

CSC 411 Fall 2018  
Machine Learning and Data Mining

Homework 3

---

Family name: Bhatia

Given name: Pooja

## solution 1

(a)

$$L_{\delta}(y, t) = H_{\delta}(y - t)$$

$$H_{\delta}(a) = \begin{cases} \frac{1}{2} a^2 & \text{for } |a| \leq \delta \\ \delta (|a| - \frac{1}{2} \delta) & \text{for } |a| > \delta \end{cases}$$

Now,

$$L_{\delta}(y, t) = \begin{cases} \frac{1}{2} (y - t)^2 & \text{if } |y - t| \leq \delta \\ \delta (y - t) - \frac{1}{2} \delta^2 & \text{otherwise} \end{cases}$$

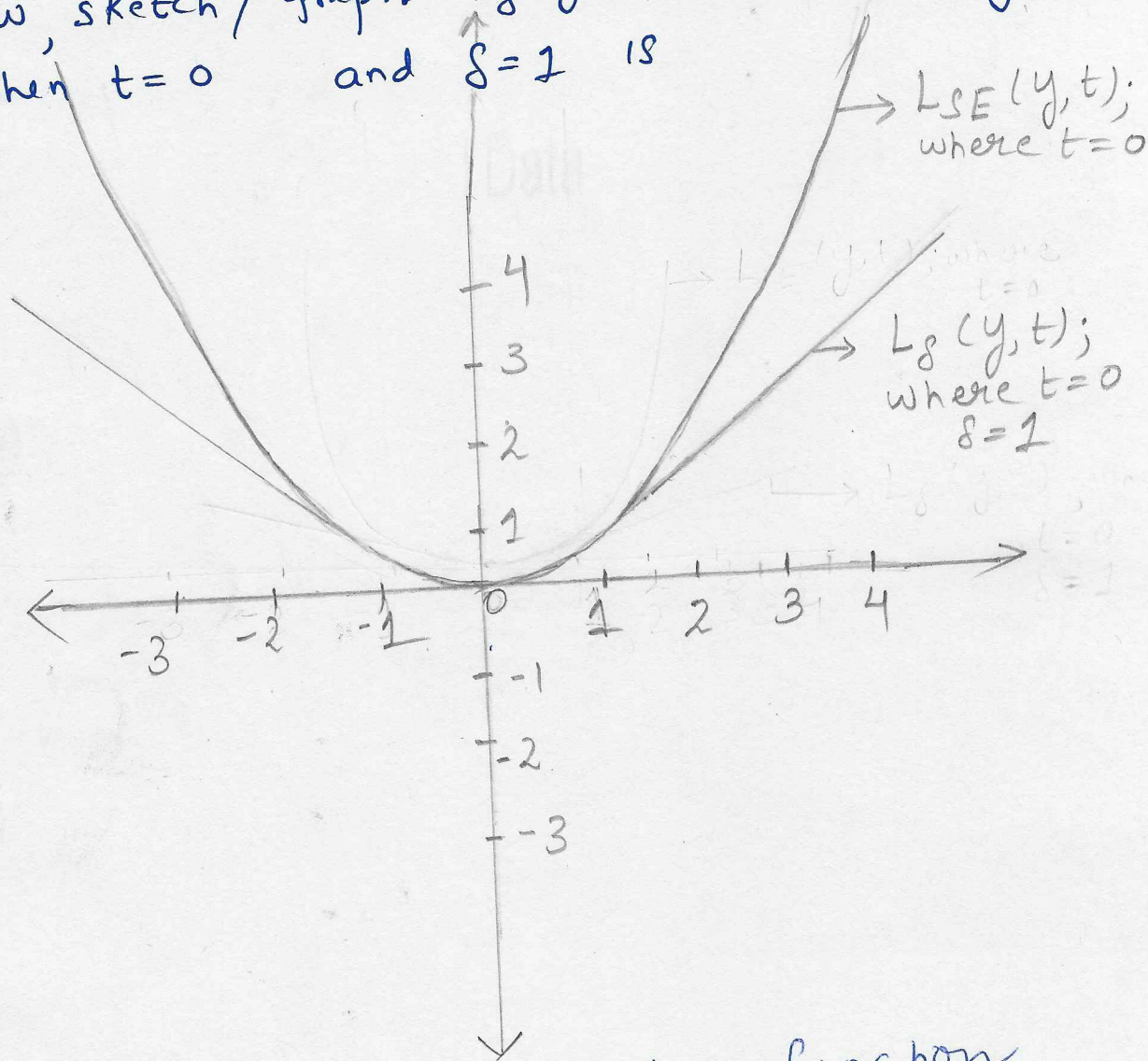
Now taking  $t = 0$ 

$$L_{\delta}(y) = \begin{cases} \frac{1}{2} y^2 & \text{if } |y| \leq \delta \\ \delta (y) - \frac{1}{2} \delta^2 & \text{otherwise.} \end{cases}$$

LSE i.e the squared Loss function  
 $L_{SE} = \frac{1}{2} (y - t)^2$ , now taking  $t = 0$   
 $L_{SE} = \frac{1}{2} y^2$



Now, sketch / graph  $L_S(y, t)$  and  $L_{SE}(y, t)$   
when  $t = 0$  and  $S = 2$  is



Based on the sketch Huber function  
is strongly convex function in a uniform  
neighborhood of its minimum  $a=0$

$$\text{i.e. } L_S(a) = \begin{cases} \frac{1}{2}a^2 & \text{for } |a| < S \\ S\{|a| - \frac{1}{2}S\} & \text{otherwise} \end{cases}$$

Therefore at the boundaries of the  
uniform neighborhood



the Huber loss function has a differentiable extension to an Affine function at  $a = -\delta$  and  $a = +\delta$ , making it less sensitive / more Robust to outliers than squared Loss error.

