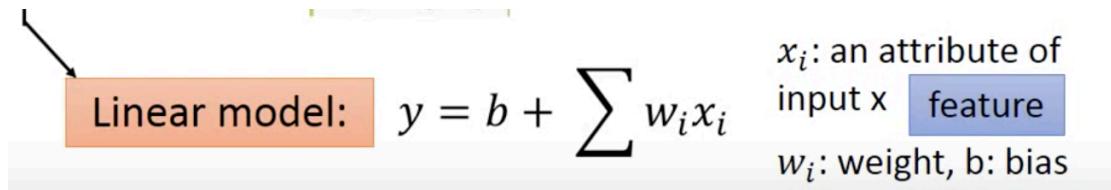
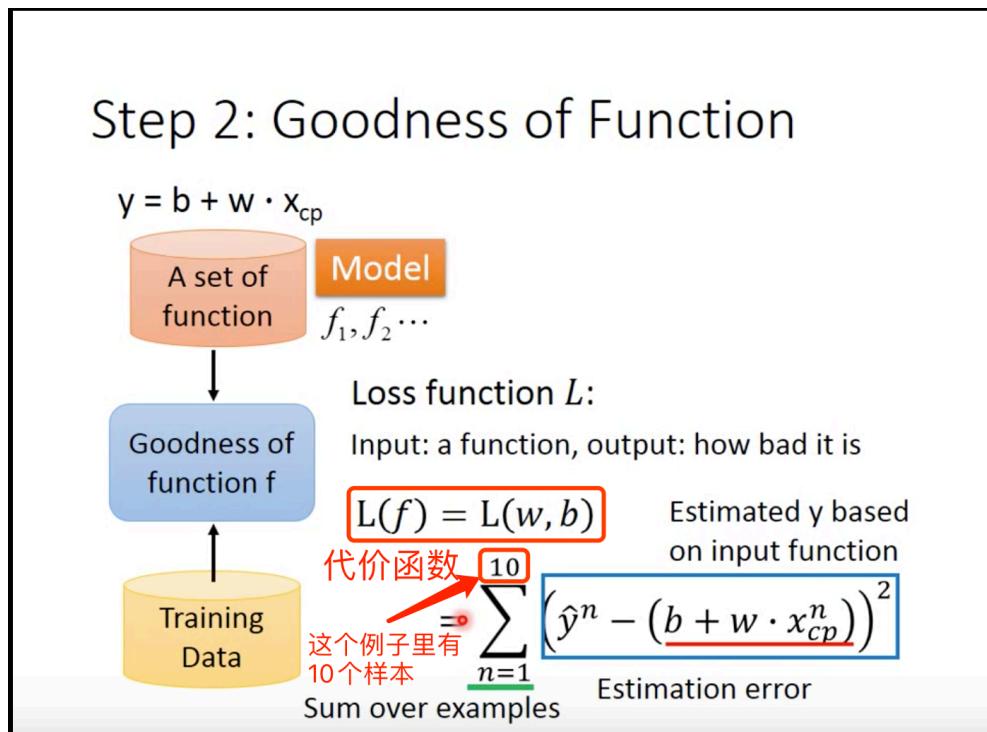


线性模型



x_i : 样本多维属性, w_i : 样本该维度的系数, b : 偏移, 其实就是求 w_i

代价函数



把每一组预测系数 w_i 和偏移 b 带入代价函数, 算误差 (这里取平方, 是一种简单的误差算法)

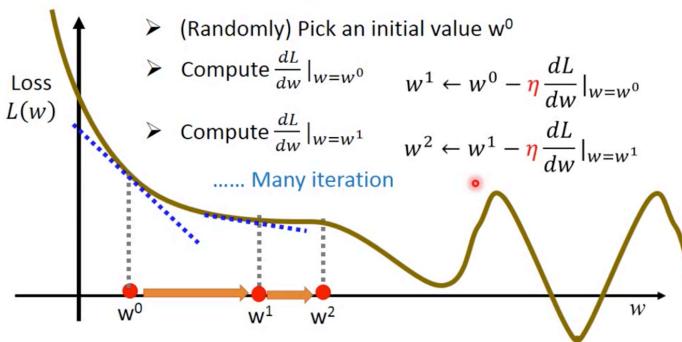
GradientDescent 梯度下降

最有效率的穷举 w_i 的方法

Step 3: Gradient Descent

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

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



先说1个参数w的例子，(不是上面的例子，上面的例子还有b) (代价函数图像L与w的关系)

