Bhavsar Jeel Alpesh
ES16BTECH11005

1) Region $R_k$ is defined as

$$R_k = \{ \bar{x} \mid y_k(\bar{x}) > y_j(\bar{x}) \quad \forall j \neq k \}$$

$$y_k(\bar{x}) = W_k^T \bar{x} + W_{k0}$$

~~If $R_k$ is convex,~~

let $\bar{x}_1, \bar{x}_2 \in R_k$

If $R_k$ is convex, $\bar{x}_0 = \lambda \bar{x}_1 + (1-\lambda) \bar{x}_2$
must also belong to $R_k$ $\forall \lambda \in [0,1]$

So $y_k(\bar{x}_0) = y_k(\lambda \bar{x}_1 + (1-\lambda) \bar{x}_2)$

Given $y_k(\bar{x})$ is linear

So $y_k(\bar{x}_0) = \lambda y_k(\bar{x}_1) + (1-\lambda) y_k(\bar{x}_2)$

Now $\lambda y_k(\bar{x}_1) > \lambda y_j(\bar{x}_1)$    —①

$(1-\lambda) y_k(\bar{x}_2) > \lambda y_j(\bar{x}_2)$    —②

So adding ① & ②,

$$y_k(\bar{x}_0) > y_j(\bar{x}_0)$$

Hence $\bar{x}_0$ belongs to $R_k$

$R_k$ is convex.

2) two class SVM

let the labels $y \in \{1, -1\}$

~~Also, for~~ $\hat{g}(\bar{x}) = w^T \bar{x} + w_0$

If $\hat{g}(\bar{x}) > 0$, $\bar{x}$ belongs to label 1

$\hat{g}(\bar{x}) < 0$, $\bar{x}$ belongs to label $-1$.

So if $(\hat{g}(\bar{x})) * y < 0$, the it is a wrong prediction

Loss function $L(y^{(i)}, \hat{g}^{(i)}(\bar{x})) = y^{(i)}(w^T \bar{x}^{(i)} + w_0) < 0$

For M errors or M miss classifications,

$$L(y, \hat{g}(x)) = \sum_{i \in M} y^{(i)}(w^T \bar{x}^{(i)} + w_0)$$

We need to find $\bar{w}$ that maximizes $L(y, \hat{g}(x))$

$$\max_{\bar{w}, w_0 \|\bar{w}\|=1} u \qquad \text{~~subject to~~} ; u > 0$$

subject to $y^{(i)}(\bar{w}^T \bar{x}^{(i)} + w_0) \geq u$

$1 \leq i \leq N$

We remove the constraint $\|\bar{w}\| = 1$ by

$$\max_{\bar{w}, w_0} u \qquad \text{subject to} \quad \frac{y^{(i)} \cdot (\bar{w}^T \bar{x}^{(i)} + w_0)}{\|\bar{w}\|} \geq u$$

$$= y^{(i)}(\bar{w}^T \bar{x}^{(i)} + w_0) \geq u \|\bar{w}\|$$
$$u > 0$$

Arbitrarily assume $\|\bar{\omega}\| = \frac{1}{\mu}$

The optimization is equivalent to

$$\min_{\bar{\omega}, \omega_0} \frac{1}{2}\|\omega\|^2 \text{ such that } y^{(i)}(\omega^T x^{(i)} + \omega_0) \geq 1$$
$$1 \leq i \leq N$$

Converting it into unconstrained problem,

$$L_p = \min_{\omega_0, \bar{\omega}} \frac{1}{2}\|\omega\|^2 - \sum_{i=1}^{N} \alpha_i \left[ y^{(i)}(\bar{\omega}\bar{x}^{(i)} + \omega_0) - 1 \right] \quad \text{①}$$

$$\nabla_\omega L_p = 0$$

i.e. $\qquad \frac{\partial L_p}{\partial \omega_0} = 0$

$$- \sum_{k=1}^{N} \alpha_i y^{(i)} = 0 \quad \text{②}$$

$$\frac{\partial L_p}{\partial \bar{\omega}} = 0$$

$$\bar{\omega} = \sum_{i=1}^{N} \alpha_i x^{(i)} y^{(i)} \quad \text{③}$$

We cannot solve ② & ③ without values of $\alpha$

Plug ③ & ② in ①,

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

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

$$- \sum_{i=1}^{N} \alpha_i \left[ y^{(i)} (\omega^T x^{(i)}) - 1 \right]$$

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

$$s.t \quad \alpha_i \geq 0 \; ; \quad \sum_{i=1}^{N} \alpha_i y^{(i)} = 0 \qquad \textcircled{A}$$