The squared Loss error is more sensitive to outliers as while taking summation the sample mean gets influenced by a large value of  $y$  in this case. i.e

$$LSE(y) = \frac{1}{2}(y - t)^2$$

$$LSE(y) = \frac{1}{2}(y)^2; \text{ when } t = 0$$

Ans

P. T. O



solution 1

(b)

$$H_\delta(a) = \begin{cases} \frac{1}{2}a^2 & \text{if } |a| \leq \delta \\ \delta(|a| - \frac{1}{2}\delta) & \text{if } |a| > \delta \end{cases}$$

 ~~$L_\delta(y, t) = H_\delta(y - t)$~~ 

$$L_\delta(y, t) = H_\delta(y - t)$$

$$H_\delta(y - t) = \begin{cases} \frac{1}{2}(y - t)^2 & \text{if } |y - t| \leq \delta \\ \delta(|y - t| - \frac{1}{2}\delta) & \text{if } |y - t| > \delta \end{cases}$$

now,  $y = w^T x + b$

$$H_\delta(y - t) = \begin{cases} \frac{1}{2}(w^T x + b - t)^2 & \text{if } |w^T x + b - t| \leq \delta \\ \delta(|w^T x + b - t| - \frac{1}{2}\delta) & \text{if } |w^T x + b - t| > \delta \end{cases}$$

→ ①

$$\frac{\partial L_s}{\partial w} = \frac{\partial (H_s(y-t))}{\partial w}$$

$$= \begin{cases} \frac{1}{2} \frac{\partial}{\partial w} (w^T x + b - t)^2, & \text{if } |w^T x + b - t| \leq \delta \\ \delta \frac{\partial}{\partial w} \left( |w^T x + b - t| - \frac{1}{2} \delta \right) & \text{if } |w^T x + b - t| > \delta \end{cases}$$

$$\frac{\partial L_s}{\partial w} = \begin{cases} (w^T x + b - t)(x) & \text{if } |w^T x + b - t| \leq \delta \\ \delta(x) & \text{if } |w^T x + b - t| > \delta \end{cases}$$

Ans.

