

22.1.2 Backpropagation algorithm for training NN

Step 1：隨機設定初始的權重 $w_1 \sim w_8$ 及偏置 $b_1 \sim b_2$ 。

1963年Vapnik提出

Step 2：輸入某個訓練資料 $i_1 \sim i_2$ ，求 $o_1 \sim o_2$ 。正確輸出是0.01和0.99，求誤差。

$$net_{h1} = 0.05 * 0.15 + 0.1 * 0.2 + 0.35 * 1 = 0.3775$$

$$net_{o1} = 0.593269992 * 0.4 + 0.596884378 * 0.45 + 1 * 0.6 = 1.105905967$$

$$out_{h1} = \frac{1}{1 + e^{-net_{h1}}} = \frac{1}{1 + e^{-0.3775}} = 0.593269992$$

$$out_{o1} = \frac{1}{1 + e^{-net_{o1}}} = 0.75136507$$

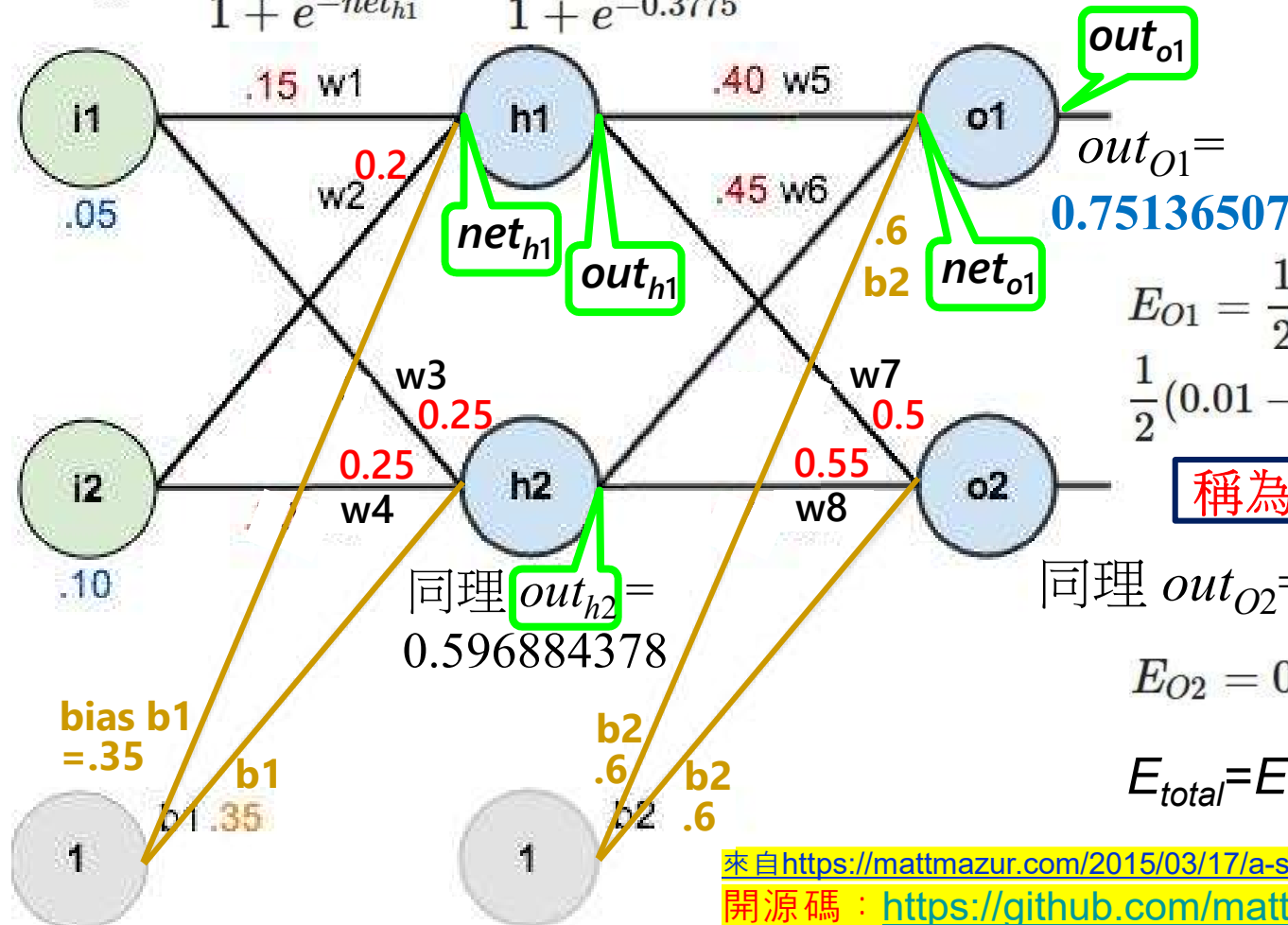
$$E_{O1} = \frac{1}{2} (target - output)^2 = \frac{1}{2} (0.01 - 0.75136507)^2 = 0.274811083$$

稱為least mean square error

同理 $out_{o2} = 0.772928465$

$$E_{O2} = 0.023560026$$

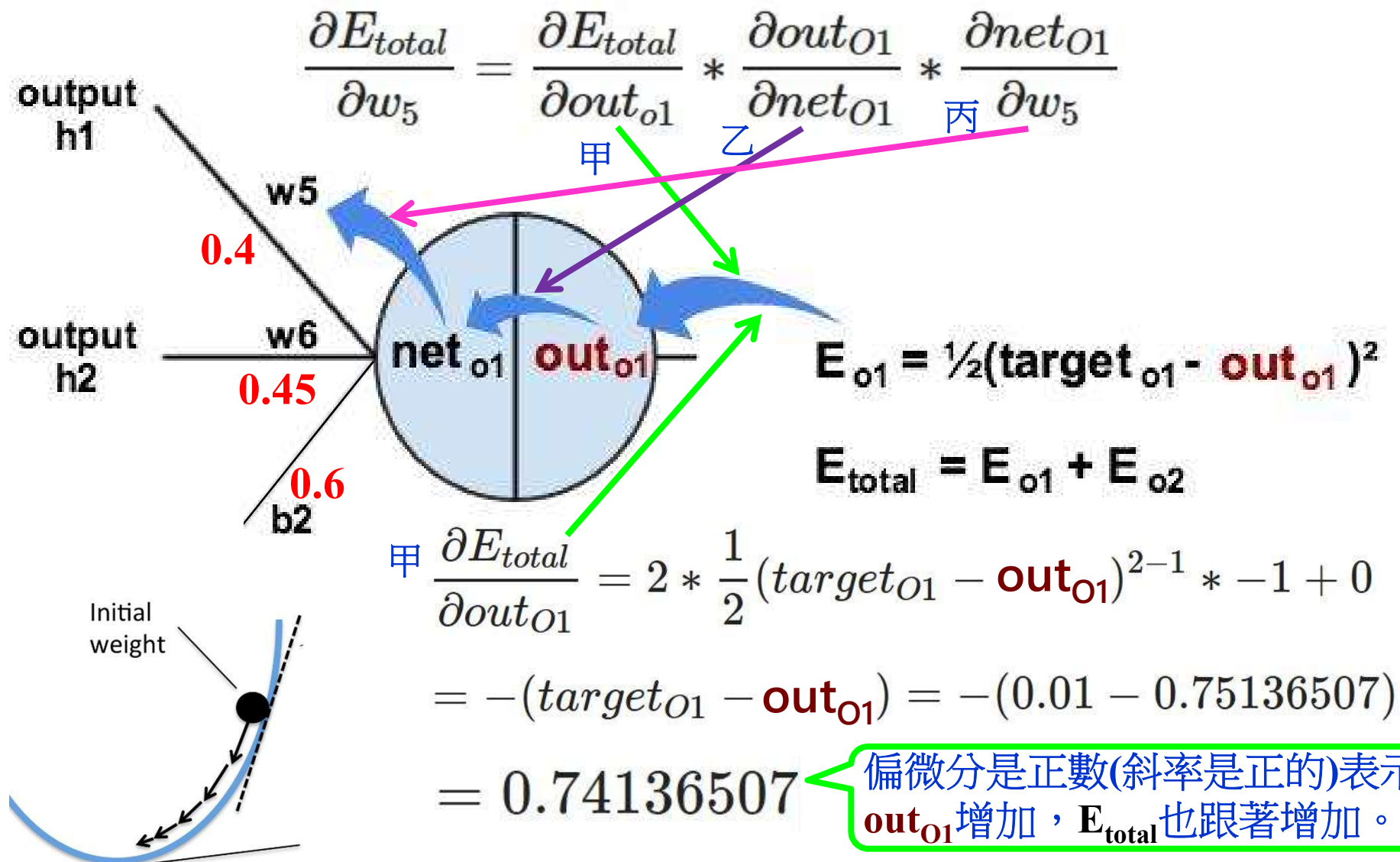
$$E_{total} = E_{O1} + E_{O2} = 0.298371109$$