In addition to constraints in $\textcircled{A}$, optimal $\alpha_i$ must satisfy $\alpha_i\left[ y^{(i)} (\omega^T x^{(i)} + \omega_0) - 1 \right] = 0$
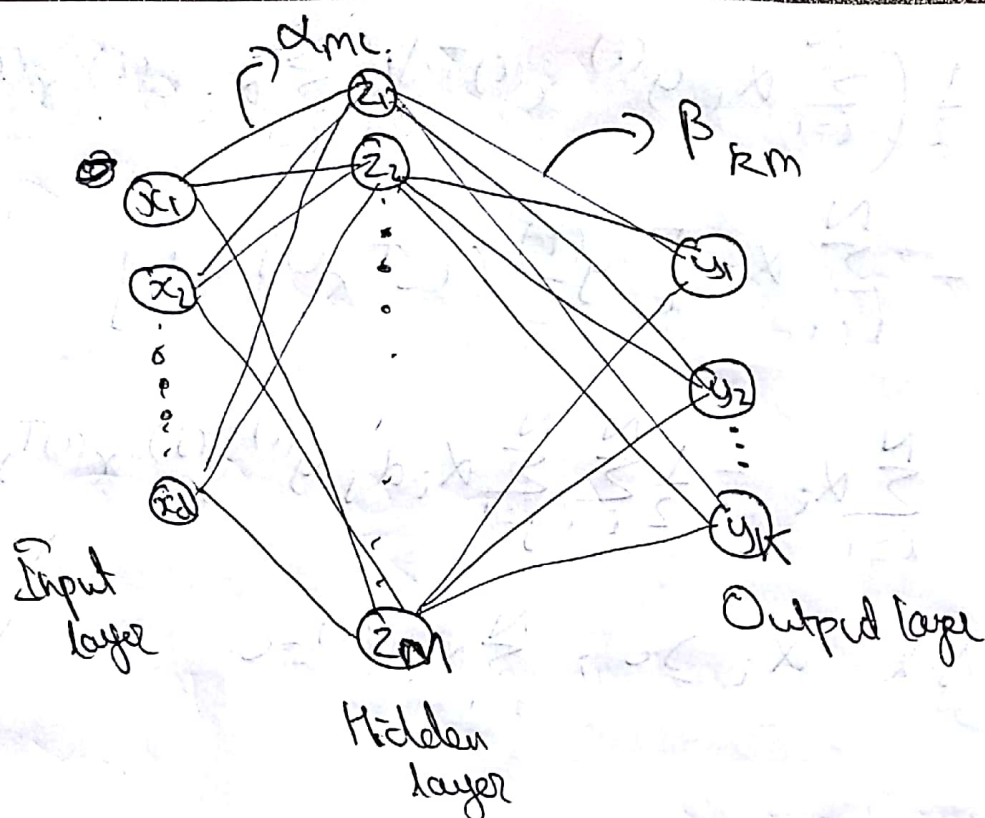
If $\alpha_i = 0$,

$x^{(i)}$ cannot help us find $\omega$

If $\alpha_i > 0$,

$$y^{(i)} (\omega^T x^{(i)} + \omega_0) - 1 = 0$$

$x^{(i)}$ will determine the separating hyperplane

3)



$\rightarrow \alpha_{mc}$

$\rightarrow \beta_{km}$

Input layer

Hidden layer

Output layer

$$Z_m = \sigma\left(\alpha_{mo} + \bar{\alpha}_m^T \bar{x}\right) \qquad 1 \leq m \leq M$$

Sigmoid func$^n$: $\sigma(x) = \dfrac{1}{1 + e^{-x}}$

$$\alpha_m = \left[\alpha_{m1} \cdots \alpha_{md}\right]^T$$

$\alpha_{mo}$ : bias associated with $m^{th}$ hidden node

$$g_k(\bar{x}) = g_k\left(\beta_{ko} + \bar{\beta}_k^T \bar{Z}\right)$$

where $g_k(\bar{x}) = \dfrac{e^{x_k}}{\sum\limits_{j=1}^{K} e^{x_j}}$ (softmax)

Generally,
$g(x)$ is softmax for classification
and Cross entropy for regression.

