

# Csc 2515 Fall 2018

## Homework 3

Soon Chee Loong  
999295793

Last Name: Soon  
First Name: Chee Loong  
cheeloong.soon@mail.utoronto.ca cdf markus: soon chee

### 1. Robust Regression

↳ more robust to outliers

#### Huber Loss

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

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

a) Sketch the Huber loss  $L_{\delta}(y, t)$

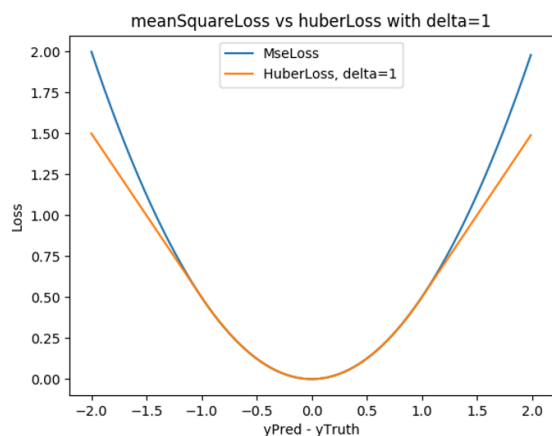
and squared error loss  $LSE(y, t) = \frac{1}{2}(y - t)^2$  for  $t = 0$ .

Based on sketch, why would you expect Huber Loss to be more robust to outliers?

Huber Loss trades off between sensitivity quadratic loss for strongly convex optimization with the robustness of absolute value loss towards outliers.

$\delta$  is a hyperparameter,  $\delta \geq 0$  that becomes more robust as we reduce it since it changes to absolute value loss quicker.

I used  $t = 0$ ,  $\delta = 1$ .



It is more robust as it is less sensitive to outliers that are far away.

b) Just as with linear regression, assume linear model

$$y = w^T \cdot x + b, \quad Y = \tilde{X} \cdot \tilde{w} + b, \quad \tilde{X} \text{ is design matrix}$$

Calculate partial derivatives  $\frac{\partial L_{\delta}}{\partial w}$  and  $\frac{\partial L_{\delta}}{\partial b}$ .

Hint: Find formula for derivative  $H'_{\delta}(a)$ , then give answers in terms of  $H'_{\delta}(y - t)$ .

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

$$= \begin{cases} \frac{1}{2}a & \text{if } |a| \leq \delta \\ \delta(a - \frac{\delta}{2}) & \text{if } |a| > \delta, a \geq 0 \\ \delta(-a - \frac{\delta}{2}) & \text{if } |a| > \delta, a < 0 \end{cases}$$

b) Just as with linear regression, assume linear model

$$y = w^T \cdot x + b, \quad Y = \tilde{X} \cdot \tilde{w} + b, \quad \tilde{X} \text{ is design matrix}$$

Calculate partial derivatives  $\frac{\partial L_S}{\partial w}$  and  $\frac{\partial L_S}{\partial b}$ .

Hint: Find formula for derivative  $H_S'(a)$ , then give answers in terms of  $H_S'(y-t)$ .

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

$$= \begin{cases} \frac{1}{2} a & \text{if } |a| \leq \delta \\ \delta(a - \frac{\delta}{2}) & \text{if } |a| > \delta, a \geq 0 \\ \delta(-a - \frac{\delta}{2}) & \text{if } |a| > \delta, a < 0 \end{cases}$$

$$a = y - t$$

$$y = \tilde{w}^T \cdot \tilde{x} + b = \sum_{i=1}^d w_i x_i + b, \quad d = \text{dimensionality of } \tilde{x}, \tilde{w}, \tilde{x}, \tilde{w} \in \mathbb{R}^d$$

