

1 (a) Given:

$$E_k = \frac{1}{2} (y_{dk} - \hat{y}_k)^2 + \frac{\lambda}{n} \sum_{k=1}^l w_{jk}^2 \quad \text{Required! } \Delta w_{jk}$$

Solution:

$$\Delta w_{jk} = \alpha \cdot \frac{\partial E_k}{\partial w} \cdot x_k = \alpha \cdot \delta_k \cdot x_k \quad \text{Where } x_k = \text{Input (Assumed to be Independent variable of the model)}$$

One is assumed for simplicity

$$\text{Then } \frac{\partial E_k}{\partial w} = \frac{\partial \left( \frac{1}{2} (y_{dk} - \hat{y}_k)^2 + \frac{\lambda}{n} \sum_{k=1}^l w_{jk}^2 \right)}{\partial w}$$

$$\approx \hat{y}_k = \sum_k x_k w_k \approx x_k w_{jk}$$

$$\frac{\partial E_k}{\partial w} = \frac{1}{2} \cdot 2 (y_{dk} - \hat{y}_k) \cdot x_k + \lambda w_{jk}$$

$$\frac{\partial E_k}{\partial w} = (y_{dk} - \hat{y}_k) \cdot x_k + \lambda w_{jk}$$

$$\therefore \Delta w_{jk} = \alpha \cdot ((y_{dk} - \hat{y}_k) \cdot x_k + \lambda w_{jk}) \quad \text{for Ridge (L}^2\text{)}$$

1 (b) Required!  $\Delta w_{jk}$  Given!  $E_k = \frac{1}{2} (y_{dk} - \hat{y}_k)^2 + \frac{\lambda}{n} \sum_{k=1}^l |w_{jk}|$

Solution:

$$\frac{\partial E_k}{\partial w} = \frac{\partial \left( \frac{1}{2} (y_{dk} - \hat{y}_k)^2 + \frac{\lambda}{n} \sum_{k=1}^l |w_{jk}| \right)}{\partial w} \quad \Rightarrow \hat{y}_k \approx x_k w_{jk} \quad \text{(Assumed for simplicity)}$$

$$\frac{\partial E_k}{\partial w} = (y_{dk} - \hat{y}_k) x_k + \lambda$$

$$\therefore \Delta w_{jk} = \alpha \cdot ((y_{dk} - \hat{y}_k) \cdot x_k + \lambda) \quad \text{for LASSO constraint.}$$

1. (c) Required!  $\Delta w_{jk} = ?$  Given!  $E_k = \frac{1}{2} (y_{d,k} - \hat{y}_k)^2 + \frac{\lambda_1}{n} \sum_{k=1}^l |w_{jk}| + \frac{\lambda_2}{n} \sum_{k=1}^l w_{jk}^2$

$$\frac{\partial E_k}{\partial w_{jk}} = \frac{\partial \left( \frac{1}{2} (y_{d,k} - \hat{y}_k)^2 + \lambda_1 w_{jk} + \lambda_2 w_{jk}^2 \right)}{\partial w_{jk}}$$

Assumption for simplicity;  
 $\hat{y}_k = x_k w_{jk}$

$$\frac{\partial E_k}{\partial w_{jk}} = (y_{d,k} - \hat{y}_k) x_k + \lambda_1 + 2\lambda_2 w_{jk}$$

$$\therefore \Delta w_{jk} = \alpha \cdot ((y_{d,k} - \hat{y}_k) x_k + \lambda_1 + 2\lambda_2 w_{jk})$$

1d) Required!  $\Delta w_{jk} = ?$  Given!  $E_k = \frac{1}{2} (y_{d,k} - \hat{y}_k)^2 + \frac{\lambda_1}{n} \sum_{k=1}^l |w_{jk}|^p$

$$\frac{\partial E_k}{\partial w_{jk}} = \frac{\partial \left( \frac{1}{2} (y_{d,k} - \hat{y}_k)^2 + \lambda_1 |w_{jk}|^p \right)}{\partial w_{jk}}$$

$$\frac{\partial E_k}{\partial w_{jk}} = (y_{d,k} - \hat{y}_k) x_k + \lambda_1 \cdot p \cdot w_{jk}^{p-1}$$

$$\Delta w_{jk} = \alpha \cdot ((y_{d,k} - \hat{y}_k) x_k + p \cdot \lambda_1 \cdot w_{jk}^{p-1})$$

2. Given! (a)  $f(\text{net}) = a \cdot \tanh(b \cdot \text{net}) = a \left[ \frac{e^{b \cdot \text{net}} - 1}{e^{b \cdot \text{net}} + 1} \right] = \frac{2a}{1 + e^{-b \cdot \text{net}}} - a$

Solution!

