S. Suprabath Reddy
EE15BTECH11026

1) For k-class linear discriminant classifier,

$$\hat{y}_k(\vec{x}) = W_k^T \vec{x} + W_{k0}$$

If $\hat{y}_k(\vec{x}) > \hat{y}_j(\vec{x})$ for $j \neq k$,

then output is class $k$.

Consider two points $\vec{x}_A$ and $\vec{x}_B$ which are lying in the decision region $R_k$.

Any point $\hat{\vec{x}}$ that lies on line connecting $\vec{x}_A$ and $\vec{x}_B$

is $\hat{\vec{x}} = \lambda \vec{x}_A + (1-\lambda) \vec{x}_B$ where $0 \leq \lambda \leq 1$.

From the linearity of discriminant function.

$$\hat{y}_k(\hat{\vec{x}}) = \lambda y_k(\vec{x}_A) + (1-\lambda) y_k(\vec{x}_B) \quad \text{—①}$$

Since $\vec{x}_A$ and $\vec{x}_B$ lie inside $R_k$

$$\hat{y}_k(\vec{x}_A) > \hat{y}_j(\vec{x}_A) \quad \& \quad \hat{y}_k(\vec{x}_B) > \hat{y}_j(\vec{x}_B)$$

① $\Rightarrow \hat{y}_k(\hat{\vec{x}}) > \lambda(\hat{y}_j(\vec{x}_A)) + (1-\lambda) y_j(\vec{x}_B)$

$$\hat{y}_k(\hat{\vec{x}}) \geq \hat{y}_j(\hat{\vec{x}})$$

∴ $\hat{\vec{x}}$ lies inside $R_k$.

∴ $R_k$ is convex

2)

maximin of $u$ subject to.

$$y^{(i)} \cdot \frac{(\vec{w}^T \vec{x}^{(i)} + w_o)}{\|\vec{w}\|} \geq u$$

$$\Rightarrow y^{(i)}(\vec{w}^T \vec{x}^{(i)} + w_o) \geq u.\|\vec{w}\|$$

If $\|\vec{w}\| = \frac{1}{u}$

$\Rightarrow$ Maximization of $u$ $\Rightarrow$ minimization of $\|\vec{w}\|^2$

such that $y^{(i)}(\vec{w}^T \vec{x}^{(i)} + w_o) \geq 1$

$$L_p = \min_{w_o, \vec{w}} \frac{1}{2}\|\vec{w}\|^2 - \sum_{i=1}^{N} \alpha_i \left[ y^{(i)}(\vec{w}^T \vec{x}^{(i)} + w_o) - 1 \right]$$
$$\text{where } \alpha_i \geq 0 \qquad \text{①}$$

We want to min $L_p$ w.r.t. $\vec{w}$ and $\vec{w}_o$.

Set $\nabla_{\vec{w}} L_p = 0$ $\Rightarrow$ $\nabla_{\vec{w}} \left( \frac{1}{2} \vec{w}^T \vec{w} - \sum_{i=1}^{N} \alpha_i \left( y^{(i)}(\vec{w}^T \vec{x}^{(i)} + w_o) - 1 \right) \right) = 0$

$$\Rightarrow \vec{w} = \sum_{i=1}^{N} \alpha_i y^{(i)} \cdot \vec{x}^{(i)} \qquad \text{②}$$

Set $\frac{\partial L_p}{\partial w_o} = 0$ $\Rightarrow$ $\sum_{i=1}^{N} \alpha_i \cdot y^{(i)} = 0 \qquad \text{③}$

From ② & ③, solve fo $\alpha_i$ by substituting in ①.

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

such that $\alpha_i \geq 0$ and $\sum_{i=1}^{N} \alpha_i y^{(i)} = 0$

In addition to above condition, our optimal solutions $\alpha_i's$ must satisfy

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

**3)** Output at the $m^{th}$ layer

$$Z_m = \sigma(\alpha_{mo} + \vec{\alpha}_m^T \cdot \vec{x})$$

where $\sigma(x) = \dfrac{1}{1 + e^{-x}}$ ; $\vec{\alpha}_m = [\alpha_{m1} \cdots \alpha_{md}]^T$

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

$$\hat{y}_k(\vec{x}) = g_k(\beta_{ko} + \vec{\beta}_k^T \cdot \vec{z})$$

where $g_k(\vec{x}) = \dfrac{e^{x_k}}{\sum\limits_{j=1}^{k} e^{x_j}}$

Parameters $= \{\theta : \alpha_{mo}, \vec{\alpha}_m, \beta_{ko}, \vec{\beta}_k\}$ ; $1 \le m \le M, 1 \le k \le K$

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

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

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

To find locally optimal parameter $\theta$, find

$$\dfrac{\partial R^{(i)}(\theta)}{\partial \beta_{km}} = \dfrac{\partial}{\partial \beta_{km}} \sum_{k'=1}^{K} (y_{k'}^{(i)} - \hat{y}_{k'}^{(i)})^2$$

$$= \dfrac{\partial}{\partial \beta_{km}} \sum_{k'=1}^{K} (y_{k'}^{(i)} - g_{k'}(\beta_{k'o} + \vec{\beta}_{k'}^T \cdot \vec{z}^{(i)}))^2$$

$$\not\Rightarrow \dfrac{\partial R^{(i)}(\theta)}{\partial \beta_{km}} = 2(y_{k'}^{(i)} \cdots g_k \cdots \beta_k$$

$$\rightarrow \dfrac{\partial R^{(i)}(\theta)}{\partial \beta_{km}} = 2(y_k^{(i)} - g_k(\beta_{ko} + \vec{\beta}_k^T \cdot \vec{z}^{(i)})).$$
$$(-g_k' \cdot (\beta_{ko} + \vec{\beta}_k^T \cdot \vec{z}^{(i)})) \cdot z_m^{(i)}$$

$$= \delta_k^{(i)} \cdot z_m^{(i)} \quad \text{where} \quad \delta_k^{(i)} = -2(y_k^{(i)} - g_k(\beta_{ko} + \vec{\beta}_k^T \cdot \vec{z}^{(i)})).$$
$$g_k'(\beta_{ko} + \vec{\beta}_k^T \cdot \vec{z}^{(i)})$$

$$\rightarrow \frac{\partial R^{(i)}(\theta)}{\partial \alpha_{m\lambda}} = S_m^{(i)} \cdot x_\lambda^{(i)} \quad -②$$