现在其实就是找图像最低点

1. 初始w⁰位置，求该点导数 (参考教材，偏导)

2. 偏导>0说明左侧下降，右侧上升

3. w⁰该往左偏移 (使代价最小)

4. 增量为 导数 dL/dw 值 * 学习率 η (自定义)

$$w^+ = -\left(\frac{\partial L}{\partial w} \times \eta\right)$$

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

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

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

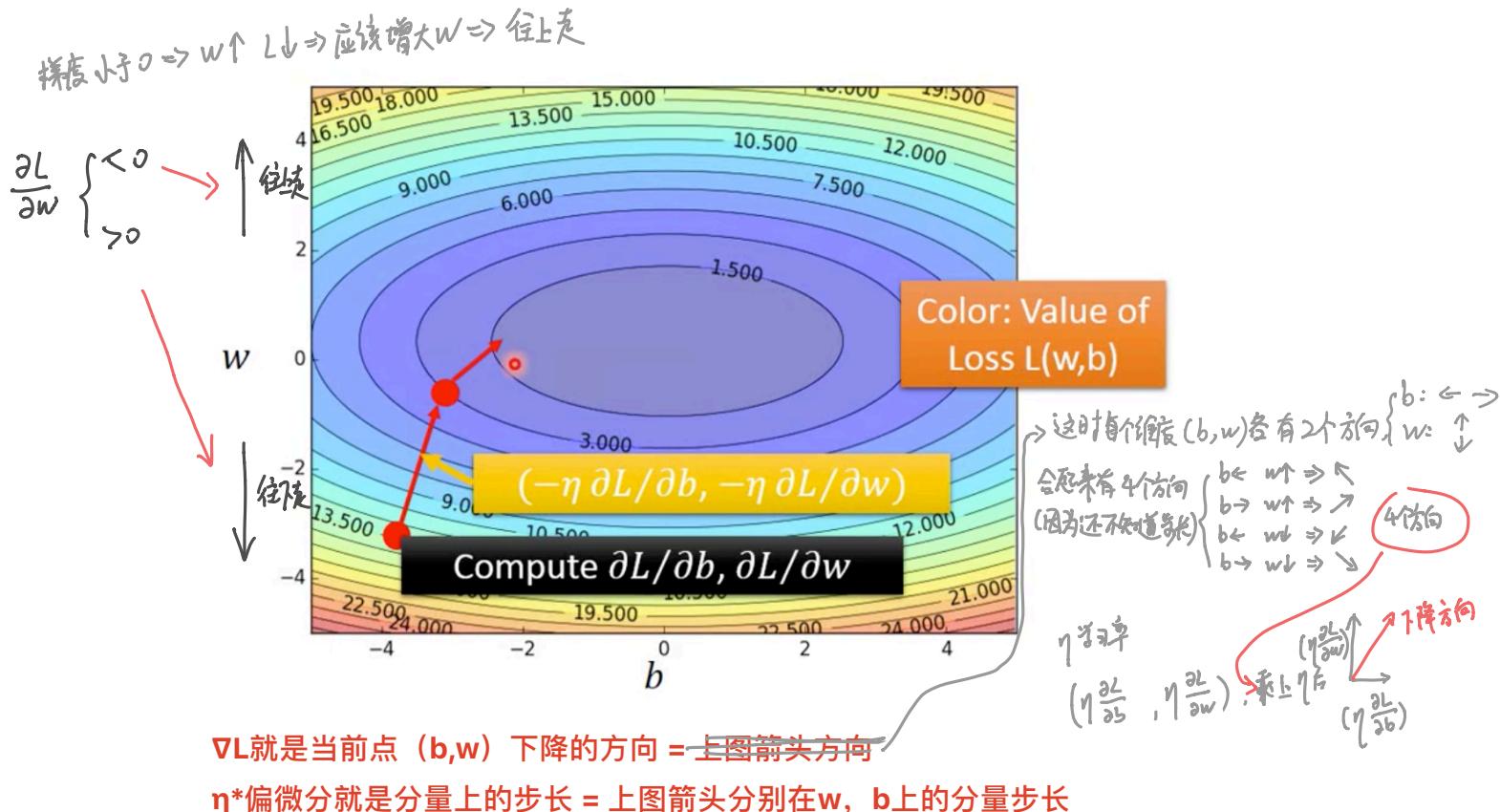
$$w^2 \leftarrow w^1 - \eta \frac{\partial L}{\partial w} |_{w=w^0, b=b^0} \quad b^2 \leftarrow b^1 - \eta \frac{\partial L}{\partial b} |_{w=w^0, b=b^0}$$

