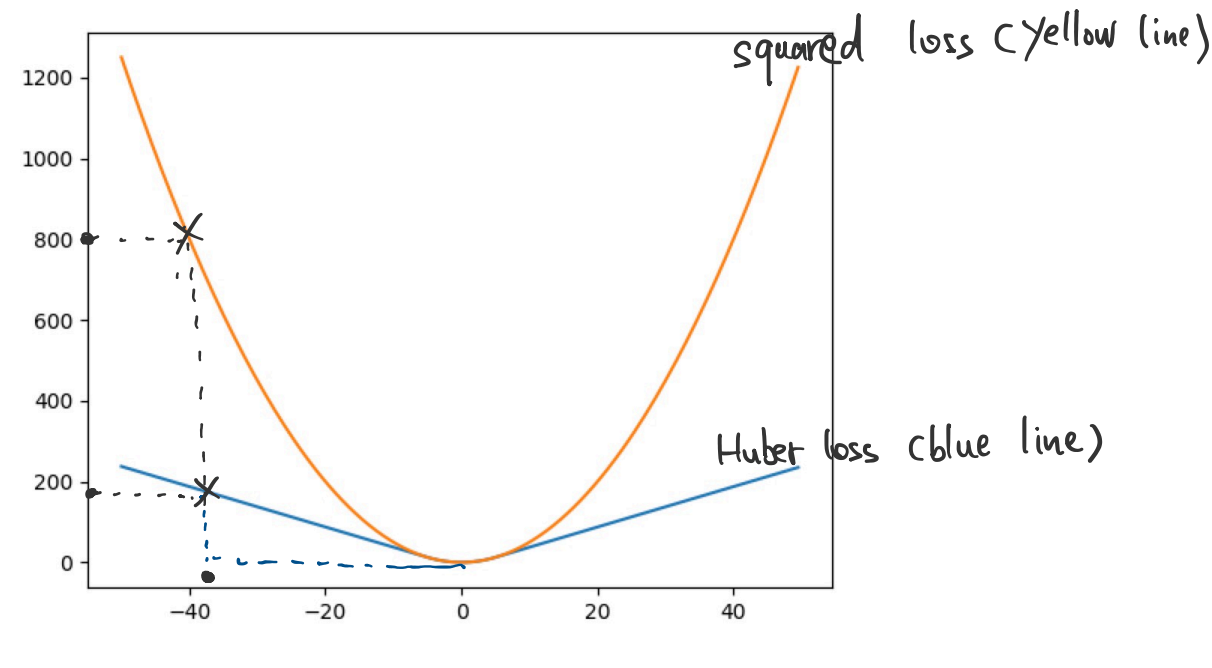


Q1 (a)



Ans: within Huber loss, we have a delta value to classify a, when $|a| > \delta$. It's effect is linear. Compare with squared loss. Whether it is correct value or outliers. It will has exponential growth rate. As we can see on the graph., assume we have outliers = -40. The y-value of huber loss is significant smaller than the y-value of squared loss.

Q1 (b)

$$\therefore H'_\delta(a) = \begin{cases} a & \text{if } |a| \leq \delta \\ \delta & \text{if } a > \delta \\ -\delta & \text{if } a < -\delta \end{cases}$$

$$\therefore H'_\delta(y-t) = \begin{cases} y-t & \text{if } |y-t| \leq \delta \\ \delta & \text{if } y-t > \delta \\ -\delta & \text{if } y-t < -\delta \end{cases}$$

$$\frac{\partial L_\delta}{\partial w} \begin{cases} ① \frac{\partial L_\delta}{\partial w} = \frac{dL_\delta}{dy} \frac{\partial y}{\partial w} = \frac{d}{dy} \left[\frac{1}{2} (y-t)^2 \right] \cdot x_i \\ \quad = (y-t) x_i \quad \text{if } |y-t| \leq \delta \\ ② \frac{\partial L_\delta}{\partial w} = \frac{dL_\delta}{dy} \frac{\partial y}{\partial w} = \frac{d}{dy} \left(\delta (y-t) - \frac{1}{2} \delta^2 \right) x_i \\ \quad = \delta x_i \quad \text{if } y-t > \delta \\ ③ \frac{\partial L_\delta}{\partial w} = \frac{dL_\delta}{dy} \frac{\partial y}{\partial w} = \frac{d}{dy} \left(-\delta (y-t) - \frac{1}{2} \delta^2 \right) x_i \\ \quad = -\delta x_i \quad \text{if } y-t < -\delta \end{cases}$$