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

From Chain Rule,

$$\frac{dL_S}{d\tilde{w}} = \frac{dH_S(a)}{da} = \left( \frac{dH_S(a)}{da} \right) \left( \frac{da}{d\tilde{w}} \right) = \left( \frac{dH_S(a)}{da} \right) \left( \frac{da}{dy} \right) \left( \frac{dy}{d\tilde{w}} \right)$$

$$\frac{dL_S}{db} = \frac{dH_S(a)}{da} = \left( \frac{dH_S(a)}{da} \right) \left( \frac{da}{db} \right) = \left( \frac{dH_S(a)}{da} \right) \left( \frac{da}{dy} \right) \left( \frac{dy}{db} \right)$$

$$H_S'(a) = \frac{dH_S(a)}{da} = \begin{cases} a & , \text{ if } |a| \leq \delta \\ \delta & , \text{ if } |a| > \delta, a > 0 \\ -\delta & , \text{ if } |a| > \delta, a < 0 \end{cases}$$

$$a' = \frac{da}{dy} = 1$$

$$\frac{dy}{d\tilde{w}} = \begin{bmatrix} \frac{dy}{dw_1} \\ \vdots \\ \frac{dy}{dw_d} \end{bmatrix}_{(d \times 1)} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}_{(d \times 1)} = \tilde{x}$$

$$\frac{dy}{db} = 1$$

$\therefore$

$$\frac{dL_S}{d\tilde{w}} = \left( \frac{dH_S(a)}{da} \right) \left( \frac{da}{dy} \right) \left( \frac{dy}{d\tilde{w}} \right)$$

$$= \begin{cases} a(1)(\tilde{x}) & , \text{ if } |a| \leq \delta \\ \delta(1)(\tilde{x}) & , \text{ if } |a| > \delta, a > 0 \\ -\delta(1)(\tilde{x}) & , \text{ if } |a| > \delta, a < 0 \end{cases}$$

$$= \begin{cases} a(\tilde{x}) & , \text{ if } |a| \leq \delta \\ \delta(\tilde{x}) & , \text{ if } |a| > \delta, a > 0 \\ -\delta(\tilde{x}) & , \text{ if } |a| > \delta, a < 0 \end{cases}$$

$$dL_S = \left( \frac{dH_S(a)}{da} \right) \left( \frac{da}{dy} \right) \left( \frac{dy}{d\tilde{w}} \right)$$

$\therefore$ 

$$\frac{dL_S}{d\vec{w}} = \left( \frac{dH_S(a)}{da} \right) \left( \frac{da}{dy} \right) \left( \frac{dy}{d\vec{w}} \right)$$

$$= \begin{cases} a(1) (\vec{x}) & , \text{ if } |a| \leq \delta \\ \delta(1) (\vec{x}) & , \text{ if } |a| > \delta, a > 0 \\ -\delta(1) (\vec{x}) & , \text{ if } |a| > \delta, a < 0 \end{cases}$$

$$= \begin{cases} a(\vec{x}) & , \text{ if } |a| \leq \delta \\ \delta(\vec{x}) & , \text{ if } |a| > \delta, a > 0 \\ -\delta(\vec{x}) & , \text{ if } |a| > \delta, a < 0 \end{cases}$$

$$\frac{dL_S}{db} = \left( \frac{dH_S(a)}{da} \right) \left( \frac{da}{dy} \right) \left( \frac{dy}{db} \right)$$

$$= \begin{cases} a(1) (1) & , \text{ if } |a| \leq \delta \\ \delta(1) (1) & , \text{ if } |a| > \delta, a > 0 \\ -\delta(1) (1) & , \text{ if } |a| > \delta, a < 0 \end{cases}$$

$$= \begin{cases} a & , \text{ if } |a| \leq \delta \\ \delta & , \text{ if } |a| > \delta, a > 0 \\ -\delta & , \text{ if } |a| > \delta, a < 0 \end{cases} \Rightarrow \text{piecewise derivative}$$

$$= \begin{cases} a & , \text{ if } |a| \leq \delta \\ \delta & , \text{ if } |a| > \delta, a > 0 \\ -\delta & , \text{ if } |a| > \delta, a < 0 \end{cases} \Rightarrow \text{piecewise derivative}$$

note: Vectorized gradients were implemented directly in code as the average between the gradients  $\forall$  training data

c) Write python code to perform full batch gradient descent.