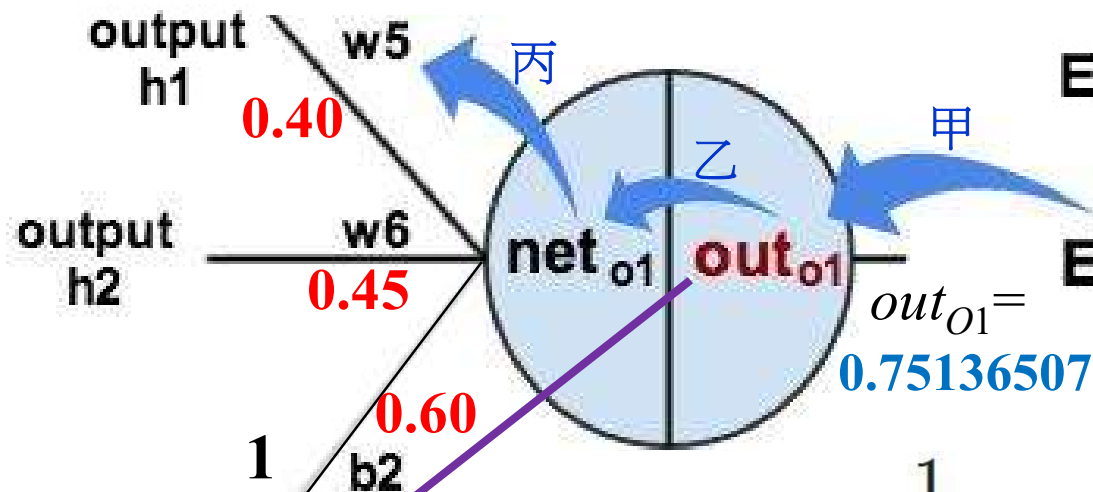


來自<https://mattmazur.com/2015/03/17/a-step-by-step-backpropagation-example/>
開源碼：<https://github.com/mattm/simple-neural-network>

Step 3：更新網路中的每一個權重，使得最終的輸出更接近**target**。

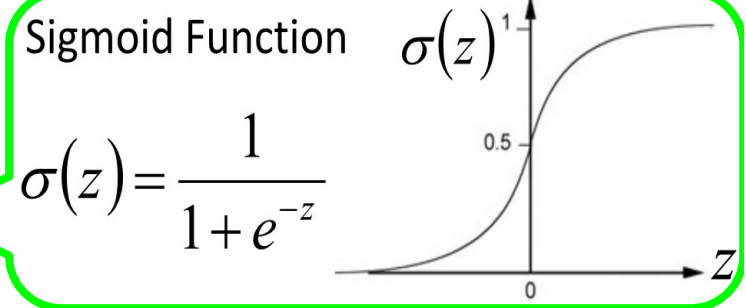
先看 w_5 ，我們想知道 w_5 的改變會影響總誤差多少，也就是 $\frac{\partial E_{total}}{\partial w_5}$ 。
根據**Chain rule**，我們可以得到：





$$E_{\text{total}} = E_{o1} + E_{o2}$$

$$E_{o1} = \frac{1}{2}(\text{target}_{o1} - \text{out}_{o1})^2$$



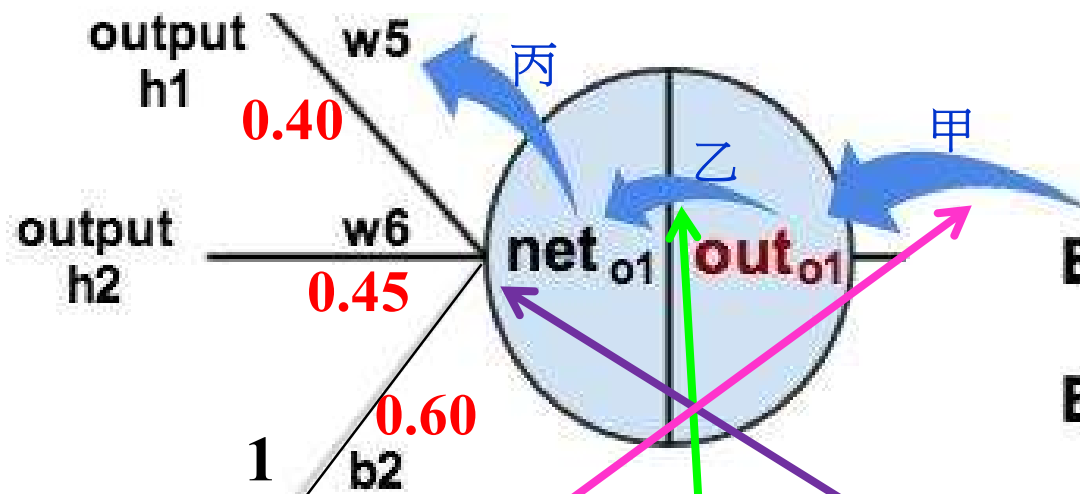
$$\text{out}_{o1} = \frac{1}{1 + e^{-\text{net}_{o1}}}$$

