

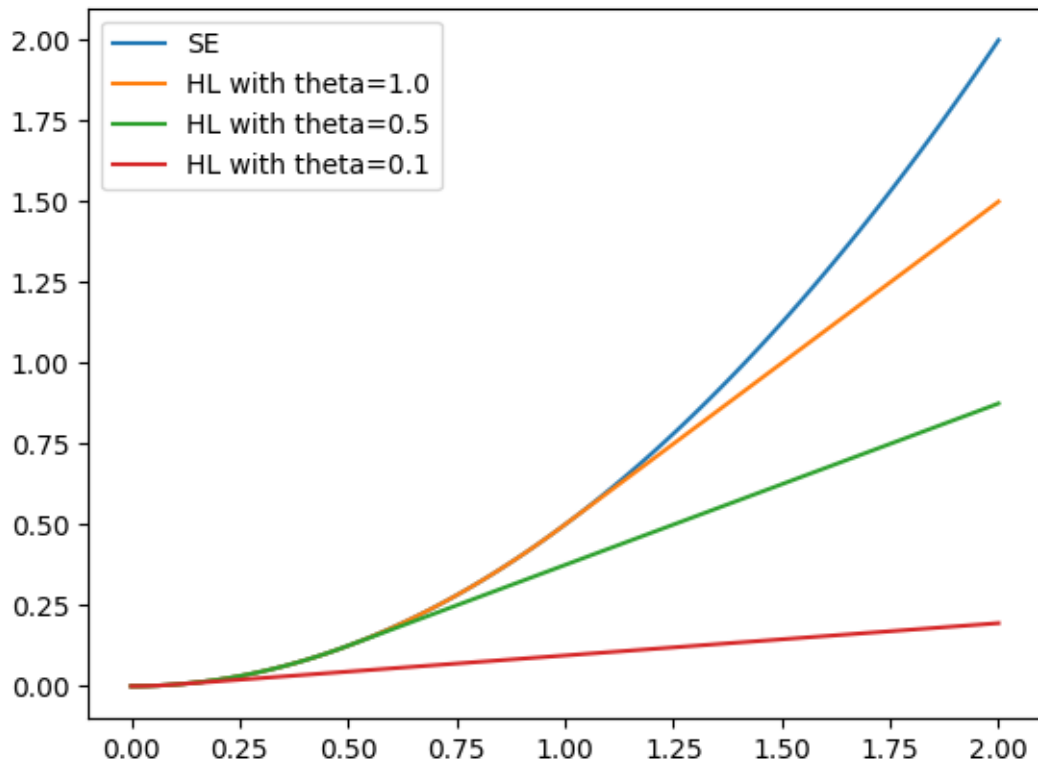
# Homework 2

Zihan Zhao

1001103708

1

(a)



Compared with squared error loss, when the residue ( $y-t$ ) increases, Huber loss is the same as squared error loss. But when it reaches over the threshold  $\delta$ , i.e. the loss is at outliers, Huber loss becomes linearly increasing by the slope of  $\delta$ . It becomes less sensitive to outliers than squared error loss. So the optimal weights can be determined more quickly by using Huber loss gradient descent. Moreover, as  $\delta$  decreases, the loss is further away from squared error loss. The degree of insensitivity becomes larger. Therefore it is robust regression.

(b)

Now determines  $\frac{dL_\delta}{dw}$ :

$$\frac{dL_\delta}{dw} = \frac{dH_\delta(a)}{da} \frac{da}{dy} \frac{dy}{dw}$$

When  $|y - t| \leq \delta$ :

$$\begin{aligned} &= \frac{d\frac{1}{2}a^2}{da} \frac{da}{dy} \frac{dy}{dw} \\ &= a * 1 * x = ax = (y - t)x \\ &= (w^\top x + b - t)x \end{aligned}$$

When  $|y - t| > \delta$ :

$$\begin{aligned} &= \frac{d\delta(|a| - \frac{1}{2}\delta)}{da} \frac{da}{dy} \frac{dy}{dw} \\ &= \begin{cases} \delta x, & y - t > \delta \\ -\delta x, & y - t < -\delta \end{cases} \end{aligned}$$