$$\frac{\partial L_\delta}{\partial b} \begin{cases} ④ \frac{\partial L_\delta}{\partial b} = \frac{dL_\delta}{dy} \frac{\partial y}{\partial b} = \frac{d}{dy} \left[\frac{1}{2} (y-t)^2 \right] \cdot 1 \\ \quad = (y-t) \quad \text{if } |y-t| \leq \delta \\ ⑤ \frac{\partial L_\delta}{\partial b} = \frac{dL_\delta}{dy} \frac{\partial y}{\partial b} = \frac{d}{dy} \left(\delta (y-t) - \frac{1}{2} \delta^2 \right) \cdot 1 \\ \quad = \delta \quad \text{if } y-t > \delta \\ ⑥ \frac{\partial L_\delta}{\partial b} = \frac{dL_\delta}{dy} \frac{\partial y}{\partial b} = \frac{d}{dy} \left(-\delta (y-t) - \frac{1}{2} \delta^2 \right) \cdot 1 \\ \quad = -\delta \quad \text{if } y-t < -\delta \end{cases}$$

$$\therefore \text{we have } \frac{\partial L_\delta}{\partial w} = \begin{cases} (y-t) x_i & \text{if } |y-t| \leq \delta \\ \delta x_i & \text{if } y-t > \delta \\ -\delta x_i & \text{if } y-t < -\delta \end{cases}$$

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

Q2 (a)

$$\therefore W^* = \arg \min \underbrace{\frac{1}{2} \sum_{i=1}^N a^{(i)} (y^{(i)} - w^T x^{(i)})^2}_{①} + \underbrace{\frac{\lambda}{2} \|w\|^2}_{②}$$

\therefore define a vector $v = y - Xw$, therefore the first term can be rewrite as $\frac{1}{2} \sum_{j=1}^N v_j^2 a_j$. Because we also know that A is a diagonal matrix where $A_{ii} = a^{(i)}$. Then $[Av]_j = a_j^{(i)} v_j$, and inner product $\langle v, Av \rangle = v^T A v = \sum_{j=1}^N v_j^2 a_j^{(i)}$.

\therefore Moreover, we can rewrite the loss function as form of vector:

$$\begin{aligned} \therefore L(w) &= \frac{1}{2} (y - Xw)^T A (y - Xw) + \frac{\lambda}{2} \|w\|^2 \\ &= \frac{1}{2} [y^T - (Xw)^T] A (y - Xw) + \frac{\lambda}{2} \|w\|^2 \\ &= \frac{1}{2} [y^T - w^T X^T] A (y - Xw) + \frac{\lambda}{2} \|w\|^2 \\ &= \frac{1}{2} [(y^T A - w^T X^T A) (y - Xw) + \lambda w^T w] \\ &= \frac{1}{2} [y^T A y - y^T A X w - w^T X^T A y + w^T X^T A X w + \lambda w^T w] \\ &= \frac{1}{2} [y^T A y - (y^T A X w)^T - w^T X^T A y + w^T X^T A X w + \lambda w^T w] \\ &= \frac{1}{2} [y^T A y - w^T X^T A y - w^T X^T A y + w^T X^T A X w + \lambda w^T w] \\ &= \frac{1}{2} y^T A y - w^T X^T A y + \frac{1}{2} w^T X^T A X w + \frac{1}{2} \lambda w^T w \end{aligned}$$

$$\therefore \nabla_w w^T w = 2w$$

$$\therefore \nabla_w w^T A w = 2A w, \text{ we know that } X^T A X \text{ is symmetric.}$$

$$\text{since } X_{1 \times n}^T A_{n \times n} X_{n \times 1} = \text{matrix of } 1 \times 1.$$