$$\begin{aligned} \text{乙 } \frac{\partial \text{out}_{o1}}{\partial \text{net}_{o1}} &= \frac{e^{-\text{net}_{o1}}}{(1 + e^{-\text{net}_{o1}})^2} = \frac{1}{1 + e^{-\text{net}_{o1}}} \cdot \frac{e^{-\text{net}_{o1}}}{1 + e^{-\text{net}_{o1}}} = \text{out}_{o1} \cdot \left(\frac{1 + e^{-\text{net}_{o1}}}{1 + e^{-\text{net}_{o1}}} - \frac{1}{1 + e^{-\text{net}_{o1}}} \right) \\ &= \text{out}_{o1}(1 - \text{out}_{o1}) = 0.75136507(1 - 0.75136507) = 0.186815602 \end{aligned}$$

$$\text{丙 } \text{net}_{o1} = w_5 * \text{out}_{h1} + w_6 * \text{out}_{h2} + b_2 * 1$$

$$\frac{\partial \text{net}_{o1}}{\partial w_5} = \text{out}_{h1} = 0.593269992$$

$$\begin{aligned} \frac{\partial E_{\text{total}}}{\partial w_5} &= \frac{\partial E_{\text{total}}}{\partial \text{out}_{o1}} * \frac{\partial \text{out}_{o1}}{\partial \text{net}_{o1}} * \frac{\partial \text{net}_{o1}}{\partial w_5} \\ &= \text{甲 } -(\text{target}_{o1} - \text{out}_{o1}) * \text{乙 } \text{out}_{o1}(1 - \text{out}_{o1}) * \text{丙 } \text{out}_{h1} \\ &= \text{甲 } 0.74136507 * \text{乙 } 0.186815602 * \text{丙 } 0.593269992 = \mathbf{0.082167041} \end{aligned}$$



$$E_{o1} = \frac{1}{2}(\text{target}_{o1} - \text{out}_{o1})^2$$

$$E_{\text{total}} = E_{o1} + E_{o2}$$

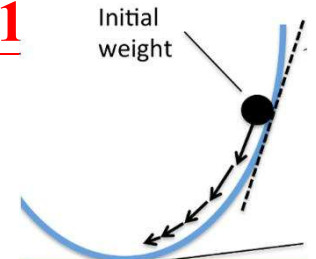
$$\begin{aligned} \frac{\partial E_{\text{total}}}{\partial w_5} &= \frac{\partial E_{\text{total}}}{\partial \text{out}_{o1}} * \frac{\partial \text{out}_{o1}}{\partial \text{net}_{o1}} * \frac{\partial \text{net}_{o1}}{\partial w_5} \\ &= \text{甲} - (\text{target}_{o1} - \text{out}_{o1}) * \text{乙} \text{out}_{o1} (1 - \text{out}_{o1}) * \text{丙} \text{out}_{h1} \\ &= \underline{0.74136507} * \underline{0.186815602} * \underline{0.593269992} = \underline{0.082167041} \end{aligned}$$

偏微分是正數，即斜率是正的，表示 w_5 增加，總誤差 E_{total} 也會跟著增加。因此我們應減少 w_5 之值，使總誤差 E_{total} 降低：

$$w_5^+ = w_5 - \eta * \frac{\partial E_{\text{total}}}{\partial w_5} = \underline{0.4} - 0.5 * 0.082167041 = \underline{0.35891648}$$

learning rate

這個規則稱為**LMS(Least Mean Squares)**更新規則，或**Widrow-Hoff learning rule**。



Gradient Descent(梯度下降法)

Gradient direction:

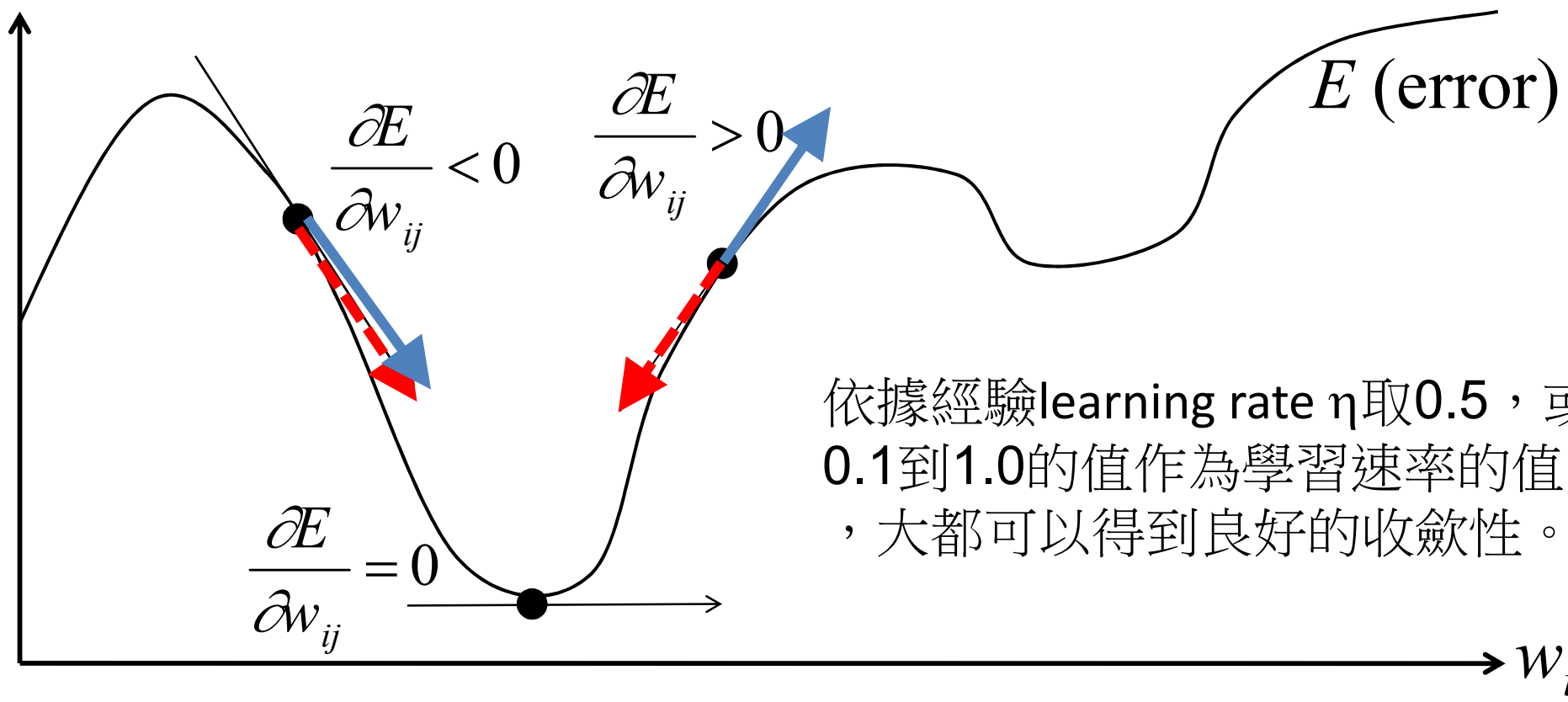


Descent direction:

