

PYTHON機器學習入門

UNIT 4 : REGRESSION

授課教師：江尚瑀

A decorative wavy line in a light green color runs vertically along the left side of the slide, separating a light beige area from a dark teal background.

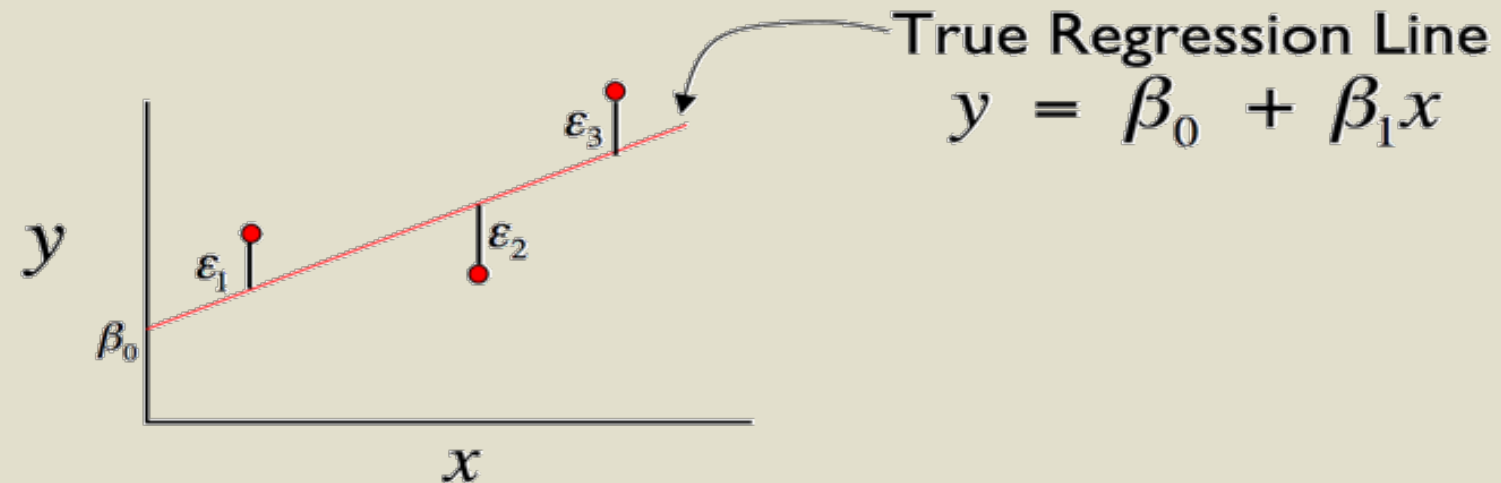
SUPERVISED LEARNING- LINEAR REGRESSION

REGRESSION 迴歸

- Definition: There exists parameters β_0 , β_1 , and σ^2 , such that for any fixed value of the independent variable x , the dependent variable is related to x through the model equation

$$y = \beta_0 + \beta_1 x + \varepsilon$$

- ε is a rv assumed to be $N(0, \sigma^2)$



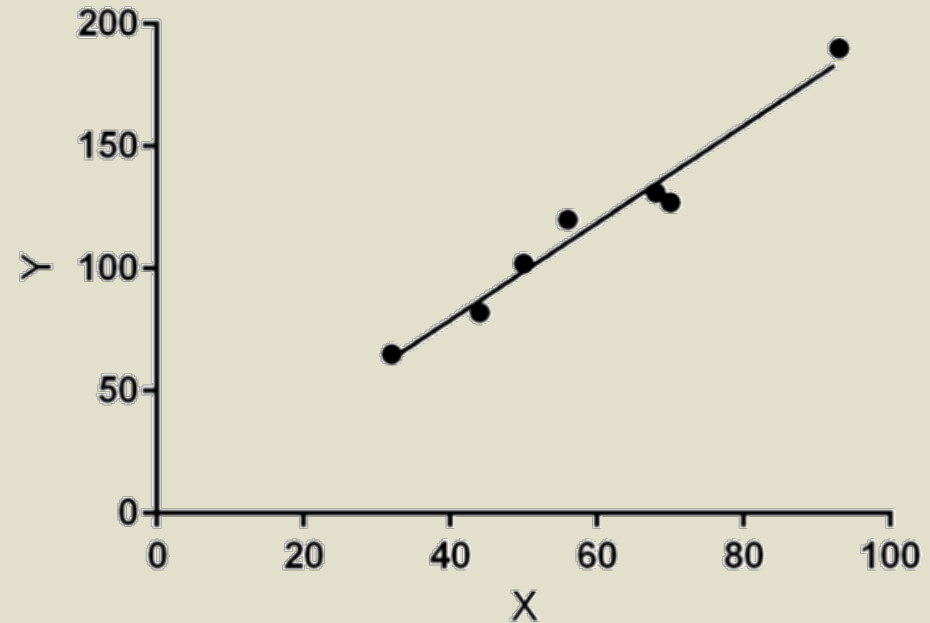
SINGLE FEATURE REGRESSION

- Single feature problem

$$\hat{Y} = \beta_0 + \beta_1 X$$

House Size (X)	House Price (Y)
50	102
70	127
32	65
68	131
93	190
44	82
56	120

house size is the single feature -> X
house price is the label -> Y



$$\hat{Y} = -0.5243 + 1.987X$$

$$\hat{Y}(85) = 168.37$$

MULTIREGRESSION

- Multi feature problem $\hat{Y} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3$

House Size (X_1)	Rooms (X_2)	Floor (X_3)	House Price (Y)
50	2	5	102
70	2	3	127
32	1	3	65
68	3	7	131
93	4	10	190
44	2	6	82
56	3	1	120

構建多元線性回歸模型時，隨著解釋變量數目的增多，其中某兩個解釋變量之間產生多重共線性的機會就會大增。此時就需要考慮是否將其中某個變量從模型中剔除出去，甚至是重新考慮模型的構建。

共線性: 若 X_1 和 X_2 非獨立變數, 則 β_1 也會影響 X_2

定義損失函數(LOSS FUNCTION)

好的**Regression Line** 就是誤差最小的那一條線

- 如何找出誤差最小的那條線？

$$\min L(y_i, f(x_i))$$

- $L(y_i, f(X_i))$ 叫做**Loss Function**或是**Cost Function**. **Loss**越小，就代表模型的配適性越好

不同問題或模型, 都必須明確定義出好的**Loss Function**

迴歸常用的LOSS FUNCTION

- MAE (Mean Absolute Error)

$$\frac{1}{m} \sum_{i=1}^m |y_i - \hat{y}_i|$$

- MSE (Mean Squared Error)

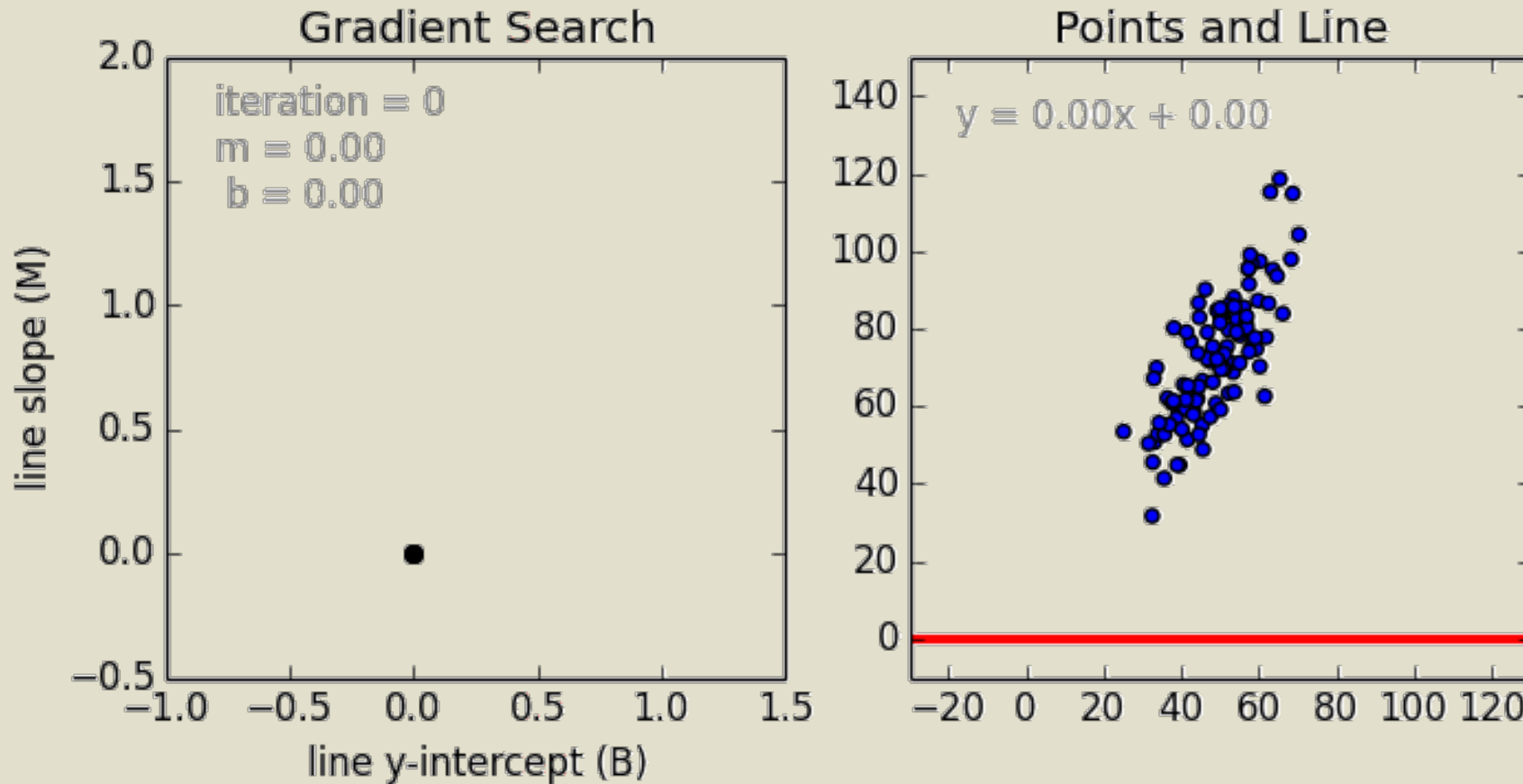
$$\frac{1}{m} \sum_{i=1}^m (y_i - \hat{y}_i)^2 \quad J(\theta) = 1/2m \sum_{i=1}^m (h(\theta)^{(i)} - y^{(i)})^2$$