Now determines  $\frac{dL_\delta}{db}$ :

$$\frac{dL_\delta}{db} = \frac{dH_\delta(a)}{da} \frac{da}{dy} \frac{dy}{db}$$

When  $|y - t| \leq \delta$ :

$$\begin{aligned} &= \frac{d\frac{1}{2}a^2}{da} \frac{da}{dy} \frac{dy}{db} \\ &= a * 1 * 1 = x \\ &= w^\top x + b - t \end{aligned}$$

When  $|y - t| > \delta$ :

$$\begin{aligned} &= \frac{d\delta(|a| - \frac{1}{2}\delta)}{da} \frac{da}{dy} \frac{dy}{dw} \\ &= \begin{cases} \delta, & y - t > \delta \\ -\delta, & y - t < -\delta \end{cases} \end{aligned}$$

(c)

Look at q1.py.

## 2

### (a)

First factor the Loss formula:

$$\begin{aligned}
L &= \frac{1}{2} \sum_{i=1}^N a^{(i)} (y^{(i)} - w^\top x^{(i)})^2 + \frac{\lambda}{2} \|w\|^2 \\
&= \frac{1}{2} A \|Y - Xw\|^2 + \frac{\lambda}{2} w^\top w \quad (\text{where } Y \text{ is } N \times 1, X \text{ is } N \times d, A \text{ is } N \times N, \text{ and } w \text{ is } d \times 1) \\
&= \frac{1}{2} (Y - Xw)^\top (A(Y - Xw)) + \frac{\lambda}{2} w^\top w \\
&= \frac{1}{2} (Y^\top AY - Y^\top AXw - (Xw)^\top AY + w^\top X^\top AXw) + \frac{\lambda}{2} w^\top w \\
&= \frac{1}{2} (Y^\top AY - Y^\top AXw - (AY)^\top Xw + w^\top (X^\top AX)w) + \frac{\lambda}{2} w^\top w \\
&= \frac{1}{2} (Y^\top AY - Y^\top AXw - Y^\top A^\top Xw + w^\top (X^\top AX)w) + \frac{\lambda}{2} w^\top w
\end{aligned}$$

Since  $A = A^\top$ ,

$$\begin{aligned}
&= \frac{1}{2} (Y^\top AY - Y^\top AXw - Y^\top AXw + w^\top (X^\top AX)w) + \frac{\lambda}{2} w^\top w \\
&= \frac{1}{2} Y^\top AY - Y^\top AXw + \frac{1}{2} w^\top (X^\top AX)w + \frac{\lambda}{2} w^\top w
\end{aligned}$$

Now take derivative of L by w:

Since  $A = A^\top$ , so  $X^\top AX$  is symmetric as well, then

$$\begin{aligned}
\frac{dL}{dw} &= 0 - Y^\top AX + \frac{1}{2} 2(X^\top AX)w + \frac{1}{2} 2\lambda w \\
&= -Y^\top AX + (X^\top AX)w + \lambda w
\end{aligned}$$

Let  $\frac{dL}{dw} = 0$ , we got

$$\begin{aligned}
-Y^\top AX + (X^\top AX)w + \lambda w &= 0 \\
(X^\top AX + \lambda I)w &= Y^\top AX \\
w &= (X^\top AX + \lambda I)^{-1} Y^\top AX
\end{aligned}$$

Since  $Y^\top AX = (AX)^\top Y = X^\top A^\top Y = X^\top AY$ ,

$$w = (X^\top AX + \lambda I)^{-1} X^\top AY$$

Done.