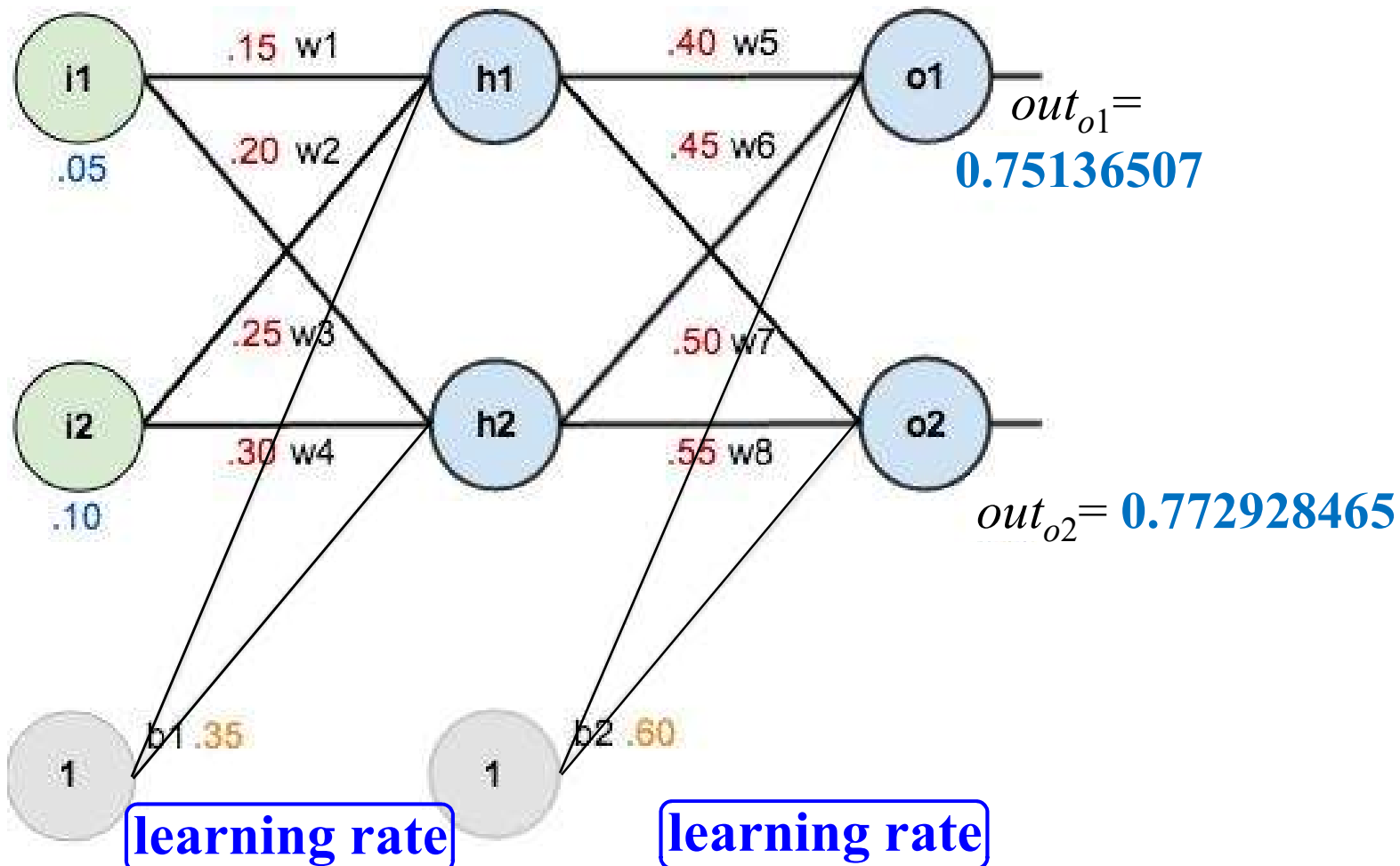


$$\Delta w_{ij} = -\eta \frac{\partial E}{\partial w_{ij}}$$

(η is learning rate)



依據經驗learning rate η 取0.5，或0.1到1.0的值作為學習速率的值，大都可以得到良好的收斂性。

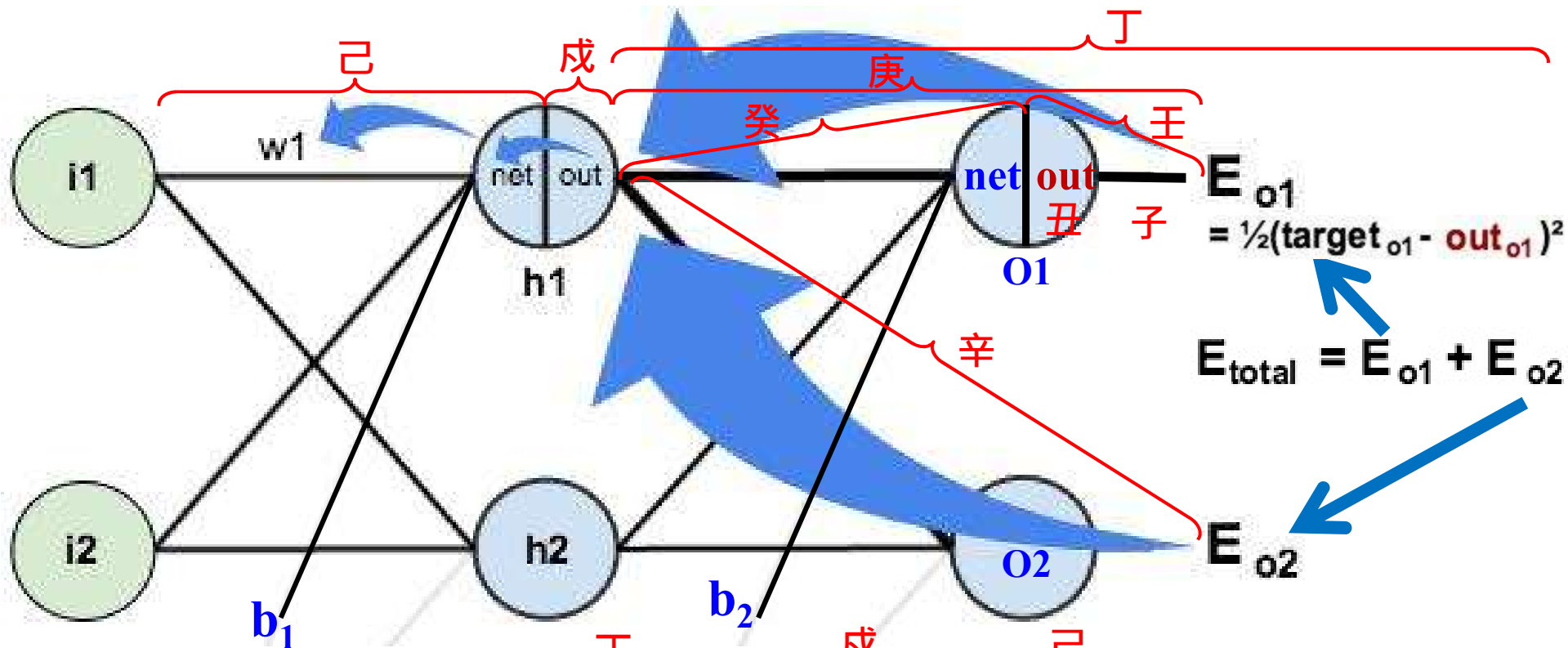


$$w_5^+ = w_5 - \eta * \frac{\partial E_{total}}{\partial w_5} = 0.4 - 0.5 * 0.082167041 = 0.35891648$$

