

(1) For  $k$  class Linear discriminant classifier,

we consider single  $k$ -class discriminant comprising of  $k$  linear functions of the form

$$y_k(x) = w_k^T x + w_{k0}$$

and assigning a point  $x$  to class  $c_k$  if

$$y_k(x) > y_j(x) \quad \forall j \neq k$$

The decision boundary is thus given by

$$y_k(x) = y_j(x) \quad \text{given by}$$

$$(w_k - w_j)^T x + (w_{k0} - w_{j0}) = 0$$

This decision region is always singly connected and convex (proof below)

To prove this,

consider two points  $x_A, x_B$  inside the decision region  $R_k$ .

Any point  $\vec{x}$  that lies on the line connecting  $\vec{x}_A$  and  $\vec{x}_B$  is of the form

$$\vec{x} = \lambda \vec{x}_A + (1-\lambda) \vec{x}_B$$

where  $0 \leq \lambda \leq 1$  From linearity of discriminant functions,

$$y_k(\vec{x}) = \lambda y_k(\vec{x}_A) + (1-\lambda) y_k(\vec{x}_B)$$

$$> \lambda y_j(\vec{x}_A) + (1-\lambda) y_j(\vec{x}_B)$$

$$> y_j(\vec{x})$$

$$\boxed{y_k(\vec{x}) > y_j(\vec{x})}$$

Proves that  $R_k$  is convex

Optimal Separating hyperplane  
Two class supervised SVM  
assume  $y \in \{-1, +1\}$

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

If we make errors or misclassification

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

find  $\vec{w}$  so that  $L(y, \hat{y}(\vec{x}))$  is maximized  
(note: - negative sign of  $L$ )

Drawbacks: Solution depends on initial choice of  $\vec{w}$

we can do better

$$\max U \quad ; \quad U > 0$$

$$\vec{w}, w_0, \|\vec{w}\| = 1$$

subject to the condition

$$y^{(i)} (\vec{w}^T \vec{x}^{(i)} + w_0) \geq U \quad 1 \leq i \leq N$$



Lets remove the  $\|w\| = 1$  constraint

Lets reformulate,

max  $\bar{u}$

$\vec{w}, w_0$

subject to condition

$$\boxed{u > 0}$$

$$(y^{(i)}) \left( \frac{w^T x^{(i)} + w_0}{\|w\|} \right) \geq u$$

$$\Rightarrow y^{(i)} (w^T x^{(i)} + w_0) \geq u \|w\| \quad \boxed{u > 0}$$

If  $\|w\| = \frac{1}{u}$  the the above optimization

is equivalent to

$$\min_{w_0, \vec{w}} \frac{1}{2} \|w\|^2$$

such that (condition)

$$y^{(i)} (w^T x^{(i)} + w_0) \geq 1$$

Using Lagrangian multipliers concept,  
This converts as an unconstrained problem.

$$L_p = \min_{\vec{\omega}, \omega_0} \frac{\|\omega\|^2}{2} - \left( \sum_{i=1}^N \alpha_i (y^i (\vec{\omega}^T x^i + \omega_0) - 1) \right)$$

Setting

$$\nabla_{\vec{\omega}} L_p = 0$$

$$\boxed{\vec{\omega} = \sum_{i=1}^N \alpha_i y^i \vec{x}^i} \quad - (1)$$

4

$$\nabla_{\omega_0} L_p = 0$$

$$\boxed{\sum_{i=1}^N \alpha_i y_i = 0} \quad - (2)$$

Putting  $\vec{\omega}$  in  $L_p$  & careful simplification  
yields the dual of  $L_p$  as

$$\text{maximize,}$$

$$L_D = \sum_{i=1}^N \alpha_i = \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i y^i \alpha_j y^j (x^i)^T x^j \quad L(3)$$

s.t

$$\alpha_i \geq 0 \quad \forall i, \quad \sum_{i=1}^N \alpha_i y_i = 0 \quad - (4)$$

