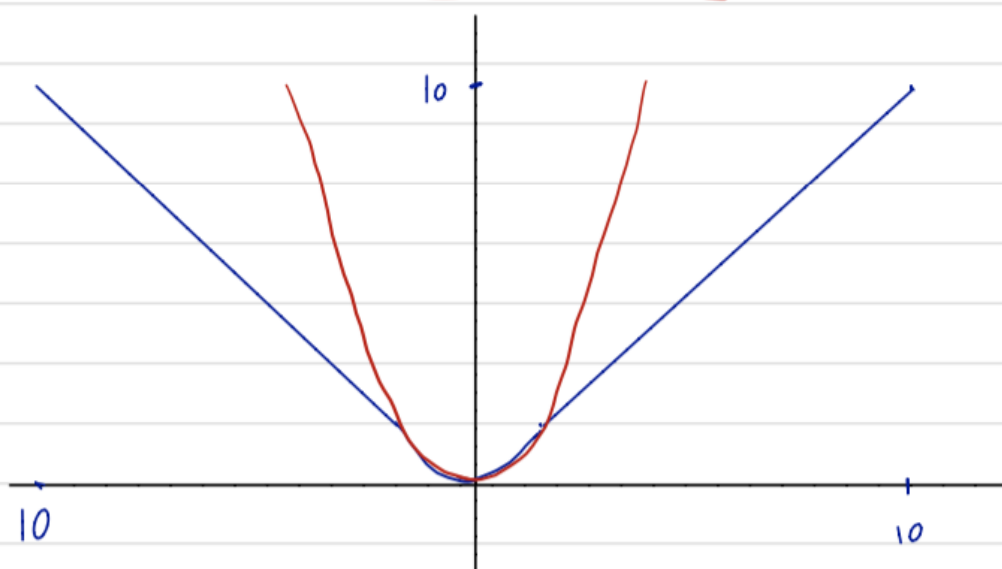


1. as huber loss.

Squared error loss.



Huber loss should be more robust to outliers. Because outside of range  $[-\delta, \delta]$ , it increases or decreases linearly. However, squared error has an increasing rate of change. The error of prediction affects the squared loss more than huber loss.

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

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

$$\therefore \frac{dL_\delta}{dw} = \frac{dL_\delta}{da} \times \frac{\partial a}{\partial w}$$

$$\frac{dL_\delta}{db} = \frac{dL_\delta}{da} \times \frac{\partial a}{\partial b}$$

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

$$a = w \cdot x + b \quad \therefore \frac{\partial a}{\partial w} = x$$

$$= \begin{cases} y-t & |y-t| \leq \delta \\ -\delta & y-t < -\delta \\ \delta & y-t > \delta \end{cases}$$

$$\therefore \frac{dL_\delta}{dw} = x \cdot H'_\delta(y-t)$$

$$= \begin{cases} (y-t)x & |y-t| \leq \delta \\ -\delta x & y-t < -\delta \\ \delta x & y-t > \delta \end{cases}$$

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

$$z = \dots + \frac{\lambda}{2} \vec{w}^T \vec{w}$$

derivative  $\Rightarrow$

$$0 = -\frac{1}{2} \sum_{i=1}^N a^{(i)} \cdot 2 \cdot \vec{x}^{(i)} \cdot (y^{(i)} - \vec{w}^T \vec{x}^{(i)}) + \frac{\lambda}{2} \vec{w} \cdot 2$$

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

$$= - \sum_{i=1}^N \vec{x}^{(i)} \cdot a^{(i)} y^{(i)} + \sum_{i=1}^N \vec{x}^{(i)} a^{(i)} \cdot \vec{w}^T \cdot \vec{x}^{(i)} + \lambda \vec{w}$$

$$\therefore \text{w.r.t } \vec{w} \text{ vector } \quad 0 = - \sum_{i=1}^N \vec{x}^{(i)} a^{(i)} \cdot y^{(i)} + \sum_{i=1}^N \vec{x}^{(i)} \cdot a^{(i)} \cdot (\vec{x}^{(i)})^T \vec{w}$$

$$\boxed{\sum_{i=1}^N \vec{x}^{(i)} \cdot a^{(i)} \cdot \vec{x}^{(i)T} \vec{w} + \lambda \vec{w} = \sum_{i=1}^N \vec{x}^{(i)} \cdot a^{(i)} \cdot y^{(i)}}$$

