

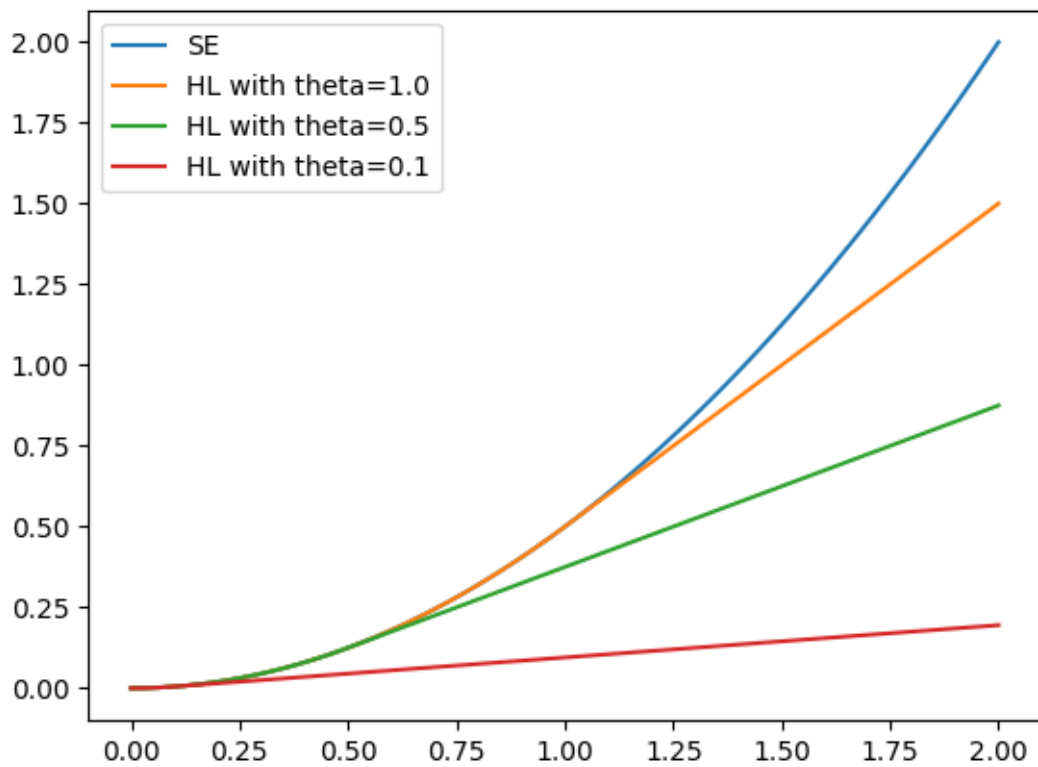
Homework 2

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 δ , i.e. the loss is at outliers, Huber loss becomes linearly increasing by the slope of δ . 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. 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)

Look at q2.py.