同理更新

$$\begin{aligned} w_6^+ &= 0.408666186 \\ w_7^+ &= 0.511301270 \\ w_8^+ &= 0.561370121 \end{aligned}$$

在用到 w_5 、 w_6 、 w_7 和 w_8 的地方，仍使用舊的權重，而不是逐一更新後的權重。



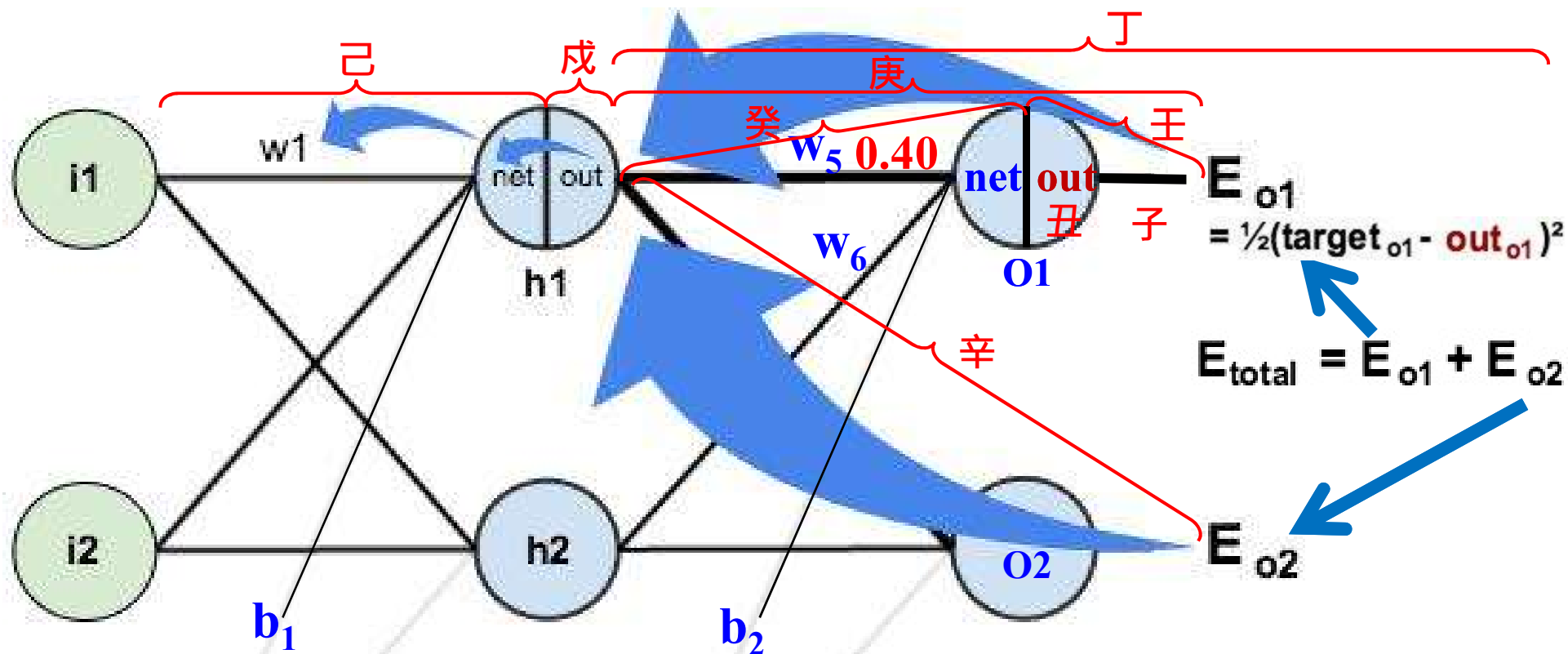
再看 w_1 如何更新: $\frac{\partial E_{total}}{\partial w_1} = \frac{\partial E_{total}}{\partial out_{h1}} * \frac{\partial out_{h1}}{\partial net_{h1}} * \frac{\partial net_{h1}}{\partial w_1}$ //根據 Chain rule

$$\frac{\partial E_{total}}{\partial out_{h1}} = \frac{\partial E_{o1}}{\partial out_{h1}} + \frac{\partial E_{o2}}{\partial out_{h1}}$$

注意 w_1 有兩條路徑影響 E_{total}

$$\frac{\partial E_{o1}}{\partial out_{h1}} = \frac{\partial E_{o1}}{\partial net_{o1}} * \frac{\partial net_{o1}}{\partial out_{h1}}$$

$$\frac{\partial E_{o1}}{\partial net_{o1}} = \frac{\partial E_{o1}}{\partial out_{o1}} * \frac{\partial out_{o1}}{\partial net_{o1}} = \underline{0.74136507 * 0.186815602} = 0.138498562$$



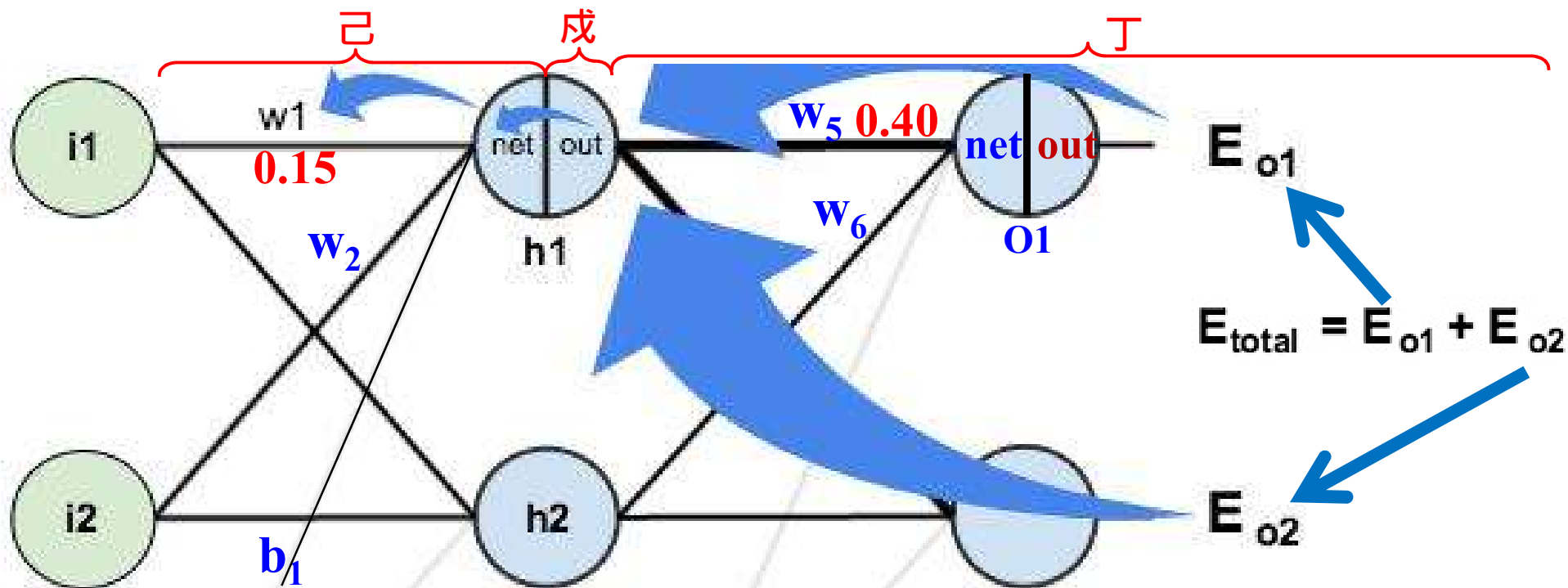
$$net_{o1} = w_5 * out_{h1} + w_6 * out_{h2} + b_2 * 1$$