- RMSE (Root Mean Squared Error)

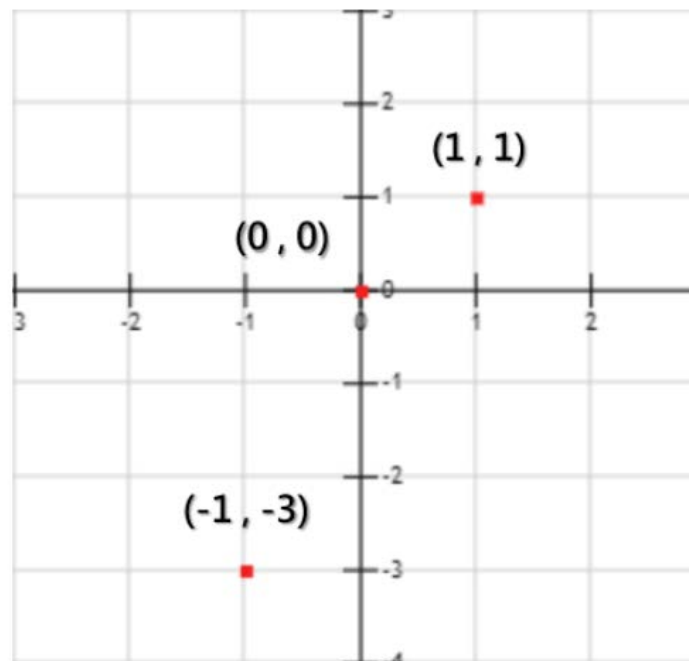
$$\sqrt{\frac{1}{m} \sum_{i=1}^m (y_i - \hat{y}_i)^2}$$

如何找出最佳參數使得**LOSS** 最小

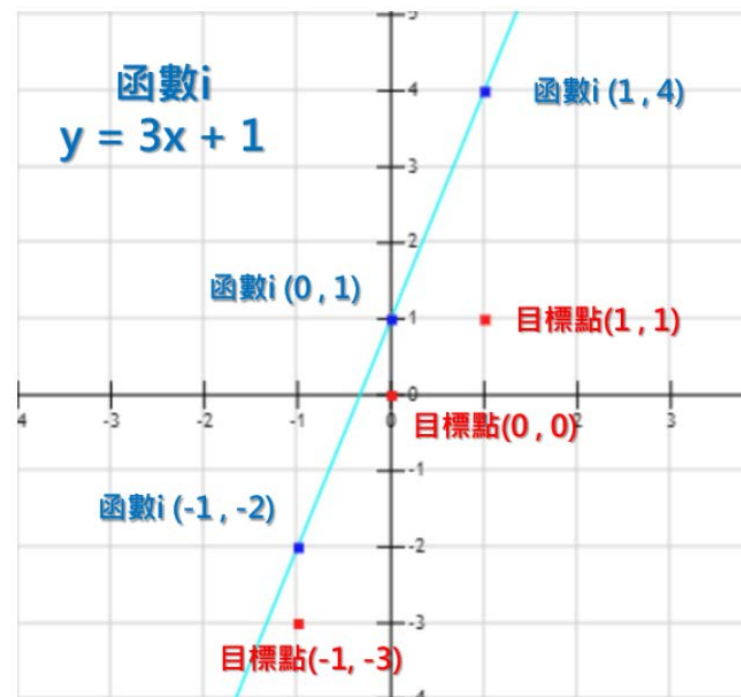
可以利用 **Gradient Descent** 方法找出最適參數 $B_0, B_1, B_2 \dots$, 使得**Loss** 最小



如何找出最佳參數使得**LOSS** 最小



三個目標點



目標點 X值	目標點 Y值	函數i: $Y = 3X + 1$, X使用目標點X值帶入得到	距離 (目標點Y值-函數i 輸出值) 取平方
-1	-3	-2	1
0	0	1	1
1	1	4	9
總距離: $1 + 1 + 9 = 11$			

如何找出最佳參數使得**LOSS** 最小

- 某目標點距離 = (某目標點的**Y**值 - 函數**i**對應某目標點**X**值的輸出值)²
- 總距離 = (每個目標點的**Y**值 - 函數**i**對應每個目標點**X**值的輸出值)²

The diagram illustrates the components of the loss function $L = (y - wx - b)^2$. The equation is written in large red characters. Annotations include:

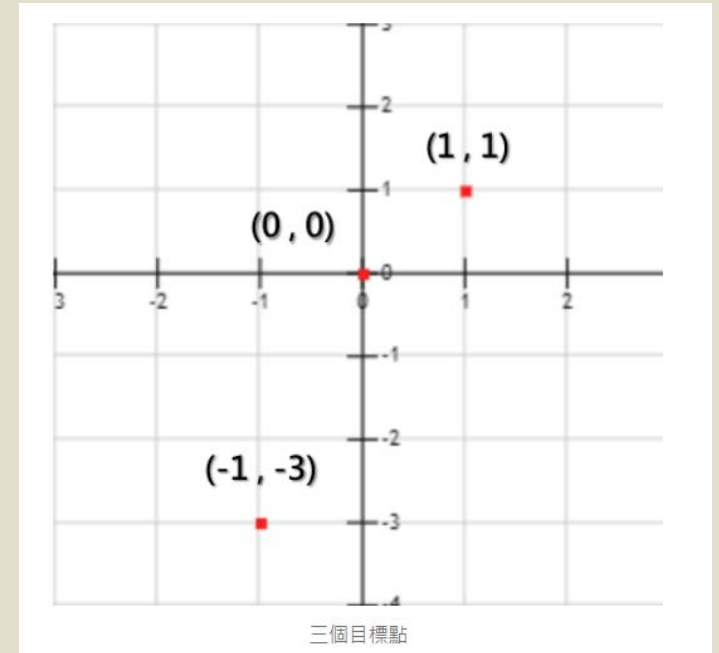
- A green box labeled "用3替換的" (Replaced by 3) pointing to the coefficient w .
- A green box labeled "用1替換的" (Replaced by 1) pointing to the bias term b .
- A blue box labeled "目標點的Y值" (Target point Y value) pointing to the variable y .
- A blue box labeled "目標點的X值" (Target point X value) pointing to the variable x .

結論**1**:損失函數**L**可以幫助我們找到接近目標的新函數

如何找出最佳參數使得**LOSS** 最小

- 損失函數**L**可以幫助我們找到接近目標的新函數

$$L = (y - wx - b)^2$$

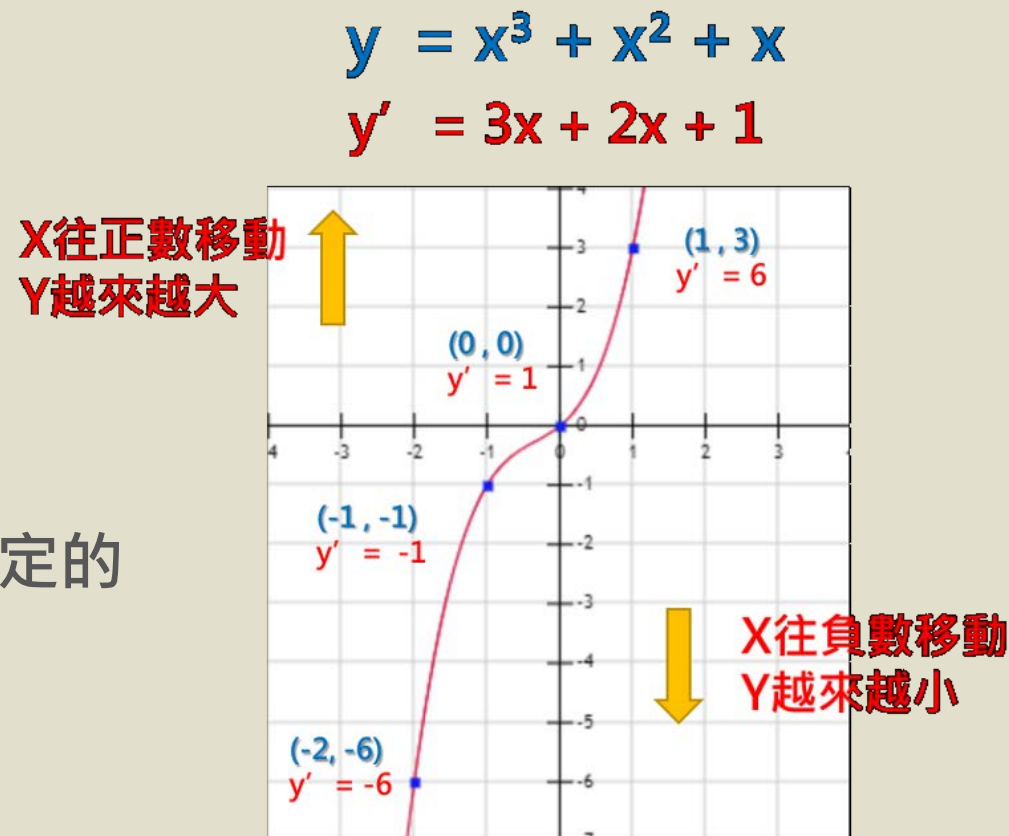


- 微分**某變數後產生的函數，可以指出**原函數在每點的變化**。

如何找出最佳參數使得**LOSS** 最小

- 微分某變數後產生的函數，可以指出原函數在每點的變化
 - 數值移動的方向的判定方法：
將**x**帶入微分後的函數**y'**
 - 我們會將帶入**x**後的**y'**產生的數值
稱之為梯度