$f'(\text{net}) \Rightarrow$  from;  $\frac{\partial \left( \frac{u}{v} \right)}{\partial x} = \frac{u' \cdot v - v' \cdot u}{v^2}$   $u$  and  $v$  are functions of  $x$ .

$$f'(\text{net}) = \frac{a \cdot b \cdot e^{b \cdot \text{net}} (e^{b \cdot \text{net}} + 1) - (b \cdot e^{b \cdot \text{net}}) (a e^{b \cdot \text{net}} - a)}{(e^{b \cdot \text{net}} + 1)^2}$$

$$2a) \quad f'_{\text{net}} = \frac{a \cdot b e^{b \cdot \text{net}} + a b e^{b \cdot \text{net}} - a \cdot b e^{b \cdot \text{net}} + a \cdot b e^{b \cdot \text{net}}}{(e^{b \cdot \text{net}} + 1)^2}$$

$$f'_{\text{net}} = \frac{2ab e^{b \cdot \text{net}}}{(e^{b \cdot \text{net}} + 1)^2}$$

$$f'_{\text{net}} = \frac{2a e^{b \cdot \text{net}}}{e^{b \cdot \text{net}} + 1} \cdot \frac{b}{e^{b \cdot \text{net}} + 1}$$

$$f'_{\text{net}} = \left( \frac{2a e^{b \cdot \text{net}}}{e^{b \cdot \text{net}} + 1} \cdot \frac{e^{-b \cdot \text{net}}}{e^{-b \cdot \text{net}}} \right) \left( \frac{b}{e^{b \cdot \text{net}} + 1} \right)$$

$$f'_{\text{net}} = \left( \frac{2a}{1 + e^{-b \cdot \text{net}}} \right) \left( \frac{b}{e^{b \cdot \text{net}} + 1} \right)$$

$$f'_{\text{net}} = \left( \frac{2a}{1 + e^{-b \cdot \text{net}}} - a + a \right) \left( \frac{b}{e^{b \cdot \text{net}} + 1} \right)$$

$$\therefore f'_{\text{net}} = (f_{\text{net}} + a) \left( \frac{b}{e^{b \cdot \text{net}} + 1} \right)$$

$$2b) \quad f(-\infty) = a \left[ \frac{e^{b(-\infty)} - 1}{e^{b(-\infty)} + 1} \right] = a \left[ \frac{0 - 1}{0 + 1} \right] = -a$$

$$f(0) = a \left[ \frac{e^{b(0)} - 1}{e^{b(0)} + 1} \right] = a \left[ \frac{1 - 1}{1 + 1} \right] = 0$$

$$f(\infty) = a \left[ \frac{e^{b(\infty)} - 1}{e^{b(\infty)} + 1} \right] = a \left[ \frac{\infty - 1}{\infty + 1} \right] = a$$

$$2(b) \quad f'(net) = \frac{2abe^{b \cdot net}}{(e^{b \cdot net} + 1)^2}$$

$$f'(-\infty) = \frac{2abe^{b(-\infty)}}{(e^{b(-\infty)} + 1)^2} = 0, \quad f'(0) = \frac{2abe^{b(0)}}{(e^{b(0)} + 1)^2} = \frac{2ab}{2^2} = \frac{1}{2}ab$$

$$f'(\infty) = \frac{2abe^{b(\infty)}}{(e^{b(\infty)} + 1)^2} = \frac{\infty}{\infty} = 1$$

$$\therefore f'(-\infty) = 0, \quad f'(0) = \frac{1}{2}ab \quad \text{and} \quad f'(\infty) = 1$$

$$2(b) \quad f''(net) = ?$$

$$\text{from } f'(net) = \frac{2abe^{b \cdot net}}{(e^{b \cdot net} + 1)^2} = \frac{2abe^{b \cdot net}}{e^{2b \cdot net} + 2e^{b \cdot net} + 1}$$

