

Machine Learning Mid-Term Exam (Homework #2)

Indra Imanuel Gunawan / 20195178

1.) a.) $E_k = \frac{1}{2} (y_{d,k} - \hat{y}_k)^2 + \frac{\lambda}{n} \sum_{j=1}^p w_{jk}^2$

λ = regularization parameter

Derive Δw_{jk} for this E_k

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

$$= 2 \cdot \frac{\lambda}{n} \sum_{j=1}^p w_{jk}$$

$$\Delta w_{jk} = \alpha \cdot \frac{\partial E_k}{\partial w_{jk}} \cdot x_k$$

$$= \alpha \cdot 2 \cdot \frac{\lambda}{n} w_{jk} \cdot x_k$$

b.) Derive Δw_{jk} for $E_k = \frac{1}{2} (y_{d,k} - \hat{y}_k)^2 + \frac{\lambda}{n} \sum_{j=1}^p |w_{jk}|$

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

$$= \frac{\lambda}{n}$$

$$\Delta w_{jk} = \alpha \cdot \frac{\partial E_k}{\partial w_{jk}} \cdot x_k$$

$$= \alpha \cdot \frac{\lambda}{n} x_k$$

c.) Derive ΔW_{jk} when $E_k = \frac{1}{2} (y_{dk} - \hat{y}_k)^2 + \frac{\lambda_1}{n} \sum_{i=1}^l |w_{jk}| + \frac{\lambda_2}{n} \sum_{i=1}^l w_{jk}^2$

$$\frac{\partial E_k}{\partial w_{jk}} = \frac{\partial \left(\frac{1}{2} (y_{dk} - \hat{y}_k)^2 + \frac{\lambda_1}{n} \sum_{i=1}^l |w_{jk}| + \frac{\lambda_2}{n} \sum_{i=1}^l w_{jk}^2 \right)}{\partial w_{jk}}$$

$$= \frac{\lambda_1}{n} \cdot 1 + 2 \cdot \frac{\lambda_2}{n} w_{jk}$$

$$\Delta w_{jk} = \alpha \cdot \frac{\partial E_k}{\partial w_{jk}} \cdot x_k$$

$$= \alpha \cdot \left(\frac{\lambda_1}{n} + 2 \cdot \frac{\lambda_2 w_{jk}}{n} \right) \cdot x_k$$

d.) Derive ΔW_{jk} when $E_k = \frac{1}{2} (y_{dk} - \hat{y}_k)^2 + \frac{\lambda_1}{n} \sum_{i=1}^l |w_{jk}|^p$

$$\frac{\partial E_k}{\partial w_{jk}} = \frac{\partial \left(\frac{1}{2} (y_{dk} - \hat{y}_k)^2 + \frac{\lambda_1}{n} \sum_{i=1}^l |w_{jk}|^p \right)}{\partial w_{jk}}$$

$$= p \cdot \frac{\lambda_1}{n} \cdot w_{jk}^{p-1}$$

$$\Delta w_{jk} = \alpha \cdot \frac{\partial E_k}{\partial w_{jk}} \cdot x_k$$

$$= \alpha \cdot \left(p \cdot \frac{\lambda_1}{n} \cdot w_{jk}^{p-1} \right) \cdot x_k$$

4/27

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

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

$$= -2 \cdot \frac{(e^{-b \cdot \text{net}} + 1)'}{(e^{-b \cdot \text{net}} + 1)^2} \cdot a + 0$$

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

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

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

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

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

$$= 2ab \left(\frac{e^{b \cdot \text{net}} \cdot b (e^{b \cdot \text{net}} + 1)^2 - 2(e^{b \cdot \text{net}} + 1) \cdot e^{b \cdot \text{net}} \cdot b}{(e^{b \cdot \text{net}} + 1)^4} \right)$$

~~Barren~~

$$= 2ab \left(\frac{e^{b \cdot \text{net}}}{(e^{b \cdot \text{net}} + 1)^2} \right)$$

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

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

$$b.) \text{net} = -\infty$$

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

$$f'(-\infty) = \frac{2 \cdot a \cdot b \cdot e^{b \cdot (-\infty)}}{(e^{b \cdot (-\infty)} + 1)^2} = \underline{\underline{0}}$$

$$f''(-\infty) = \frac{2 \cdot a \cdot b \cdot (e^{b \cdot (-\infty)} \cdot b (e^{b \cdot (-\infty)} + 1)^2 - 2(e^{b \cdot (-\infty)} + 1) \cdot e^{b \cdot (-\infty)} \cdot b)}{(e^{b \cdot (-\infty)} + 1)^4} = \frac{0}{(0 + 1)^4} = \underline{\underline{0}}$$

$$\text{net} = 0$$

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

$$f'(0) = \frac{2 \cdot a \cdot b \cdot e^{b(0)}}{(e^{b \cdot 0} + 1)^2} = \frac{2 \cdot a \cdot b}{2 \cdot 2} = \underline{\underline{\frac{a \cdot b}{2}}}$$

