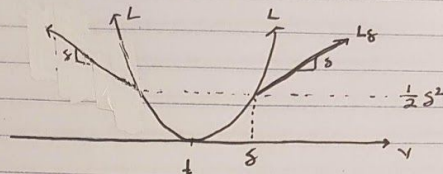


CSC411 Homework 3

CSC411 HW3 Feb 5, 2019

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



After $(y-t) > \delta$, the loss only increases linearly, and thus it influences the classifier less.

$$b) \rightarrow \frac{\partial L_S}{\partial \mathbf{w}} = \frac{\partial H_S}{\partial y} \cdot \frac{\partial y}{\partial \mathbf{w}} \quad \begin{aligned} y &= \mathbf{w}_1 x_1 + \dots + \mathbf{w}_n x_n \\ \frac{\partial y}{\partial \mathbf{w}} &= [x_1 \ x_2 \ \dots \ x_n] \\ \frac{\partial y}{\partial \mathbf{w}} &= \mathbf{x}^T \end{aligned}$$

Letting $a = y - t$

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

$$\frac{\partial H_S}{\partial a} = \begin{cases} a & |a| \leq \delta \\ \delta & |a| > \delta, a > \delta \\ -\delta & |a| > \delta, a < -\delta \end{cases}$$

$$\therefore \frac{\partial L_S}{\partial \mathbf{w}} = \begin{cases} (y-t) \mathbf{x}^T & |y-t| \leq \delta \\ \delta \mathbf{x}^T & |y-t| > \delta, y-t > \delta \\ -\delta \mathbf{x}^T & |y-t| > \delta, y-t < -\delta \end{cases}$$

$$\frac{\partial L_S}{\partial b} = \frac{\partial H_S(a)}{\partial a} \frac{\partial a}{\partial y} \frac{\partial y}{\partial b} = 1$$

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

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

$$y = \begin{bmatrix} y^1 \\ \vdots \\ y^n \end{bmatrix} \quad Xw = \begin{bmatrix} w^T x^1 \\ \vdots \\ w^T x^n \end{bmatrix}$$

$$\left(\begin{bmatrix} a^1 \\ \vdots \\ a^n \end{bmatrix} \begin{bmatrix} y_1 - w^T x^1 \\ \vdots \\ y_n - w^T x^n \end{bmatrix} \right)^T \begin{bmatrix} y_1 - w^T x^1 \\ \vdots \\ y_n - w^T x^n \end{bmatrix} = (A(y-Xw))^T (y-Xw) = (y-Xw)^T A^T (y-Xw) = (y-Xw)^T \tilde{A} (y-Xw)$$

$$\therefore C = \frac{1}{2} (y-Xw)^T A (y-Xw) + \frac{\lambda}{2} w^T w$$

$$\text{Let } B = \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix}$$

$$B^T A B = (b_1 \dots b_n) \begin{bmatrix} a^1 & \dots & a^n \end{bmatrix} \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix} = a^1 b_1^2 + \dots + a^n b_n^2$$

$$\frac{d(B^T A B)}{dB} = (2a^1 b_1 \quad \dots \quad 2a^n b_n)$$

$$= 2 B^T A$$

Notice: if $A = I$

$$\frac{d(B^T B)}{dB} = 2 B^T$$

$$Xw = \begin{bmatrix} w^T x^1 \\ \vdots \\ w^T x^n \end{bmatrix} \quad \frac{d(Xw)}{dw} = X$$

$$\frac{dC}{dw} = \frac{1}{2} \cdot \frac{d((y-Xw)^T A (y-Xw))}{d(y-Xw)} \cdot \frac{d(y-Xw)}{dw} + \frac{\lambda}{2} \frac{d(w^T w)}{dw}$$

