② 2- Class Support Vector Machine

Let the seperating hyperplane be given by the equation

$$\text{sign} \{W^T x^{(i)} + W_0\} = \hat{y} \qquad \text{for } i \in N \quad \text{i.e, there are } N \text{ samples}$$

Let the labels be $y = \{-1, 1\}$ and $W, W_0$ are such that

when
$$W^T x^{(i)} + W_0 < 0 \Rightarrow \hat{y} = -1$$
$$W^T x^{(i)} + W_0 > 0 \Rightarrow \hat{y} = +1$$

$\Rightarrow$ Error occurs when $y^{(i)} (W^T x^{(i)} + W_0) < 0$

Let the loss function be $L(y, \hat{y})$

$$L(y, \hat{y}) = \sum_{i \in N} y^{(i)} \left( \vec{w}^T x^{(i)} + W_0 \right)$$

maximizing $L(y, \hat{y})$ reduces the errors

$W^*$ S·T $L(y, \hat{y})$ is minimized

But we need an initial $W$ & $W_0$ for this

Note:- $\dfrac{W^T x^{(i)} + W_0}{\|W\|}$ is the distance of the points $x^{(i)}$ from $W$ plane

∴ Let $y^{(i)} \left( W^T x^{(i)} + W_0 \right) \geqslant \mu \qquad \mu > 0$

$\Rightarrow \left( \begin{array}{c} \text{Gnd} \\ \text{truth} \end{array} \times \begin{array}{c} \text{notion of distance} \\ \text{with sign} \end{array} \right) \geq \begin{array}{c} \text{pure distance} \\ \text{measure} \end{array}$

Letting the distance to be thresholded to be greater than

$\mu$ reduces the chances of errors

$\Rightarrow \max_{W, W_0, \|W\|=1} \mu \left\{ \text{Subject to} \quad y^{(i)} \left( W^T x^i + W_0 \right) \geqslant \mu \atop \mu \geqslant 0 \right\}$

the contraint $\quad y^{(i)} \left[ W^T x^{(i)} + W_0 \right] \geqslant \mu$

can be changed to $\quad \dfrac{y^{(i)} \left[ W^T x^{(i)} + W_0 \right]}{\|W\|} \geqslant \mu$

Since $\mu$ is arbitrary

$\qquad$ & Let $\|W\| = \dfrac{1}{\mu}$ since arbitrary

$\Rightarrow \left\{ y^{(i)} \left[ W^T x^{(i)} + W_0 \right] \geqslant 1 \right\}$

& $\max_{W, W_0, \|W\|=1} \mu$ changes to $\quad \left[ \max_{W_0, \vec{W}} \frac{1}{2} \|\vec{W}\|^2 \right\}$

Note:- $\quad \max_{W, W_0} \|W\| , \quad \max_{W, W_0} \|W\|^2$ & $\max_{W, W_0} \frac{1}{2} \|W\|^2$

$\qquad$ yields the same results.

$\therefore$ Finally

$\qquad \max_{W, W_0} \frac{1}{2} \|W\|^2 \quad \text{S.T} \quad y^{(i)} \left[ W^T x^{(i)} + W_0 \right] \geqslant 1$

by introducing Lagrange multipliers, we can change it into an unconstrained problem

$$L_p = \left[ \frac{\|w\|^2}{2} - \sum_{i=1}^{N} \alpha_i \left\{ y^{(i)} \left[ w^T x^{(i)} + w_0 \right) - 1 \right\} \right] \quad \alpha_i \geq 0.$$

minimizing $L_p$ gives us the required weights

$\Rightarrow \nabla_{w^*} L_p = 0$ gives us $w^*$

$$\nabla_w L_p = \quad w - \nabla \sum_{i=1}^{N} \alpha_i \left\{ y^i \left( w^T x^{(i)} + w_0 \right) - 1 \right\} = 0$$

$$\Rightarrow \quad w - \sum_{i=1}^{N} \alpha_i y^{(i)} x^{(i)} = 0$$

$$\Rightarrow \quad \boxed{w^* = \sum_{i=1}^{N} \alpha_i y^{(i)} x^{(i)}} \quad —①$$

$$\nabla_{w_0} L_p = 0 \quad \Rightarrow \quad \nabla_{w_0} \sum_{i=1}^{N} \alpha_i \left( y^{(i)} \left[ w^T x^{(i)} + w_0 \right] - 1 \right) = 0$$