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

$$\lhook ③$$

Final updated weights

$$\left. \begin{array}{l} \beta_{km}^{(\lambda+1)} = \beta_{km}^{(\lambda)} - \gamma_\lambda \sum_{i=1}^{\gamma} \frac{\partial R^{(i)}(\theta)}{\partial \beta_{km}^{(\gamma)}} \\[4mm] \alpha_{m\lambda}^{(\lambda+1)} = \alpha_{m\lambda}^{(\lambda)} - \gamma_\lambda \cdot \sum_{i=1}^{N} \frac{\partial R^{(i)}(\theta)}{\partial \alpha_{m\lambda}^{(\lambda)}} \end{array} \right\} \quad -④$$

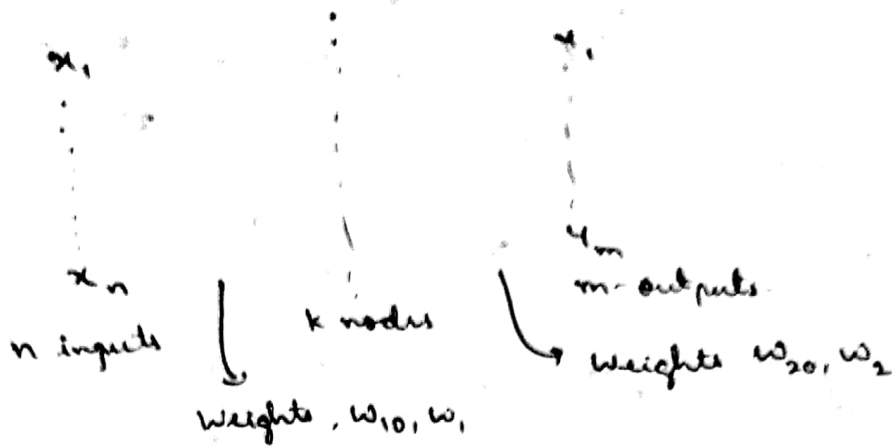From ① & ②, $\delta_k^{(i)}$ ; $S_m^{(i)}$ are errors from current model at output & hidden layer units.

③ is a back propagation eqn.

With the updated weight values in ④, we can implement back-propagation algorithm.

In forward pass, the current weights are fixed and predicted values are computed.

In backward pass, errors $\delta_k^{(i)}$ are computed, then using ③, to give errors $S_m^{(i)}$.

Both sets of errors are used to computed gradients.

4) Assume the below given neural network

$x_1$

$\vdots$

$y_1$

$\vdots$

$x_n$

$y_m$

n inputs

k nodes

m - outputs

weights $w_{10}, w_1$

weights $w_{20}, w_2$

$$h = w_1 x + w_{10} \quad — ①$$

$$h_a = \frac{1}{1 + e^{-x}}$$

$$z = w_2 h_a + w_{20} \quad — ②$$

Loss function = loss entropy function

$$L = - \sum_{i=1}^{M} y^{(i)} \cdot \log(\hat{y}^{(i)})$$

$$\rightarrow \quad \frac{\partial L}{\partial w_2^{(i)}} = - \left( \left( \sum_{j \neq i} \frac{y^{(j)}}{\hat{y}^{(j)}} \left( -\hat{y}^{(i)} \cdot \hat{y}^{(j)} \right) + \frac{y^{(i)}}{\hat{y}^{(i)}} \cdot \hat{y}^{(i)} (1 - \hat{y}^{(i)}) \right) \cdot h_a \right)$$

$$= - h_a \left( - \sum_{j=1}^{M} y^{(j)} \times \hat{y}^{(i)} + y^{(i)} \right)$$

$$\sum_{j=1}^{M} y^{(i)} = 1$$

$$\therefore \quad \frac{\partial L}{\partial w_2^{(i)}} = h_a \left( \hat{y}^{(i)} - y^{(i)} \right)$$

$$w_2^{(i)} \text{ updated} = w_2^{(i)} - h_a \left( \hat{y}^{(i)} - y^{(i)} \right)$$

$$\frac{\partial L}{\partial w_1^{(i)}} = \frac{\partial L}{\partial h_a^{(i)}} \cdot \frac{\partial h_a^{(i)}}{\partial h^{(i)}} \cdot \frac{\partial h^{(i)}}{\partial w_1^{(i)}}$$

$$= \frac{\partial L}{\partial h_a^{(i)}} \cdot h^{(i)} \cdot (1 - h^{(i)}) \cdot x \qquad (\text{From ① & ②}).$$

Scanned by CamScanner

$$= \frac{\partial L}{\partial z} \left( w_2^{(i)} \cdot h^{(i)} \cdot \left(1 - h^{(i)}\right) \right) \cdot x$$

$\hookrightarrow$ all weights from $i^{th}$ layer to output layer.

$$= \left(\hat{y} - y\right) \times \underline{w_2^{(i)} \cdot h^{(i)} \cdot \left(1 - h^{(i)}\right) \cdot x}$$

be $p$.

$$\therefore w_1^{(i)} \text{ updated} = w_1^{(i)} + p \times x$$