$\frac{\partial L_s}{\partial w}$  in terms of  $\gamma$

$$\frac{\partial L_s}{\partial w} = \begin{cases} (y-t) \cdot (x) & \text{if } |y-t| \leq \delta \\ (\delta) \cdot (x) & \text{if } |y-t| > \delta. \end{cases}$$

Ans



Similarly from equation ①

$$\frac{\partial L_S}{\partial b} = \frac{\partial (H_S(y-t))}{\partial b}$$

$$= \begin{cases} \frac{1}{2} \frac{\partial}{\partial b} (w^T x + b - t)^2 & \text{if } |w^T x + b - t| \leq \delta \\ \delta \frac{\partial}{\partial b} (|w^T x + b - t| - \frac{1}{2} \delta) & \text{if } |w^T x + b - t| > \delta \end{cases}$$

$$\frac{\partial L_S}{\partial b} = \begin{cases} (w^T x + b - t) & \text{if } |w^T x + b - t| \leq \delta \\ \delta & \text{if } |w^T x + b - t| > \delta \end{cases}$$

Ans.

$\frac{\partial L_S}{\partial b}$  in terms of  $y$

$$\frac{\partial L_S}{\partial b} = \begin{cases} (y - t) & \text{if } |y - t| \leq \delta \\ \delta & \text{if } |y - t| > \delta \end{cases}$$

Ans



Solution 2(a)

Given Training data  $\{ (x^{(1)}, y^{(1)}) \dots$   
 $\dots (x^{(N)}, y^{(N)}) \}$  and positive  
 weights  $a^{(1)} \dots a^{(N)}$

Now, writing/representing them in Matrix  
 Notation

$$Y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(N)} \end{bmatrix}_{N \times 1} \quad X = \begin{bmatrix} 1 & \dots & x^{(1)(D)} \\ \vdots & & \vdots \\ 1 & x^{(N)(1)} & \dots & x^{(N)(D)} \end{bmatrix}_{N \times (D+1)}$$

$$A = \begin{bmatrix} a^{(1)} & 0 & \dots & 0 \\ \vdots & a^{(2)} & & \\ 0 & \dots & & a^{(N)} \end{bmatrix}_{N \times N} \quad W = \begin{bmatrix} b \\ w^{(1)} \\ \vdots \\ w^{(D)} \end{bmatrix}_{(D+1) \times 1}$$

$$y = w^T x + b \quad (\text{in case of linear regression})$$



Now, The objective in this case would be

$$E(w) = \frac{1}{2} \sum_{i=1}^N a^{(i)} (y^{(i)} - w^T x^{(i)})^2 + \frac{\lambda}{2} \|w\|^2$$

Now, representing the same in Matrix Notation

$$E(w) = \frac{A}{2} (Y - XW)^T (Y - XW) + \frac{\lambda}{2} w^T w \quad \rightarrow \textcircled{1}$$

Now, we need to minimize the square loss function along with regularizer term in above equation

Now, taking derivative w.r.t  $w$  of equation  $\textcircled{1}$  and equating it to zero

$$\begin{aligned} \nabla E(w) &= -X^T A (Y - XW) + \lambda w = 0 \\ &= -X^T A Y + X^T A X W + \lambda w = 0 \end{aligned}$$

$$\begin{aligned} \Rightarrow X^T A Y &= X^T A X W + \lambda w \\ X^T A Y &= (X^T A X + \lambda I) w \end{aligned}$$

(where  $I$  is the identity matrix)



Taking inverse on both sides of equation

$$X^T A Y = (X^T A X + \lambda I) W$$

$$\boxed{W^* = (X^T A X + \lambda I)^{-1} X^T A Y}$$

∴  
Solution to the weighted least square problem

$$W^* = \operatorname{argmin} \frac{1}{2} \sum_{i=1}^N a^{(i)} (y^{(i)} - W^T x^{(i)})^2 + \frac{\lambda}{2} \|W\|^2$$

is

$$W^* = (X^T A X + \lambda I)^{-1} X^T A Y$$

Ans



## Solution Q2

(d)

When  $z \rightarrow 0$

Then both the training loss and test loss are approaching the same value.

In the graph you will witness the plots of ~~test losses~~ train-losses and test-losses is merging.

When  $z \rightarrow \infty$

Then <sup>gap between</sup> the training loss and test loss is increasing.

They are approaching different values.

In the graph you will witness the plots train-losses and test-losses is diverging