Initialize  $w$  and  $b$  to zeros

Assume dataset is given as design matrix  $X$  and target vector  $y$ .

No for loops over training examples or input dimensions.

Hint: use `np.where`

Since gradients for Huber Loss are piecewise,

`np.where()` is useful for figuring out which piece each training instance is in.

Submitted as `q1.py`

## 2. Locally Weighted Regression

a) Given  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$  and positive weights  $a^{(1)}, \dots, a^{(n)}$ .

Show that the solution to the weighted least square problem

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

is given by formula

$$\vec{w}^* = (\vec{X}^T \vec{A} \vec{X} + \lambda \mathbf{I})^{-1} \vec{X}^T \vec{A} \vec{y}$$

$X$  = Design Matrix

## 2. Locally Weighted Regression

a) Given  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\}$  and positive weights  $a^{(1)}, \dots, a^{(N)}$ .

Show that the solution to the weighted least square problem

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

is given by formula

$$\vec{w}^* = (\vec{X}^T \vec{A} \vec{X} + \lambda \mathbf{I})^{-1} \vec{X}^T \vec{A} \vec{y}$$

$\vec{X}$  = Design Matrix

$$\vec{A} = \operatorname{diag}(\vec{a}), \vec{A}_{ii} = a^{(i)}$$

Can ignore constant of  $\frac{1}{2}$  during minimization

The loss function to minimize is

$$\begin{aligned} L_1(\vec{w}) &= \frac{1}{2} \sum_{i=1}^N a^{(i)} (y^{(i)} - \vec{w}^T \vec{x}^{(i)})^2 + \frac{\lambda}{2} \|\vec{w}\|^2 \\ &= \frac{1}{2} \left( \sum_{i=1}^N a^{(i)} (y^{(i)} - \vec{w}^T \vec{x}^{(i)})^2 + \lambda \|\vec{w}\|^2 \right) \end{aligned}$$

Since  $\operatorname{minimize}(c L(\vec{w})) = \operatorname{minimize}(L(\vec{w}))$  for any loss function,  $L(\vec{w})$

multiplied by a constant,  $c$ , we can ignore the constant

of  $\frac{1}{2}$  by minimizing  $2 \times L_1(\vec{w})$  instead.

Let  $L = 2 \times L_1$

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

where

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

$$= \underset{\vec{w}}{\operatorname{argmin}} (L_1(\vec{w}))$$

$$= \underset{\vec{w}}{\operatorname{argmin}} (2 \times L_1(\vec{w}))$$

$$= \underset{\vec{w}}{\operatorname{argmin}} (L(\vec{w}))$$

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

Re-write Loss function,  $L(\vec{w})$  using Matrices instead of summations

Let

Design Matrix,  $\vec{X}$

$$\vec{X} = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & \dots & x_d^{(1)} \\ x_1^{(2)} & \dots & & x_d^{(2)} \\ \vdots & & \ddots & \vdots \\ x_1^{(N)} & x_2^{(N)} & \dots & x_d^{(N)} \end{bmatrix}_{(N \times d)} = \begin{bmatrix} (\vec{x}^{(1)})^T \\ (\vec{x}^{(2)})^T \\ \vdots \\ (\vec{x}^{(N)})^T \end{bmatrix}$$

$$\vec{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(N)} \end{bmatrix}, \quad \vec{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix}$$

