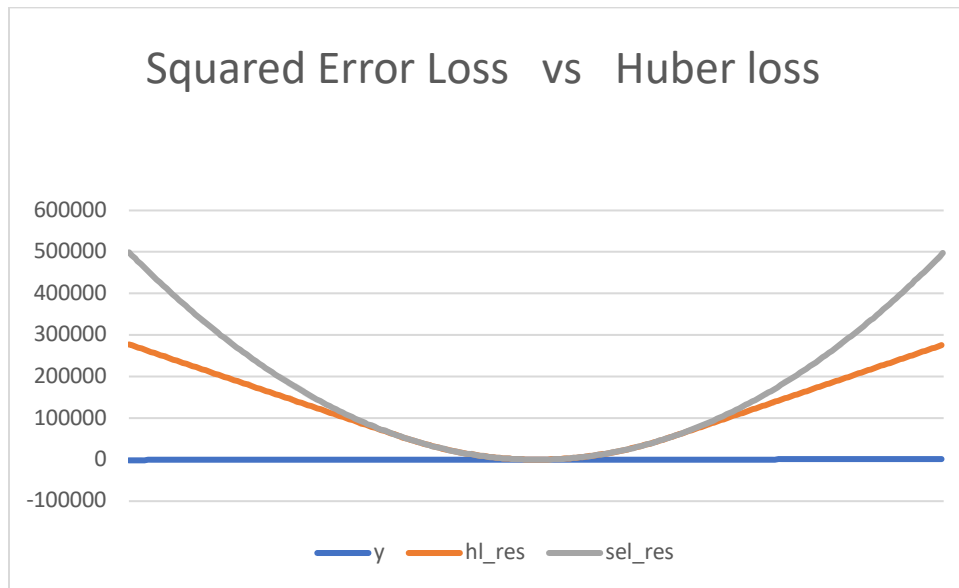


Q1

(a)

When the value  $y$  is outlier, compared to Squared Error Loss, the Huber loss is much less which is shown in below plot. ( $\delta = 333$ ,  $y$ 's interval is  $[-999, 999]$ )



(b)

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

$$\Rightarrow H'_{\delta}(a) = \begin{cases} a & \text{if } |a| \leq \delta \\ \delta & \text{if } a > \delta \\ -\delta & \text{if } a < -\delta \end{cases} \quad (1)$$

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

$$= \sum_i w_i x_i + b$$

$$\frac{\partial L_{\delta}}{\partial w} = H'_{\delta}(y-t) \cdot \frac{\partial (y-t)}{\partial w}$$

$$= H'_{\delta}(y-t) \cdot \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix}$$

where we can apply formula (1)

$$\frac{\partial L_{\delta}}{\partial b} = H'_{\delta}(y-t) \cdot \frac{\partial (y-t)}{\partial b}$$

$$= H'_{\delta}(y-t) \cdot 1$$

$$= H'_{\delta}(y-t)$$

where we can apply formula (1)

Q2  
(a)

$$\begin{aligned} J &= \frac{1}{2} \sum_{i=1}^N \alpha^{(i)} (y^{(i)} - w^T x^{(i)})^2 + \frac{\lambda}{2} \|w\|^2 \\ &= \frac{1}{2} \sum_{i=1}^N \alpha^{(i)} (y^{(i)^2} - 2y^{(i)} w^T x^{(i)} + (w^T x^{(i)})^2) + \frac{\lambda}{2} \|w\|^2 \\ &= \frac{1}{2} \sum_{i=1}^N \alpha^{(i)} y^{(i)^2} - \sum_{i=1}^N \alpha^{(i)} y^{(i)} w^T x^{(i)} + \frac{1}{2} \sum_{i=1}^N \alpha^{(i)} (w^T x^{(i)})^2 + \frac{\lambda}{2} \|w\|^2 \\ &= \frac{1}{2} y^T A y - w^T X^T A y + \frac{1}{2} w^T X^T A X w + \frac{\lambda}{2} w^T w \end{aligned}$$

$$\text{Let } \nabla_w J = 0$$

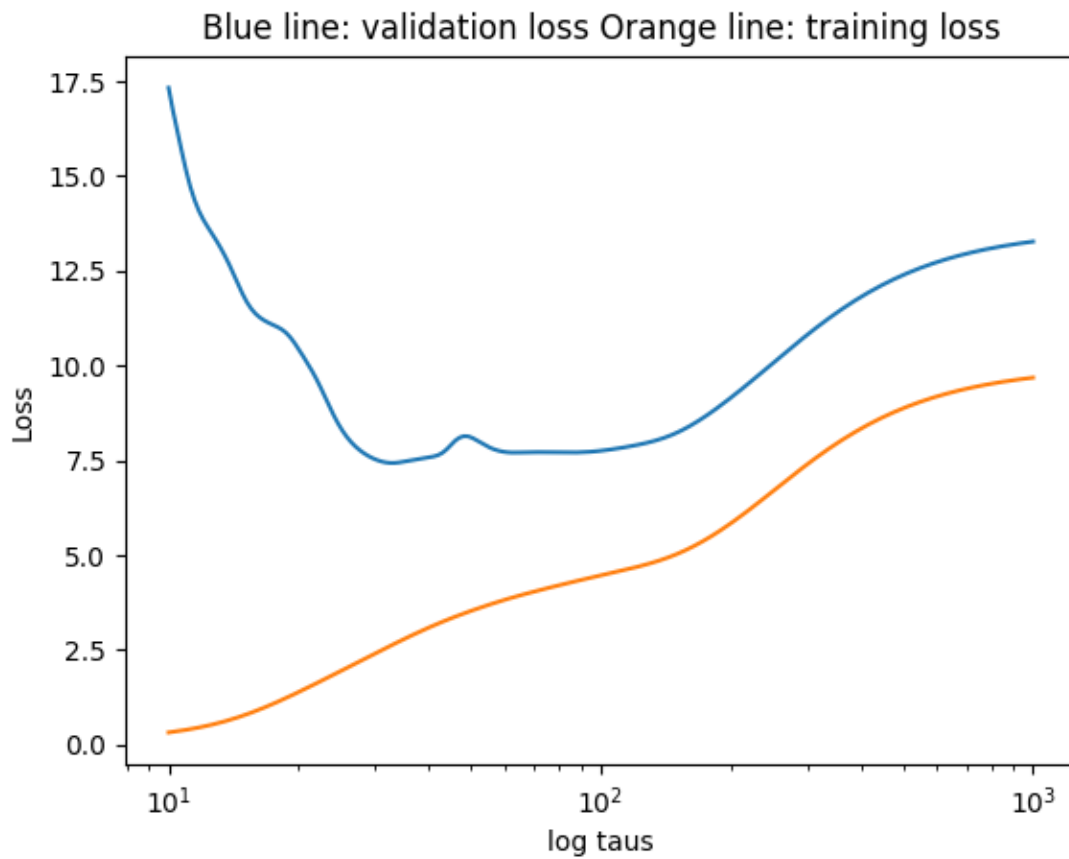
$$\Rightarrow -X^T A y + X^T A X w + \lambda w = 0$$

$$(X^T A X + \lambda I) w = X^T A y$$

$$\Rightarrow w = (X^T A X + \lambda I)^{-1} X^T A y$$

$$\Rightarrow w^* = (X^T A X + \lambda I)^{-1} X^T A y$$

(c)



(d)

When  $\gamma \rightarrow \infty$ , each training example is equally weighted, so the model is close to standard regression with  $L^2$  penalty.

When  $\gamma \rightarrow 0$ , the test point becomes dramatic driver to the model which could lead poor generalization.