In addition to this our optimal solutions must follow the KKT condition

$$\alpha_i [y_i (\omega^T x_i + \omega_0) - 1] = 0 \quad (5)$$

Thus from the conditions (4) and (5)

ie  $\alpha_i \geq 0 \forall i$        $\alpha_i [y_i (\omega^T x_i + \omega_0) - 1] = 0$

(1)  $\alpha_i \neq 0$  ie  $\alpha_i > 0$

$$y_i (\omega^T x_i + \omega_0) - 1 = 0$$

$\rightarrow \omega_0$  can be found from here  
corresponding to that non zero  $\alpha_i$

(2)  $\alpha_i \neq 0$  ie  $\alpha_i > 0$

They contribute to

weight vector since

$$\bar{\omega} = \sum_{i=1}^N \alpha_i y_i^{(1)} \bar{x}_i^{(1)}$$

These vectors are called support vectors

(3) Let us take the case discussed in

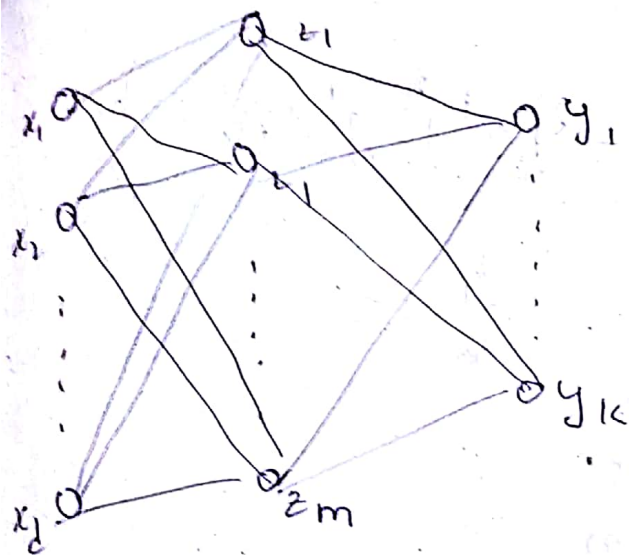
the class  
there is input layer then 1 hidden layer  
and then output layer.

Let the complete set of weights be

$$\{ \alpha_m, \vec{\alpha}_m ; m = 1, 2, \dots, M \} \quad M(d+1) \text{ weights}$$

$$\{ \beta_k, \vec{\beta}_k ; k = 1, 2, \dots, K \} \quad K(M+1) \text{ weights}$$





where

$$z_m = \sigma (d_0 m + d_m^T x) \quad m = 1, \dots, M$$

$$y_k = g_k (\beta_0 k + \beta_k^T z) \quad k = 1, \dots, K$$

Here,

$\sigma$  is Sigmoid

$g$  is Softmax function.

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

They both are activation

functions.

$$g_k^T = \frac{e^{T_k}}{\sum_{i=1}^K e^{T_i}}$$

where,

$$z = [z_1, z_2, \dots, z_m]$$

$$y = (y_1, y_2, \dots, y_k)$$

We use Sum of squared errors as our measure of fit

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

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

Since functions such as Sigmoid, Softmax etc which have been used are not convex we only have a optimal 'LOCAL' minimum

The generic approach to minimizing  $R(\theta)$  is by gradient descent called back propagation

Here is back propagation in detail for squared error loss function.

$$\text{Let } z_{mn}^{(i)} = \sigma(w_{mn} + d_{mn} z_m^{(i)})$$

$$\vec{z}^{(i)} = (z_1^{(i)}, z_2^{(i)}, \dots, z_m^{(i)})$$

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



With derivatives;

$$\frac{\partial R_i}{\partial \beta_{km}} = \left( -2(y_k^{(i)} - \hat{y}_k(x^{(i)})) g'_k(\beta_k^T z_i) \right) z_m^{(i)}$$

$\downarrow$   
 $s_k^{(i)}$

$$\frac{\partial R_i}{\partial \beta_{km}} = s_k^{(i)} z_m^{(i)}$$

$$\frac{\partial R_i}{\partial \alpha_{ml}} = \left( -\sum_{k=1}^K 2(y_k^{(i)} - \hat{y}_k(x^{(i)})) g'_k(\beta_k^T z_i) \beta_{km} \right) \frac{\partial \alpha_{ml}}{\partial x_l^{(i)}}$$