結論**2**:數值移動的方向是可以從**梯度**判定的



如何找出最佳參數使得**LOSS** 最小

$$L = (y - wx - b)^2$$

- 假設 $x = -3$ $y = -1$ ，我們該怎麼找到 w 跟 b 呢？
 - 需要分別對 w 跟 b 做微分，來得知 w 跟 b 怎麼移動才會讓 $(-3 + w - b)^2$ 越來越小
- 對 w 取微分後得到 $(-3 + w - b)^2 = 9 - 6w + 6b - 2wb + w^2 + b^2$ ， $L'w = -6 - 2b + 2w$
- 對 b 取微分後得到 $L'b : 6 - 2w + 2b$

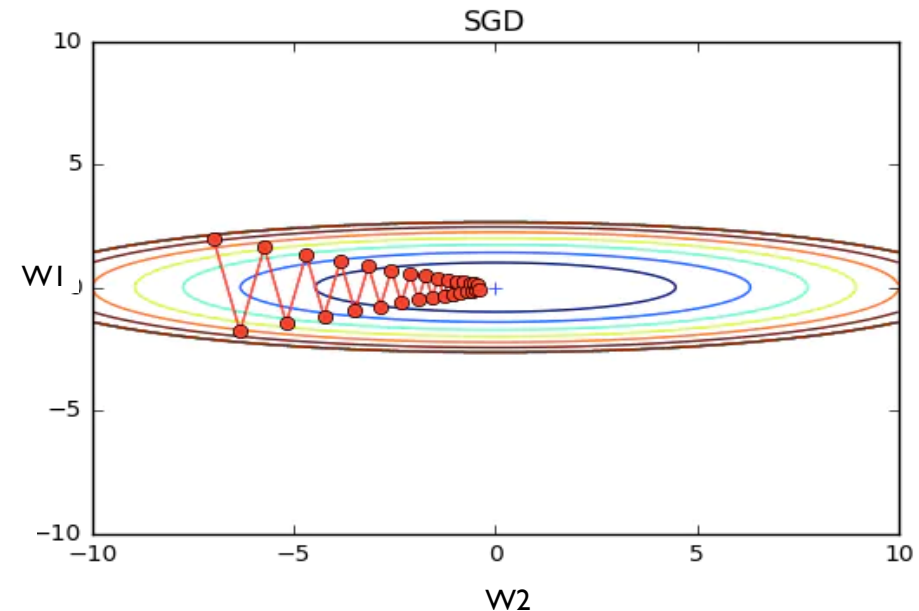
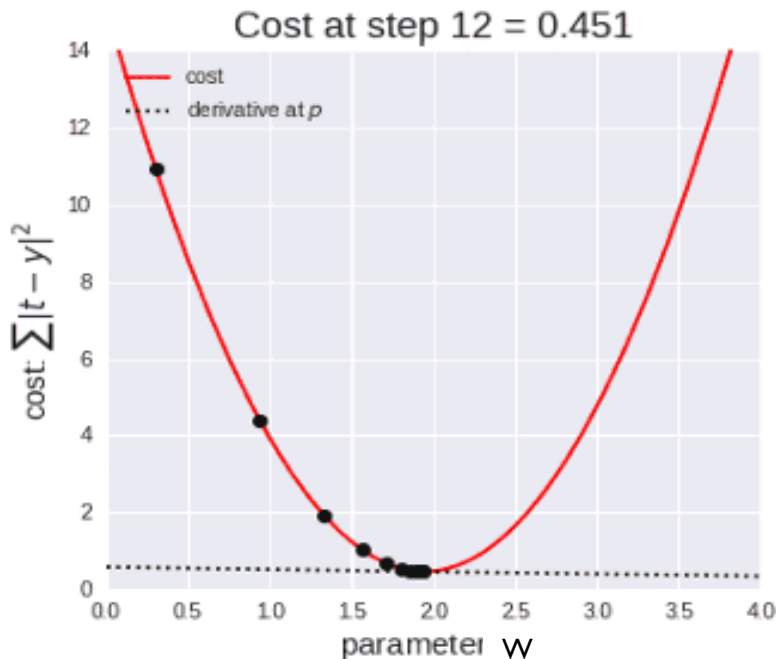
w 越往**負數**方向移動，損失函數會**越小**

b 越往**正數**方向移動，損失函數會**越大**

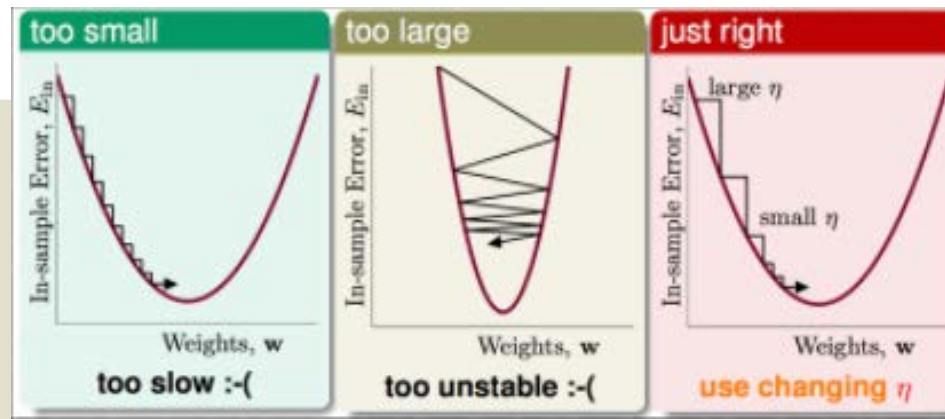
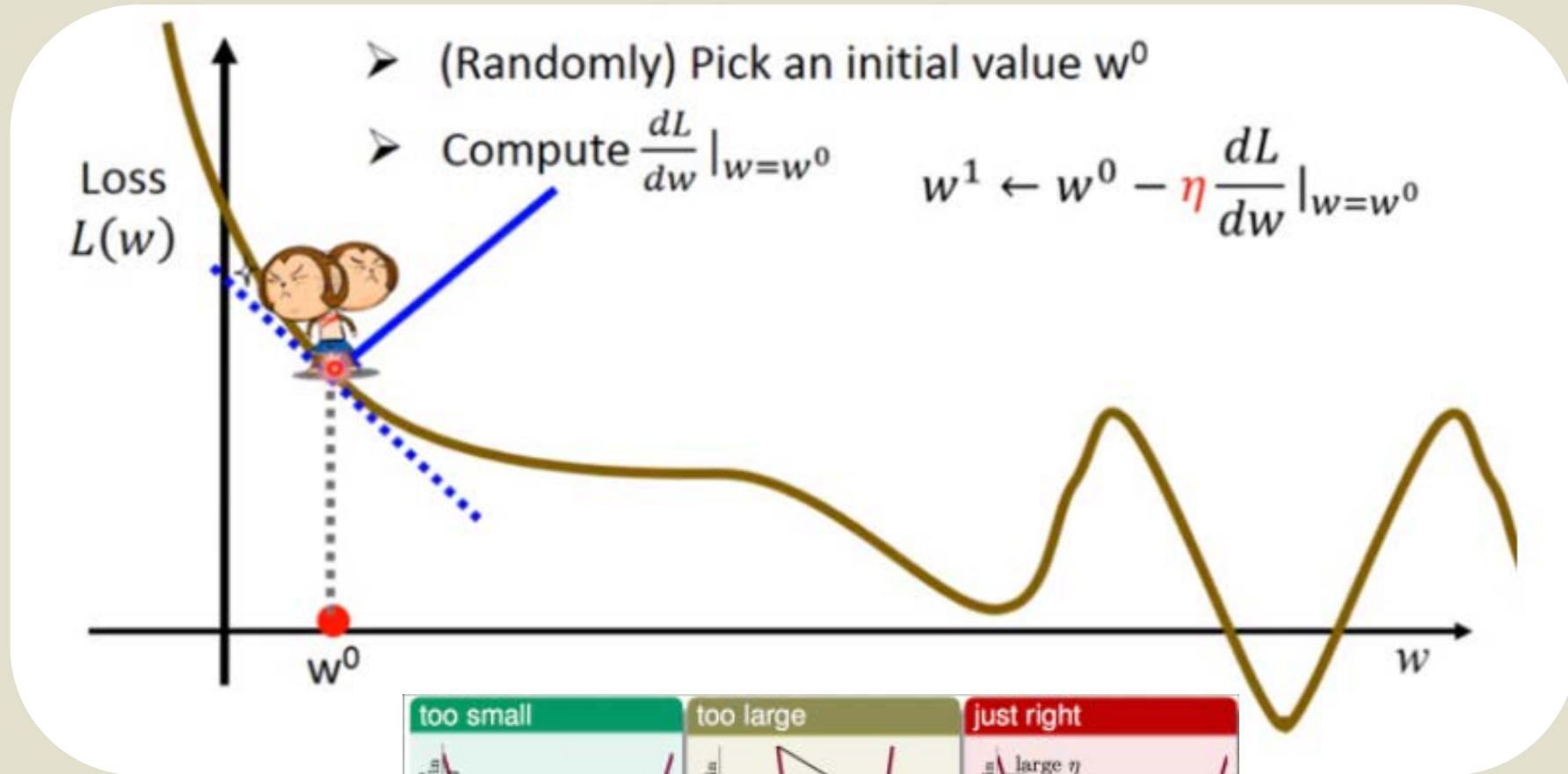
結論3：損失函數 L 移動的方向剛好跟梯度的方向**相反**