Re-write Loss function,  $L(w)$  using Matrices instead of summations

Let

Design Matrix,  $\vec{X}$

$$\vec{X} = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & \dots & x_d^{(1)} \\ x_1^{(2)} & \dots & & x_d^{(2)} \\ \vdots & & & \vdots \\ x_1^{(N)} & x_2^{(N)} & \dots & x_d^{(N)} \end{bmatrix}_{(N \times d)} = \begin{bmatrix} (\vec{x}^{(1)})^T \\ (\vec{x}^{(2)})^T \\ \vdots \\ (\vec{x}^{(N)})^T \end{bmatrix}$$

$$Y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(N)} \end{bmatrix}_{(N \times 1)}, \quad \vec{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix}_{(d \times 1)}$$

$$A = \text{diag}(\vec{a}) = \begin{bmatrix} a^{(1)} & 0 & \dots & 0 \\ 0 & a^{(2)} & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 0 & a^{(N)} \end{bmatrix}_{(N \times N)}$$

$$B = \text{diag}(\sqrt{\vec{a}}) = \begin{bmatrix} \sqrt{a^{(1)}} & 0 & \dots & 0 \\ 0 & \sqrt{a^{(2)}} & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 0 & \sqrt{a^{(N)}} \end{bmatrix}_{(N \times N)} \Rightarrow B^T \cdot B = B \cdot B^T = A$$

$$\therefore L(w) = \sum_{i=1}^N a^{(i)} (y^{(i)} - w^T x^{(i)})^2 + \lambda \|w\|^2$$

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

$$= B(Y - Xw))^2 + \lambda \|w\|_F^2, \quad \|w\|_F^2 = \vec{w}^T \cdot \vec{w}$$

$$\vec{w}^* = \underset{\vec{w}}{\text{argmin}} (B(Y - Xw))^2 + \lambda \|w\|_F^2)$$

Matrix Derivatives and proof:

Since we are dealing with Matrix Algebra, will prove some Matrix Differentiation formula before using these formula to come up with closed-form solution for Locally Weighted Regression.

$F$  independent of  $\vec{w}$ ,  $\vec{w} \in \mathbb{R}^{d \times 1}$ ,  $F \in \mathbb{R}^{m \times d}$

$$g = F \cdot \vec{w}$$

$$g = \begin{bmatrix} F_{11} & \dots & F_{1d} \\ \vdots & \ddots & \vdots \\ F_{m1} & \dots & F_{md} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} \Rightarrow g_i = \sum_{j=1}^d F_{ij} w_j, \quad i \in \{1, 2, \dots, m\}$$

$$\frac{\partial g}{\partial \vec{w}} = \begin{bmatrix} \frac{\partial g_1}{\partial w_1} & \dots & \frac{\partial g_m}{\partial w_1} \end{bmatrix} \Rightarrow \frac{\partial g_i}{\partial w_k} = \frac{\partial}{\partial w_k} \sum_{j=1}^d F_{ij} w_j = F_{ik}$$

## Matrix Derivatives and proof:

Since we are dealing with Matrix Algebra, will prove some Matrix Differentiation formula before using these formula to come up with closed-form solution for Locally Weighted Regression.

$F$  independent of  $\vec{w}$ ,  $\vec{w} \in \mathbb{R}^{d \times 1}$ ,  $F \in \mathbb{R}^{m \times d}$

$$g = F \cdot \vec{w}$$

$$g = \begin{bmatrix} F_{11} & \dots & F_{1d} \\ \vdots & & \vdots \\ F_{m1} & \dots & F_{md} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} \Rightarrow g_i = \sum_{j=1}^d F_{ij} w_j, i \in \{1, 2, \dots, m\}$$

$$\frac{\partial g}{\partial \vec{w}} = \begin{bmatrix} \frac{\partial g_1}{\partial w_1} & \dots & \frac{\partial g_m}{\partial w_1} \\ \vdots & & \vdots \\ \frac{\partial g_1}{\partial w_d} & \dots & \frac{\partial g_m}{\partial w_d} \end{bmatrix} \Rightarrow \frac{\partial g_i}{\partial w_k} = \frac{\partial \sum_{j=1}^d F_{ij} w_j}{\partial w_k} = F_{ik}$$

i comes from g  
k " " w

$$= \begin{bmatrix} F_{11} & F_{12} & \dots & F_{1d} \\ F_{21} & F_{22} & & \vdots \\ \vdots & & & \vdots \\ F_{m1} & \dots & & F_{md} \end{bmatrix} = F^T$$