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

$$\text{net} = 00$$

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

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

$$f''(s) = \text{scribble}$$

$$= \frac{s}{s} = 1$$

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

$$a.) \frac{\partial E}{\partial w} = \frac{\partial (\frac{1}{2} \sigma^2 - Y_d w + \frac{1}{2} Y_x w^2)}{\partial w}$$

$$= \text{scribble} - Y_d + Y_x w$$

$$b.) w_{k+1} = w_k + \alpha \left(-\frac{\partial E}{\partial w} \right)$$

α = learning rate

$$= w_k + \alpha (- (-Y_d + Y_x w))$$

$$= w_k + \alpha \cdot (Y_d - Y_x w)$$

c.) Optimum value for w for which $E(w)$ becomes minimal

$$\frac{\partial E}{\partial w} = 0 \quad \text{scribble}$$

$$-Y_d + Y_x w = 0$$

$$Y_x w = Y_d$$

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

4.) a) ReLU and Leaky ReLU

$$\text{ReLU}(x) = \begin{cases} x & \text{if } x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\text{Leaky ReLU}(x) = \begin{cases} x & \text{if } x \geq 0 \\ \alpha x & \text{if } x < 0 \end{cases}$$

b.) PReLU(y_i) = $\begin{cases} y_i, & \text{if } y_i > 0 \\ a_i y_i, & \text{if } y_i \leq 0 \end{cases}$

c.) SELU(x) = $\lambda \begin{cases} x, & \text{if } x > 0 \\ \alpha e^{-x}, & \text{if } x \leq 0 \end{cases}$ or SELU(x) = $\lambda \begin{cases} x & \text{if } x \geq 0 \\ \alpha (\exp(x) - 1) & \text{if } x < 0 \end{cases}$

d.) Swish(x) = $x \cdot \text{Sigmoid}(\beta x)$

3.) If a data set is linearly separable, the Perceptron will find a separating hyperplane in a finite number of updates.

Suppose $\exists w^*$ such that $y_i (x_i^T w^*) > 0 \forall (x_i, y_i) \in D$

Suppose that we rescale each data point and the w^* such that

$$\|w^*\| = 1 \quad \text{and} \quad \|x_i\| \leq 1 \quad \forall x_i \in D$$

Here's the proof:

Consider the effect of an update (w becomes $w + yx$) on the two terms $w^T w^*$ and $w^T w$. We'll use 2 facts:

- $y (x^T w) \leq 0 \rightarrow x$ is misclassified by w
- $y (x^T w^*) > 0 \rightarrow w^*$ is a separating hyperplane and classifies all points correctly.

1.) Consider the effect of an update on $w^T w^*$

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

The inequality forms from the distance from the hyperplane defined by w^* to x must be at least γ . This means, for each update, $w^T w^*$ grows at least by γ

next page \rightarrow

(2) Consider the effect of an update on $W^T W$:

$$(W + \gamma x)^T (W + \gamma x) = W^T W + \underbrace{2\gamma (W^T x)}_{< 0} + \underbrace{\gamma^2 (x^T x)}_{0 \leq \leq 1} \leq W^T W + 1$$

the inequality forms because:

-) $2\gamma (W^T x) < 0$ as we had to make an update, meaning x was misclassified
-) $0 \leq \gamma^2 (x^T x) \leq 1$ as $\gamma^2 = 1$ and all $x^T x \leq 1$ (because $\|x\| \leq 1$)

(3) Now we know that after M updates the following two inequalities must hold:

$$(1) W^T W^* \geq M\gamma$$

$$(2) W^T W \leq M$$

We can complete the proof:

$$M\gamma \leq W^T W^*$$

$$= \|W\| \cos\theta$$

$$\leq \|W\|$$

$$= \sqrt{W^T W}$$

$$\leq \sqrt{M}$$

$$\Rightarrow M\gamma \leq \sqrt{M}$$

$$\Rightarrow M^2 \gamma^2 \leq M$$

$$\Rightarrow M \leq \frac{1}{\gamma^2}$$

By (1)

by definition of inner product, where θ is the angle between W and W^*
by definition of \cos , we must have $\cos(\theta) \leq 1$
by definition of $\|W\|$

By (2)

And hence, the number of updates M is bounded from above by a constant