

Stat 600 HW3: Jivoung Park 531006793

Task 1. Prove $\beta_k^{(t+1)} = (\beta_k^{(t)} - \gamma (X^T W_k X + \lambda I)^{-1} [X^T (P_n - \mathbb{1}_{\{Y=n\}}) + \lambda \beta_k^{(t)}])$ is a damped Newton method for multi-class Ridge regression.

Sol. 1. Damped Newton's method is written as

$$\beta_n^{(t+1)} = \beta_n^{(t)} - \eta (\nabla_{xx}^2 f(\beta))^{-1} \nabla_n f(\beta)$$

In this case, $f(\beta) = (-\sum_i^n (\sum_k^{K-1} \mathbb{1}_{\{y_i=k\}} \log p_k(x_i)) + \frac{\lambda}{2} \sum_k^{K-1} \sum_j^p \beta_{k,j}^2$

$$\begin{aligned}
 &= -\sum_{i=1}^n \left(\sum_k^{k-1} \mathbb{1}_{\{y_i=k\}} \left(\log(e^{x_i^T \beta_k}) - \log\left(\sum_{l=0}^{k-1} e^{x_i^T \beta_l}\right) \right) + \frac{1}{2} \sum_k^{k-1} \sum_j^p (\beta_{k,j})^2 \right) \\
 &= -\sum_{i=1}^n \left(\sum_k^{k-1} \mathbb{1}_{\{y_i=k\}} \left(x_i^T \beta_k - \log\left(\sum_{l=0}^{k-1} e^{x_i^T \beta_l}\right) \underbrace{\sum_k^{k-1} \mathbb{1}_{\{y_i=k\}}}_{=1} \right) + \dots \right)
 \end{aligned}$$

$$= -\sum_{\mu}^n (\sum_k^{k-1} \mathbb{1}_{\{y_{\mu} = k\}} x_{\mu}^{\top} \beta_k - \log(\sum_{k=0}^{k-1} e^{x_{\mu}^{\top} \beta_k})) + \frac{\lambda}{2} \sum_k^{k-1} \sum_j^p \beta_{k,j}^2$$

① derivative

$$\frac{\partial f(\beta)}{\partial \beta_m} = - \underbrace{\sum_{i=1}^n \left(\frac{\partial}{\partial \beta_m} \sum_{k=0}^{K-1} \mathbb{1}_{\{y_i=k\}} x_i^{\top} \beta_k \right)}_{\parallel \leftarrow \frac{\partial \sigma^{\top} x}{\partial x} = a} \underbrace{- \frac{\partial}{\partial \beta_m} \log \left(\sum_{k=0}^{K-1} e^{x_i^{\top} \beta_k} \right)}_{\parallel \frac{e^{x_i^{\top} \beta_m}}{\sum_{k=0}^{K-1} e^{x_i^{\top} \beta_k}} = p_m(x_i)} + \underbrace{\frac{\lambda}{2} \frac{\partial}{\partial \beta_m} \sum_{k=0}^{K-1} \sum_j \beta_{k,j}^2}_{\parallel \frac{\partial}{\partial \beta_m} \beta_m^{\top} \beta_m = 2\beta_m} = \underbrace{\lambda \beta_m}_{\parallel}$$

$$\therefore \frac{\partial f(\beta)}{\partial \beta_k} = \underbrace{\sum_{i=1}^n x_i (p_k(x_i) - \mathbb{1}_{\{y_i=k\}})}_{= X^T (P_k - \mathbb{1}_{\{Y=k\}})} + \lambda \beta_k \quad \text{①}$$

②. hessium

$$\begin{aligned} \frac{\partial f(\beta)}{\partial \beta_k \beta_m} &= \frac{\partial}{\partial \beta_m^T} X^T (P_k - \mathbb{1}_{y=k}) + \lambda \frac{\partial}{\partial \beta_m^T} \beta_k \\ &= \sum_i^n \frac{\partial X_i (P_k(X_i) - \mathbb{1}_{y_i=k})}{\partial \beta_m^T} + \lambda \frac{\partial}{\partial \beta_m^T} \beta_k \end{aligned}$$

$$= \sum_{i=1}^n x_i \frac{\partial}{\partial \beta_m^+} p_n(x_i) + \lambda \frac{\partial}{\partial \beta_m^+} \beta_n$$

