

ANEEL DAS

1. Robust Regression.

a) Plots are attached.

• Reason why Huber Loss is better for handling outliers:

when using squared error loss  $LSE(y, t) = \frac{1}{2} (y - t)^2$ , we see that when the predicted value is far from the actual value, the loss is very large as the difference is squared.

So, when our model comes across an outlier, the loss is very large. So, the outlier gets more importance as fitting/predicting values ~~close~~ that are close to the outlier is the only way of minimizing loss.

Since, the Huber loss uses MAE  $LSE(y, t) = \text{abs}(y - t)$  for difference ~~greater~~ greater than threshold ( $\delta$ ), it is more robust as the loss ~~scales~~ scales linearly. This is why Huber loss is more robust to outliers.

1.b)  $y = w^T x + b$

To find  $\frac{\partial L_S}{\partial w}$ ,  $\frac{\partial L_S}{\partial b}$  for Huber loss.

$$L_S(y, t) = H_S(y - t)$$

$$H_S(a) = \begin{cases} \frac{1}{2} a^2 & ; |a| \leq \delta \\ \delta(|a| - \frac{1}{2}\delta) & ; |a| > \delta \end{cases}$$

where  $a = y - t$ .  
 $y \rightarrow$  predicted value.  
 $t \rightarrow$  true value.

First we find  $\frac{\partial L_S}{\partial w}$

case 1  $|a| \leq \delta$

$$\frac{\partial L_S}{\partial w} = \frac{\partial H(a)}{\partial w} = \frac{\partial H(a)}{\partial a} \times \frac{\partial a}{\partial y} \times \frac{\partial y}{\partial w}$$

$$\frac{\partial H(a)}{\partial a} = \frac{\partial}{\partial a} \left( \frac{1}{2} a^2 \right) = \frac{1}{2} \frac{\partial}{\partial a} (a^2) = \frac{1}{2} \times 2a = \underline{\underline{a}}.$$

$$\frac{\partial a}{\partial y} = \frac{\partial (y-t)}{\partial y} = \underline{\underline{1.}} \quad \text{and} \quad \frac{\partial y}{\partial w} = \frac{\partial}{\partial w} (wx + b) = \underline{\underline{x}}$$

$$\text{So, } \frac{\partial L}{\partial w} = \frac{\partial H(a)}{\partial a} \times \frac{\partial a}{\partial y} \times \frac{\partial y}{\partial w}$$

$$= a \cdot 1 \cdot x = \underline{\underline{ax}} \quad \text{So, } \frac{\partial L_s}{\partial w} = ax; |a| \leq \delta.$$

Now,

Case 2

$$|a| > \delta$$

$$\Rightarrow -\delta > a > \delta$$

$$\frac{\partial L}{\partial w} = \frac{\partial H(y-t)}{\partial w} = \frac{\partial H(a)}{\partial a} \times \frac{\partial a}{\partial y} \times \frac{\partial y}{\partial w}.$$

$$\frac{\partial H(a)}{\partial a} = \frac{\partial}{\partial a} \left( \delta |a| - \frac{1}{2} \delta^2 \right)$$

$$= \delta \cdot \frac{\partial}{\partial a} |a|$$

when  $a > \delta$ ;  $|a| = a$

$$\text{So, } \delta \frac{\partial |a|}{\partial a} = \delta \frac{\partial a}{\partial a} = \underline{\underline{\delta}}$$

when  $-\delta > a$ ;  $|a| = -a$ .

$$\text{So, } \delta \frac{\partial |a|}{\partial a} = \underline{\underline{-\delta}}$$

$$\frac{\partial a}{\partial y} = \underline{\underline{1}}$$

$$\frac{\partial y}{\partial w} = \frac{\partial}{\partial w} (wx + b) = \underline{\underline{x}}$$

So,

$$\frac{\partial L_s}{\partial w} = \begin{cases} \delta \cdot 1 \cdot x = \underline{\underline{\delta x}}; & a > \delta \\ -\delta \cdot 1 \cdot x = \underline{\underline{-\delta x}}; & -\delta > a. \end{cases}$$

