

4a) Derivation of the closed form solution for parameter θ that minimizes the loss function $J(\theta)$ in ridge regression.

$$J(\theta) = \frac{1}{2} \sum_{i=1}^N (h_{\theta}(x_i) - y_i)^2 + \frac{\lambda}{2} \sum_{j=1}^d \theta_j^2$$

The term d in regularization part denotes amount of covariables used in the model.

We can rewrite $J(\theta)$ in matrix notation and further break it down.

$$J(\theta) = \frac{1}{2} (X\theta - y)^T (X\theta - y) + \frac{\lambda}{2} (\theta^T \theta)$$

The term $\frac{\lambda}{2} \sum_{j=1}^d \theta_j^2$ represents regularization we apply on coefficients.

$$\Rightarrow \frac{1}{2} [X^T \theta^T X \theta - X^T \theta^T y - X \theta y^T + y^T y + \lambda \theta^T \theta]$$

$[X \theta]^T y = (X \theta) y^T$ the a transposed scalar is same scalar.

$$\Rightarrow \frac{1}{2} [X^T \theta^T X \theta - 2 X^T \theta^T y + y^T y + \theta^T \lambda I \theta] \quad [I = \text{Identity matrix}]$$

$$\Rightarrow \frac{1}{2} [y^T y - 2 X^T \theta^T y + \theta^T (X^T X + \lambda I) \theta]$$

θ should minimize $J(\theta)$.

By matrix differentiation rule $\frac{\partial x^T A x}{\partial x} = (A + A^T)x = 2Ax$

We can apply it in here as

$$\frac{\partial J}{\partial \theta} \Rightarrow \frac{1}{2} [0 - 2 X^T y + 2(X^T X + \lambda I) \theta] = 0$$

$$\Rightarrow -X^T y + (X^T X + \lambda I) \theta = 0$$

$$\Rightarrow (X^T X + \lambda I) \theta = X^T y$$

$$\boxed{\theta = (X^T X + \lambda I)^{-1} X^T y} \text{ is a closed form eq.}$$

5) $X = \begin{bmatrix} 1 & x_{11} & x_{j1} \\ 1 & x_{12} & x_{j2} \\ \vdots & \vdots & \vdots \\ 1 & x_{1n} & x_{jn} \end{bmatrix}_{n \times 3}$

$$X^T = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_{11} & x_{12} & \dots & x_{1n} \\ x_{j1} & x_{j2} & \dots & x_{jn} \end{bmatrix}_{3 \times n}$$

$$x_{j1} = 2x_{i1}$$

$$So \quad X = \begin{bmatrix} 1 & x_{i1} & 2x_{i1} \\ 1 & x_{i2} & 2x_{i2} \\ \vdots & \vdots & \vdots \\ 1 & x_{in} & 2x_{in} \end{bmatrix}_{n \times 3}$$

$$X^T = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_{i1} & x_{i2} & \dots & x_{in} \\ 2x_{i1} & 2x_{i2} & \dots & 2x_{in} \end{bmatrix}_{3 \times n}$$

$$\begin{bmatrix} 1 & 1 & \dots & 1 \\ x_{i1} & x_{i2} & \dots & x_{in} \\ 2x_{i1} & 2x_{i2} & \dots & 2x_{in} \end{bmatrix} \begin{bmatrix} 1 & x_{i1} & 2x_{i1} \\ 1 & x_{i2} & 2x_{i2} \\ \vdots & \vdots & \vdots \\ 1 & x_{in} & 2x_{in} \end{bmatrix}$$

$$X^T X = \begin{bmatrix} 1 & \sum_{i=1}^n x_i & 2 \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i & \sum_{i=1}^n x_i^2 & 2 \sum_{i=1}^n x_i^2 \\ 2 \sum_{i=1}^n x_i & 2 \sum_{i=1}^n x_i^2 & 4 \sum_{i=1}^n x_i^2 \end{bmatrix}$$

Row 1 and Row 2 are linearly dependent

$$\text{Rank}(X^T X) = 1 \quad (\text{using row echelon form}).$$

The columns of $X^T X$ are linearly dependent, $\text{col } 1 \times \sum_{i=1}^n x_i = \text{col } 2$ and etc.

As per invertible matrix theorem linearly dependent matrix do not have inverses.

So, $X^T X$ is non-invertible.

(6b) Laplace PDF

$$P(x) = \frac{e^{-|x|/b}}{2b}$$

hypothesis be,

$$y_i = h_\theta(x_i) + \epsilon_i$$

ϵ_i - random noise generated from Laplace distribution.

$$f(y_i | x_i; \theta, b) = \frac{1}{2b} e^{-\frac{|y_i - (\theta_0 + \theta_1 x_i)|}{b}}$$

(func of y_i given x_i and parameters for Laplace distribution)

For, Assuming that the points are independent.

$$f(y_1, \dots, y_n | x_1, \dots, x_n, \theta, b) \cdot P[y | x, \theta] = \text{Max}_{\theta} L(\theta)$$

$$\text{MLE } L(\theta) = \prod_{i=1}^n f(y_i | x_i, \theta)$$

$$\log(L(\theta)) = \sum_{i=1}^n \log(f(y_i | x_i, \theta))$$

$$\Rightarrow \sum_{i=1}^n \log\left(\frac{1}{2b} \cdot e^{-|y_i - (\theta_0 + \theta_1 x_i)|/b}\right)$$

$$\Rightarrow \sum_{i=1}^n \left[\log\left(\frac{1}{2b}\right) - \frac{|y_i - (\theta_0 + \theta_1 x_i)|}{b} \right]$$

Minimize all negative terms to maximize the MLE function.

$$(ie) \quad J(\theta) = \sum_{i=1}^n \frac{|y_i - (\theta_0 + \theta_1 x_i)|}{b}$$

$$\therefore J(\theta) = \frac{1}{b} \sum_{i=1}^n \frac{|y_i - (\theta_0 + \theta_1 x_i)|}{1} \quad \xrightarrow{\text{MAE}}$$