2个参数w和b的例子 (与上个例子类似)

1. 初始w⁰与b⁰，各自求偏导

2. 根据各自偏导，各自偏移

∇L 就是 gradient, 他是一个向量，表示各个维度上的斜率（梯度）



偏微分计算

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

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

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

发现一组 w 和 b (直线) 不能非常好的拟合样本以及预测

均方误差代价函数

$$J(\Theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

↑ 预测函数
↑ 方程估计梯度
↑ 时被 2 约去

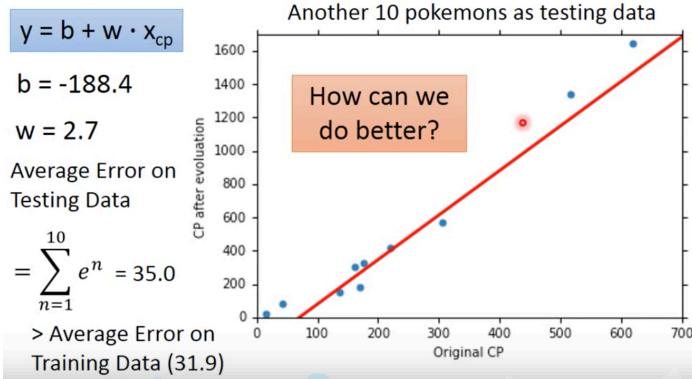
$$\nabla J(\Theta) = \left\langle \frac{\delta J}{\delta \Theta_0}, \frac{\delta J}{\delta \Theta_1} \right\rangle$$

$$\frac{\delta J}{\delta \Theta_0} = \frac{1}{m} \sum_{i=1}^m (h_\Theta(x^{(i)}) - y^{(i)})$$

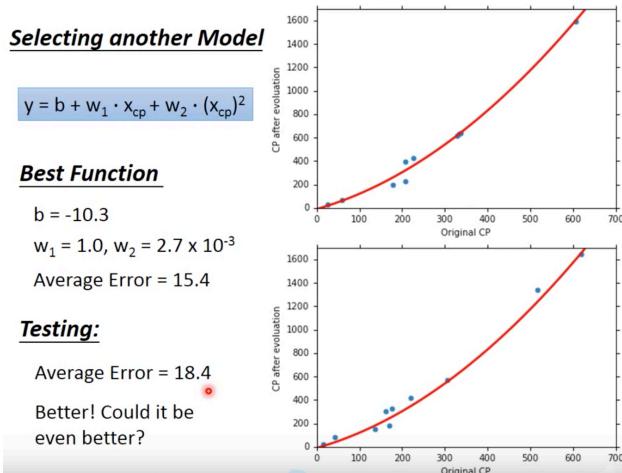
$$\frac{\delta J}{\delta \Theta_1} = \frac{1}{m} \sum_{i=1}^m (h_\Theta(x^{(i)}) - y^{(i)}) x_1^{(i)}$$