### (b)

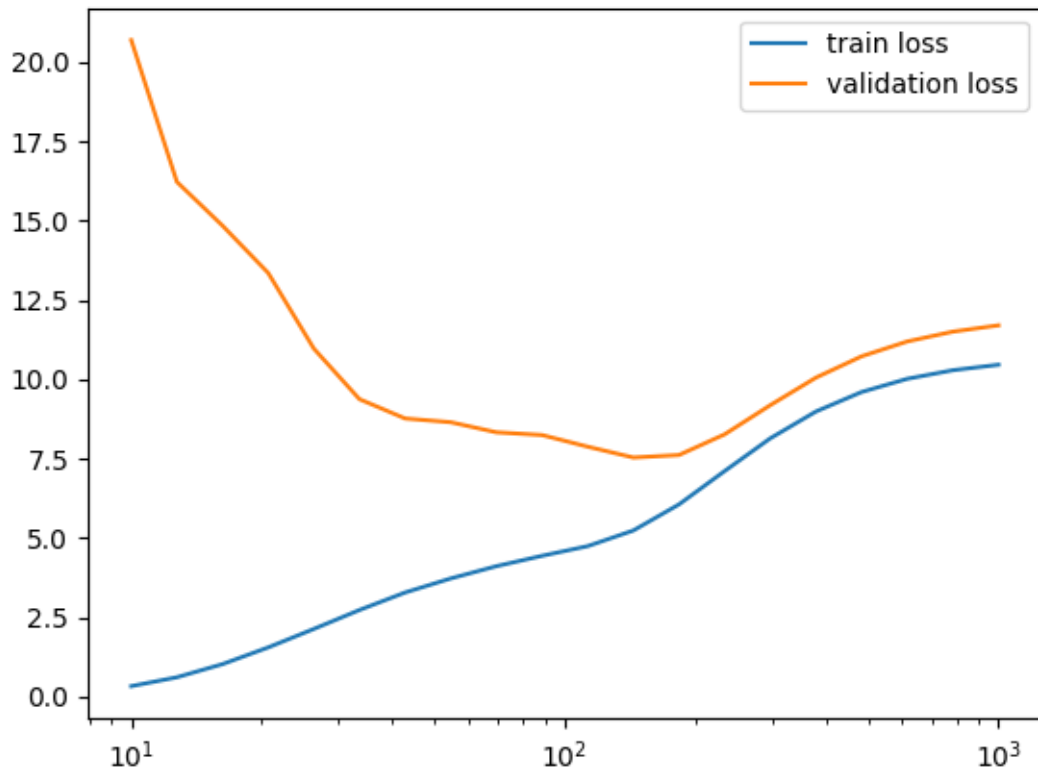
Look at q2.py.

(c)

Now take a look at the training error. When  $\tau$  is small, the divergence among the elements in  $A_{ii}$  gets large, the reweighted error loss will put larger weights on the data whose is far away from its target and put smaller weights on the data which is close to its target. So the optimmal weight will produce more precise prediction. When  $\tau$  is too large,  $a^{(i)}$  approaches to the same value  $\frac{1}{N}$ . So the loss will put same weights on each data and cannot enlarge the error where data is far from its target. The training loss is larger than when  $\tau$  is small.

Then take a look at the validation error. When  $\tau$  is small, overfitting occurs. The reason is that when training the data, the gradient descent specifies each weight for each training data and enlarges the data weight as it is far from its target. It is overfitting. Then it causes the total loss of validation data to be higher. When  $\tau$  becomes larger, weights of each data approaches to the same value, overfitting is reduced. Then validation error follows training error to increase.

Look at q2.py. The training loss and validation loss are shown in the following figure. My guess is similar with the figure.



### 3

Denote:

$$\begin{aligned} LEFT\ SIDE &= err'_t \\ &= \frac{\sum_{i=1}^N \omega'_i 1\{h_t(x^i) \neq t^i\}}{\sum_{i=1}^N \omega'_i} \\ RIGHT\ SIDE &= \frac{1}{2} \end{aligned}$$

Also from Tips

$$\begin{aligned} E &= \{i : h_t(x^i) \neq t^i\} \\ E^c &= \{i : h_t(x^i) = t^i\} \end{aligned}$$

So we can pick all  $i$  whose  $i \in E$  to rewrite LEFT SIDE:

$$LEFT\ SIDE = \frac{\sum_{i \in E} \omega'_i}{\sum_{i=1}^N \omega'_i}$$

We know the updated weight  $\omega'_i = \omega_i e^{-\alpha_t t^i h_t(x^i)}$ , plug it in LEFT SIZE:

$$= \frac{\sum_{i \in E} \omega_i e^{-\alpha_t t^i h_t(x^i)}}{\sum_{i=1}^N \omega'_i}$$