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

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

Setting to 0

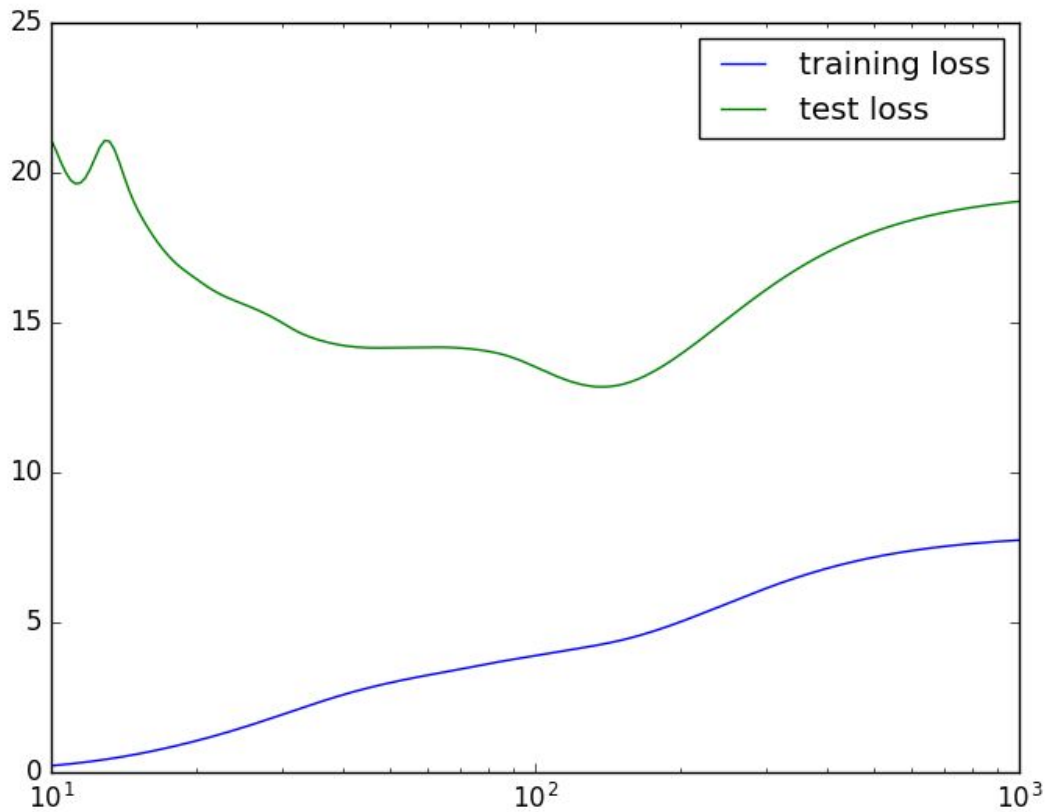
$$[0] = -(y^T - w^T X^T) A X + \lambda w^T$$

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

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

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

$$w = (X^T A X + \lambda I)^{-1} X^T A Y \quad \text{as required}$$



Training and test loss v.s. tau value

2d) As tau goes to infinity, all residuals get an equal weighting and the LRLS algorithm behaves like linear regression. Thus, the training loss increases because points further away influence the local decision. As tau goes to 0, only the closest points are considered. It is locally correct at the expense of points further away. Thus, the line is drawn better for each specific point and the training loss goes down. This is precisely the behaviour we see. A proof of behaviour as tau goes to 0 and infinity is on the next page.

Interestingly, the validation/test loss achieves a minimum somewhere in the middle. This makes sense. As tau goes to 0, the line is overfit and error increases. As tau goes to infinity, the line is locally underfit, and error increases.

Interpreting Tau

$$a^i = \frac{\exp(-b_i/2\tau^2)}{\sum_{j=1}^N \exp(-b_j/2\tau^2)} \quad b_i = \|x - x^i\|^2 \geq 0$$

$$\lim_{\tau \rightarrow \infty} a^i = \frac{\exp(0)}{\sum_{j=1}^N \exp(0)} = \frac{1}{N}$$

$$\lim_{\tau \rightarrow 0} a^i \neq \lim_{\tau \rightarrow 0} \frac{\exp(-b_i/2\tau^2) \cdot \frac{b_i}{\tau^2}}{\sum_{j=1}^N \exp(-b_j/2\tau^2) \cdot \frac{b_j}{\tau^2}} = ? \quad \times$$

$$a^i = \frac{1}{\sum_{j=1}^N \frac{\exp(-b_j/2\tau^2)}{\exp(-b_i/2\tau^2)}} = \frac{1}{\sum_{j=1}^N \exp\left(\frac{b_i - b_j}{2\tau^2}\right)}$$

$$a^i = \frac{1}{1 + \sum_{j \neq i} \exp\left(\frac{b_i - b_j}{2\tau^2}\right)}$$

$$\text{suppose } b_i > b_j : \lim_{\tau \rightarrow 0} \exp\left(\frac{b_i - b_j}{2\tau^2}\right) = \infty$$

$$\text{suppose } b_i < b_j : \lim_{\tau \rightarrow 0} \exp\left(\frac{b_i - b_j}{2\tau^2}\right) = 0$$

Clearly $a^i = 1$ or 0 .

Consider the closest point to x . Then $b_i \leq b_j \forall j$.

That point receives $a = 1$.

All other points get $a = 0$.