How's the results?
- Generalization

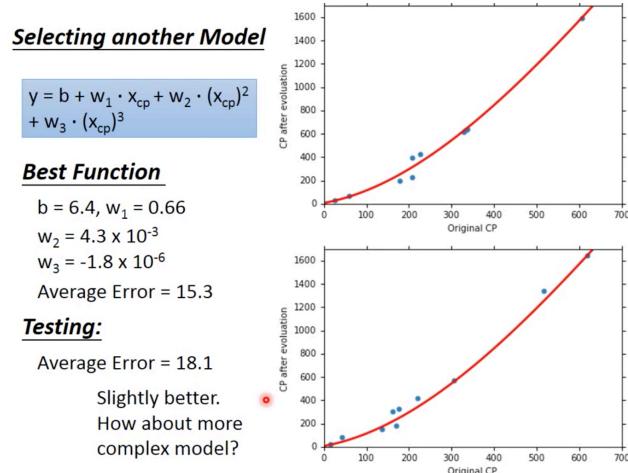
What we really care
about is the error on
new data (testing data)



引入二次项来更好的拟合，预测上好了很多18.4



引入三次项，预测上稍微又好了一点18.1



引入四次项，样本拟合不断变好，可是预测变差了28.8

Selecting another Model

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

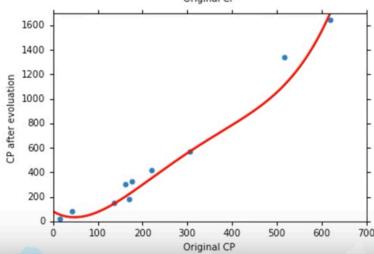
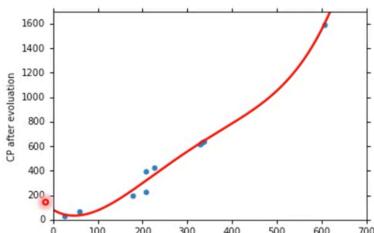
Best Function

Average Error = 14.9

Testing:

Average Error = 28.8

The results become worse ...



引入五次项，炸了

Selecting another Model

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

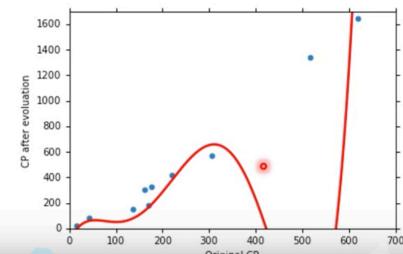
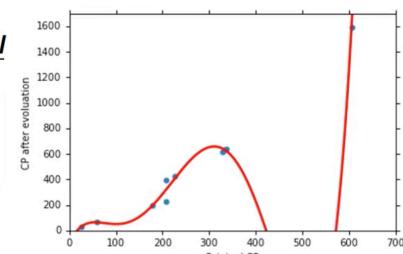
Best Function

Average Error = 12.8

Testing:

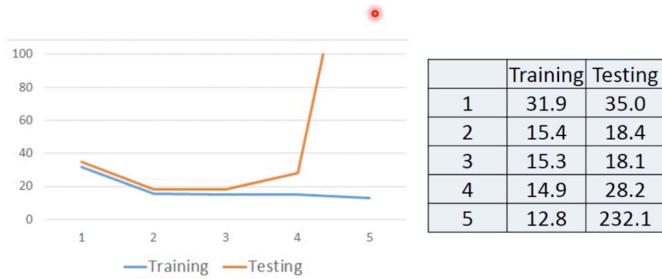
Average Error = 232.1

The results are so bad.