$$\text{癸} \frac{\partial net_{o1}}{\partial out_{h1}} = w_5 = 0.40$$

$$\text{庚} \frac{\partial E_{o1}}{\partial out_{h1}} = \text{壬} \frac{\partial E_{o1}}{\partial net_{o1}} * \text{癸} \frac{\partial net_{o1}}{\partial out_{h1}} = 0.138498562 * 0.40 = 0.055399425$$

同理可得

$$\text{丁} \frac{\partial E_{total}}{\partial out_{h1}} = \text{庚} \frac{\partial E_{o1}}{\partial out_{h1}} + \text{辛} \frac{\partial E_{o2}}{\partial out_{h1}} = 0.055399425 + (-0.019049119) = 0.036350306$$



$$out_{h1} = \frac{1}{1 + e^{net_{h1}}}$$

戌 $\frac{\partial out_{h1}}{\partial net_{h1}} = out_{h1}(1 - out_{h1}) = 0.59326999(1 - 0.59326999) = 0.241300709$

net_{h1} = $w_1 * i_1 + w_2 * i_2 + b_1 * 1$ 故 $\frac{\partial net_{h1}}{\partial w_1} = i_1 = 0.05$


丁 $\frac{\partial E_{total}}{\partial w_1} = \frac{\partial E_{total}}{\partial out_{h1}} * \frac{\partial out_{h1}}{\partial net_{h1}} * \frac{\partial net_{h1}}{\partial w_1} = 0.036350306 * 0.241300709 * 0.05 = 0.000438568$

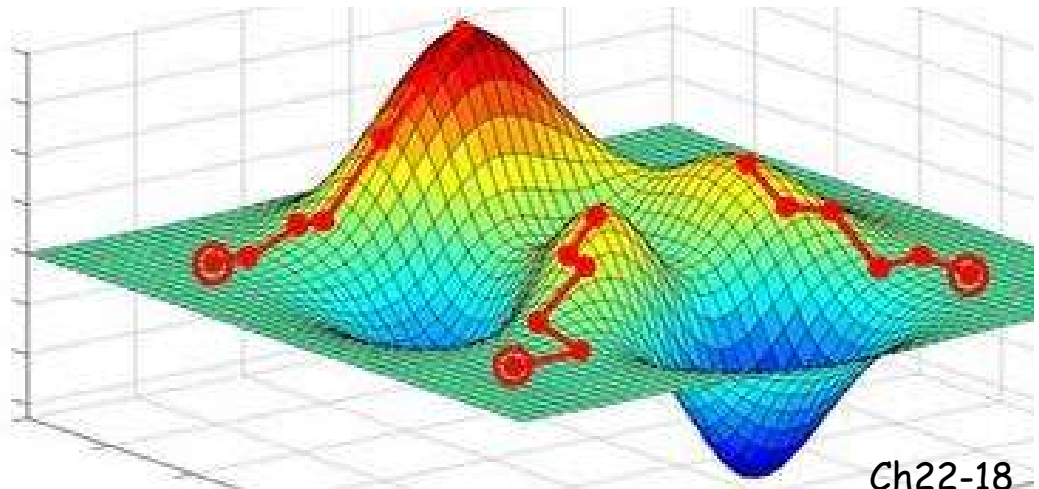
更新 w_1 : $w_1^+ = w_1 - \eta * \frac{\partial E_{total}}{\partial w_1} = 0.15 - 0.5 * 0.000438568 = 0.149780716$

同理更新：

$$w_2=0.2 \quad w_3=0.25 \quad w_4=0.3$$

$$w_2^+ = 0.19956143 \quad w_3^+ = 0.24975114 \quad w_4^+ = 0.29950229$$

- 現在已經更新所有的 w 權重。(我們沒有更動 b_1 及 b_2) 練習
- 在最初，輸入為0.05和0.1的時候，網路的誤差為**0.298371109**。  親自檢視、並執行開源碼看看。
- 經過第一次 **Backpropagation** 後，誤差降低到了**0.291027924**，看來降得不多。
- 在重複這個過程**10000**次以後，網路的誤差就降到了**0.0000351085**。 開源碼：<https://github.com/mattm/simple-neural-network>
- 這個時候，當我們把0.05和0.1再輸入進去，兩個神經元的輸出為**0.015912196**(理想值**0.01**)和**0.984065734**(理想值**0.99**)。



function BACK-PROP-LEARNING(*examples, network*) **returns** a neural network //反向傳播學習法

inputs: *examples*, a set of examples, each with input vector **x** and output vector **y**

network, a multilayer network with *L* layers, weights $w_{i,j}$, activation function *g*

