

Derivations

Aarpan Ghosh

UIN: 236002102

Damped Newton's Method for multi-class
Logistic Regression:

In the question, the multiclass logistic regression problem with ridge regularization,

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

where $p_k(x_i) = \frac{e^{x_i^T \beta_k}}{\sum_{l=0}^{K-1} e^{x_i^T \beta_l}}$ and $\lambda > 0$

We want to find the damped Newton's update step with learning rate η ,

Note that,

$$f(\beta) = \left[- \sum_{i=1}^n \left\{ \sum_{l=0}^{K-1} \mathbb{1}(y_i = l) x_i^T \beta_l - \log \left(\sum_{l=0}^{K-1} e^{x_i^T \beta_l} \right) \right\} + \frac{\lambda}{2} \sum_{k=0}^{K-1} \sum_{j=1}^P \beta_{k,j}^2 \right]$$

Observe that, $- \mathbb{1}(y_i = l) x_i^T \beta_l$ is linear in β_l for each β_l and hence convex.

Also, taking $g(z) = \log \left(\sum_{k \neq 1} e^{x_k^T \beta} + e^z \right)$

$$g'(z) = \frac{e^z}{\sum_{k \neq 1} e^{x_k^T \beta} + e^z}$$

$$\Rightarrow g''(z) = \frac{e^z \sum_{k \neq 1} e^{x_k^T \beta}}{\left(\sum_{k \neq 1} e^{x_k^T \beta} + e^z \right)^2} \geq 0$$

$\Rightarrow g(z)$ is convex.

$\Rightarrow \log \left(\sum_{k=0}^{K-1} e^{x_k^T \beta} \right) = g(x_i^T \beta_k)$ is convex.

Further, $\sum_{j=1}^P \beta_{k,j}^2 = \|\beta_k\|^2$ is convex in β_k .

So, addition of two also gives convex

So, $f(\beta)$ is convex in each β_k .

To find $\hat{\beta} = \arg \min_{\beta} f(\beta)$, it is enough

to solve for $\frac{\partial f(\beta)}{\partial \beta_k} = 0 \quad \forall k$.

So, update step for damped Newton's method with learning rate η is,

$$\beta_{1k}^{(t+1)} = \beta_k^{(t)} - \eta \left[\frac{\partial^2 f(\beta)}{\partial \beta_k^2} \bigg|_{\beta_k = \beta_k^{(t)}} \right]^{-1} \left(\frac{\partial f(\beta)}{\partial \beta_k} \bigg|_{\beta = \beta_k^{(t)}} \right)$$

Gradient:

From the above expression $f(\beta)$,

$$\frac{\partial f(\beta)}{\partial \beta_k} = \left[- \sum_{i=1}^n \left\{ 1(y_i=k) x_i - \frac{1}{\sum_{l=0}^{K-1} e^{x_i^T \beta_l}} \cdot e^{x_i^T \beta_k} x_i \right\} + \lambda \beta_k \right]$$

$$= \left[- \sum_{i=1}^n x_i \{ 1(y_i=k) - p_k(x_i) \} + \lambda \beta_k \right]$$

$$= \left[\sum_{i=1}^n x^T e_i e_i^T \{ p_k - 1(y=k) \} + \lambda \beta_k \right]$$

$$= \left[x^T I_n \{ p_k - 1(y=k) \} + \lambda \beta_k \right]$$

$$= \left[x^T \{ p_k - 1(y=k) \} + \lambda \beta_k \right]$$

[Here, e_i is the i th canonical basis vector for \mathbb{R}^n and I_n is the identity matrix]

$$\Rightarrow \frac{\partial f(\beta)}{\partial \beta_k} \Big|_{\beta = \beta_k^{(k)}} = \left[x^T \{ p_k^T - 1(y=k) \} + \lambda \beta_k^{(k)} \right]$$

Hessian:

$$\text{Now, } \frac{\partial^2 f(\beta)}{\partial \beta_k^2} = \frac{\partial}{\partial \beta_k} \frac{\partial f(\beta)}{\partial \beta_k}$$

$$= \frac{\partial}{\partial \beta_k} \left[- \sum_{i=1}^n \left\{ 1(y_i=k) x_i - \frac{1}{\sum_{l=0}^{K-1} e^{x_i^T \beta_l}} \cdot e^{x_i^T \beta_k} x_i \right\} + \lambda \beta_k \right]$$

$$= \left[\sum_{i=1}^n x_i \left\{ \frac{\left(\sum_{l=0}^{k-1} e^{x_i^T \beta_l} \right) e^{x_i^T \beta_k} - e^{x_i^T \beta_k} \cdot e^{x_i^T \beta_k}}{\left(\sum_{l=0}^{k-1} e^{x_i^T \beta_l} \right)^2 + \lambda I_k} \right\} \right]_{k^T}$$

$$= \left[\sum_{i=1}^n x^T e_i \{ p_k(x_i) - p_k(x_i^2) \} e_i^T x + \lambda I_k \right]$$

$$= \left[x^T \left\{ \sum_{i=1}^n e_i p_k(x_i) (1 - p_k(x_i)^2) e_i^T \right\} x + \lambda I_k \right]$$

$$= [x^T W_k x + \lambda I_k]$$

[Since, it is given, $\sum_{i=1}^n e_i p_k(x_i) (1 - p_k(x_i)^2) e_i^T = W_k$,

$$\Rightarrow \frac{\partial^2 f(\beta)}{\partial \beta_k^2} \Big|_{\beta = \beta_k^{(k)}} = [x^T W_k^{(k)} x + \lambda I_k]$$

So, the final update eqn becomes,

$$\beta_k^{(k+1)} = \beta_k^{(k)} - \eta \left[\frac{x^T W_k^{(k)} x + \lambda I_k}{x^T \{ p_k^{(k)} - \mathbb{1}(\psi=k) \} + \lambda \beta_k^{(k)}} \right]$$

(proved).