Now, we find  $\frac{\partial L_s}{\partial b}$

Case 1:  $|a| \leq \delta$

$$\frac{\partial L(y,t)}{\partial b} = \frac{\partial H(y-t)}{\partial b} = \frac{\partial H(a)}{\partial b} = \frac{\partial H(a)}{\partial a} \times \frac{\partial a}{\partial y} \times \frac{\partial y}{\partial b}.$$

$$= \frac{\partial H(a)}{\partial a} \times \frac{\partial a}{\partial y} \times \frac{\partial y}{\partial b} \rightarrow \frac{\partial (wx+b)}{\partial b} = \underline{\underline{1}}$$

$$\frac{\partial}{\partial a} \left( \frac{1}{2} a^2 \right) = \frac{1}{2} \frac{\partial (a^2)}{\partial a} = \underline{\underline{a}}$$

$$\frac{\partial (y-t)}{\partial y} = \underline{\underline{1}}$$

So,  $\frac{\partial L_S}{\partial b} = a ; |a| \leq \delta.$

Case 2

$|a| > \delta$   
 $\Rightarrow -\delta > a > \delta$

$$\frac{\partial L}{\partial b} = \frac{\partial H(a)}{\partial a} \times \frac{\partial a}{\partial y} \times \frac{\partial y}{\partial b}.$$

As calculated earlier.

$$\frac{\partial H(a)}{\partial a} = \begin{cases} \delta ; a > \delta \\ -\delta ; -\delta > a \end{cases}$$

and  $\frac{\partial a}{\partial y} = \underline{\underline{1}}.$

$$\frac{\partial y}{\partial b} = \frac{\partial (wx+b)}{\partial b} = \underline{\underline{1}}$$

So,

$$\frac{\partial L}{\partial b} = \begin{cases} \delta ; a > \delta \\ -\delta ; -\delta > a \end{cases}$$

1.c) Code provided as q1.py.

## Q2> Locally weighted Regression

Given:  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\} \rightarrow$  observation pairs.  
 where  $y$  is actual value.  
 $x$  is input value.

$a^{(1)}, \dots, a^{(n)} \rightarrow$  the weights.

So, each weight multiplies the corresponding input pair.

$$A = \begin{bmatrix} a_{11} & & & \\ & a_{22} & & \\ & & \ddots & \\ & & & a_{nn} \end{bmatrix}$$

diagonal matrix of all the weights.

and  $w^* = \underset{\text{weighted square loss}}{\operatorname{argmin}} \underbrace{\frac{1}{2} \sum_{i=1}^N a^{(i)} (y^{(i)} - w^T x^{(i)})^2}_{\text{weighted square loss}} + \underbrace{\frac{\lambda}{2} \|w\|^2}_{\text{Regularization term}}$

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

Let us take  $r = y - Xw$

So, the loss can be written as.

$$= \frac{1}{2} \sum_{i=1}^N r_i^2 a^{(i)} + \frac{\lambda}{2} \|w\|^2$$

Now, the  $i^{\text{th}}$  element of  $Ar$  is  $a^{(i)} \cdot r_i$ . So,

$$r \cdot Ar = r^T \cdot A \cdot r$$

So, loss =

$$\frac{1}{2} r^T A \cdot r + \frac{\lambda}{2} \|w\|^2$$

$$= \frac{1}{2} (y - Xw)^T \cdot A (y - Xw) + \frac{\lambda}{2} \|w\|^2 \quad ; \text{ putting } r = y - Xw$$

$$L = \frac{1}{2} (y^T A y - 2w^T X^T A y + w^T X^T A X w + \lambda w^T w)$$

Now,  $w$  is of the form  $(w_1, w_2, \dots, w_n)$

So,  $w^T \cdot w = w_1^2 + w_2^2 + \dots + w_n^2$

$$= \sum_{i=1}^n w_i^2$$

$$\text{So, } \frac{d}{dw} \sum_{i=1}^n w_i^2 = \frac{d}{dw} (w^2) = \underline{\underline{2w}}$$



$$\text{So, } \frac{dL}{d\omega} = \frac{1}{2} \frac{d}{d\omega} (y^T A y - 2\omega^T x^T A y + \omega^T x^T A X \omega + \lambda \omega^T \omega)$$

$$= \frac{1}{2} \left\{ \frac{d}{d\omega} (y^T A y) - 2 \frac{d}{d\omega} (\omega^T x^T A y) + \frac{d}{d\omega} (\omega^T x^T A X \omega) + \lambda \frac{d}{d\omega} (\omega^T \omega) \right\}$$

$$= \frac{1}{2} \left\{ -2x^T A y + 2 \cdot x^T A X \cdot \omega + 2\lambda \omega \right\}$$

$$\frac{dL}{d\omega} = -x^T A y + x^T A X \omega + \lambda \omega.$$

$$\text{for } \frac{dL}{d\omega} = 0. \Rightarrow \omega = \omega^*$$

$$0 = -x^T A y + x^T A X \omega^* + \lambda \omega^*$$

$$\Rightarrow x^T A y = x^T A X \omega^* + \lambda \omega^*$$

$$\Rightarrow x^T A y = (x^T A X + \lambda I) \omega^*$$

$$\text{So, } \underline{\underline{(x^T A X + \lambda I)^{-1} \cdot x^T A y = \omega^*}} \xleftarrow{\text{Ans.}}$$

Q2b) Code in q2.py.