SGD(STOCHASTIC GRADIENT DESCENT)

- Choose an initial vector of parameters **W** and learning rate **η**
- Repeat until an approximate minimum is obtained:
 - Randomly shuffle samples in the training set.
 - For each sample do $W \leftarrow W - \eta \frac{\partial L}{\partial W}$

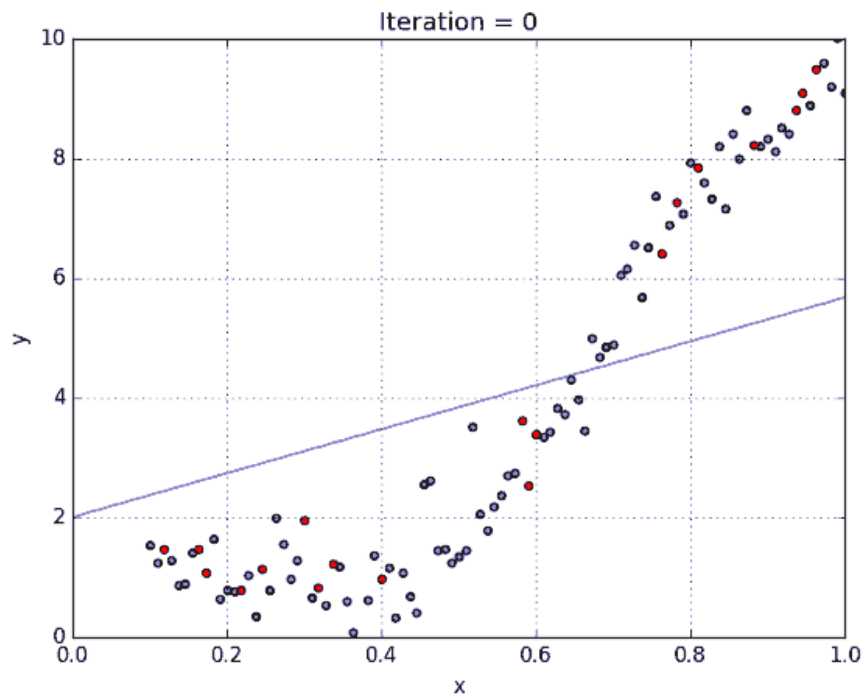


SGD(STOCHASTIC GRADIENT DESCENT)



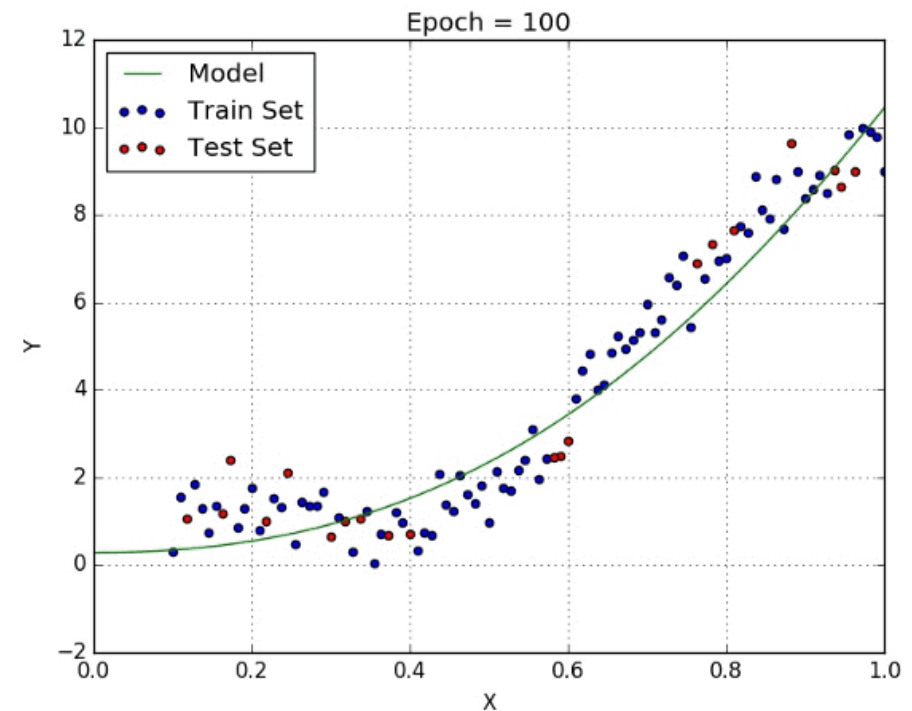
POLYNOMIAL LINEAR REGRESSION

Linear Regression



$$\hat{y} = xw + b = w_1x_1 + b$$

Polynomial regression



$$\hat{y} = xw + b = w_1x_1 + w_2x_1^2 + b$$

模型泛化程度

- 泛化(**generalize**)：指**ML**模型「對未知資料集」的預測能力。

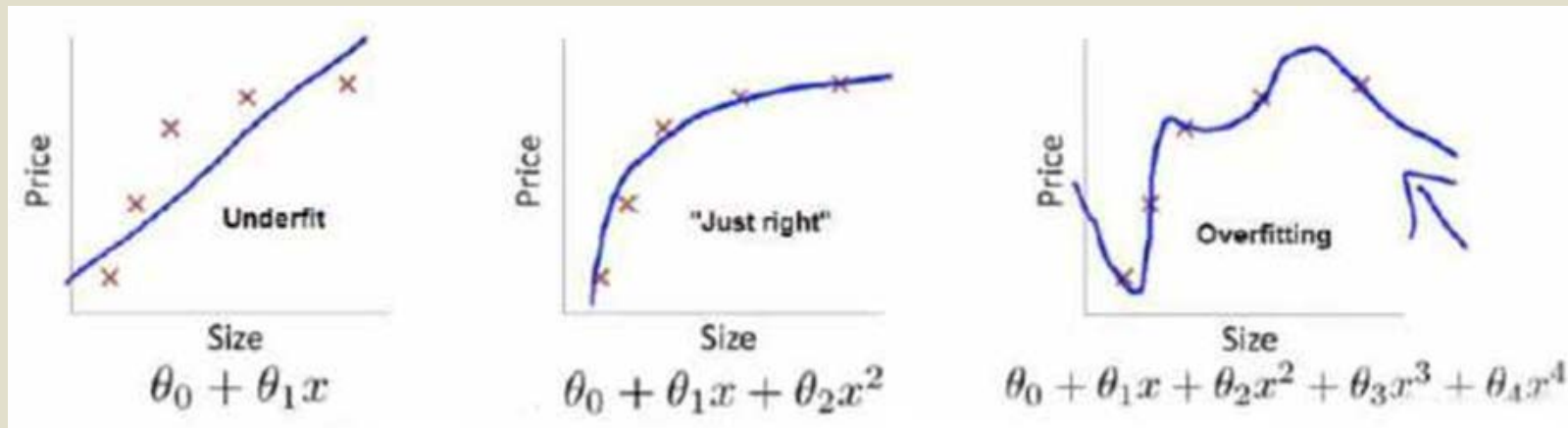
如果你的 **model** 在 **training dataset** 表現不錯，但是在 **testing dataset** 表現卻很差，那就是 **overfitting**

Overfitting

- 常常發生在 **model** 很複雜、有很多參數的時候或是 **dataset** 裡有很多 **noise** 或 **outlier**。使得表現在 **training set** 的準確率很高，但是在 **testing set** 的準確率卻很低。

Underfitting

- 通常發生在 **model** 太簡單的時候，其表現就算在 **training set** 上的錯誤率也很高。



OVERFITTING ? UNDERFITTING ?