$$\Rightarrow \quad \boxed{\sum_{i=1}^{N} \alpha_i y^{(i)} = 0} \quad —②$$

Plugging ① and ② in primal form yields the dual form

$$L_D = \frac{1}{2} \left\| \sum_{i=1}^{N} \alpha_i y^{(i)} x^{(i)} \right\|^2 - \sum_{i=1}^{N} \alpha_i \left\{ y^{(i)} \left\{ \sum_{j=1}^{N} \alpha_j y^{(j)} x^{(j)T} x^i + w_0 \right\} - 1 \right\}$$

$$= \sum_{i=1}^{N} \alpha_i - \sum_{i=1}^{N} \alpha_i y^{(i)} w_0 - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i y^{(i)} x^{(i)T} \alpha_j y^{(j)} x^{(j)}$$

$$\sum_{i=1}^{N} \alpha_i y^{(i)} = 0 \quad \Rightarrow \quad \sum_{i=1}^{N} \alpha_i y^{(i)} w_0 = 0.$$

$$\Rightarrow L_D = \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \left( \alpha_i \alpha_j \, y^{(i)} y^{(j)} \, x^{(i)^T} x^{(j)} \right)$$

$$L_D = \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \left\| \sum_{i=1}^{N} \alpha_i \, y^{(i)} x^{(i)} \right\|^2$$

$$\boxed{L_D = \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \left\| \, diag(\alpha) \cdot Y \cdot X \, \right\|^2}$$

$$\left\{ \begin{array}{l} diag(\alpha) : \text{diagonal matrix with} \\ \text{diagonal elements as } \alpha \end{array} \right\}$$
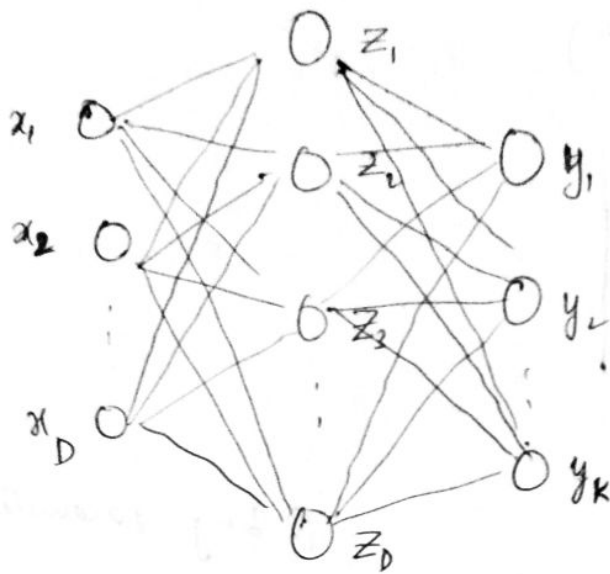
Convex optimizing $L_D$ yields values for $\alpha$

Subst. $\alpha$ in $w = \sum_{i=1}^{N} \alpha_i \, y^{(i)} x^{(i)}$ yields $w^*$

Additionally from KKT condition

$$\alpha_i \left[ y^{(i)} \left( \vec{w}^T x^{(i)} + w_0 \right) - 1 \right] = 0$$

subst. $\left( \begin{array}{c} \text{non zero} \\ \alpha_i \end{array} \right)$ in KKT gives $w_0$.

(3) Let each data point be of dimension $D$ & number of hidden nodes layer be $M$. Let output be of dimension $K$



Each node is connected to every other node in the immediate next layer. Let $A_{11}, A_2, \dots A_{1M}$ be the weights associated with $x_1$ to $\vec{z}$ or in general

$$\vec{A_m} = \{A_{d1}, A_{d2}, \dots A_{dM}\}$$ be weights from $\vec{x_d}$ to $\vec{z}$

in precise $A_{dm}$ be weight from $x_d$ to $z_m$. Let its bias be $A_{om}$. Be more like $\{A_{01}, A_{02}, \dots A_{om}\}$ be bias' for each hidden node.

$\Rightarrow$ Let the activation function be sigmoid.

