

Q1 $b = -34.084$

$w = \begin{bmatrix} 0.6482 \\ -0.0632 \end{bmatrix}$

Test error = 21.4859

Predicted response at $x = [100, 100] = 24.4581$

Q2

$$J = \min_{w, b} \sum_{i=1}^n c_i (y_i - w^T x_i - b)^2$$

$$= \min_{w, b} \sum_{i=1}^n c_i (y_i - \theta^T x_i)^2$$

$\theta = [1, \dots, x_p]$
 $\theta = [b, w_1, \dots, w_p]^T$

$$= (Y - X\theta)^T C (Y - X\theta)$$

where $Y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$ $X = \begin{bmatrix} x_1^{(1)} & \dots & x_1^{(p+1)} \\ \vdots & & \vdots \\ x_n^{(1)} & \dots & x_n^{(p+1)} \end{bmatrix}$ $C = \begin{bmatrix} c_1 & & 0 \\ & c_2 & \\ 0 & & \ddots \\ & & & c_n \end{bmatrix}$

$$\theta = \begin{bmatrix} b \\ w_1 \\ \vdots \\ w_p \end{bmatrix}$$

$$J = (Y^T - \theta^T X^T) (CY - CX\theta)$$

$$J = Y^T CY - \underbrace{Y^T CX\theta}_{\text{same}} - \theta^T X^T CY + \theta^T X^T CX\theta$$

$$\nabla J = -2 X^T CY + 2 X^T CX\theta = 0$$

$$\therefore \theta = (X^T CX)^{-1} X^T CY$$

Q3a) For P of robust regression, the MM algorithm has

$$\rho_t(r) = \rho(r_{t,i}) - \underbrace{\frac{1}{2} r_{t,i} \psi(r_{t,i})}_{\substack{\text{as constant} \\ \text{wrt } \theta}} + \underbrace{\frac{1}{2} \frac{\psi(r_{t,i})}{r_{t,i}}}_{C_{t,i}} \underbrace{r^2}_{\substack{\text{as constant} \\ \text{wrt } \theta}} = \underbrace{C_{t,i}}_{\substack{\text{as constant} \\ \text{wrt } \theta}} (Y - X\theta)^T (Y - X\theta)$$

~~$\rho_t(r)$~~ = $C_{t,i}$

$$\frac{\partial \rho_t(r)}{\partial \theta} = \frac{\partial}{\partial \theta} (Y - X\theta)^T C_t (Y - X\theta) = 0$$

where $C_t = \begin{bmatrix} C_{t,1} & & 0 \\ & \ddots & \\ 0 & & C_{t,n} \end{bmatrix}$

$$\theta = (X^T C_t X)^{-1} X^T C_t Y$$

Hence the same form as weighted linear reg.
Also, the weights (ie $C_{t,i} \approx f(r_{t,i})$) are some function of the residuals. Therefore it is called iteratively reweighted least squares.

As, we can see above, the weight θ

$$w_i = \frac{\psi(r_{t,i})}{2 r_{t,i}}$$

& $\psi(r_{t,i})/r_{t,i}$ is non-increasing for $r > 0$

ie as $r \uparrow$ $w_i \downarrow$ or remains same.

ie More weights is given to inliers than outliers

whereas in case of OLS,

$$w_i = 1$$

is same weight for inliers as well as outliers.

& hence the algo achieves robustness.

Q4 a) $T_+(\theta) = T(\theta^{(t)}) + \nabla T(\theta^{(t)})^T (\theta - \theta^{(t)}) + \frac{1}{2} (\theta - \theta^{(t)})^T B (\theta - \theta^{(t)})$

Using Taylor series expansion,

$$T(\theta) = T(\theta^{(t)}) + \nabla T(\theta^{(t)})^T (\theta - \theta^{(t)}) + \frac{1}{2} (\theta - \theta^{(t)})^T \nabla^2 T(\tilde{\theta}) (\theta - \theta^{(t)})$$

~~where~~

$$\therefore T_+(\theta) - T(\theta) = \frac{1}{2} (\theta - \theta^{(t)})^T (B - \nabla^2 T(\tilde{\theta})) (\theta - \theta^{(t)})$$

$$\because B - \nabla^2 T(\tilde{\theta}) \succeq 0$$

$$\therefore x^T (B - \nabla^2 T(\tilde{\theta})) x \geq 0 \quad \forall x$$

$$\text{Let } x = \theta - \theta^{(t)}$$

Then,

$$(\theta - \theta^{(t)})^T (B - \nabla^2 T(\tilde{\theta})) (\theta - \theta^{(t)}) \geq 0$$

$$\text{or } T_+(\theta) - T(\theta) \geq 0$$

or $T_+(\theta)$ is majorizing function for T .