$$\sum_{i=1}^N a^{(i)} \cdot \vec{x}^{(i)} \cdot y^{(i)} = a^{(1)} \cdot \vec{x}^{(1)} \cdot y^{(1)} + a^{(2)} \cdot \vec{x}^{(2)} \cdot y^{(2)} \dots a^{(N)} \cdot \vec{x}^{(N)} \cdot y^{(N)}$$

$$\Rightarrow \sum_{i=1}^N a^{(i)} \cdot x^{(i)} y^{(i)} = X^T \cdot A \cdot Y$$

$$X^T \quad N \times M \quad M \times 1$$

$$\Rightarrow N \times 1$$

$$A \Rightarrow M \times M$$

$$\sum_{i=1}^N \vec{x}^{(i)} \cdot a^{(i)} \cdot \vec{x}^{(i)T} \cdot \vec{w} = X \cdot A \cdot X^T \cdot \vec{w}$$

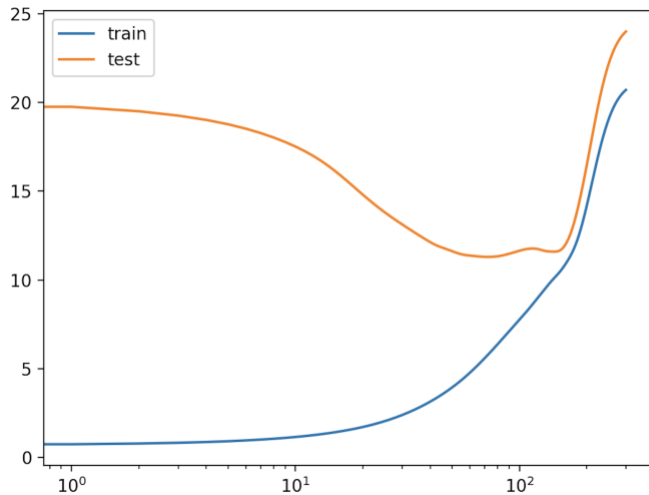
$$\begin{matrix} N \times M & \times & M \times M & \times & M \times N & \times & N \times 1 \\ X & & A & & X^T & & w \end{matrix}$$

$$\therefore 0 = X^T A y \cdot (X^T A X w - X^T I w)$$

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

$$\therefore 0 = X^T A y \cdot (X^T A X + \lambda I) w \Rightarrow$$

### Question 2 C)



Question:

Based on our understanding of overfitting and underfitting, how would you expect the training error and the validation error to vary as a function of  $\tau$ ? (I.e., what do you expect the curves to look like?)

Now run the experiment. Randomly hold out 30% of the dataset as a validation set. Compute the average loss for different values of  $\tau$  in the range [10,1000] on both the training set and the validation set. Plot the training and validation losses as a function of  $\tau$  (using a log scale for  $\tau$ ). Was your guess correct?

Answer:

First, we should know the  $\tau$  is used to decide the weight for each point. When the  $\tau$  is small, the training error is small, which also means the algorithm pays more attention to the closest point to the datum. The algorithm would be overfit to the training dataset. Therefore, the validation error should be huge when the  $\tau$  is small. Conversely, when the  $\tau$  increases, the algorithm will consider the weight of each point fairly. And it will turn to be underfit the features of training dataset. Therefore, the training error will increase and validation error should be decrease. However, if the  $\tau$  is so huge. Both validation error and training error should be increase.

My guess is almost correct. For the small value of  $\tau$ , the algorithm had low training loss error but high validation loss error. Because the algorithm overfitted the training data and lost more generalization. When the  $\tau$  value increased over 100, the algorithm underfit the training data too much, and it cannot generate a linear regression model to fit the training dataset. Therefore, the validation losses also increase. When the  $\tau$  value is about 50 (seen from the graph above), the algorithm reached a suitable situation, it didn't overfit the training dataset, and the generalization increased. The weight for each point is perfect for us to predict the result. So we can see from the graph above, the validation error reached the minimum value.