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

$$Z_{dm} = \sigma\left(A_{mo} + \vec{A_{dM}}^T x_d\right)$$

$$\vec{Z_m} = \sigma\left(\vec{A_{mo}} + \vec{A_M}^T \vec{x}\right)$$

for the last layer, assigning weights and biases similarly, we g have

$$\hat{y}_k = g_k\left(\beta_{k0} + \vec{\beta_k}^{\,T} Z\right)$$

$$\boxed{g_k(\vec{x}) = \frac{e^{x_k}}{\sum\limits_{i=1}^{K} e^{x_i}}}$$

let the cost function be :

O being parameters

$$R(\theta) = \sum_{i=1}^{N} \|\, \vec{y}^{(i)} - \hat{\vec{y}}(x^{(i)})\|^2$$

$$R^{(i)}(\theta) = \sum_{k=1}^{K} \left( y_k^{(i)} - \hat{y}_k(\vec{x}^{(i)})\right)^2$$

find $\dfrac{\partial R^{(i)}(\theta)}{\partial A_{dm}} = 0$     $\dfrac{\partial R^{(i)}(\theta)}{\partial \beta_{mk}} = 0$     for locally optimal params $\theta$

$$\frac{\partial}{\partial \beta_{mk}} \sum \left( y_k^{(i)} - \hat{y}_k(x^{(i)})\right)^2 = 0$$

$$\rightarrow 2\left( y_k^{(i)} - \hat{y}_k(x^{(i)})\right)\left(-\frac{\partial}{\partial \beta_{mk}}\left(\hat{y}_k(x^{(i)})\right)\right) = 0$$

$$\hat{y}_k^{(i)} = g_k\left(\beta_{k0} + \vec{\beta}_{mk}\, x^{(i)}\right)$$

$$\rightarrow \frac{\partial R^{(i)}(\theta)}{\partial \beta_{mk}} = 0 \rightarrow -2\left( y_i - \hat{y}_i\right)\left( g'\left(\beta_{k0} + \vec{\beta}_{mk} Z^{(i)}\right)\right)\left( Z_m^{(i)}\right)$$

$$= \delta_k^{(i)} Z_m^{(i)}$$

$$\delta_k^{(i)} = -2\left( y - \hat{y}\right)\left(\frac{\partial}{\partial \beta_{km}} \hat{y}\right)$$

$$(11)^y \quad \frac{\partial R^{(i)}(\theta)}{\partial A_{dm}} = \frac{\partial}{\partial A_{dm}}\left( \sum_{i=1}^{K} (y-\hat{y})^2 \right)$$

$$= -2(y-\hat{y})\left( \frac{\partial}{\partial A_{dm}}(\hat{y}) \right)$$

$$= \sum_{k=1}^{K} -2(y_i - \hat{y}_i)\left( g'(\beta_0 + \beta_m Z)\right)\left( \frac{\partial}{\partial A_{dm}} Z \right)(\beta_{km})$$

$$= \sum_{k=1}^{K} \beta_{km}\, \delta_k^{(i)}\, \sigma'\left( \alpha_{mo} + \alpha_m^T x^{(i)} \right) x^{(i)}$$

$$\therefore \boxed{\begin{array}{l} \dfrac{\partial R^{(i)}(\theta)}{\partial A_{dm}} = S_m^{(i)}\, \vec{x}^{(i)} \\[2mm] S_m^{(i)} = \sum_{k=1}^{K}\left( \delta_k^{(i)} \beta_{km} \right) \sigma'\left( A_{mo} + A_m^T x^{(i)} \right) \end{array}}$$

Similar to Newton-Raphson approach, we can find optimal $\vec{\beta}$ and $\vec{A}$ by back propagation with some learning Rate

Learning Rate:- By what fraction of the gradient we're correcting our weights and biases

Back propagation:- for each epoch, we tune our weights by sending back information of gradients of cost functions at current datapoint and correct our current weights.

$$\left\{ \beta_{mk\varkappa}^{r+1} = \beta_{mk}^{r} - (Lr) \sum_{i=1}^{N} \frac{\partial R^{(i)}(\theta)}{\partial \beta_{ml}} \right\}$$

$$\left\{ A_{dm}^{r+1} = A_{dm}^{r} - (Lr) \sum_{i=1}^{N} \frac{\partial R^{(i)}(\theta)}{\partial A_{dm}} \right\}$$

$$\left\{ \begin{array}{l} \lambda = \text{epoch number} \\ Lr = \text{learning rate} \end{array} \right\}$$

④    Given the cost function is crossentropy loss function

$$R^{(i)}(\theta) = -\sum_{c=1}^{M} y_{\sigma c} \log(P_{\sigma,c})$$

$$y_{\sigma,c} = \left\{ \begin{array}{ll} 1 & \text{if } x \in \text{class } c \\ 0 & \text{else} \end{array} \right.$$