$\downarrow$   
 $s_{mi}$

Given these derivatives,

A gradient descent update at the  $(r+1)^{th}$  iteration

has the form

$$\beta_{km}^{(r+1)} = \beta_{km}^{(r)} - \gamma_r \sum_{i=1}^N \frac{\partial R_i}{\partial \beta_{km}^{(r)}}, \quad \text{where } \gamma_r \text{ is the learning rate.}$$

$$d_{ml}^{(r+1)} = d_{ml}^{(r)} - \gamma_r \sum_{i=1}^N \frac{\partial R_i}{\partial d_{ml}^{(r)}}$$

At each iteration the value of weights <sup>4 bias</sup> are being updated, thus this is known as back propagation since gradient <sup>of cost fn.</sup> from previous

value is used to update the weights to new value,

The expression for bias is

$$\beta_{0k}^{(r+1)} = \beta_{0k}^{(r)} - \eta_r \sum_{i=1}^N s_k^{(i)} \quad \text{Since } \frac{\partial R^{(i)}(\theta)}{\partial \beta_{0k}} = -s_k^{(i)}$$

$$d_{0m}^{(r+1)} = d_{0m}^{(r)} - \eta_r \sum_{i=1}^N s_m^{(i)} \quad \text{Since } \frac{\partial R^{(i)}(\theta)}{\partial d_{0m}} = s_m^{(i)}$$

This is similar to Newton Raphson method of finding  $f(x) = 0$  here our  $f(x)$  is  $\nabla R(\theta)$ .  
By back propagating with some learning rate.

Learning Rate :- Converging rate is by what fraction of gradient we are correcting error and updating the weights (or) bias.

(5)

For the above qn,

If the remaining set up being same as qn 4.

$$R(\theta) = - \sum_{i=1}^N \sum_{k=1}^K y_k^{(i)} \log [\hat{y}_k^{(i)}]$$
$$= + \sum_{i=1}^N R^i(\theta)$$

$$\text{where } R^i(\theta) = - \sum_{k=1}^K y_k^{(i)} \log [\hat{y}_k(\vec{x}^i)]$$

with derivatives

$$\frac{\partial R^i(\theta)}{\partial \beta_{km}} = \left[ \frac{-y_k^{(i)}}{\hat{y}_k(\vec{x}^i)} g'_k(\beta_{0k} + \beta_k^T \vec{z}^i) \right] z_m^i$$
$$= \delta_k^{(i)} z_m^i$$

$$\frac{\partial R^i(\theta)}{\partial \alpha_{ml}} = \left[ - \sum_{k=1}^K \left[ \frac{y_k^{(i)}}{\hat{y}_k(\vec{x}^i)} \right] g'_k(\beta_{0k} + \beta_k^T \vec{z}^i) \beta_{km} \cdot \sigma'(\alpha_{0m} + \alpha_m^T \vec{x}^i) \right] x_l^i$$
$$= \delta_m^{(i)} x_l^i$$

A gradient descent update at the  $(r+1)^{\text{th}}$  iteration

has the form



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

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

$$\beta_{ok}^{(r+1)} = \beta_{ok}^{(r)} - \gamma_r \sum_{i=1}^N \frac{\partial R^{(i)}(\theta)}{\partial \beta_{ok}^{(r)}}, \text{ where } \frac{\partial R^{(i)}(\theta)}{\partial \beta_{ok}^{(r)}} = s_k^{(i)}$$

$$d_{om}^{(r+1)} = d_{om}^{(r)} - \gamma_r \sum_{i=1}^N \frac{\partial R^{(i)}(\theta)}{\partial d_{om}^{(r)}}, \text{ where } \frac{\partial R^{(i)}(\theta)}{\partial d_{om}^{(r)}} = s_m^{(i)}$$

At each iteration the value of weights & bias are being updated from the gradient of cost function at previous value is used to update the weights to new value justifying the name back propagation.