Question 3: to prove  $err' = \frac{\sum_{i=1}^N w_i' I(h(x_i) \neq t^{(i)})}{\sum_{i=1}^N w_i'} = \frac{1}{2}$

given: target set  $t^{(i)} \in \{-1, 1\}$   
 $h(x_i)$  classifier return the value  $\{-1, 1\}$

We set  $W_n = \sum_{i=1}^N w_i$        $W_E = \sum_{i \in E} w_i$        $I(h(x_i) \neq t_i)$   
 $W_{CE} = \sum_{i \in CE} w_i$        $I(h(x_i) = t_i)$

$i \in E$   $h(x_i) t_i = -1$  means the classifier did wrong  
 $i \in E_c$   $h(x_i) t_i = 1$  means the classifier did right

①  $err_t = \frac{\sum_{i=1}^N w_i I(h(x_i) \neq t^{(i)})}{\sum_{i=1}^N w_i}$

$= \frac{\sum_{i \in E} w_i}{\sum_{i=1}^N w_i} = \frac{W_E}{W_E + W_{CE}}$

②  $err_t' = \frac{\sum_{i=1}^N w_i' I(h(x_i) \neq t^{(i)})}{\sum_{i=1}^N w_i'} = \frac{\sum_{i=1}^N w_i' I(E)}{\sum_{i=1}^N w_i'} = \frac{\sum_{i \in E} w_i'}{\sum_{i \in E} w_i' + \sum_{i \in E_c} w_i'}$

③  $w_i' = w_i \cdot \exp \{ -\alpha t^{(i)} h(x^{(i)}) \}$

we separate  $w_i'$  equation for  $i \in E$  and  $i \in E_c$

$\therefore w_i' = \begin{cases} \text{a) } i \in E \\ \text{b) } i \in E_c \end{cases}$

when  $i \in E$   $t(i) = h_t(x(i)) = 1$

$\therefore w_i = w_i \cdot \exp(\alpha)$

$\alpha = \frac{1}{2} \left( \log \frac{1 - e_{nt}}{e_{nt}} \right)$

$\therefore w_i = w_i \cdot \exp \left( \frac{1}{2} \log \frac{1 - e_{nt}}{e_{nt}} \right) = e^{\frac{1}{2}} = \sqrt{e}$

$\therefore w_i = w_i \sqrt{\frac{1 - e_{nt}}{e_{nt}}}$

using the same method, when  $i \in E_c$

$i \in E_c$   $w_i = w_i \cdot \exp(-\alpha)$

$= w_i \exp \left( -\frac{1}{2} \log \frac{1 - e_{nt}}{e_{nt}} \right)$

$= w_i \exp \left( \frac{1}{2} \log \frac{e_{nt}}{1 - e_{nt}} \right) = w_i \cdot \sqrt{\frac{e_{nt}}{1 - e_{nt}}}$

Continuous with (2)

$$e_{nt}' = \frac{\sum_{i \in E} w_i'}{\sum_{i \in E} w_i' + \sum_{i \in E_c} w_i'} = \frac{\sum_{i \in E} w_i \sqrt{\frac{1 - e_{nt}}{e_{nt}}}}{\sum_{i \in E} w_i \sqrt{\frac{1 - e_{nt}}{e_{nt}}} + \sum_{i \in E_c} w_i \sqrt{\frac{e_{nt}}{1 - e_{nt}}}}$$

from (1) we know  $e_{nt} = \frac{W_E}{W_E + W_{E_c}} \therefore 1 - e_{nt} = \frac{W_{E_c}}{W_E + W_{E_c}}$

$\therefore e_{nt}' = \frac{W_E \cdot \sqrt{\frac{W_{E_c}}{W_E}}}{W_E \cdot \sqrt{\frac{W_{E_c}}{W_E}} + W_{E_c} \cdot \sqrt{\frac{W_E}{W_{E_c}}}}$

$= \frac{W_{E_c} \cdot W_E}{W_E \cdot \sqrt{\frac{W_{E_c}}{W_E}} + W_{E_c} \cdot \sqrt{\frac{W_E}{W_{E_c}}}}$

$= \frac{W_{E_c} \cdot W_E}{\sqrt{W_E \cdot W_{E_c}}}$

$\sqrt{\frac{W_{E_c} \cdot W_E^2}{W_E}} + \sqrt{\frac{W_E \cdot W_{E_c}^2}{W_{E_c}}}$

$= \frac{\sqrt{W_{E_c} \cdot W_E}}{\sqrt{W_{E_c} \cdot W_E}}$

$\therefore = \frac{1}{2}$