$\Delta[j]$ 的計算
次序：↘

local variables: Δ , a vector of errors, indexed by network node

repeat

for each weight $w_{i,j}$ **in** *network* **do**

$w_{i,j} \leftarrow$ a small random number // 隨機設定初始權重

for each example (**x**, **y**) **in** *examples* **do** // 輸入某個訓練資料

/ Propagate the inputs forward to compute the outputs */*

for each node *i* in the input layer **do**

$a_i \leftarrow x_i$ // 讀取輸入層資料

for $\ell = 2$ **to** *L* **do**

for each node *j* in layer ℓ **do**

$$in_j \leftarrow \sum_i w_{i,j} a_i$$

$$a_j \leftarrow g(in_j)$$

// 計算中間層資料，
// 正向傳遞預測的值

/ Propagate deltas backward from output layer to input layer */*

for each node *j* in the output layer **do**

$$\Delta[j] \leftarrow g'(in_j) \times (y_j - a_j)$$

// 計算輸出層誤差

for $\ell = L - 1$ **to** 1 **do**

for each node *i* in layer ℓ **do**

$$\Delta[i] \leftarrow g'(in_i) \sum_j w_{i,j} \Delta[j]$$

// 由後往前計算中間層誤差，backpropate errors

// 反向傳遞誤差的值

/ Update every weight in network using deltas */*

$$a_j = g\left(\sum_{i=0}^n w_{i,j} a_i\right)$$

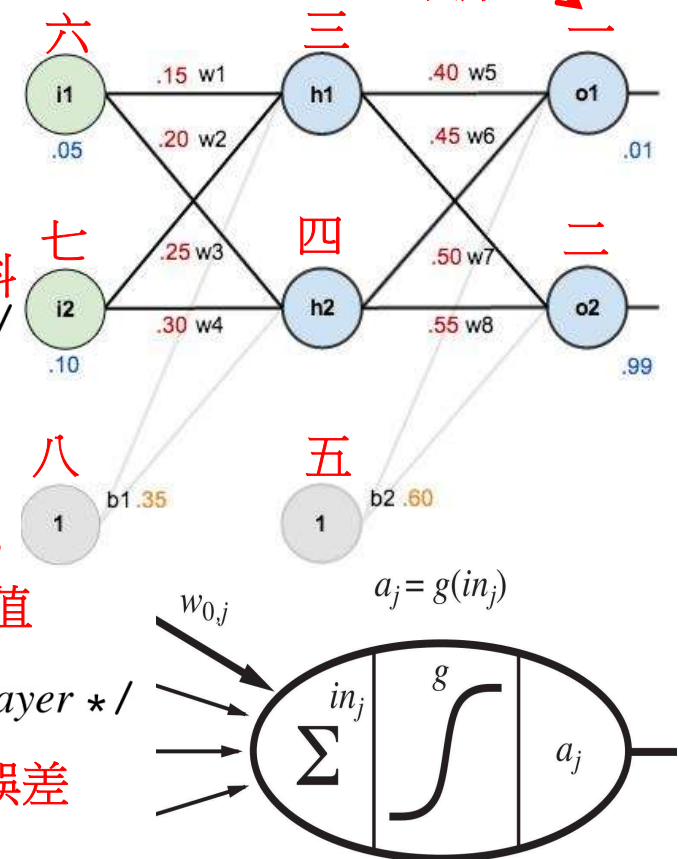
for each weight $w_{i,j}$ **in** *network* **do**

$$w_{i,j} \leftarrow w_{i,j} + \alpha \times a_i \times \Delta[j]$$

// 更新網路中的每一個權重， α 是 learning rate

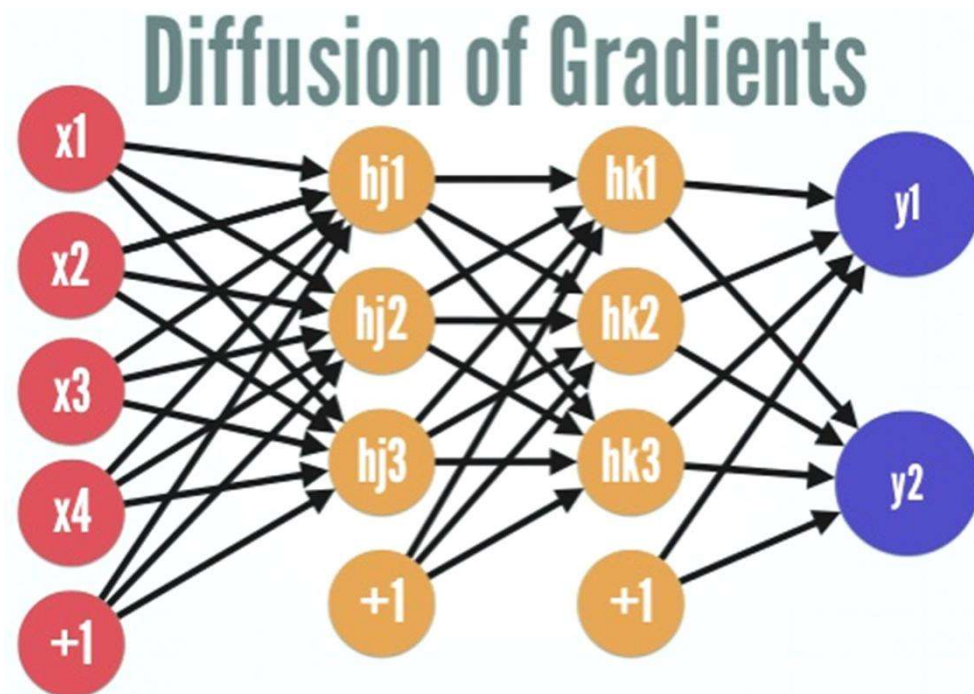
until some stopping criterion is satisfied // 結束條件：總誤差夠小、權值變化很小、執行逾時。

return *network*

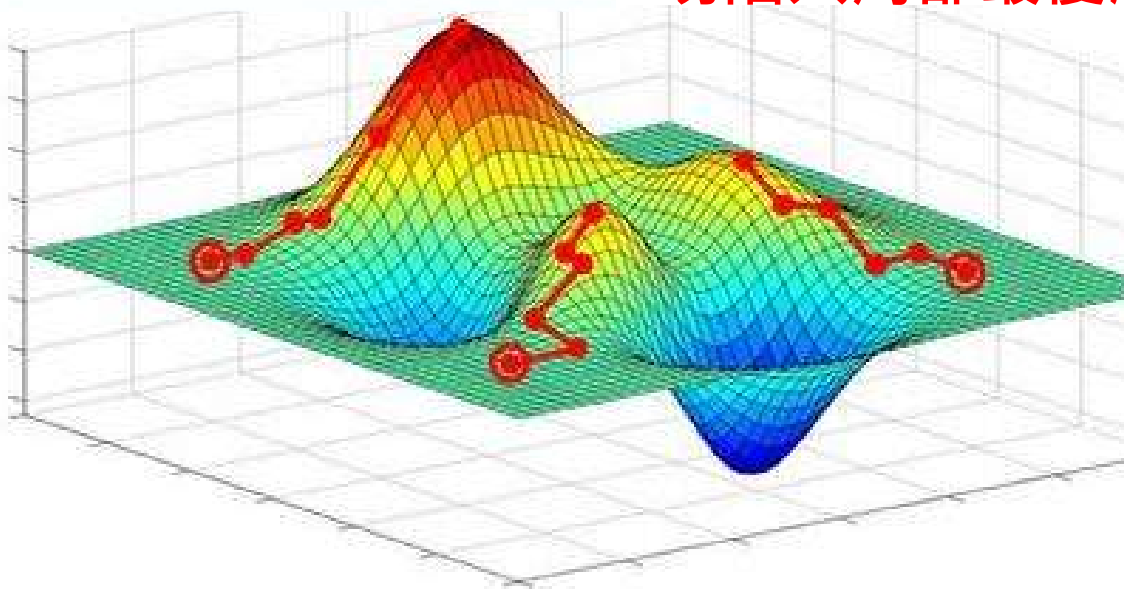


1963年Vapnik 提出

BP演算法的問題

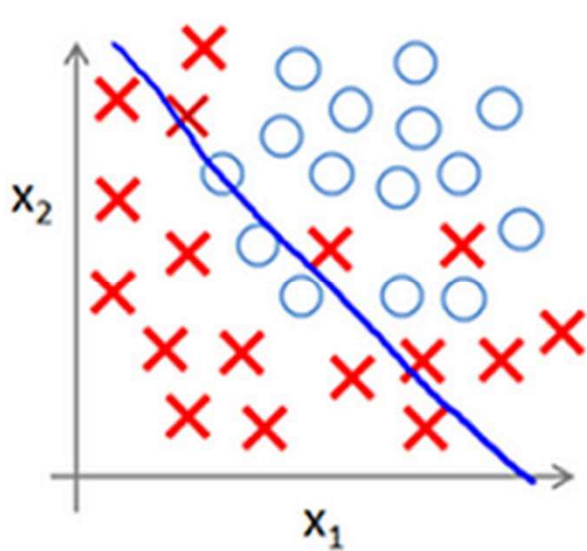


◆ gradient diffusion (梯度擴散，又稱 vanishing gradient)：誤差傳播到前面的layer將使梯度的幅度會急劇減小，導致前幾層神經元的權重更新非常緩慢，變成了前幾層相對固定，只能改變最後幾層的結果，很容易陷入局部最優解。