$$= \sum_i x_i \frac{\partial \beta_m}{\partial \beta_m^*} \cdot \frac{\partial f_k(x_i)}{\partial \beta_m} + \lambda \frac{\partial \beta_m}{\partial \beta_m^*} \cdot \frac{\partial \beta}{\partial \beta_m}$$

$$\left(\begin{array}{l} = \text{map}: \mathbb{R}^p \rightarrow (\mathbb{R}^p)^* \\ x_i \mapsto x_i^t \end{array} \right) (\because \text{dual map's derivative is transposed identity map})$$

$$\textcircled{2} \left(\begin{aligned} &= \frac{\mathbb{1}_{\{m=k\}} X_i \exp(X_i^T \beta_k) \cdot \sum_{\ell=0}^{k-1} \exp(X_i^T \beta_\ell) - \exp(X_i^T \beta_k) \cdot X_i \exp(X_i^T \beta_m)}{\left(\sum_{\ell=0}^{k-1} \exp(X_i^T \beta_\ell) \right)^2} \\ &= \mathbb{1}_{\{m=k\}} X_i p_k(X_i) - X_i p_k(X_i) p_m(X_i) = X_i p_k(X_i) (\mathbb{1}_{\{m=k\}} - p_m(X_i)) \end{aligned} \right.$$

plugging $\Theta = X_i p_k(X_i) (\mathbb{1}_{\{m=k\}} - p_m(X_i))$ into original eq.

$$\begin{aligned} \frac{\partial f}{\partial \beta_k \beta_k} &= \sum_i^n X_i \underbrace{\left(X_i p_k(X_i) (\mathbb{1}_{\{m=k\}} - p_m(X_i)) \right)^2}_{\text{scalar.}} + \lambda \underbrace{\left(\frac{\partial \beta_k}{\partial \beta_m} \right)^2}_{\substack{\text{take out} \\ = (\mathbb{1}_{\{m=k\}} \mathbb{I})^2 = \mathbb{1}_{\{m=k\}} \mathbb{I}}} \\ &= \sum_i^n p_k(X_i) (\mathbb{1}_{\{m=k\}} - p_m(X_i)) \cdot X_i X_i^T + \lambda \mathbb{1}_{\{m=k\}} \mathbb{I}. \end{aligned}$$

Now, note $\nabla_{k \times k}^2 f = \frac{\partial f}{\partial \beta_k \beta_k} = \sum_i^n p_k(X_i) (1 - p_k(X_i)) \cdot X_i X_i^T + \lambda \mathbb{I}.$

$$= \sum_i^n W_{k,i} X_i X_i^T + \lambda \mathbb{I} = X^T W_k X + \lambda \mathbb{I}$$

since $X^T W_k X = \begin{pmatrix} X_{11} & \dots & X_{n1} \\ \vdots & & \vdots \\ X_{1p} & \dots & X_{np} \end{pmatrix} \begin{pmatrix} p(X_{11})(1-p(X_{11})) \\ \vdots \\ p(X_{1p})(1-p(X_{1p})) \end{pmatrix} X$

$$= \begin{pmatrix} X_{11} & \dots & X_{1p} \\ \vdots & & \vdots \\ X_{n1} & \dots & X_{np} \end{pmatrix} \begin{pmatrix} W_{11}^k \\ \vdots \\ W_{1p}^k \end{pmatrix}$$

$$= \sum_k^n X_{ik} W_{ik}^k X_{kj} = \sum_i^n W_{ii}^k X_i X_i^T \quad \square \quad \textcircled{2}$$

$\therefore \begin{cases} \nabla_k f = X^T (P_k - \mathbb{1}_{\{Y=k\}}) + \lambda \beta_k \\ \nabla_{k \times k}^2 f = X^T W_k X + \lambda \mathbb{I}. \end{cases}$

$\therefore \beta_k^{(t+1)} = \beta_k^{(t)} - \gamma (\nabla_{k \times k}^2 f)^{-1} \nabla_k f \quad (\text{damped Newton})$
 $= \beta_k^{(t)} - \gamma (X^T W_k X + \lambda \mathbb{I})^{-1} \cdot (X^T (P_k - \mathbb{1}_{\{Y=k\}}) + \lambda \beta_k)$
 (target form).

\square (Task 1)