$$f''(net) = \frac{2ab^2e^{b \cdot net}(e^{2b \cdot net} + 2e^{b \cdot net} + 1) - 2abe^{b \cdot net}(2be^{2b \cdot net} + 2be^{b \cdot net})}{(e^{b \cdot net} + 1)^4}$$

$$f''(net) = \frac{2ab^2e^{3b \cdot net} + 4ab^2e^{2b \cdot net} + 2ab^2e^{b \cdot net} - 4ab^2e^{3b \cdot net} - 4ab^2e^{2b \cdot net}}{(e^{b \cdot net} + 1)^4}$$

$$f''(-\infty) = \frac{0}{(0+1)^4} = \frac{0}{1} = 0 \quad f''(0) = \frac{2ab^2 + 4ab^2 + 2ab^2 - 4ab^2 - 4ab^2}{2^4}$$

$$f''(0) = 0$$

$$f''(\infty) = \frac{\infty}{(\infty+1)^4} = \frac{\infty}{\infty} = 1$$

$$\therefore f''(-\infty) = 0, \quad f''(0) = 0 \quad \text{and} \quad f''(\infty) = 1$$

3. Given!  $E(w) = \frac{1}{2} \sigma^2 - Y_d w + \frac{1}{2} Y_x w^2$

a.  $\frac{\partial E}{\partial w} = ?$

$$\Rightarrow \frac{\partial E}{\partial w} = \frac{\partial \left( \frac{1}{2} \sigma^2 - Y_d w + \frac{1}{2} Y_x w^2 \right)}{\partial w}$$

$$\frac{\partial E}{\partial w} = 0 - Y_d + \frac{1}{2} (2 Y_x w)$$

$$\frac{\partial E}{\partial w} = -Y_d + Y_x w = Y_x w - Y_d$$

$$\therefore \frac{\partial E}{\partial w} = Y_x w - Y_d$$

b) Update equation for  $w_{k+1}$

Then,  $w_{k+1} = w_k + \alpha \left( -\frac{\partial E}{\partial w} \right)$  but  $\frac{\partial E}{\partial w} = Y_x w - Y_d$

$$w_{k+1} = w_k + \alpha \left( - (Y_x w - Y_d) \right)$$

$$w_{k+1} = w_k + \alpha (Y_d - Y_x w)$$

$$\therefore w_{k+1} = w_k + \alpha (Y_d - Y_x w)$$

Where  $\alpha$  is the learning rate.

c) Optimum value of  $w$ ; To get minimal value gradient/slope must be equal to 0 (zero) in this case

$$\frac{\partial E}{\partial w} = 0; \text{ Then, } 0 = Y_x w - Y_d$$

$$* w = \frac{Y_d}{Y_x}$$

$\therefore$  At  $w = \frac{Y_d}{Y_x}$   $E(w)$  becomes minimal.



4. Weight update;

@ ReLU and Leaky ReLU;

Given!

$w_{13} = 0.5, w_{14} = 0.9, w_{24} = 1.0, w_{25} = 0.4, w_{35} = 1.2, w_{45} = 1.1, \theta_3 = 0.8, \theta_4 = -0.1$   
and  $\theta_5 = 0.3$

$$\text{func: } g(x) = \begin{cases} x & \text{if } x \geq 1 \\ 0 & \text{otherwise.} \end{cases}$$

$$y_3 = \text{ReLU}(x_1 w_{13} + x_2 w_{23} - \theta_3) = g(1 \times 0.5 + 1 \times 0.9 - 0.8) = 0.1$$

$$y_4 = \text{ReLU}(x_1 w_{14} + x_2 w_{24} - \theta_4) = g(1 \times 0.9 + 1 \times 1 + 0.1) = 2.0$$

$$y_5 = \text{ReLU}(y_3 w_{35} + y_4 w_{45} - \theta_5) = g(0.1 \times 1.2 + 2.0 \times 1.1 - 0.3) = 2.02$$