$$\therefore, \frac{d(F \cdot \vec{w})}{d\vec{w}} = F^T \Rightarrow \text{Take transpose of what's on left of } \vec{w}$$

$$\text{Let } h = \vec{w}^T \cdot F^T$$

$$\frac{\partial h}{\partial \vec{w}} = \frac{\partial \vec{w}^T \cdot F^T}{\partial \vec{w}}$$

$$\text{Let } p = h^T = F \cdot w$$

$$\Rightarrow \frac{\partial h}{\partial \vec{w}} = \frac{\partial h^T}{\partial \vec{w}} \text{ since } h \in \mathbb{R} \Rightarrow \text{scalars are always symmetric}$$

$$\frac{\partial h^T}{\partial w} = F^T \text{ as proven above}$$

$$\therefore, \frac{\partial \vec{w}^T \cdot F^T}{\partial \vec{w}} = F^T \Rightarrow \text{just rewrite what's on right of } \vec{w}$$

Quadratic Form,

independent of  $x$ ,  $A \in \mathbb{R}^{n \times n}$ ,  $x \in \mathbb{R}^n$

$$f = \vec{w}^T \cdot H \cdot \vec{w}, H \text{ independent of } w, H \in \mathbb{R}^{d \times d}, w \in \mathbb{R}^d$$

$$f = \sum_{j=1}^d \left( \sum_{i=1}^d w_i \cdot H_{ij} \right) x_j \in \mathbb{R}$$

$$\frac{\partial f}{\partial \vec{w}} = \begin{bmatrix} \frac{\partial f}{\partial w_1} \\ \vdots \\ \frac{\partial f}{\partial w_d} \end{bmatrix}, \frac{\partial f}{\partial w_k} = \frac{\partial}{\partial w_k} \left( \sum_{j=1}^d \left( \sum_{i=1}^d w_i \cdot H_{ij} \right) x_j \right)$$

$$= \sum_{i=1}^d H_{ik} w_i + \sum_{j=1}^d H_{kj} w_j, \text{Product Rule}$$

$$\therefore, \frac{\partial \vec{w}^T \cdot \vec{F}^T}{\partial \vec{w}} = \vec{F}^T \Rightarrow \text{just rewrite what's on right of } \vec{w}$$

Quadratic Form,

independent of  $\vec{x}$ ,  $A \in \mathbb{R}^{n \times n}$ ,  $\vec{x} \in \mathbb{R}^n$

$$f = \vec{w}^T \cdot H \cdot \vec{w}, \quad H \text{ independent of } \vec{w}, \quad H \in \mathbb{R}^{d \times d}, \quad \vec{w} \in \mathbb{R}^d$$

$$f = \sum_{j=1}^d \left( \sum_{i=1}^d w_i \cdot H_{ij} \right) x_j \in \mathbb{R}$$

$$\frac{\partial f}{\partial \vec{w}} = \begin{bmatrix} \frac{\partial f}{\partial w_1} \\ \vdots \\ \frac{\partial f}{\partial w_d} \end{bmatrix}, \quad \frac{\partial f}{\partial w_k} = \frac{\partial}{\partial w_k} \left( \sum_{j=1}^d \left( \sum_{i=1}^d w_i \cdot H_{ij} \right) x_j \right)$$

$$= \sum_{i=1}^d H_{ik} w_i + \sum_{j=1}^d H_{kj} w_j, \quad \text{Product Rule}$$

$$= \begin{bmatrix} \sum_{i=1}^d H_{i1} w_i + \sum_{j=1}^d H_{1j} w_j \\ \vdots \\ \sum_{i=1}^d H_{id} w_i + \sum_{j=1}^d H_{dj} w_j \end{bmatrix}$$

$$= \begin{bmatrix} \sum_{i=1}^d H_{i1} w_i \\ \vdots \\ \sum_{i=1}^d H_{id} w_i \end{bmatrix} + \begin{bmatrix} \sum_{j=1}^d H_{1j} w_j \\ \vdots \\ \sum_{j=1}^d H_{dj} w_j \end{bmatrix}$$

$$= H^T \cdot \vec{w} + H \cdot \vec{w}$$

$$= (H^T + H) \cdot \vec{w}$$

$$\therefore, \frac{\partial (\vec{w}^T \cdot H \cdot \vec{w})}{\partial \vec{w}} = (H^T + H) \cdot \vec{w}, \quad \text{if } H^T = H \Rightarrow 2H \cdot \vec{w}$$

