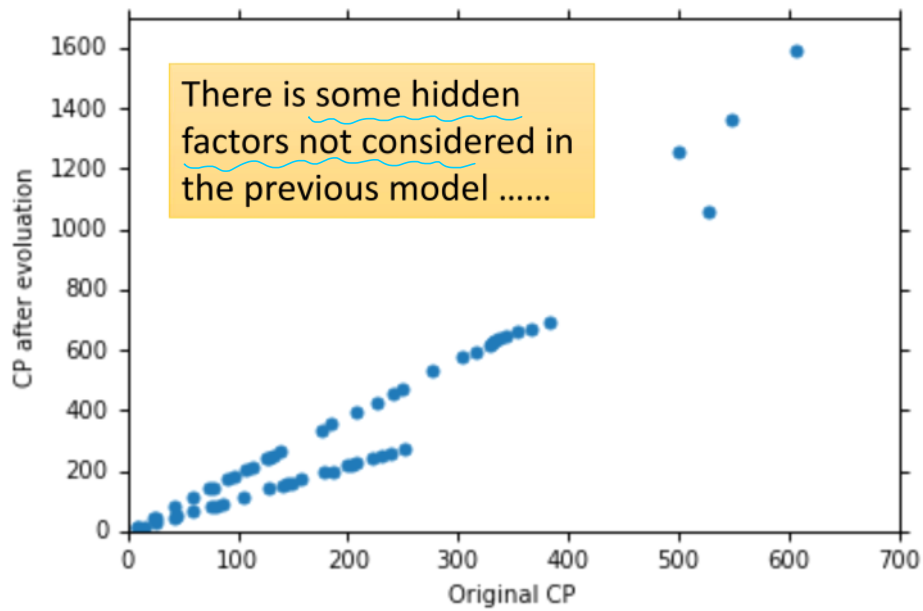
A more complex model does not always lead to better performance on **testing data**.

This is **Overfitting**. *并不* Select suitable model *合适的 model*

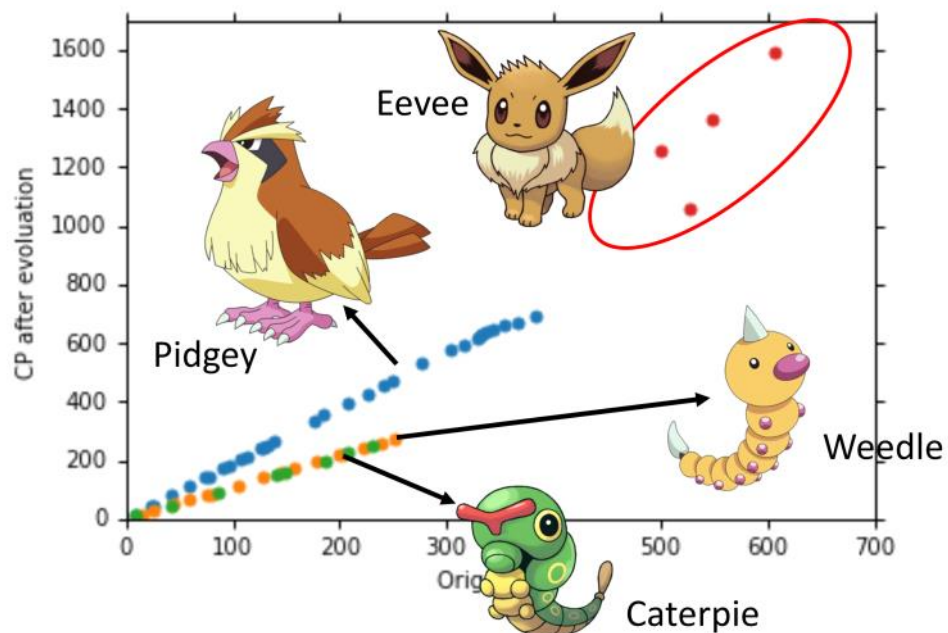
*过拟合 → 模型越复杂, 但 error 不是最优*  
*test 的 error 很小, 但 training 的 error 很大*

test的error很少  
但 training的error较大

# Let's collect more data



## What are the hidden factors?





# Back to step 1: Redesign the Model

重新设计 function

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

Linear model?