- Underfitting

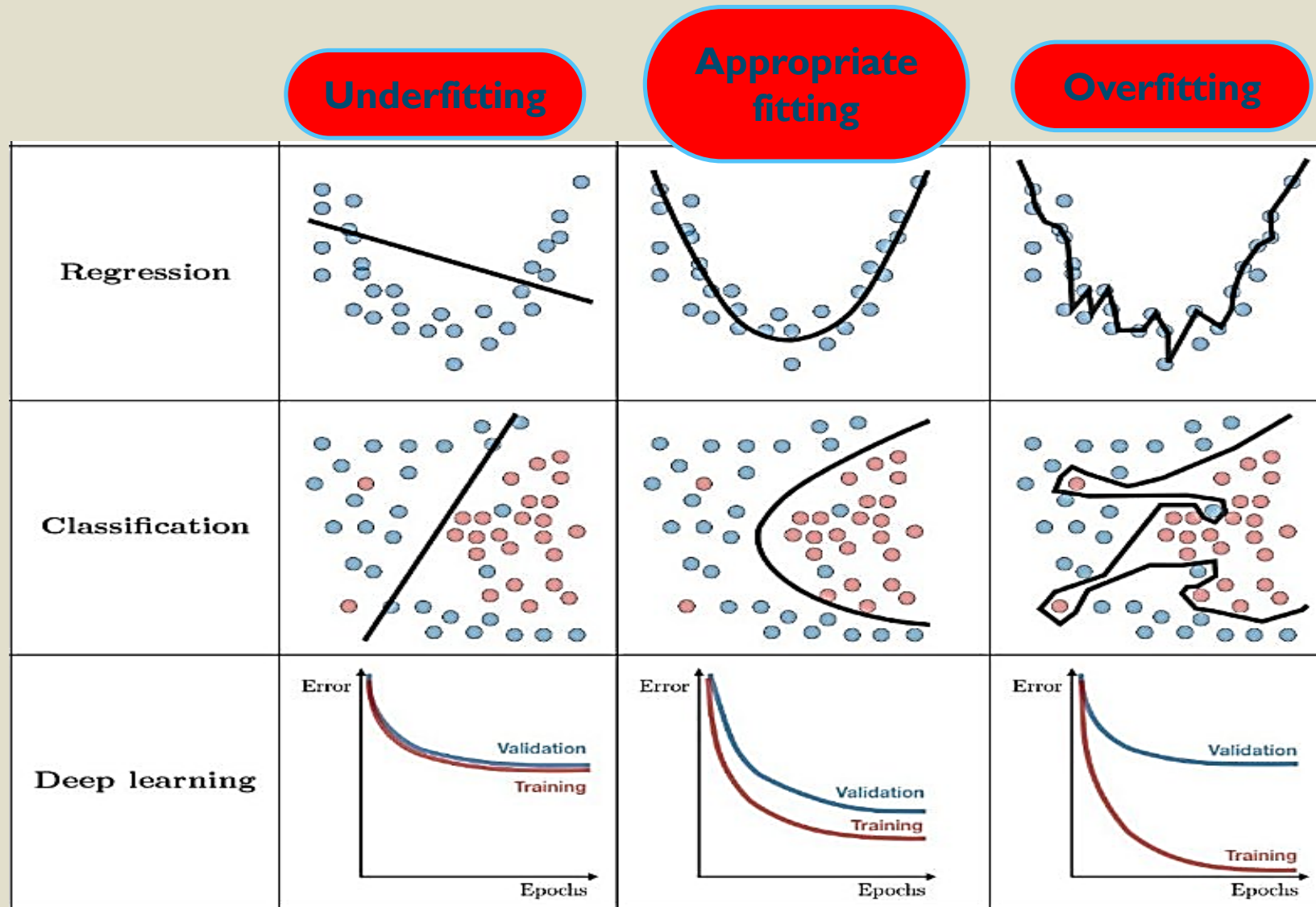
- is the case where the model has“ not learned enough” from the training data, resulting in low generalization and unreliable predictions.

- Overfitting

- is the case where the overall cost is really small, but the generalization of the model is unreliable. This is due to the model learning “too much” from the training data set.

	Underfitting	Just right	Overfitting
Symptoms	<ul style="list-style-type: none">- High training error- Training error close to test error- High bias	<ul style="list-style-type: none">- Training error slightly lower than test error	<ul style="list-style-type: none">- Low training error- Training error much lower than test error- High variance
Remedies	<ul style="list-style-type: none">- Complexify model- Add more features- Train longer		<ul style="list-style-type: none">- Regularize- Get more data

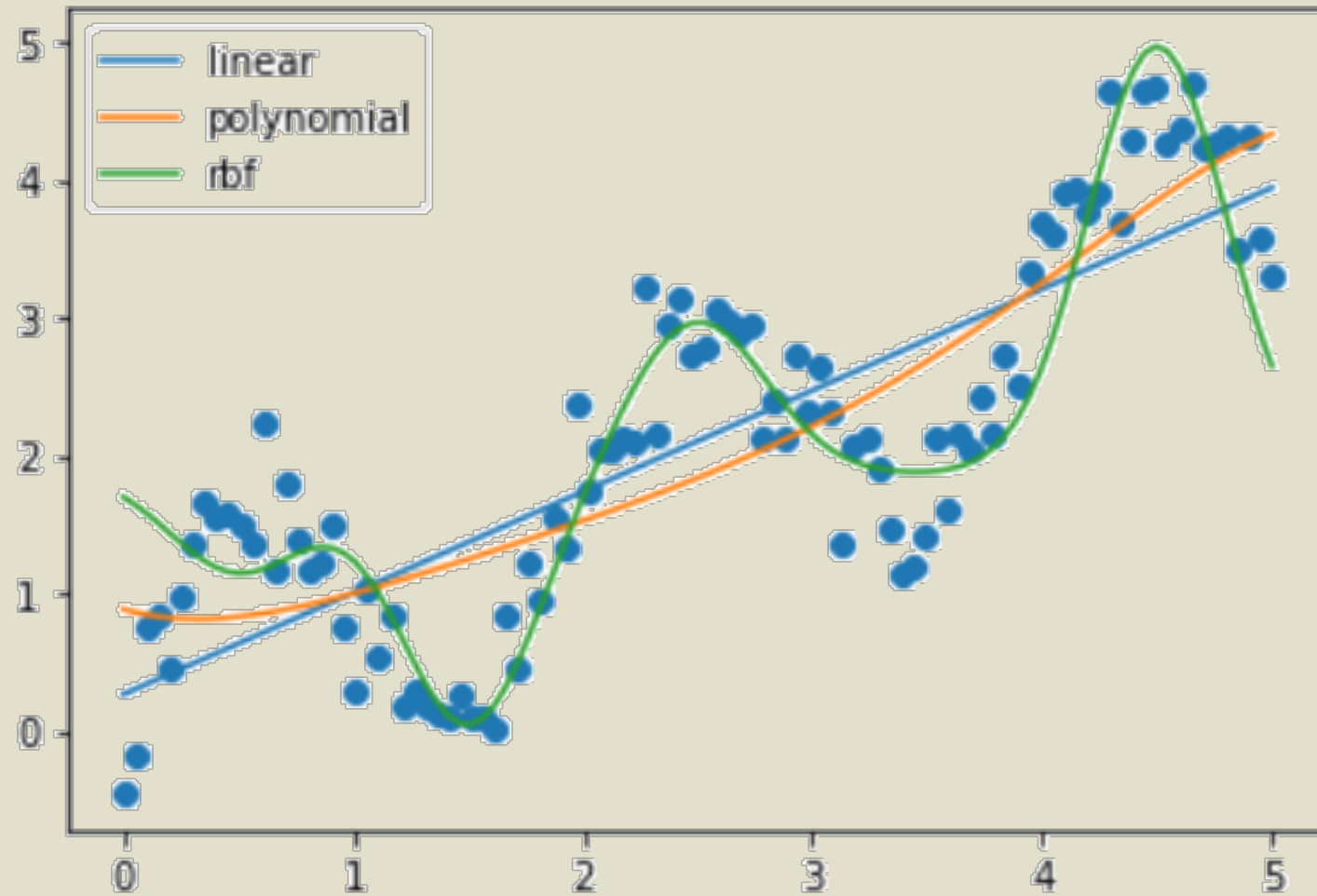
OVERFITTING ? UNDERFITTING ?



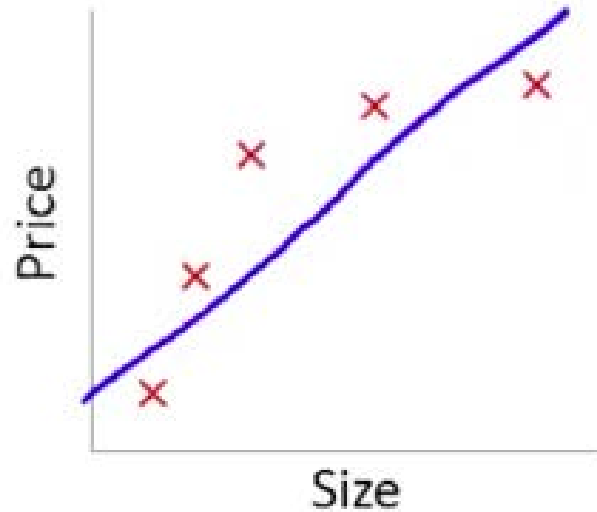
實作

- [Loss Function.ipynb](#)
- [Linear Regression_1.ipynb](#)
- [Linear Regression_2.ipynb](#)
- [Linear Regression_3.ipynb](#)