Closed-Form Solution for Locally Weighted Regression

Take derivative and set to 0 to solve for  $\vec{w}^*$

$$L(\vec{w}) = (\vec{B} (\vec{y} - \vec{x} \vec{w}))^2 + \lambda \|\vec{w}\|_F^2$$

$$= (\vec{B} (\vec{y} - \vec{x} \vec{w}))^T (\vec{B} (\vec{y} - \vec{x} \vec{w})) + \lambda \vec{w}^T \vec{w}$$

$$= (\underbrace{(\vec{y} - \vec{x} \vec{w})^T}_{(1 \times m)} \underbrace{\vec{B}^T}_{(m \times 1)}) \cdot (\underbrace{\vec{B} \vec{y}}_{(m \times 1)} - \underbrace{\vec{B} \vec{x} \vec{w}}_{(m \times 1)}) + \lambda \vec{w}^T \vec{w}$$

$$= ((\vec{y}^T - \vec{w}^T \vec{x}^T) \vec{B}^T) \cdot (\vec{B} \vec{y} - \vec{B} \vec{x} \vec{w}) + \lambda \vec{w}^T \vec{w}$$

$$= (\vec{y}^T \vec{B}^T - \vec{w}^T \vec{x}^T \vec{B}^T) (\vec{B} \vec{y} - \vec{B} \vec{x} \vec{w}) + \lambda \vec{w}^T \vec{w}$$

$$= \vec{y}^T \vec{B}^T \vec{B} \vec{y} - \vec{y}^T \vec{B}^T \vec{B} \vec{x} \vec{w} - \vec{w}^T \vec{x}^T \vec{B}^T \vec{B} \vec{y} + \vec{w}^T \vec{x}^T \vec{B}^T \vec{B} \vec{x} \vec{w} + \lambda \vec{w}^T \vec{w}$$

## Closed-Form Solution for Locally Weighted Regression

Take derivative and set to 0 to solve for  $w^*$

$$L(\tilde{w}) = (\tilde{B} (Y - \tilde{X}\tilde{w}))^2 + \lambda \|\tilde{w}\|_F^2$$

$$= (\tilde{B} (Y - \tilde{X}\tilde{w}))^T (\tilde{B} (Y - \tilde{X}\tilde{w})) + \lambda \tilde{w}^T \tilde{w}$$

$$= (\underbrace{(Y - \tilde{X}\tilde{w})^T}_{(1 \times m)} \tilde{B}^T) \cdot \underbrace{(\tilde{B} Y - \tilde{B} \tilde{X} \tilde{w})}_{(m \times 1)} + \lambda \tilde{w}^T \tilde{w}$$

$$= ((Y^T - \tilde{w}^T \tilde{X}^T) \tilde{B}^T) \cdot (\tilde{B} Y - \tilde{B} \tilde{X} \tilde{w}) + \lambda \tilde{w}^T \tilde{w}$$

$$= (\tilde{Y}^T \tilde{B}^T - \tilde{w}^T \tilde{X}^T \tilde{B}^T) (\tilde{B} Y - \tilde{B} \tilde{X} \tilde{w}) + \lambda \tilde{w}^T \tilde{w}$$

$$= \tilde{Y}^T \tilde{B}^T \tilde{B} Y - \tilde{Y}^T \tilde{B}^T \tilde{B} \tilde{X} \tilde{w} - \tilde{w}^T \tilde{X}^T \tilde{B}^T \tilde{B} Y + \tilde{w}^T \tilde{X}^T \tilde{B}^T \tilde{B} \tilde{X} \tilde{w} + \lambda \tilde{w}^T \tilde{w}$$

