# Stat 600: Homework 3

For Multi-class logistic regression, we know that,

$$P(y_i = k \mid x_i) = P_k(x_i; \beta), \quad \sum_{k=0}^{K-1} P_k(x_i; \beta) = 1$$

$$P_k(x_i; \beta) = \frac{e^{x_i^T \beta_k}}{\sum_{l=0}^{K-1} e^{x_i^T \beta_l}}$$

$$f(\beta) = \left[ -\sum_{i=1}^{n} \left\{ \sum_{k=0}^{K-1} 1(y_i = k) \log P_k(x_i; \beta) \right\} + \frac{\lambda}{2} \sum_{k=0}^{K-1} \sum_{j=1}^{P} \beta_{kj}^2 \right] \leftarrow ①$$

Taking derivatives of different parts and then combining them to form a gradient,

$$\frac{\partial}{\partial \beta_k} \log P_k(x_i; \beta) = \frac{1}{P_k(x_i; \beta)} \cdot \frac{\partial P_k(x_i; \beta)}{\partial \beta_k} \leftarrow ②$$

$$\frac{\partial}{\partial \beta_k} P_k(x_i; \beta) = \frac{\partial}{\partial \beta_k} \left( \frac{e^{x_i^T \beta_k}}{\sum_{l=0}^{K-1} e^{x_i^T \beta_l}} \right)$$

$$\frac{\partial}{\partial \beta_k} P_k(x_i; \beta) = P_k(x_i; \beta)(1 - P_k(x_i; \beta)) x_i \leftarrow ③$$

When $\ell \neq k$,

$$\frac{\partial}{\partial \beta_\ell} \, p_k(x_i;\beta) = - p_k(x_i;\beta)\, p_\ell(x_i;\beta)\, x_i \quad \leftarrow \textcircled{4}$$

Therefore, from $\textcircled{2}$, $\textcircled{3}$, and $\textcircled{4}$,

$$\frac{\partial}{\partial \beta_k}\left(-1(y_i=k)\log p_k(x_i;\beta)\right)$$

$$= -\left(1(y_i=k) - p_k(x_i;\beta)\right)x_i \quad \leftarrow \textcircled{5}$$

$$\frac{\partial}{\partial \beta_k}\, \frac{\lambda}{2}\sum_{k=0}^{K-1}\sum_{j=1}^{P}\beta_{kj}^2 = \lambda\beta_k \quad \leftarrow \textcircled{6}$$

Therefore from $\textcircled{1}$, $\textcircled{5}$ and $\textcircled{6}$, we get the gradient term,

$$\nabla_{\beta_k}\, f(\beta) = \underbrace{-\sum_{i=1}^{n}\left(1(y_i=k) - p_k(x_i;\beta)\right)x_i}_{\text{likelihood term}} + \underbrace{\lambda\beta_k}_{\text{Regularization term}} \quad , \leftarrow \textcircled{7}$$

For $k = l$,

$$\frac{\partial p_k (x_i ; \beta)}{\partial \beta_k} = p_k (x_i ; \beta)(1 - p_k (x_i ; \beta)) x_i$$

$$\frac{\partial^2 p_k (x_i ; \beta)}{\partial \beta_k^2} = p_k (x_i ; \beta)(1 - p_k (x_i ; \beta)) x_i x_i^T \leftarrow \text{⑧}$$

For $k \neq l$,

$$\frac{\partial p_k (x_i ; \beta)}{\partial \beta_l} = - p_k (x_i ; \beta) p_l (x_i ; \beta) x_i$$

$$\frac{\partial^2 p_k (x_i ; \beta)}{\partial \beta_k \, \partial \beta_l} = - p_k (x_i ; \beta) p_l (x_i ; \beta) x_i x_i^T \leftarrow \text{⑨}$$

The full Hessian matrix is obtained by summing over all samples $i$:

$$H_{kl} = \sum_{i=1}^{n} H_{kl}^{(i)}$$

The second derivative of ⑤ gives us $\lambda I_p$ where $I_p$ is a $p \times p$ identity matrix.

Thus,

$$H_{kk} = \sum_{i=1}^{n} P_k(x_i;\beta)\left(1 - P_k(x_i;\beta)\right) x_i x_i^T + \lambda I_p$$

$$H_{k\ell} = -\sum_{i=1}^{n} P_k(x_i;\beta) P_\ell(x_i;\beta) x_i x_i^T \qquad (k \neq \ell)$$

Here, we approximate by assuming that the off-diagonal interactions between different classes are small and can be ignored.

Thus, we approximate the Hessian for class $k$ by:

$$H_k \approx \sum_{i=1}^{n} P_k(x_i;\beta)\left(1 - P_k(x_i;\beta)\right) x_i x_i^T + \lambda I_p$$

The matrix form of this approximation is written as:

$$H_k = X^T W_k X + \lambda I_p \qquad \textcircled{10}$$

The Damped Newton's update rule is given by,

$$\beta_k^{(t+1)} = \beta_k^{(t)} - \eta H^{-1} \nabla f(\beta_k) \leftarrow \text{①}$$

where

$\eta$ = learning rate (eta)
$\lambda$ = ridge parameter

From ⑦, ⑩ and ⑪, we get the desired equation,

$$\beta_k^{(t+1)} = \beta_k^{(t)} - \eta \left( X^T W_k X + \lambda I \right)^{-1} \left[ X^T \{ P_k - \mathbb{1}(y=k) \} + \lambda \beta_k^{(t)} \right]$$

$$\text{where } k = 0, 1, \ldots, K-1$$