$$\tilde{J}_t(\theta) = J(\theta^{(t)}) + \nabla J(\theta^{(t)})^T \theta - \nabla J(\theta^{(t)})^T \theta^{(t)} + \frac{1}{2} (\theta^T B \theta - \theta^{(t)T} B \theta^{(t)} - \theta^{(t)T} B \theta + (\theta^{(t)})^T B \theta^{(t)})$$

$$\nabla \tilde{J}_t(\theta) = \nabla J(\theta^{(t)}) + \frac{1}{2} (2B\theta - 2B\theta^{(t)})$$

$$0 = \nabla J(\theta^{(t)}) + B(\theta - \theta^{(t)})$$

$$\therefore \theta = \theta^{(t)} - B^{-1}(\nabla J(\theta^{(t)}))$$

b) $J(\theta) = -\mathcal{L}(\theta) + \lambda \|\theta\|^2$

$$= \sum \log(1 + \exp(-y_i \theta^T \tilde{x}_i)) + \lambda \|\theta\|^2$$

$$J(\theta) = \frac{1}{n} \sum \log(1 + \exp(-y_i \theta^T \tilde{x}_i)) + \lambda \theta^T \theta$$

$$\frac{\partial J(\theta)}{\partial \theta_j} = \sum_{i=1}^n \frac{1}{1 + \exp(-y_i \theta^T \tilde{x}_i)} \times -y_i \tilde{x}_i + 2\lambda \theta$$

$$\frac{\partial J(\theta)}{\partial \theta} = \sum (1 - \sigma) (-y_i \tilde{x}_i) + 2\lambda \theta$$

$$\frac{\partial \theta}{\partial J(\theta)} = \sum \tilde{x}_i \left(\frac{e^{\theta^T \tilde{x}_i}}{1 + e^{\theta^T \tilde{x}_i}} - y_i \right) + 2\lambda \theta$$

$$\frac{\partial^2 J(\theta)}{\partial \theta^2} = \sum \tilde{x}_i \tilde{x}_i^T \left(\frac{e^{\theta^T \tilde{x}_i}}{(1 + e^{\theta^T \tilde{x}_i})^2} \right) + 2\lambda \mathbf{I}$$

$$\max \left(\frac{\partial^2 J(\theta)}{\partial \theta^2} \right) = \sum \tilde{x}_i \tilde{x}_i^T \underbrace{\{\sigma - (1 - \sigma)\}}_{\max = 1/4} + 2\lambda \mathbf{I}$$

$$\therefore \max \left(\frac{\nabla^2 J(\theta)}{\cancel{\theta}} \right) = \sum \tilde{x}_i \tilde{x}_i^T \times \frac{1}{4} + 2\lambda \mathbb{I} \quad \rightarrow (p \times (p+1))$$

$$\max \left(\frac{\nabla^2 J(\theta)}{\cancel{\theta}} \right) = \underbrace{\frac{\tilde{X}^T \tilde{X}}{4} + 2\lambda \mathbb{I}}$$

Let this be B , then $B - \nabla^2 J(\theta)$ will always be PSD.

$$\& \theta^{(t+1)} = \theta^{(t)} - \left(\frac{\tilde{X}^T \tilde{X}}{4} + 2\lambda \mathbb{I} \right)^{-1} \nabla J(\theta^{(t)})$$

$$\theta^{(t+1)} = \theta^{(t)} - \left(\frac{(\tilde{X}^T \tilde{X})^{-1}}{4} + 2\lambda \mathbb{I} \right) \nabla J(\theta^{(t)})$$

$$\text{where } \nabla J(\theta^{(t)}) = \sum \tilde{x}_i (1 - \sigma_i - y_i) + 2\lambda \theta$$

$$\sigma_i = \frac{1}{1 + \exp(-y_i \theta^T \tilde{x}_i)}$$

~~Then~~

The Computational advantage of mm algo compared to Newton's method:-

- 1) We have calculate just $\tilde{X}^T \tilde{X}$ instead of Hessian.
- 2) Inverse of a symmetric matrix is computationally less intensive than that of a non symmetric one. because ~~(AT)~~ its inverse will also be symmetric.