(or)   more precisely

for    $R_k^{(i)}(\theta)$    i.e.,   $x^{(i)} \longrightarrow y^{(k)}$ (or) class $k$

$$R_k^{(i)}(\theta) = -y_k \log\left(P_r\{x^{(i)} \in C_k\}\right)$$

generay   $\hat{y} : P_r\{x^{(i)} \in C_k\}$

$$\boxed{R^{(i)}(\theta) = \sum_{k=1}^{K} -y_{i,k} \log\left(\hat{y}\right)_{i,k}}$$

for optimal solution consider:

$$\frac{\partial R^{(i)}(\theta)}{\partial A_{dm}} = \sum_{k=1}^{K} \frac{\partial}{\partial A_{dm}} \left( -y_{i,k} \left( \log \left( \hat{y}_{i,k} \right) \right) \right)$$

$$\frac{\partial R^{(i)}(\theta)}{\partial A_{dm}} = \sum_{k=1}^{K} \frac{-y_{i,k}}{\hat{y}_{i,k}} \cdot \frac{\partial}{\partial A_{dm}} \left( \hat{y}_{i,k} \right)$$

$$\frac{\partial R^{(i)}(\theta)}{\partial A_{dm}} = \sum_{k=1}^{K} \left( \frac{-y_{i,k}}{\hat{y}_{i,k}} \right) \left( g' \left( \beta_{ok} + \beta_{mk} z_k \right) \right) \left( \sigma' \left( A_{od} + A_{dm} x^{(i)} \right) \right) x^{(i)}$$

$$\frac{\partial R^{(i)}(\theta)}{\partial \beta_{mk}} = \sum_{k=1}^{K} \frac{-y_{i,k}}{\hat{y}_{i,k}} g' \left( \beta_{ok} + \beta_{mk} z_k \right) z_k$$

$$\boxed{\frac{\partial R^{(i)}(\theta)}{\partial \beta_{mk}} = \delta_k^{(i)} z_k}$$

$$\boxed{\delta_k^{(i)} = \sum_{k=1}^{K} \frac{-y_{i,k}}{\hat{y}_{i,k}} g'\left( \beta_{ok} + \beta_{mk} z_k \right)}$$

$$\frac{\partial R^{(i)}(\theta)}{\partial A_{dm}} = \left\{ \delta_k^{(i)} \cdot \sigma'\left( A_{od} + A_{dm} x^{(i)} \right) \right\} x^{(i)} \cdot \beta_{mk}$$

$$\boxed{\frac{\partial R^{(i)}(\theta)}{\partial A_{dm}} = \delta_m^{(i)} x^{(i)}}$$

$$\boxed{\delta_{(m)}^{(i)} = \sum_{k=1}^{K} \frac{-y_{i,k}}{\hat{y}_{i,k}} g''(\beta) \sigma'(A) \beta_{mk}}$$

① Let the decision hyperplane b/w class $j$ & $k$ be

$$W_k^T x^i + W_{k0} = \hat{y}_k$$

Let $\quad W^T x^i + W_0 = \hat{y}_j > 0 \quad$ if $x_i \in j$

$$\hat{y}_k > 0 \qquad x_i \in k$$

now consider two points in $j$

$$x_1 \ \& \ x_2$$



any point $x_\lambda$ b/w $x_1, x_2$ is given by

$$x_\lambda = \lambda x_2 + (1-\lambda) x_1 \qquad \boxed{\lambda \geq 0} \quad \underline{\lambda \in [0,1]}$$

$$W^T x_\lambda + W_0 < 0 \quad \text{if } \underbrace{\text{the separating hyperplane}}$$
$$\underbrace{\text{is convex}}$$

$$W^T (\lambda x_2 + (1-\lambda) x_1) + W_0$$

$$= \lambda W^T x_2 + W_0 + W^T (1-\lambda) x_1$$

Now for j-k boundary, consider.

$$x^i \in C_k \quad \text{if} \quad \hat{y}_k(x^{(i)}) - \hat{y}_j(x^{(i)}) > 0$$

$$(w_k - w_j)^T x^{(i)} + (w_{ok} - w_{oj}) > 0.$$

$$y_k(\hat{x}) = \lambda y_k(x_1) + (1-\lambda) y_k(x_2)$$

if $x_1, x_2 \in C_k$

$$y_k(x_1), \ y_k(x_2) > 0$$

$$\Rightarrow \ \lambda y_k(x_1) + (1-\lambda) y_k(x_2) > 0$$

$$\Rightarrow \ \underline{\underline{\hat{x} \in C_k}}$$

$$\Rightarrow \boxed{y_k \text{ is convex}}$$