Q2c) As ' $\tau$ ' increases, the bandwidth of points that should be considered to make a prediction increases.

If ' $\tau$ ' is very small, you consider only the points which are in the close vicinity of our test point. So, we end up fitting only those selected set of points. This leads to overfitting.

If ' $\tau$ ' is large, you consider a large set of points when making a prediction. So, you end up with a general curve that fits all of the data but fails to capture irregularities. So, we end up underfitting.

Our training loss should increase as we are adding more points as ' $\tau$ ' increases. Our validation loss would decrease ~~until~~ <sup>as</sup> the model considers nearby points that follow a similar trend to the test point. Then, it would increase when other points are considered

as ' $\gamma$ ' increases.

Code for q2 is in q2.py.

Q3) AdaBoost

target labels  $\rightarrow \{-1, +1\}$

$$h_t \leftarrow \underset{h \in H}{\operatorname{argmin}} \sum_{i=1}^N w_i \mathbb{I} \{ h(x^{(i)}) \neq t^{(i)} \}$$

where  $t^{(i)}$  is the true label and  $h(x^{(i)})$  is the predicted value.

The Error is given by (I will be using ' $E_t$ ' instead of ' $\operatorname{err}_t$ ' to make it easier to read)

$$E_t = \frac{\sum_{i=1}^N w_i \mathbb{I} \{ h_t(x^{(i)}) \neq t^{(i)} \}}{\sum_{i=1}^N w_i}$$

=  $\frac{\text{sum of weights of all wrong samples}}{\text{sum of weights of all samples}}$ .

$$\text{Given} \rightarrow \alpha_t = \frac{1}{2} \log \left( \frac{1 - E_t}{E_t} \right). \quad \text{--- (1)}$$

and.

$$w_i' \leftarrow w_i \cdot e^{-\alpha_t \cdot h_t(x^{(i)}) \cdot t^{(i)}}. \quad \text{--- (2)}$$

Now,  $\begin{cases} h_t(x^{(i)}) \cdot t^{(i)} = 1 & \text{when } h_t(x^{(i)}) = t^{(i)} \text{ (Correct prediction)} \\ h_t(x^{(i)}) \cdot t^{(i)} = -1 & \text{when } h_t(x^{(i)}) \neq t^{(i)} \text{ (Wrong prediction)} \end{cases}$

Putting these in (2);

$$w_i' = \begin{cases} w_i e^{-\alpha_t} & \text{; correct prediction} \\ w_i e^{\alpha_t} & \text{; wrong prediction.} \end{cases} \quad \text{--- (3)}$$

Substituting the value of  $\alpha$  in (3).

$$w_i' = \begin{cases} w_i \cdot e^{-\frac{1}{2} \log \left( \frac{1-\epsilon_t}{\epsilon_t} \right)} & ; \text{ correct prediction} \\ w_i \cdot e^{\frac{1}{2} \log \left( \frac{1-\epsilon_t}{\epsilon_t} \right)} & ; \text{ wrong prediction} \end{cases}$$

simplifying, we get,

$$\Rightarrow w_i' = \begin{cases} w_i \sqrt{\frac{\epsilon_t}{1-\epsilon_t}} & ; \text{ correct pred.} \\ w_i \sqrt{\frac{1-\epsilon_t}{\epsilon_t}} & ; \text{ wrong pred / incorrect pred.} \end{cases} \quad \text{--- (6)}$$

we know,  $\frac{\sum_{\text{Incorrect}} w_i}{\sum_{\text{all}} w_i} = \epsilon_t \Rightarrow \frac{\sum_{\text{correct}} w_i}{\sum_{\text{all}} w_i} = 1 - \epsilon_t$

So,  $\sum_{i \in \text{incorrect}} w_i = \epsilon_t \cdot \sum_{i \in \text{all}} w_i \quad \text{--- (4)}$

$\sum_{i \in \text{correct}} w_i = (1 - \epsilon_t) \cdot \sum_{i \in \text{all}} w_i \quad \text{--- (5)}$