Then, error;

$$e = y_{d5} - y_5 = 0 - 2.02 = -2.02$$

Then, error gradient for neuron 5 in output layer;

$$\text{func: } g'(x) = \frac{\partial(g(x))}{\partial x} = \begin{cases} 1 & \text{if } x \geq 1 \\ 0 & \text{otherwise.} \end{cases}$$

$$\delta_5 = g'(x) \cdot e \quad \text{But } x = y_5 = 2.02$$

$$\delta_5 = g'(2.02) \cdot -2.02$$

$$\text{for } g'(2.02) = 1$$

$$\text{Hence } \delta_5 = -2.02$$

#### \* Weight corrections

$$\Delta w_{35} = \alpha y_5 \cdot \delta_5 = 0.1 \times 2.02 \times 0.1 = \underline{0.0202}$$

$$\Delta w_{45} = \alpha y_4 \cdot \delta_5 = 0.1 \times 2 \times 2.02 = \underline{0.404}$$

$$\Delta \theta_5 = \alpha(-1) \cdot \delta_5 = 0.1 \times (-1) \times 2.02 = \underline{-0.202}$$

#### \* Error gradient for neurons 3 and 4 in the hidden layer

$$\delta_3 = g'(y_3) \cdot \delta_5 \cdot w_{35}$$

$$\text{for } g'(y_3) = g(0.1) = 0$$

$$\text{hence, } \delta_3 = 0$$

$$\delta_4 = g'(y_4) \cdot \delta_5 \cdot w_{45}$$

$$\delta_4 = g'(2.0) \cdot \delta_5 \cdot w_{45}$$

$$\delta_4 = 1 \times 2.02 \times 1 = 2.02$$

$$\delta_4 = 2.42$$

#### \* Then, Weight corrections;

$$\Delta w_{13} = \alpha x_1 \cdot \delta_3 = 0.1 \times 1 \times 0 = 0$$

$$\Delta w_{23} = \alpha x_2 \cdot \delta_3 = 0.1 \times 1 \times 0 = 0$$

$$\Delta \theta_3 = \alpha(-1) \cdot \delta_3 = 0.1 \times (-1) \times 0 = 0$$

$$\Delta w_{14} = \alpha x_1 \cdot \delta_4 = 0.1 \times 1 \times 2.42 = 0.242$$

$$\Delta w_{24} = \alpha x_2 \cdot \delta_4 = 0.1 \times 1 \times 2.42 = 0.242$$

$$\Delta \theta_4 = \alpha(-1) \cdot \delta_4 = 0.1 \times (-1) \times 2.42 = -0.242$$

#### \* Update weights and biases (Threshold).

$$w_{13} = w_{13} + \Delta w_{13} = 0.5 + 0 = 0.5$$

$$w_{14} = w_{14} + \Delta w_{14} = 0.9 + 0.242 = 1.142$$

$$w_{23} = w_{23} + \Delta w_{23} = 0.4 + 0 = 0.4$$

$$w_{24} = w_{24} + \Delta w_{24} = 1.0 + 0.242 = 1.242$$

$$w_{35} = w_{35} + \Delta w_{35} = -1.2 + 0.0202 = -1.1808$$

$$w_{45} = w_{45} + \Delta w_{45} = 1.1 + 0.404 = 1.504$$

$$\theta_3 = \theta_3 + \Delta \theta_3 = 0.8 + 0 = 0.8$$

$$\theta_4 = \theta_4 + \Delta \theta_4 = -0.1 + 0.242 = 0.142$$

$$\theta_5 = \theta_5 + \Delta \theta_5 = 0.3 + (-0.202) = 0.098$$

#### 4⑥ PReLU

Given!  $w_{13}=0.5, w_{14}=0.3, w_{24}=1.0, w_{31}=0.4, w_{35}=1.2, w_{45}=1.1, \theta_3=0.8, \theta_4=-0.1$   
and  $\theta_5=0.3$

PReLU activation function:  $g(x) = \begin{cases} x & \text{if } x > 0 \\ a_i x & \text{if otherwise} \end{cases}$

But in PReLU activation function:

$a_i$  = Coefficient controlling the slope

for  $a_i$  is learnable parameter.

$$\text{hence, } \frac{\partial g(x_i)}{\partial a_i} = \begin{cases} 0, & \text{if } x_i > 0 \\ x_i, & \text{if } x_i \leq 0 \end{cases}$$