Take Matrix derivative set to 0, applying closed form derivatives derived above

$$\frac{dL(\tilde{w})}{d\tilde{w}} = -\tilde{X}^T \tilde{B}^T \tilde{B} Y - \tilde{X}^T \tilde{B}^T \tilde{B} \tilde{X} \tilde{w} + 2\tilde{X}^T \tilde{B}^T \tilde{B} \tilde{X} \tilde{w} + \lambda 2\tilde{w} = 0$$

$$\Rightarrow (2\tilde{X}^T \tilde{B}^T \tilde{B} \tilde{X} + \lambda \cdot I) \tilde{w} = 2\tilde{X}^T \tilde{B}^T \tilde{B} Y$$

$$\Rightarrow (\tilde{X}^T \tilde{B}^T \tilde{B} \tilde{X} + \lambda \cdot I) \tilde{w} = \tilde{X}^T \tilde{B}^T \tilde{B} Y$$

$$\Rightarrow \tilde{w} = (\tilde{X}^T \tilde{B}^T \tilde{B} \tilde{X} + \lambda I)^{-1} \tilde{X}^T \tilde{B}^T \tilde{B} Y$$

Since  $\tilde{B}^T \tilde{B} = \tilde{A}$ ,

$$\Rightarrow \tilde{w} = (\tilde{X}^T \tilde{A} \tilde{X} + \lambda I)^{-1} \tilde{X}^T \tilde{A} Y \quad \therefore, \text{ proven}$$

b) Locally reweighted least squares combines ideas from kNN and linear regression.

For each new test example  $\tilde{x}$ , we compute distance-based weights for each training example

$$a^{(i)} = \frac{e^{-\frac{\|x - x^{(i)}\|^2}{2\tau^2}}}{\sum_j e^{-\frac{\|x - x^{(j)}\|^2}{2\tau^2}}}$$

$$W^* = \arg \min \frac{1}{2} \sum_{i=1}^N a^{(i)} (y^{(i)} - w^T \cdot x^{(i)})^2 + \frac{\lambda}{2} \|w\|^2$$

$$\hat{y} = x^T \cdot w^*$$

Complete implementation in `q2.py`.

Notes for implementation.

→ Do not invert any matrix. Use `np.linalg.solve`

$$\rightarrow \frac{\exp(A_i)}{\sum_j \exp(A_j)} = \frac{\exp(A_i - B)}{\sum_j \exp(A_j - B)} = \frac{\exp(A_i - \max_j(A_j))}{\sum_j \exp(A_j - \max_j(A_j))} \rightarrow \text{more numerically stable}$$

⇒ Use `scipy.misc.logsumexp`.



b) Locally reweighted least squares combines ideas from kNN and linear regression. For each new test example  $\tilde{x}$ , we compute distance-based weights for each training example

$$a^{(i)} = \frac{e^{-\frac{\|x - x^{(i)}\|^2}{2\tau^2}}}{\sum_j e^{-\frac{\|x - x^{(j)}\|^2}{2\tau^2}}}$$

$$W^* = \arg\min \frac{1}{2} \sum_{i=1}^N a^{(i)} (y^{(i)} - w^T \cdot x^{(i)})^2 + \frac{\lambda}{2} \|w\|^2$$

$$\hat{y} = x^T \cdot w^*$$

Complete implementation in `q2.py`.

Notes for implementation.

→ Do not invert any matrix. Use `np.linalg.solve`

$$\rightarrow \frac{\exp(A_i)}{\sum_j \exp(A_j)} = \frac{\exp(A_i - B)}{\sum_j \exp(A_j - B)} = \frac{\exp(A_i - \max_j(A_j))}{\sum_j \exp(A_j - \max_j(A_j))} \rightarrow \text{more numerically stable}$$

⇒ Use `scipy.misc.logsumexp`.

$$\text{Let } A_i = -\frac{1}{2\tau^2} \|x - x^{(i)}\|^2 \in \mathbb{R}^{n \times n}$$