过拟合（越复杂的模型不一定能很好的预测）

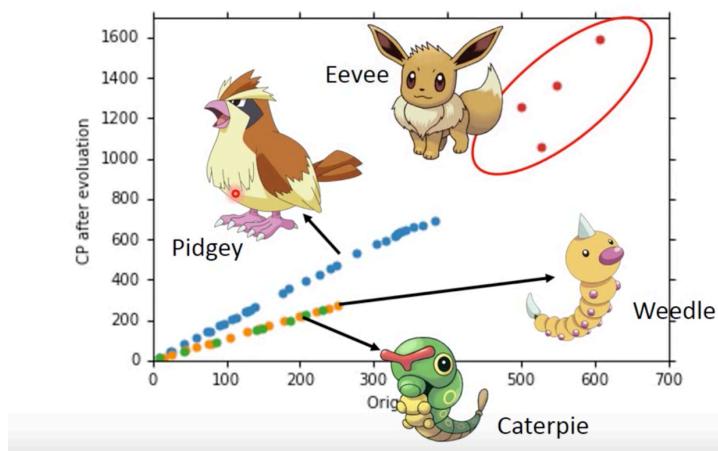
Model Selection



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

修改模型，寻找隐藏因素

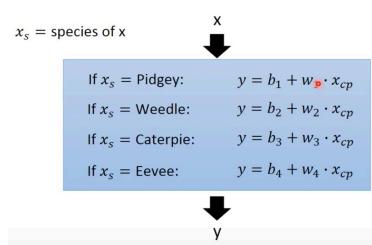
What are the hidden factors?



这个例子中，隐藏因素为物种

用if else 分物种分类处理

Back to step 1:
Redesign the Model



转成数学表达式，就是让那一项系数变1，不需要的项系数变0

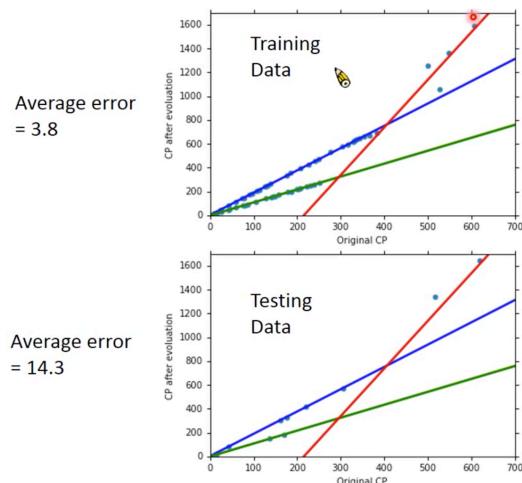
Back to step 1:
Redesign the Model

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

Linear model?

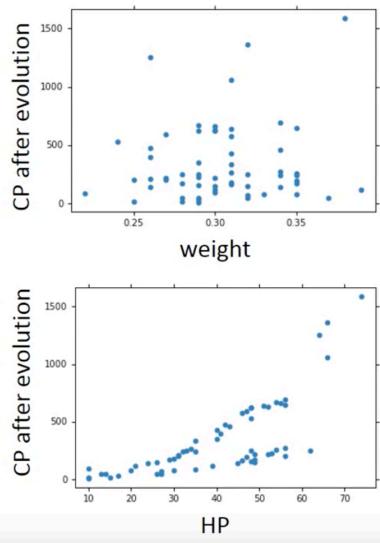
$$\begin{aligned} y &= b_1 \cdot 1 \\ &+ w_1 \cdot 1 \quad \bullet \quad \delta(x_s = \text{Pidgey}) \\ &+ b_2 \cdot 0 \\ &+ w_2 \cdot 0 \quad \left\{ \begin{array}{ll} =1 & \text{If } x_s = \text{Pidgey} \\ =0 & \text{otherwise} \end{array} \right. \\ &+ b_3 \cdot 0 \\ &+ w_3 \cdot 0 \quad \text{If } x_s = \text{Pidgey} \\ &+ b_4 \cdot \delta(x_s = \text{Eevee}) \\ &+ w_4 \cdot \delta(x_s = \text{Eevee}) x_{cp} \end{aligned}$$