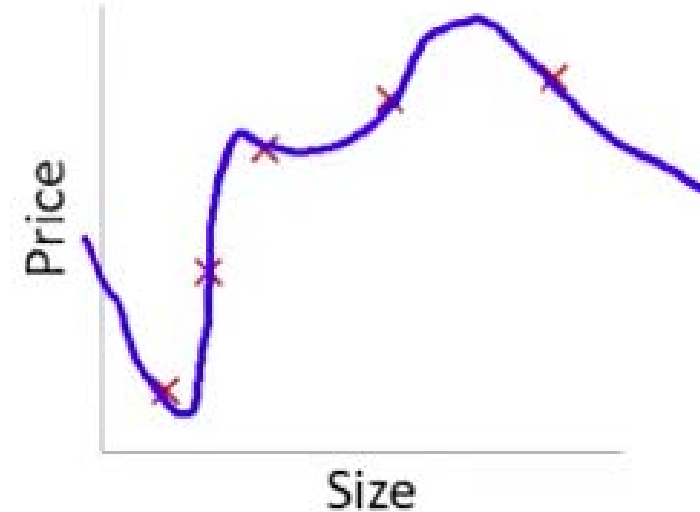
PROBLEM



PROBLEM



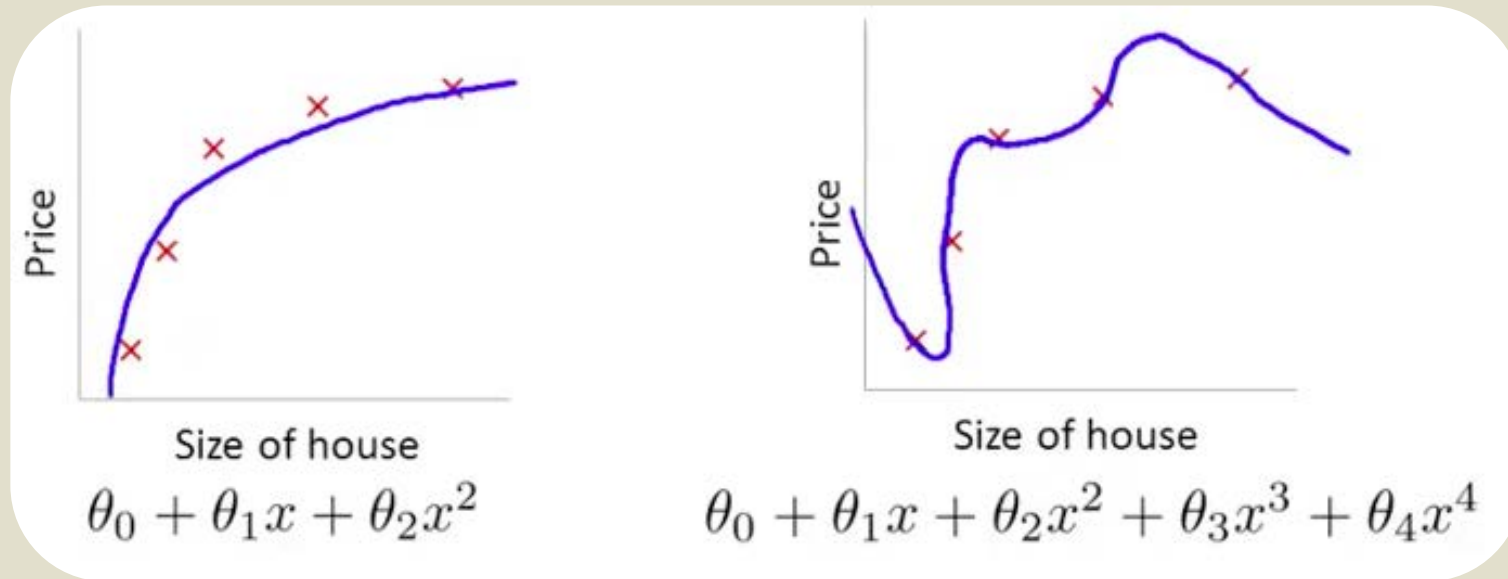
$\rightarrow \theta_0 + \theta_1 x$
"Underfit" "High bias"



$\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$
"Overfit" "High variance"

如何解決OVERFITTING

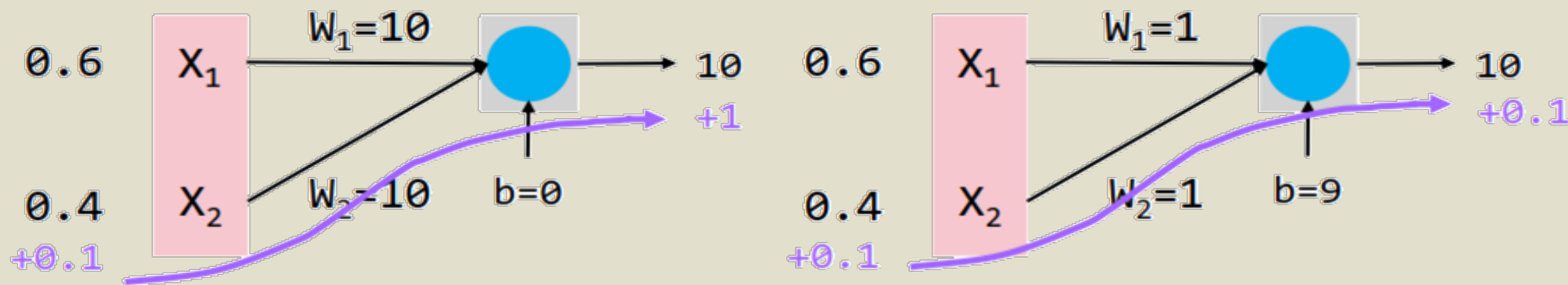
- If we have too many features, the learned hypothesis may fit the training set very well ($J(\theta)=0$), but fail to generalize to new examples(Predictions on new examples)



1. 降低**features**的數量：人工選擇、**model selection algorithm**
2. **Regularization**：維持現有的**features**，但是降低部分不重要**feature**的影響力。

REGULARIZATION

- 限制 weights 的大小讓 output 曲線比較平滑
- 為什麼要限制呢？



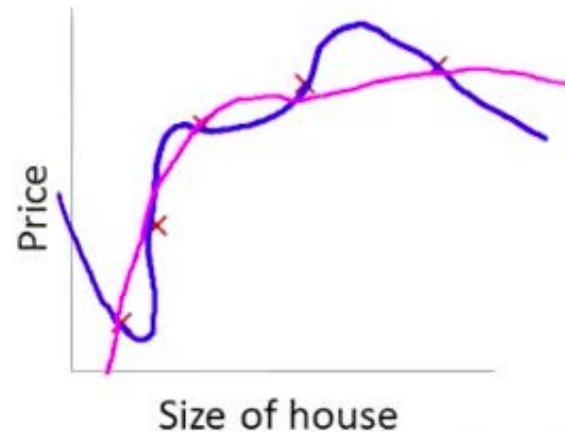
w_i 較小 $\rightarrow \Delta x_i$ 對 \hat{y} 造成的影響 ($\Delta \hat{y}$) 較小
 \rightarrow 對 input 變化比較不敏感 \rightarrow generalization 好

REGULARIZATION

$$J(\theta) = 1/2m \sum_{i=1}^m (h(\theta)^{(i)} - y^{(i)})^2$$

$$\rightarrow \min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + 1000 \underline{\theta_3^2} + 1000 \underline{\theta_4^2}$$

$\underline{\theta_3 \approx 0} \quad \underline{\theta_4 \approx 0}$



$$\underline{\theta_0 + \theta_1 x + \theta_2 x^2} + \cancel{\theta_3 x^3} + \cancel{\theta_4 x^4}$$

$\uparrow \quad \uparrow$

REGULARIZATION

$$J(\theta) = 1/2m \sum_{i=1}^m (h(\theta)^{(i)} - y^{(i)})^2$$



$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

REGULARIZATION

- 為了限制**weights** 的大小，以避免落入**Overfitting**狀態，我們將 **$J(f)$** 加入損失函數中，我們叫做**Regularization**

$$\min L(y_i, f(x_i)) + \lambda J(f)$$

$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

λ 是用來調整regularization 的比重
小心顧此失彼(降低weights 的大小而犧牲模型準確性)

L1 AND L2 REGULARIZERS

- Regularization又分兩種：
 - 第一種是L1正則化 (Lasso)，第二種是L2正則化 (Ridge)