So,

$$y_k = g\left(\sum_{j=1}^m x_{jk}(p) \cdot w_{jk}(p) - \theta_k\right)$$

$$y_j = g\left(\sum_{k=1}^n x_{jk}(p) \cdot w_{jk}(p) - \theta_j\right)$$

Weight Training/Update:

$$\delta_k(p) = y_k g'(y_k) \cdot e_k(p) \quad \text{Where } e_k(p) = y_{d,k}(p) - y_k(p)$$

Weight Correction:

$$\delta_j(p) = g'(y_j) \cdot \sum_{k=1}^l \delta_k(p) w_{jk}(p)$$

$$\Delta w_{jk}(p) = \alpha \cdot y_j(p) \cdot \delta_k(p)$$

$$\Delta w_{jk}(p) = \alpha \cdot x_j(p) \cdot \delta_j(p)$$

$$\therefore w_{jk}(p+1) = w_{jk}(p) + \Delta w_{jk}(p)$$

$$\therefore w_{ij}(p) = w_{ij}(p) + \Delta w_{ij}(p)$$



5. What we always want;

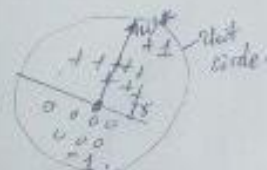
for each data point

$$\{x : w^T x + b = 0\}$$

For generalisation;

$$w^T x = 0$$

Example: Given;



$w^*$  Exist hyper-plane defined by  $(w^*)$  with  $\|w^*\| = 1$ .

$r$  - Distance of hyperplane to the closest datapoint.

$$\text{So, } y = +1 \Rightarrow w^T x > 0$$

$$y = -1 \Rightarrow w^T x < 0 \quad \left. \begin{array}{l} y = +1 \Rightarrow w^T x > 0 \\ y = -1 \Rightarrow w^T x < 0 \end{array} \right\} y w^T x > 0$$

Hence; Update rule;

$$y = +1 \Rightarrow w \leftarrow w + x$$

$$y = -1 \Rightarrow w \leftarrow w - x$$

or If the above holds then  $\frac{1}{2} n$  mistakes are made at most.

Proof needed to show  $y w^T x > 0$  will be

true for every data point by separating hyper-plane and classifies all data point correctly.

For each that  $\forall (x, y) \in D$

$$\alpha y w^* x > 0$$

Scalar value

$$w^* = \alpha w$$

$\|w^*\| = 1$  Mean all datapoints will be in the within the radius of 1 unit.

$$\forall; \|x_i\| \leq 1$$

So,  $x_i = x$  so that the maximum  $\max_j |x_j|$  one should have exactly norm 1.

Then; On update of ~~xx~~  $w^T w^*$

$$(w + yx)^T w^* = w^T w^* + y(x^T w^*) \geq w^T w^* + \gamma$$

This means;  $y(x^T w^*) = |x^T w^*| \geq \gamma$

On single update step;  $w^T w^*$  grows by  $\approx \gamma$  (at least)

Then; Update on  $w^T w$

$$(w + yx)^T (w + yx) = w^T w + 2y(w^T x) + y^2(x^T x) \leq w^T w + 1$$

$2y(w^T x) < 0$  means  $x$  is misclassified, update is required.

$$0 \leq y^2(x^T x) \leq 1 \text{ as } y^2 = 1 \text{ and } x^T x \leq 1 \text{ because } \|x\| \leq 1.$$

Means  $w^T w$  grows at most 1.  
for each update.

Therefore; after  $K$  updates

$w^T w \leq K$  and  $w^T w^* \geq K\gamma$  these must be true

$$K\gamma \leq w^T w^*$$

$$= \|w\| \cos(\theta) \text{ inner product, } \theta \text{ is angle between } w \text{ and } w^*$$

$$\leq \|w\| \quad \cos(\theta) \text{ must have } \cos \theta \leq 1.$$

$$= \sqrt{w^T w} \text{ by definition of } \|w\|$$

$$\leq \sqrt{K} \text{ because } w^T w \leq K$$

$$\text{Then; } K\gamma \leq \sqrt{K}$$

$$K^2 \gamma^2 \leq K$$

$$\therefore K \leq \frac{1}{\gamma^2} \quad K \text{ is bounded from by a constant } \frac{1}{\gamma^2}$$