每种物种一根预测线



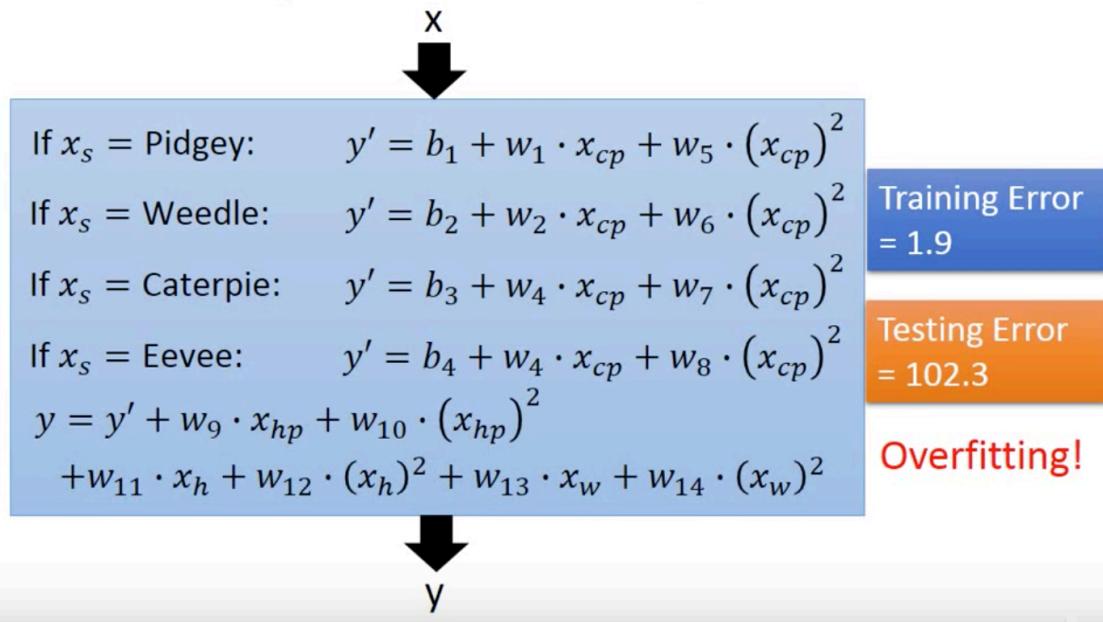
到底预测值与 哪个因素有关？

Are there any other hidden factors?



寻找一个足够复杂的模型，包括所有因素

Back to step 1:
Redesign the Model Again



使代价函数更加平滑

Back to step 2: Regularization

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

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

The functions with smaller w_i are better

$$+ \lambda \sum (w_i)^2$$

- Why smooth functions are preferred?

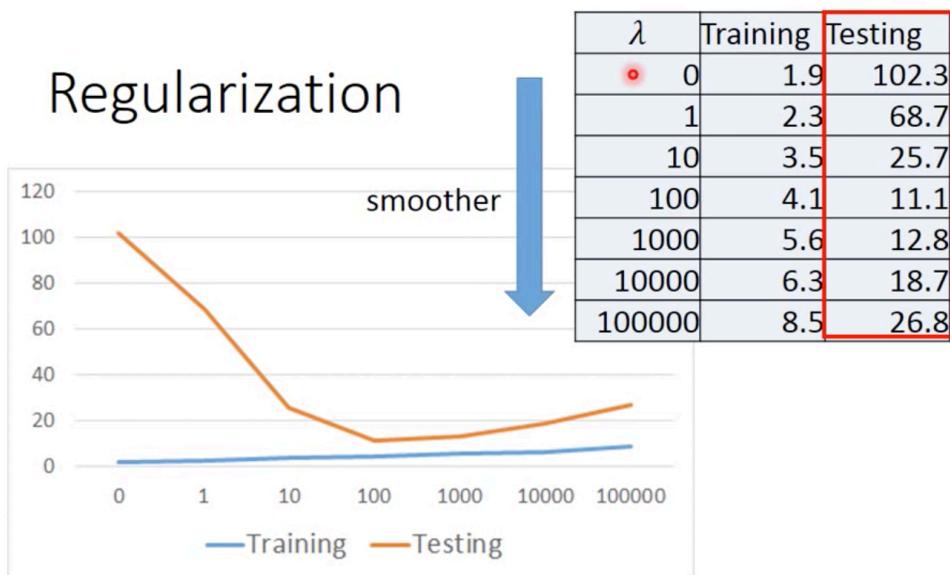
$$y = b + \sum w_i x_i + w_i \Delta x_i + \Delta x_i$$

- If some noises corrupt input x_i when testing

A smoother function has less influence.

可以看到，不加入红框项时，此时 y 随 Δx 变化= $w_i * \Delta x$ ， w_i 越大=> y 变化越大

平滑后，训练样本误差变大，预测误差变小



- Training error: larger λ , considering the training error less