In  $E$ ,  $t^i$  and  $h_t(x^i)$  should be  $1/-1$  or  $-1/1$ , so  $t^i h_t(x^i) = -1$ . Then:

$$= \frac{\sum_{i \in E} \omega_i e^{\alpha_t}}{\sum_{i=1}^N \omega'_i}$$

Plug  $\alpha_t = \frac{1}{2} \log \frac{1-err_t}{err_t}$  in LEFT SIZE:

$$\begin{aligned} &= \frac{\sum_{i \in E} \omega_i e^{\frac{1}{2} \log \frac{1-err_t}{err_t}}}{\sum_{i=1}^N \omega'_i} \\ &= \frac{\sum_{i \in E} \omega_i \left(\frac{1-err_t}{err_t}\right)^{\frac{1}{2}}}{\sum_{i=1}^N \omega'_i} \end{aligned}$$

The denominator is converted in the same way. Notice in  $E^c$ ,  $t^i h_t(x^i) = 1$ ,

$$\begin{aligned} &= \frac{\sum_{i \in E} \omega_i \left(\frac{1-err_t}{err_t}\right)^{\frac{1}{2}}}{\sum_{i=1}^E \omega'_i + \sum_{j=1}^{E^c} \omega'_j} \\ &= \frac{\sum_{i \in E} \omega_i \left(\frac{1-err_t}{err_t}\right)^{\frac{1}{2}}}{\sum_{i \in E} \omega_i e^{\frac{1}{2} \log \frac{1-err_t}{err_t}} + \sum_{j \in E^c} \omega_j e^{-\frac{1}{2} \log \frac{1-err_t}{err_t}}} \\ &= \frac{\sum_{i \in E} \omega_i \left(\frac{1-err_t}{err_t}\right)^{\frac{1}{2}}}{\sum_{i \in E} \omega_i \left(\frac{1-err_t}{err_t}\right)^{\frac{1}{2}} + \sum_{j \in E^c} \omega_j \left(\frac{1-err_t}{err_t}\right)^{-\frac{1}{2}}} \\ &= \frac{\left(\frac{1-err_t}{err_t}\right)^{\frac{1}{2}} \sum_{i \in E} \omega_i}{\left(\frac{1-err_t}{err_t}\right)^{\frac{1}{2}} \sum_{i \in E} \omega_i + \left(\frac{1-err_t}{err_t}\right)^{-\frac{1}{2}} \sum_{j \in E^c} \omega_j} \\ &= \frac{\sum_{i \in E} \omega_i}{\sum_{i \in E} \omega_i + \left(\frac{1-err_t}{err_t}\right)^{-1} \sum_{j \in E^c} \omega_j} \end{aligned}$$

$err_t$  is the similar as  $err'_t$ , that is  $err_t = \frac{\sum_{i \in E} \omega_i}{\sum_{i=1}^N \omega_i}$ . Then,

$$\begin{aligned}
&= \frac{\sum_{i \in E} \omega_i}{\sum_{i \in E} \omega_i + \left( \frac{1 - \frac{\sum_{i \in E} \omega_i}{\sum_{i=1}^N \omega_i}}{\frac{\sum_{i \in E} \omega_i}{\sum_{i=1}^N \omega_i}} \right)^{-1} \sum_{j \in E^c} \omega_j} \\
&= \frac{\sum_{i \in E} \omega_i}{\sum_{i \in E} \omega_i + \left( \frac{\sum_{i=1}^N \omega_i - \sum_{i \in E} \omega_i}{\sum_{i \in E} \omega_i} \right)^{-1} \sum_{j \in E^c} \omega_j} \\
&= \frac{\sum_{i \in E} \omega_i}{\sum_{i \in E} \omega_i + \left( \frac{\sum_{i \in E} \omega_i}{\sum_{i=1}^N \omega_i - \sum_{i \in E} \omega_i} \right) \sum_{j \in E^c} \omega_j} \\
&= \frac{\sum_{i \in E} \omega_i}{\sum_{i \in E} \omega_i + \frac{\sum_{i \in E} \omega_i}{\sum_{j \in E^c} \omega_j} \sum_{j \in E^c} \omega_j} \\
&= \frac{\sum_{i \in E} \omega_i}{\sum_{i \in E} \omega_i + \cancel{\frac{\sum_{i \in E} \omega_i}{\sum_{j \in E^c} \omega_j}} \cancel{\sum_{j \in E^c} \omega_j}} \\
&= \frac{\sum_{i \in E} \omega_i}{2 \sum_{i \in E} \omega_i} \\
&= \frac{1}{2} = RIGHT\ SIDE
\end{aligned}$$

Done.