$$\lambda J(f)$$

L1 norm (Lasso Regression) → can also do feature selection

$$L_1 = \sum_{i=1}^N |W_i|$$

Sum of absolute values

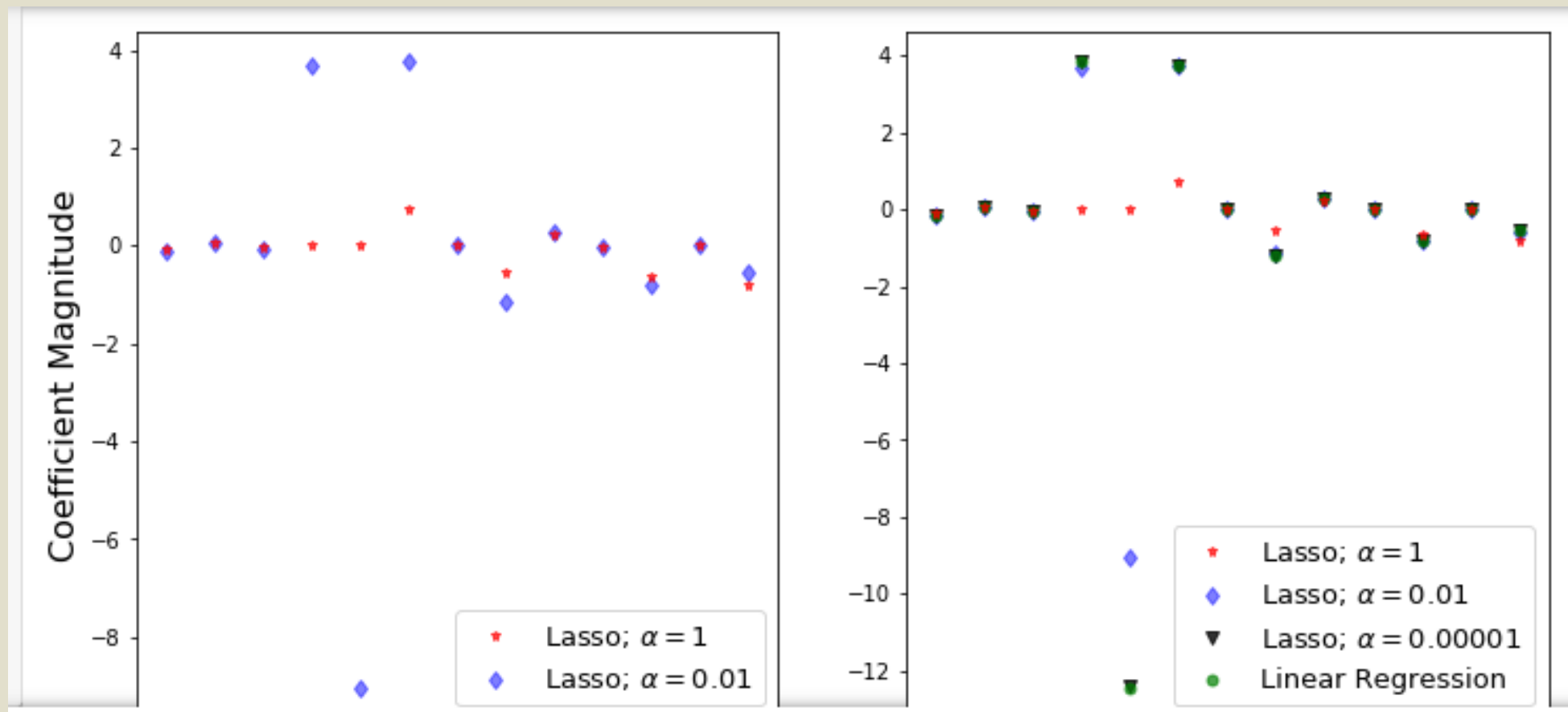
L2 norm (Ridge Regression) → reduces the coefficients close to zero

$$L_2 = \sqrt{\sum_{i=1}^N |W_i|^2}$$

Root mean square of absolute values

實作

- `lasso_regression.ipynb`



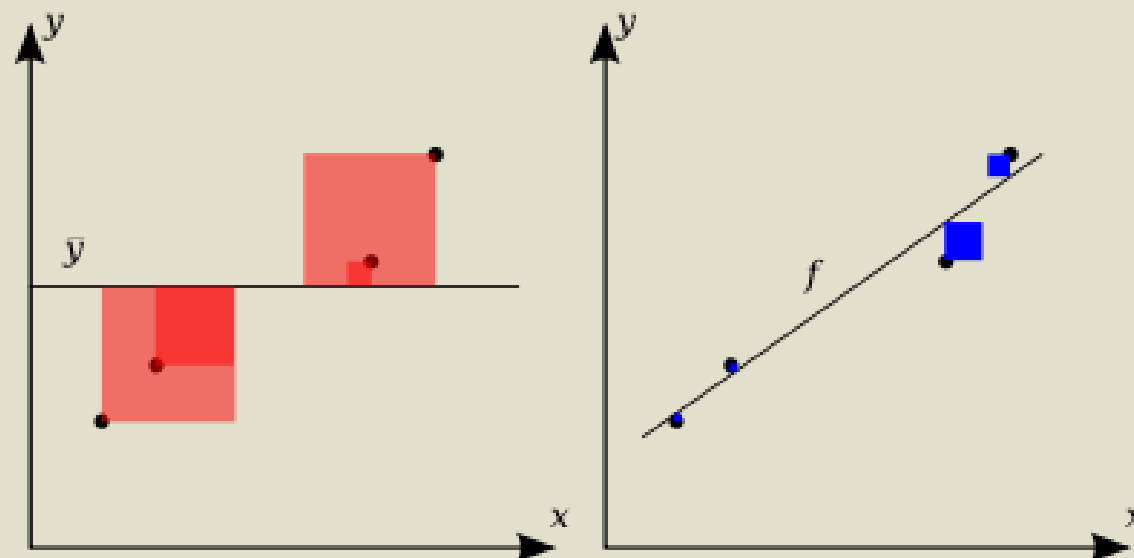
R^2

$$R^2 \equiv 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}$$

$$SS_{\text{res}} = \sum_i (y_i - f_i)^2 = \sum_i e_i^2$$

$$SS_{\text{tot}} = \sum_i (y_i - \bar{y})^2$$

$$R^2 = 1 - \frac{MSE(\hat{y}, y)}{Var(y)}$$



結果是1 → 模型無錯誤

結果是0 → 模型跟瞎猜差不多

結果是0~1之間的數 → 即模型的好壞程度

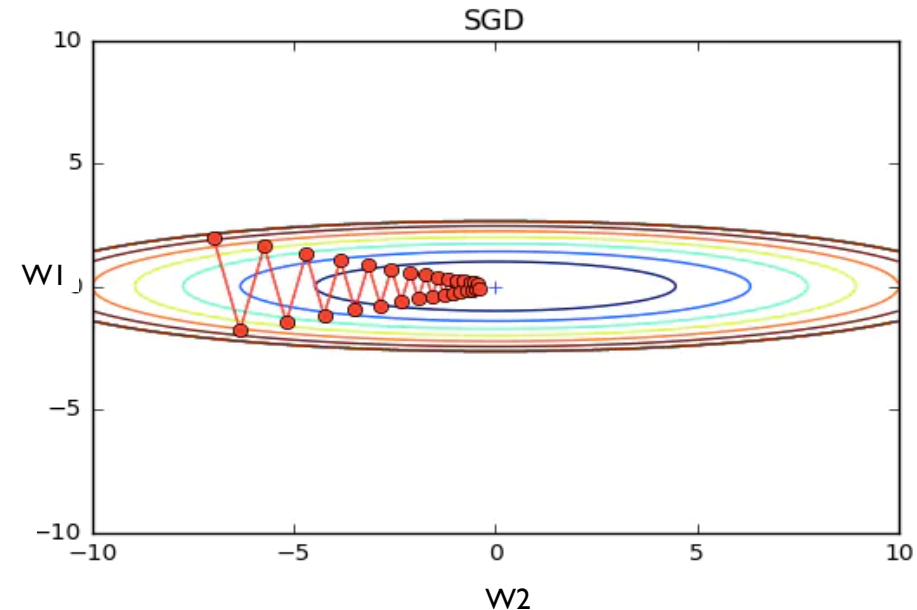
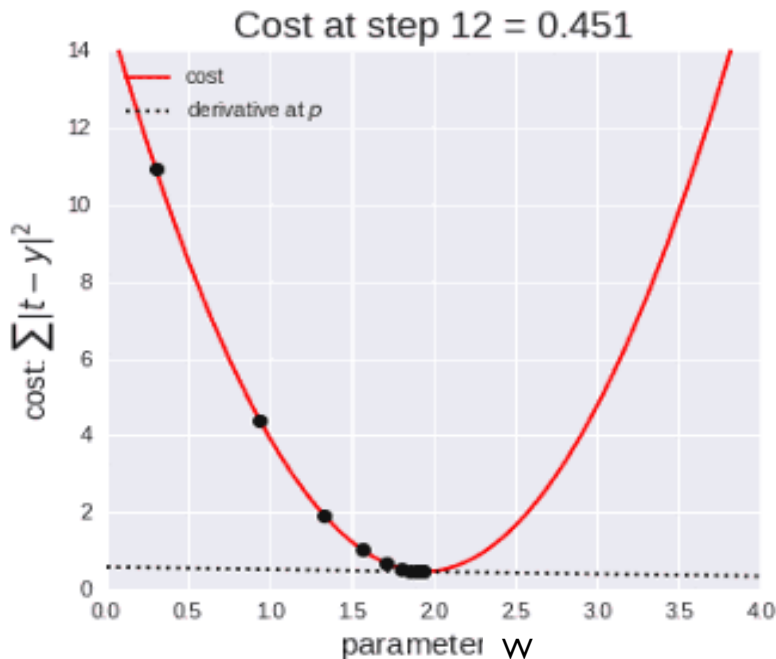
結果是負數 → 我們的模型還不如瞎猜

A decorative wavy line in a light green color runs vertically along the left side of the slide, separating a light beige area from a dark teal background.