$x_s$  = species of x 4个物种



If $x_s$ = Pidgley:	$y = b_1 + w_1 \cdot x_{cp}$
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:	$y = b_4 + w_4 \cdot x_{cp}$

不同的物种有不同的y



## Back to step 1: Redesign the Model

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

Linear model?

$b_1 \ b_2 \ b_3 \ b_4$   
 $w_1 \ w_2 \ w_3 \ w_4$

$$y = b_1 \cdot \delta(x_s = \text{Pidgley})$$

$$+ w_1 \cdot \delta(x_s = \text{Pidgley}) x_{cp}$$

$$+ b_2 \cdot \delta(x_s = \text{Weedle})$$

$$+ w_2 \cdot \delta(x_s = \text{Weedle}) x_{cp}$$

$$+ b_3 \cdot \delta(x_s = \text{Caterpie})$$

$$+ w_3 \cdot \delta(x_s = \text{Caterpie}) x_{cp}$$

$$+ b_4 \cdot \delta(x_s = \text{Eevee})$$

$$+ w_4 \cdot \delta(x_s = \text{Eevee}) x_{cp}$$

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

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

# Back to step 1: Redesign the Model

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

Linear model?

$$y = b_1 \cdot \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} + w_1 x_{cp}$$

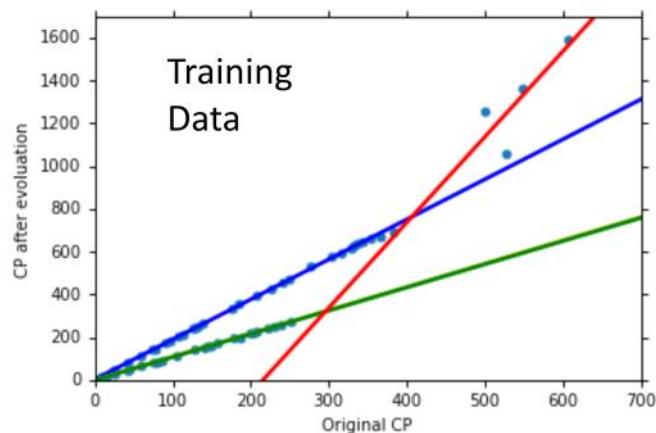
$$\delta(x_s = \text{Pidgey})$$

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

If  $x_s = \text{Pidgey}$

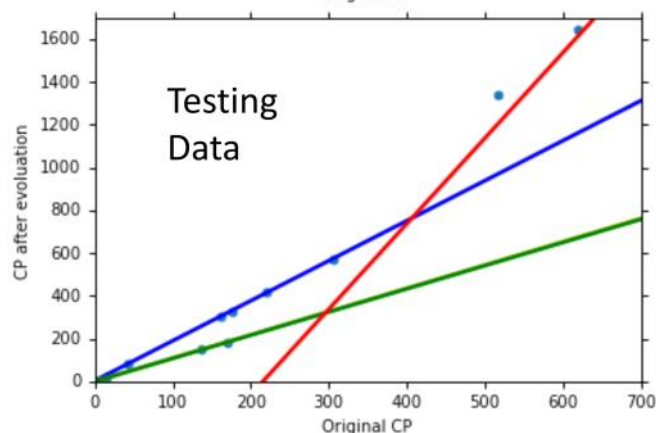
$$y = b_1 + w_1 \cdot x_{cp}$$

Average error  
= 3.8



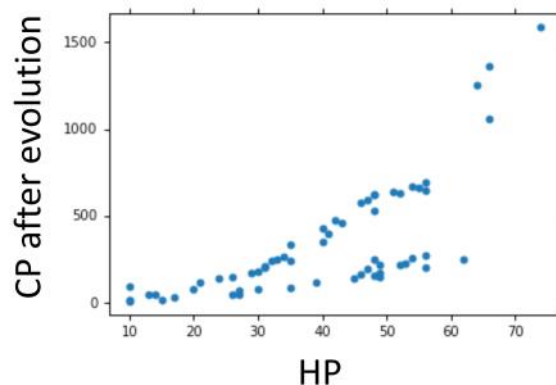
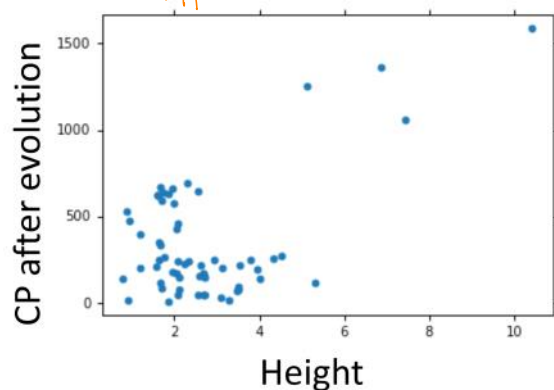
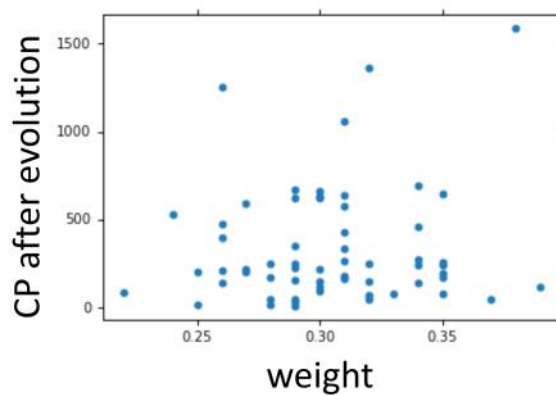
不同种类  
有不同的 model

Average error  
= 14.3



Are there any other hidden factors?

$y \sim$  weight  
height  
HP } 是否有关



Back to step 1:  
Redesign the Model Again

将模型变复杂一点

不是



If  $x_s = \text{Pidgey}$ :

$$y' = b_1 + w_1 \cdot x_{cp} + w_5 \cdot (x_{cp})^2$$

If  $x_s = \text{Weedle}$ :

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

