HW2
1. derivation of optimal W for a binary SVM
2. derivation of optimal W's for a 3-layer MLP
3. define a loss-function for SVM

## derivation of optimal W for a binary SVM

$$margin = \rho = \frac{2}{\|W\|}$$

$$max\rho \Leftrightarrow max\rho^2 \Leftrightarrow min\frac{1}{2}\|W\|^2$$

$$X_i^T W + b \geq +1, y_i = +1$$
$$X_i^T W + b \geq -1, y_i = -1$$

$$minJ(W) = min\frac{1}{2}\|W\|^2$$

$$s.t. \quad yi(X_i^T W + b) \geq 1, y_i = 1,2,...,n$$

$$L(W, b, \alpha) = \frac{1}{2}\|W\|^2 - \sum_{i=1}^{n} \alpha_i [y_i(X_i^T W + b) - 1]$$

$$maxL(W, b, \alpha) = +\infty$$

$$maxL(W, b, \alpha) = J(W) = \frac{1}{2}\|W\|^2$$

$$minmaxL(W, b, \alpha)$$

$$\nabla_W L(W, b, \alpha) = W - \sum_{i=1}^{n} \alpha_i y_i X_i = 0 \Rightarrow W = \sum_{i=1}^{n} \alpha_i y_i X_i$$

$$\nabla_b L(W, b, \alpha) = -\sum_{i=1}^{n} \alpha_i y_i = 0 \Rightarrow \sum_{i=1}^{n} \alpha_i y_i = 0$$

$$minL(W, b, \alpha) = -\frac{1}{2}\sum_{i=1}^{n}\sum_{j=1}^{n} a_i a_j y_i y_j X_i^T X_i + \sum_{i=1}^{n} a_i$$

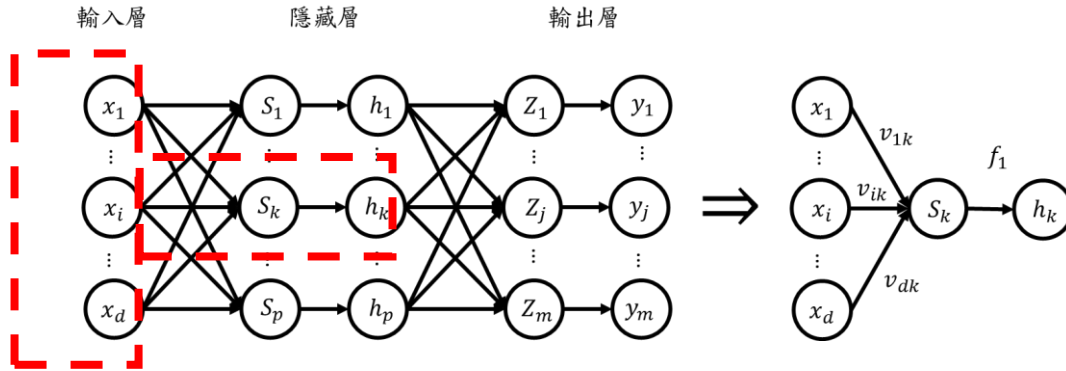$$min = \frac{1}{2}\sum_{i=1}^{n}\sum_{j=1}^{n} a_i a_j y_i y_j X_i^T X_i - \sum_{i=1}^{n} a_i$$

$$\sum_{i=1}^{n} \alpha_i y_i = 0$$

$$\alpha_i \geq 0, i = 1,2,...,n$$
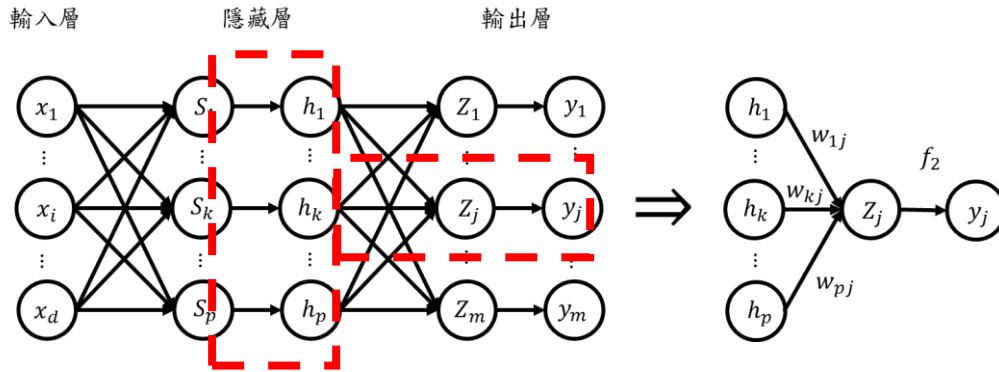
# 1. derivation of optimal W's for a 3-layer MLP