$$\therefore \nabla_w w^T x = x$$

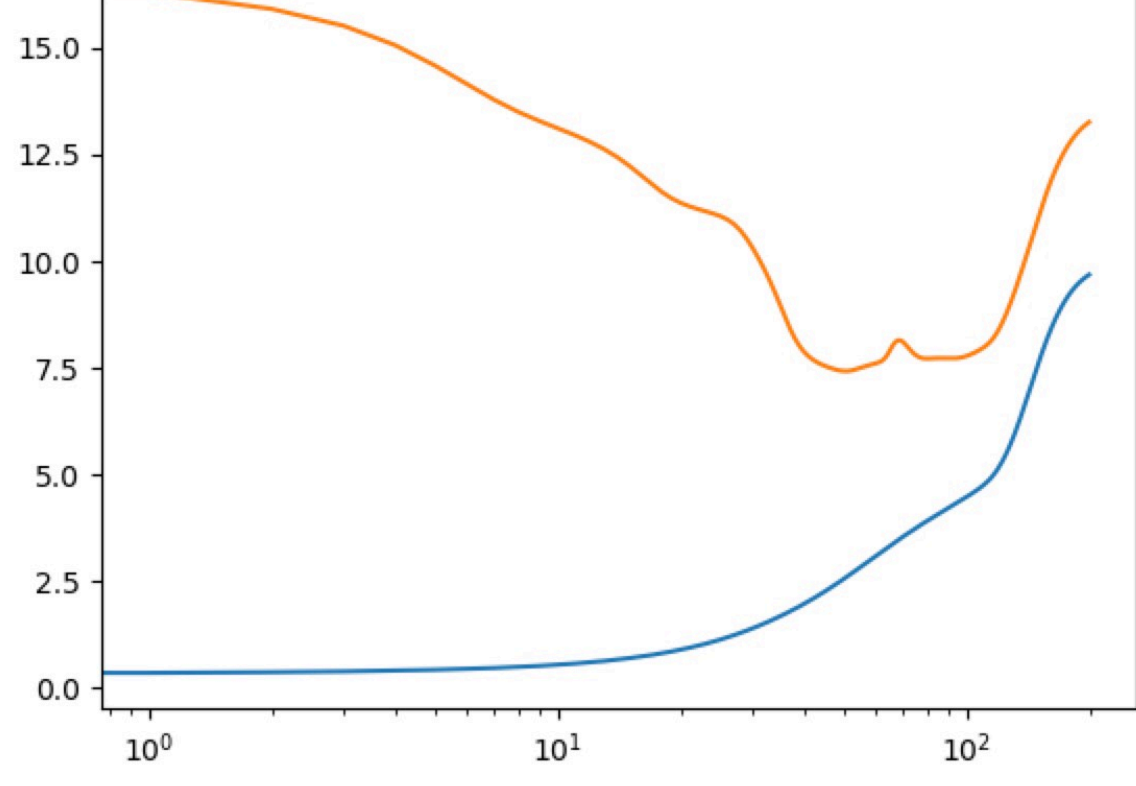
$$\begin{aligned} \therefore \text{Compute } \nabla_w L(w) &= \nabla_w \frac{1}{2} y^T A y - \nabla_w w^T X^T A y + \nabla_w \frac{1}{2} w^T X^T A X w + \nabla_w \frac{1}{2} \lambda w^T w \\ &= 0 - X^T A y + X^T A X w + \lambda w \\ &= -X^T A y + X^T A X w + \lambda w \end{aligned}$$

when we setting $\nabla_w L(w) = 0$, we have optimal w^*

$$\begin{aligned} -X^T A y + X^T A X w^* + \lambda w^* &= 0 \\ X^T A X w^* + \lambda w^* &= X^T A y \\ (X^T A X + I \lambda) w^* &= X^T A y \\ w^* &= (X^T A X + I \lambda)^{-1} X^T A y \end{aligned}$$

proved!

Q2 (c):



Q2 (D) when $\tau \rightarrow \infty$, we will have $a^{(i)} \rightarrow \frac{1}{n}$, therefore when we take it back to w^* function. It looks like squared error loss which can only predices linear model.

Therefore when $\tau \rightarrow \infty$, it is under fit.

when $\tau \rightarrow 0$, we will have $a^{(i)} \rightarrow \text{large value}$.

Therefore it will over-fit the training sample. which means it will increases loss of validation set. So when

$\tau \rightarrow 0$, it is over fit.