~~and finally~~

finally,

$$\sum_{i \in \text{all}} w_i' = \sum_{i \in \text{incorrect}} w_i' + \sum_{i \in \text{correct}} w_i'$$

putting (6) in above equation.

$$= w_i \cdot I \left\{ h(x^{(i)}) \neq t^{(i)} \right\} \sqrt{\frac{1-\epsilon_t}{\epsilon_t}}$$

$$= \sum_{i \in \text{incorrect}} w_i \sqrt{\frac{1-\epsilon_t}{\epsilon_t}} + \sum_{i \in \text{correct}} w_i \sqrt{\frac{\epsilon_t}{1-\epsilon_t}}$$

putting (4) & (5) in above equation.

$$\sum_{i \in \text{all}} w_i' = \epsilon_t \cdot \sqrt{\frac{1-\epsilon_t}{\epsilon_t}} \cdot \sum_{i \in \text{all}} w_i + (1-\epsilon_t) \cdot \sqrt{\frac{\epsilon_t}{1-\epsilon_t}} \cdot \sum_{i \in \text{all}} w_i$$

$$\sum_{i \in \text{all}} w_i' = \sqrt{\epsilon_t(1-\epsilon_t)} \sum_{i \in \text{all}} w_i + \sqrt{\epsilon_t(1-\epsilon_t)} \cdot \sum_{i \in \text{all}} w_i$$

$$\boxed{\sum_{i \in \text{all}} w_i' = 2 \sqrt{\epsilon_t(1-\epsilon_t)} \cdot \sum_{i \in \text{all}} w_i}$$

$$\sum_{i \in \text{incorrect}} w_i' = \sum_{i \in \text{incorrect}} w_i \cdot \sqrt{\frac{1-\epsilon_t}{\epsilon_t}}$$

$$= \epsilon_t \cdot \sum_{i \in \text{all}} w_i \cdot \sqrt{\frac{1-\epsilon_t}{\epsilon_t}} = \sqrt{\epsilon_t(1-\epsilon_t)} \cdot \sum_{i \in \text{all}} w_i'$$

$$\text{So, } \boxed{\sum_{i \in \text{incorrect}} w_i' = \sqrt{\epsilon_t(1-\epsilon_t)} \sum_{i \in \text{all}} w_i}$$

finally,  $\epsilon_t' = \frac{\sum_{i \in \text{incorrect}} w_i'}{\sum_{i \in \text{all}} w_i'}$

plugging all the values as found above:

$$\epsilon_t' = \frac{\sqrt{\epsilon_t(1-\epsilon_t)} \cdot \sum_{i \in \text{all}} w_i}{2 \sqrt{\epsilon_t(1-\epsilon_t)} \cdot \sum_{i \in \text{all}} w_i} = \frac{1}{2}$$