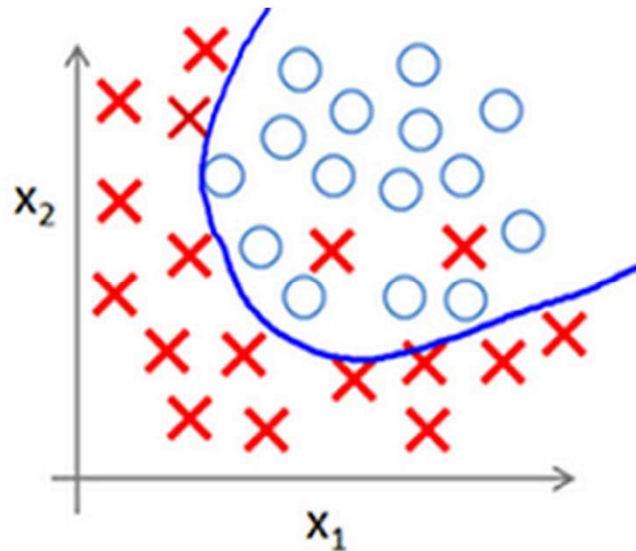


BP演算法的問題

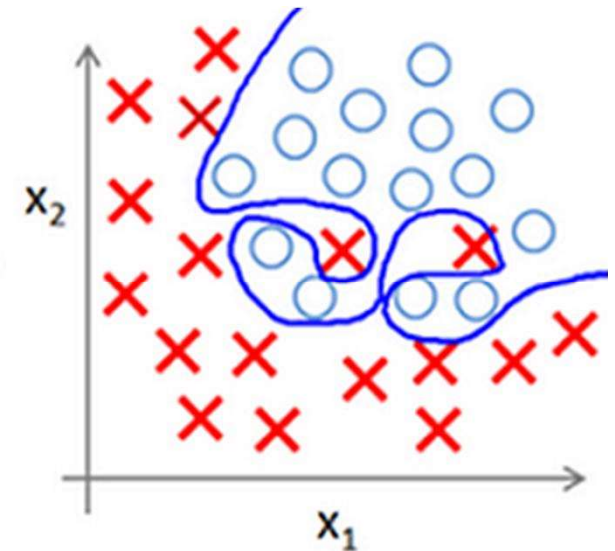
◆ **overfitting (過度擬合)**：訓練進行的太徹底，把訓練資料的所有特徵幾乎都學到了。只針對訓練資料集合的局部特徵，在測試時反而結果不佳。例如：訓練天鵝(或非天鵝)的特徵，知道了天鵝是有翅膀的，天鵝的嘴巴是長的彎的，脖子是長的有點曲度，體型像一個"2"。但很不巧訓練的天鵝全是白色的，於是機器會認為白的(局部特徵)才是天鵝，以後看到黑的天鵝就會認為那不是天鵝。



underfitting



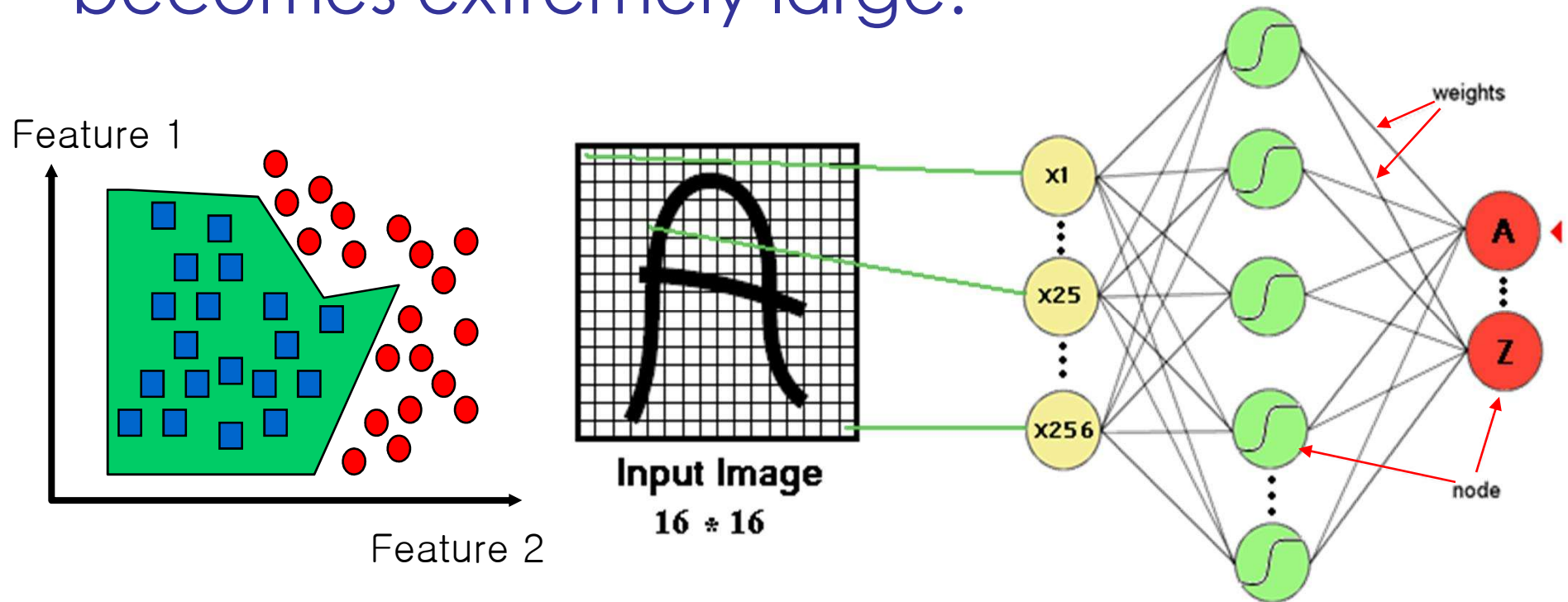
balanced



overfitting

Drawbacks of **Multi-layer** networks

- ☯ The number of trainable parameters becomes extremely large.



- ☯ Fully connected network of sufficient size can produce outputs that are invariant with respect to such variations.
 - ☹ But training time? Network size? Free parameters?