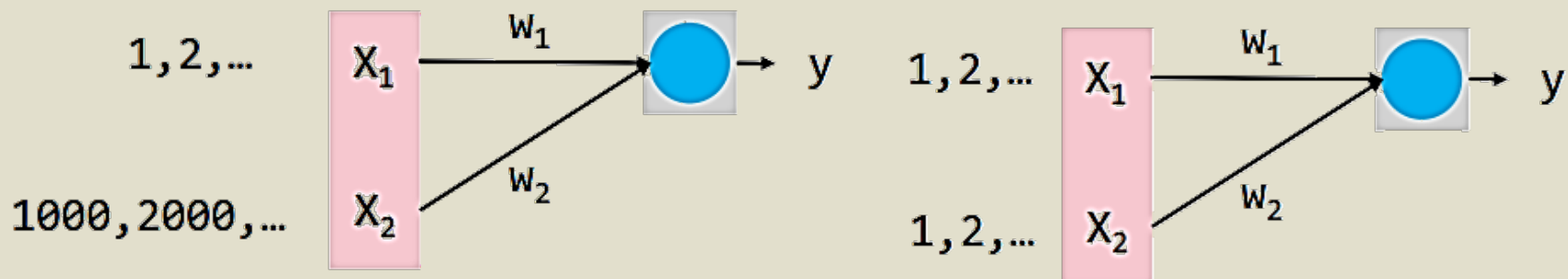
APPENDIX GRADIENT DESCENT ALGORITHM

SGD(STOCHASTIC GRADIENT DESCENT)

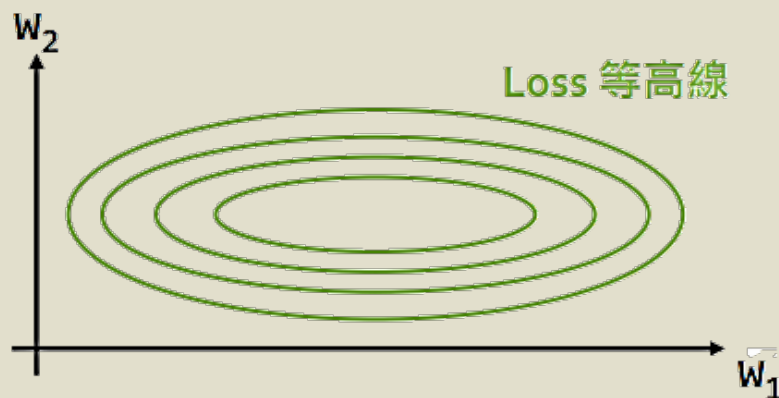
- Choose an initial vector of parameters **W** and learning rate **η**
- Repeat until an approximate minimum is obtained:
 - Randomly shuffle samples in the training set.
 - For each sample do $W \leftarrow W - \eta \frac{\partial L}{\partial W}$



資料 RE-SCALE 對 WEIGHT 影響



W_2 的修正(ΔW)對於 loss 的影響比較大



如果不re-scale, 值域特大的特徵會對Loss貢獻比較多
修正W2對LOSS影響比較大

有做re-scale 收斂速度也比較快

GRADIENTS

$$Y = XW^T + b$$

$$X = \begin{bmatrix} x_1^1 & x_2^1 \\ x_1^2 & x_2^2 \end{bmatrix} \quad w = [w_1 \quad w_2] \quad b = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}$$

$$\hat{Y} = XW^T + b$$

$$\hat{Y} = \begin{bmatrix} x_1^1 w_1 + x_2^1 w_2 + b_1 \\ x_1^2 w_1 + x_2^2 w_2 + b_2 \end{bmatrix} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \end{bmatrix} \quad \text{np.dot(X,w.T)+b}$$

$$X = \begin{bmatrix} 1 & x_1^1 & x_2^1 \\ 1 & x_1^2 & x_2^2 \end{bmatrix} \quad 2 \times 3$$

Bias 這一項其實可以直接在X加上一欄“1”

np.c_[np.ones((len(X),1)),X]

$$w = [w_0 \quad w_1 \quad w_2]$$

$$\begin{bmatrix} w_0 + x_1^i w_1 + x_2^i w_2 \\ w_0 + x_1^i w_1 + x_2^i w_2 \end{bmatrix}$$

np.dot(X,w.T)

GRADIENTS

Loss function: MSE

$$\text{loss}(w) = 1/m \sum_{i=1}^m (\hat{y}^i - y^i)^2$$

gradient

$$\begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} - \eta \begin{bmatrix} \frac{\partial}{\partial w_0} 1/m \sum_{i=1}^m (w_0 + x_1^i w_1 + x_2^i w_2 - y_i)^2 \\ \frac{\partial}{\partial w_1} 1/m \sum_{i=1}^m (w_0 + x_1^i w_1 + x_2^i w_2 - y_i)^2 \\ \frac{\partial}{\partial w_2} 1/m \sum_{i=1}^m (w_0 + x_1^i w_1 + x_2^i w_2 - y_i)^2 \end{bmatrix} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix} - \eta \begin{bmatrix} 2/m \sum_{i=1}^m (w_0 + x_1^i w_1 + x_2^i w_2 - y_i) 1 \\ 2/n \sum_{i=1}^m (w_0 + x_1^i w_1 + x_2^i w_2 - y_i) x_1^i \\ 2/n \sum_{i=1}^m (w_0 + x_1^i w_1 + x_2^i w_2 - y_i) x_2^i \end{bmatrix}$$

$$w_j := w_j - \eta(2/m \sum_{i=1}^m (\hat{y}^i - y^i) \cdot x_j^i)$$

$$w = w - \eta(2/m) \frac{\partial \text{Loss}}{\partial w} = w - \eta(2/m)(X^T \cdot \text{Loss})$$

$$w = w - l * 2/m * \text{np.dot}(X.T, \text{loss})$$

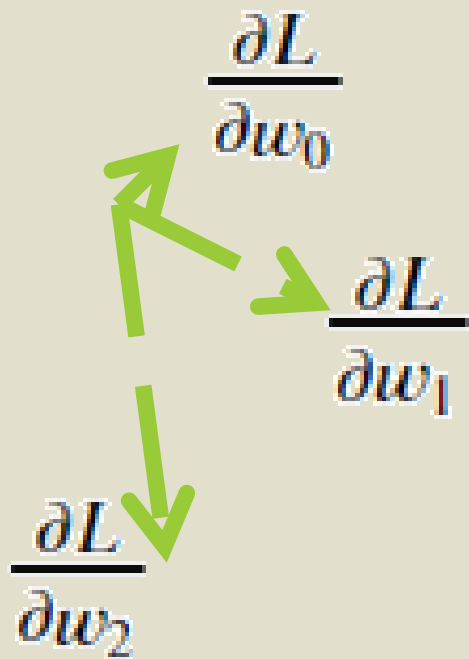
```
loss=np.dot(X,w.T)-y
l/m*np.sum(np.square(loss))
```

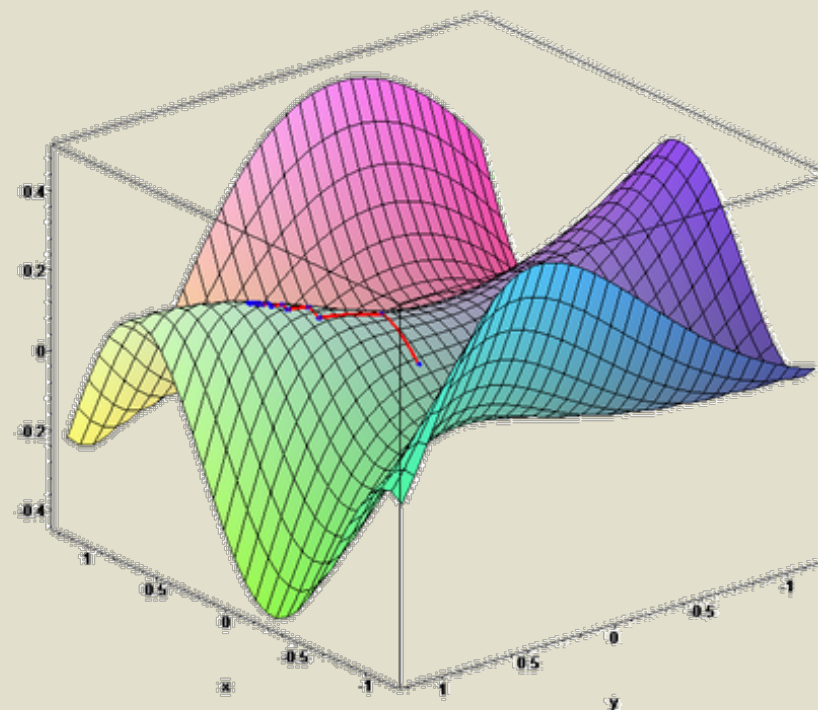
$$\begin{bmatrix} 1 & 1 \\ x_1^1 & x_1^2 \\ x_2^1 & x_2^2 \end{bmatrix}_{3 \times 2} \cdot \begin{bmatrix} \text{Loss}_1 \\ \text{Loss}_2 \end{bmatrix}_{2 \times 1} = \begin{bmatrix} \frac{\partial L}{\partial w_0} \\ \frac{\partial L}{\partial w_1} \\ \frac{\partial L}{\partial w_2} \end{bmatrix}_{3 \times 1}$$

X^T Loss Gradient

GRADIENT

- Gradient is a derivative of a function at a certain point.
- 每一變數的偏微分值，代表函數中該變數在某一點所看到坡度。不同變數 w_i 代表不同面向，而偏微分值的大小是傾斜程度


$$\frac{\partial L}{\partial w_0}$$
$$\frac{\partial L}{\partial w_1}$$
$$\frac{\partial L}{\partial w_2}$$



Gradient Vector 描述了每一個人自己對於谷底位置的認知強度

GRADIENT實作

- 調整Hyperparameters，例如 x_start 、epochs、 lr ，測試逼近的過程。