$$\text{So, } \boxed{\epsilon_t' = \frac{1}{2}} \quad (\text{hence proved}).$$



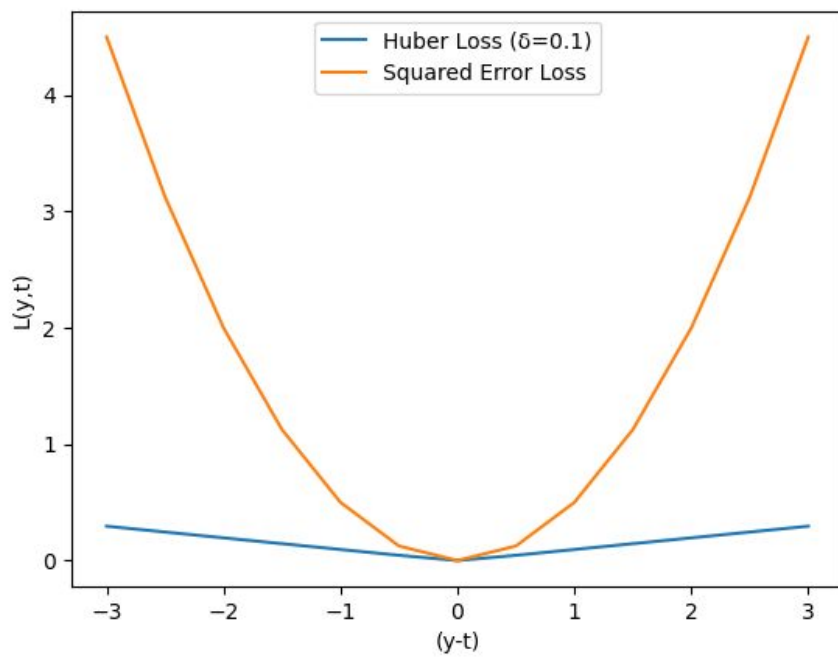
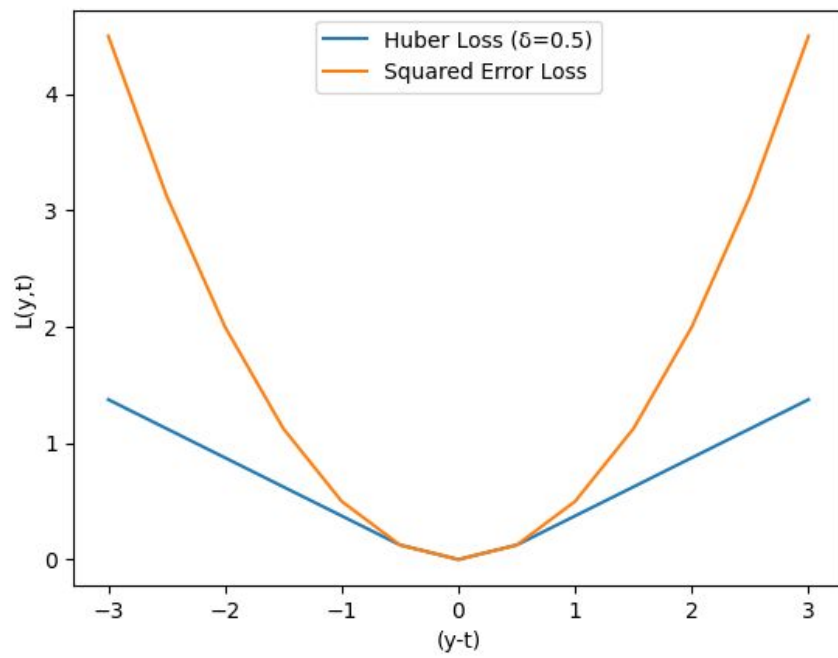
## Interpretation

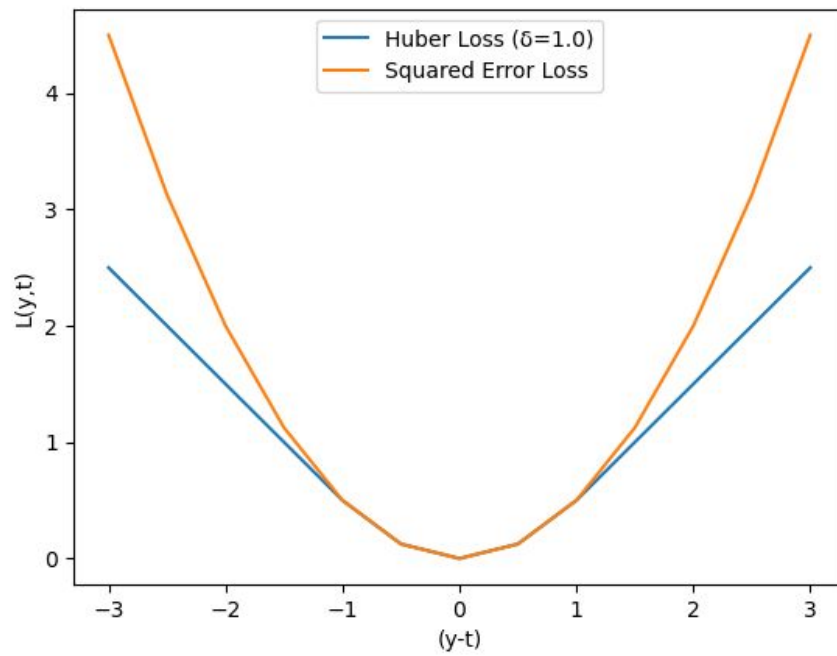
$$\epsilon_t' = \frac{\sum_{i \in \text{incorrect}} w_i'}{\sum_{i \in \text{all}} w_i'} = \frac{1}{2}$$

This means that the sum of all correct <sup>sample</sup> weights is equal to the sum of all incorrect sample weights. So, if we have weights  $w_i$  in the  $t^{\text{th}}$  time step, we will scale the weights at the  $t+1^{\text{th}}$  time step, so, that ~~weights~~ sum of all incorrect samples and correct samples are equal.

---

Q1.a. Plots :





Q2c. Plot: You can also generate the plot by running the code in q2.py