$$\{(x^{(i)}, y^{(i)})\}, i = 1, \dots, n, \qquad x_i \in R^d, y_i \in R^m$$

前向傳遞(Forward propagation)：



$$S_k = \sum_{i=0}^{d} v_{ik} x_i$$

$$h_k = f_1(S_k)$$



$$Z_j = \sum_{k=0}^{p} w_{kj} h_k$$

$$\hat{y}_j = f_2(Z_j)$$

# 2. <u>derivation of optimal W's for a 3-layer MLP</u>

運算法則

$$y = XW + b$$

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

輸入 X 乘以權重 W 得到 y，再通過啟動函數得到輸出（O）。 在這裡，啟動函數是 sigmoid 函數

$$\sigma(y) = \frac{1}{1 + e^{-y}}$$

E 是 loss 函數值，這裡是輸出值（output）與真實值（target）的歐式距離

$$E = \frac{1}{2}(O_0^1 - t)^2$$

E 的大小是評價感知器模型好壞的指標之一，w 權重是描述這個感知器模型的參數，通過計算 E 來優化感知器模型，即優化 w 的值

$w_{jk}^I$ 表示第 I 層，第 j 個輸入連結第 k 個輸出的權值 w。 以下先對一個權重（值）w 求得感知器模型的梯度

$$\frac{\partial E}{\partial w_{j0}^1} = (O_0^1 - t)\frac{\partial O_0^1}{\partial w_{j0}^1}$$

$$\frac{\partial E}{\partial w_{j0}^1} = (O_0^1 - t)\frac{\partial \sigma(x_0^1)}{\partial w_{j0}^1}$$

$$\frac{\partial E}{\partial w_{j0}^1} = (O_0^1 - t)\frac{\partial \sigma(x_0^1)}{\partial x_0^1}\frac{\partial x_0^1}{\partial w_{j0}^1}$$

$$\frac{\partial E}{\partial w_{j0}^1} = (O_0^1 - t)\sigma(x_0^1)(1 - \sigma(x_0^1))\frac{\partial x_0^1}{\partial w_{j0}^1}$$

$$\frac{\partial E}{\partial w_{j0}^1} = (O_0^1 - t)O_0^1(1 - O_0^1)x_j^0\frac{\partial x_0^1}{\partial w_{j0}^1}$$

$$\frac{\partial E}{\partial w_{j0}^1} = (O_0^1 - t)O_0^1(1 - O_0^1)x_j^0$$

現在把單個輸出的感知器模型推廣成多輸出感知器模型

$$\frac{\partial E}{\partial w_{jk}^1} = (O_k^1 - t_k)\frac{\partial O_k^1}{\partial w_{jk}^1}$$

$$\frac{\partial E}{\partial w_{jk}^1} = (O_k^1 - t_k)\sigma(x_k^1)(1 - \sigma(x_k^1))\frac{\partial x_k^1}{\partial w_{jk}^1}$$

$$\frac{\partial E}{\partial w_{jk}^1} = (O_k^1 - t_k)O_k^1(1 - O_k^1)x_j^0$$

# define a loss-function for SVM

用來評估模型的預測值與真實值不一致的程度，也是神經網絡中優化的目標函數，神經網絡訓練或者優化的過程就是最小化損失函數的過程，損失函數越小，說明模型的預測值就越接近真實值，模型的健壯性也就越好。

Hinge Loss function：

$$L(y, f(x)) = \max(0, 1 - yf(x))$$