If  $x_s = \text{Caterpie}$ :

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

If  $x_s = \text{Eevee}$ :

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

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

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



y

Training Error  
= 1.9

Testing Error  
= 102.3

Overfitting!

train error << test error

## Back to step 2: Regularization

model  $y = b + \sum w_i x_i$

The functions with smaller  $w_i$  are better

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

Regularization 正则化

➤ Smaller  $w_i$  means ... smoother

参数越接近 0 → 平滑的 function

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

输入变化 → Output 变化较小

对 input 的变化不敏感

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

输入变化

➤ We believe smoother function is more likely to be correct

为什么喜欢平滑的 function?

Do you have to apply regularization on bias?

## Back to step 2: Regularization

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

The functions with smaller  $w_i$  are better

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

➤ Why smooth functions are preferred?

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

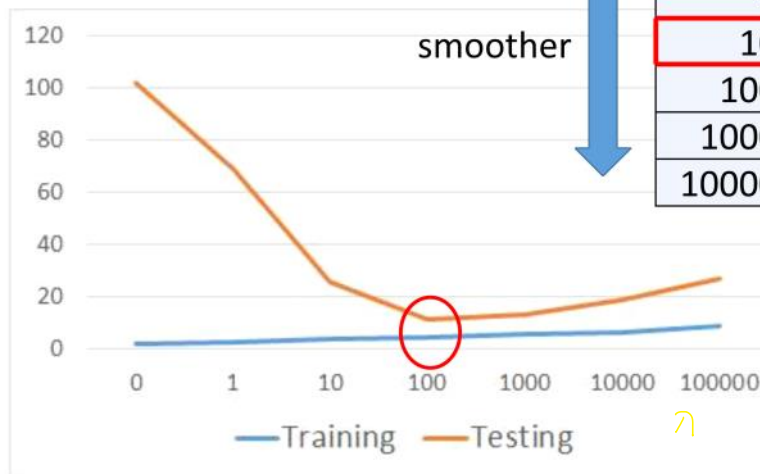
$+w_i \Delta x_i$   $+ \Delta x_i$

➤ If some noises corrupt input  $x_i$  when testing

噪音

A smoother function has less influence.

## 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

How smooth?

Select  $\lambda$  obtaining the best model

- ① Training error: larger  $\lambda$ , considering the training error less  
 $\rightarrow$  更倾向于考虑  $w$  本来值, 减少考虑 error
- ③ 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  
 $\rightarrow$  验证, 校验

## 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