We are assuming $g_k(x)$ to be sigmoid
function in this problem

Parameters $\left[\theta : a_{n_0}, \vec{a}_m, \beta_{k_0}, \vec{\beta}_k\right]$

$$1 \leq m \leq M$$
$$1 \leq k \leq K$$
$$\vec{a}_m \in \mathbb{R}^q \quad , \quad \vec{\beta}_k \in \mathbb{R}^M$$

Cost func$^n$: $P(\theta) = \sum_{i=1}^{N} \sum_{k=1}^{K} \left( y_k^{(i)} - \hat{g}_{(k)}(x^{(i)}) \right)^2$

$$= \sum_{i=1}^{N} R^{(i)}(\theta)$$

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

~~For finding optimal~~

For finding locally optimal parameters $\theta$,

$$\frac{\partial R^{(i)}(\theta)}{\partial \beta_{km}} = \frac{\partial}{\partial \beta_{km}} \sum_{k'=1}^{K} \left( y_{k'}^{(i)} - g_{k'}^{(i)} \right)^2$$

$$= \frac{\partial}{\partial \beta_{km}} \sum_{k'=1}^{K} \left( y_{k'}^{(i)} - g_{k'}(\beta_{k'0} + \vec{\beta}_{k'}^T z^{(i)}) \right)^2$$

$$\frac{\partial R^{(i)}(\theta)}{\partial \beta_{km}} = 2 \underbrace{\left( y_k^{(i)} - g_k(\beta_{k0} + \vec{\beta}_k^T z^{(i)}) \right) \left( -g_k'(\beta_{k0} + \beta_k^T z^{(i)}) z_m^{(i)} \right)}_{\delta_k^{(i)}}$$

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

$$\frac{\partial R^{(r)}(\theta)}{\partial \alpha_{m\ell}} = S_m^{(i)} x_\ell^{(i)}$$

where $S_m^{(i)} = \left(\sum_{k=1}^{K} \delta_k^{(i)} \beta_{km}\right) \sigma'\left(\alpha_{m0} + \alpha_m^T \bar{x}^{(i)}\right)$

$$\sigma'(x) = \frac{\partial}{\partial x} \sigma(x)$$

Now update the parameters with gradient descent.

$$\beta_{km}^{(r+1)} = \beta_{km}^{(r)} - \gamma_r \sum_{i=1}^{N} \frac{\partial R^{(i)}(\theta)}{\partial \beta_{km}^{(r)}}$$

$$\alpha_{m\ell}^{(r+1)} = \alpha_{m\ell}^{(r)} - \gamma_r \sum_{i=1}^{N} \frac{\partial R^{(i)}(\theta)}{\partial \alpha_{m\ell}^{(r)}}$$

$\gamma_r \rightarrow$ learning rate

& The convergence of $\alpha$ & $\beta$ largely depends on the learning rate $\gamma_r$.

Subsequently we can find $\hat{z}$ & then $\hat{y}$

4) For cross entropy loss function,

$$R(\theta) = -\sum_{i=1}^{N} \sum_{j=1}^{k} y_j^{(i)} \log\left(\hat{y}_j(x^{(i)})\right)$$

Follow the same procedure as in Question 3
solution until

$$R(\theta) = \sum_{i=1}^{N} R^{(i)}(\theta)$$

$$R^{(i)}\theta = -\sum_{k=1}^{k} y_k^{(i)} \log \hat{y}_k(x^{(i)})$$

$$\frac{\partial R^{(i)}(\theta)}{\partial \beta_{km}} = -\sum_{k'=1}^{k} \cdot y_{k'}^{(i)} \frac{\partial}{\partial \beta_{km}}\left(\log \hat{y}_{k'}(x^{(i)})\right)$$

$$= -\sum_{k'=1}^{k} y_{k'}^{(i)} \frac{1}{\hat{y}_{k'}(x^{(i)})} g'(\beta_{k'0} + \beta_{k'}^T z) z_m^{(i)}$$

$$= -y_k^{(i)} \frac{1}{\hat{y}_k(x^{(i)})} g'(\beta_{k0} + \beta_k^T z) z_m^{(i)}$$

$$\underbrace{\qquad\qquad\qquad}_{\delta_k^{(i)}}$$

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

Similarly we find

$$\frac{\partial R^{(i)}(\theta)}{\partial \alpha_{m\ell}} = \delta_m^{(i)} x_\ell^{(i)} = \left(\sum_{k=1}^{k} \delta_k^{(i)} \beta_{km}\right) \sigma'(\alpha_{m0} + \alpha_m^T x^{(i)})$$

Now update $\beta$ & $\alpha$ with gradient descent
And subsequently find $z$ & $\hat{y}$ as shown in Ans 3.