⇒ pairwise distance between  $\forall$  pair of points  
 →  $A_i$  in hint function above

$$a^{(i)} = \frac{e^{A_i}}{\sum_j e^{A_j}}, \quad \log(a^{(i)}) = A_i - \log\left(\sum_j e^{A_j}\right) \\ = A_i - \text{scipy.misc.logsumexp}(A)$$

⇒ use  $\log(a^{(i)})$  as weights instead

$$\tilde{A} = \text{np.diagflat}(a)$$

→  $\tilde{A}$  as weight for each training point

$$\tilde{w} = (\tilde{x}^T \tilde{A} \tilde{x} + \lambda \mathbf{I})^{-1} \tilde{x}^T \tilde{A} \tilde{y}$$

$$\Rightarrow (\tilde{x}^T \tilde{A} \tilde{x} + \lambda \mathbf{I}) \tilde{w} = \tilde{x}^T \tilde{A} \tilde{y}$$

$a x = b$  in linear algebra

$$\Rightarrow \text{np.linalg.solve}(\tilde{x}^T \tilde{A} \tilde{x} + \lambda \mathbf{I}, \tilde{x}^T \tilde{A} \tilde{y})$$

Submitted under `q2.py`

c) Randomly hold out 30% of dataset as validation set.

Compute the average loss for different values of  $\tau$  in the range  $[10, 1000]$  on both the training set and validation set.

Plot the training and validation losses as a function of  $\tau$ .

(using a log scale for  $\tau$ )

Submitted under `q2.py`

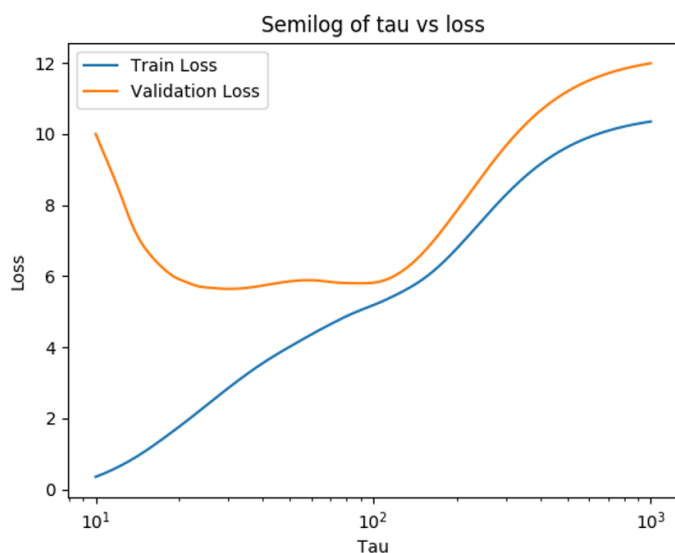
c) Randomly hold out 30% of dataset as validation set.

Compute the average loss for different values of  $\tau$  in the range  $[10, 1000]$  on both the training set and validation set.

Plot the training and validation losses as a function of  $\tau$ .

(using a log scale for  $\tau$ )

Submitted under q2.py



d) How would you expect this algorithm to behave as  $\tau \rightarrow \infty$ ?

When  $\tau \rightarrow 0$ ? Is this what actually happened?

As  $\tau \rightarrow \infty \Rightarrow$  all training point has equal weight

$\Rightarrow$  algorithm fits a single global unweighted linear regression for each training point.

$\Rightarrow$  underfitting  $\Rightarrow$  high bias

$\tau \rightarrow 0 \Rightarrow$  closest training point has most weight,

algorithm fits separate local linear regression for every training point.

$\Rightarrow$  overfitting to training point  $\Rightarrow$  high variance

Yes, this occurs in the plot shown in Q2. c) above.

As  $\tau \rightarrow 0$ , training loss is minimized, but test loss is high

As we sweep  $\tau$ , train loss should increase as no longer overfitting

Test loss shows a nice minimum due to

bias-variance decomposition of sweeping through  $\tau$ )