

LINE PLUS AD PLATFORM: Additional Derivation

1. Mathematical derivation of cyclical coordinate descent algorithm of Logistic Regression with NO PENALTY (IRLS)

Let $\mathbf{y} = (y_1, \dots, y_n)'$, $\mathbf{X} = (\mathbf{1}, \mathbf{x}_1, \dots, \mathbf{x}_p)$, $\mathbf{x}_j = (x_{1j}, \dots, x_{nj})'$, $\boldsymbol{\beta} = (\beta_0, \dots, \beta_p)$, $\sigma(x) = [1 + \exp(-x)]^{-1}$, $p_i = \sigma(\boldsymbol{\beta}' \mathbf{x}_i)$. Then, the log likelihood of logistic regression is

$$L(\boldsymbol{\beta} | \mathbf{x}) = \prod_{i=1}^n p_i^{y_i} (1 - p_i)^{1-y_i} \quad (1)$$

Maximizing (1) is equivalent to minimize

$$\mathcal{L} = -\log L(\boldsymbol{\beta} | \mathbf{x}) \quad (2)$$

If \mathcal{L} in (2) is convex function, we could apply Newton's method to get numerical solution of (2). Let's verify this step by step

Deriving $\nabla_{\boldsymbol{\beta}} \mathcal{L}$

$$\begin{aligned} \bullet \quad \frac{\partial}{\partial \beta_j} \log p &= \frac{\partial}{\partial \beta_j} \left(\log \frac{1}{1 + \exp(-\boldsymbol{\beta}' \mathbf{x})} \right) = \frac{x_j \exp(-\boldsymbol{\beta}' \mathbf{x})}{1 + \exp(-\boldsymbol{\beta}' \mathbf{x})} = \frac{x_j}{1 + \exp(\boldsymbol{\beta}' \mathbf{x})} = x_j(1 - p) \\ \bullet \quad \frac{\partial}{\partial \beta_j} \log(1 - p) &= \frac{\partial}{\partial \beta_j} \left\{ -\boldsymbol{\beta}' \mathbf{x} - \log(1 + \exp(-\boldsymbol{\beta}' \mathbf{x})) \right\} = -x_j + x_j(1 - p) = -px_j \end{aligned}$$

$$\begin{aligned} \therefore \frac{\partial}{\partial \beta_j} \mathcal{L} &= - \sum_{i=1}^n \{ y_i x_{ij} (1 - p_i) + (1 - y_i) (-p_i x_{ij}) \} \\ &= - \sum_{i=1}^n \{ y_i x_{ij} - y_i x_{ij} p_i - p_i x_{ij} + y_i p_i x_{ij} \} \\ &= \sum_{i=1}^n x_{ij} (p_i - y_i) \end{aligned}$$

$$\therefore \nabla_{\boldsymbol{\beta}} \mathcal{L} = \mathbf{X}' (\mathbf{p} - \mathbf{y}) \quad (3)$$

Because gradient w.r.t. $\boldsymbol{\beta}$ is not linear in $\boldsymbol{\beta}$, we need to apply Newton's method.

Deriving $\nabla_{\boldsymbol{\beta}}^2 \mathcal{L}$ (Hessian Matrix)

If we show that the Hessian matrix is positive semi-definite, \mathcal{L} is convex function and we can apply Newton's method to get numerical solution in (3). Note that

$$\frac{\partial^2}{\partial \beta_j \partial \beta_k} \mathcal{L} = \sum_{i=1}^n x_{ij} \frac{\partial}{\partial \beta_k} p_i$$

Because $\partial \log p = \frac{1}{p} \partial p \iff \partial p = p \partial \log p = p x_j (1 - p)$, it is obvious that $\frac{\partial}{\partial \beta_k} p_i = x_{ik} p_i (1 - p_i)$. Therefore,

$$\begin{aligned} \sum_{i=1}^n x_{ij} \frac{\partial}{\partial \beta_k} p_i &= \sum_{i=1}^n x_{ij} x_{ik} p_i (1 - p_i) \\ &= \mathbf{x}_j' \mathbf{W} \mathbf{x}_k \end{aligned}$$

$$\text{where } \mathbf{W} = \begin{bmatrix} p_1(1-p_1) & & & \\ & \ddots & & \\ & & \ddots & \\ & & & p_n(1-p_n) \end{bmatrix}$$

$$\therefore \nabla_{\beta}^2 \mathcal{L} = \mathbf{X}' \mathbf{W} \mathbf{X} = \left(\mathbf{W}^{1/2} \mathbf{X} \right)' \left(\mathbf{W}^{1/2} \mathbf{X} \right) \quad (4)$$

The eigen values of $\nabla_{\beta}^2 \mathcal{L}$ are non-negative, so $\nabla_{\beta}^2 \mathcal{L}$ is p.s.d. matrix. Therefore, \mathcal{L} is convex function w.r.t. β .

Newton's method

The general iterative equation for getting numerical solution via Newton's method is

$$\beta_{t+1} = \beta_t - \mathbf{H}^{-1} \mathbf{g}$$

where \mathbf{H} is Hessian matrix and \mathbf{g} is gradient w.r.t. β . Using (3), (4) results, we get

$$\begin{aligned} \beta_{t+1} &= \beta_t - \mathbf{H}^{-1} \mathbf{g} \\ &= \beta_t - \left(\mathbf{X}' \mathbf{W} \mathbf{X} \right)^{-1} \mathbf{X}' (\mathbf{p} - \mathbf{y}) \\ &= \left(\mathbf{X}' \mathbf{W} \mathbf{X} \right)^{-1} \mathbf{X}' \mathbf{W} (\mathbf{X} \beta_t - \mathbf{W}^{-1} (\mathbf{p} - \mathbf{y})) \\ &= \left(\mathbf{X}' \mathbf{W} \mathbf{X} \right)^{-1} \mathbf{X}' \mathbf{W} \mathbf{z}_t \\ &\text{where } \mathbf{z}_t = \mathbf{X} \beta_t - \mathbf{W}^{-1} (\mathbf{p} - \mathbf{y}) \end{aligned} \quad (5)$$

From (5), we can see that β_{t+1} is solution of weighted least squares, where we minimize following quantity.

$$\operatorname{argmin}_{\beta} \sum_{i=1}^n b_i (z_i - \beta' \mathbf{x}_i)^2, \quad b_i = p_i(1-p_i) \quad (6)$$

In conclusion, to get mle, minimizing quantity (2) is equivalent to minimizing quantity (6). This algorithm is called Iteratively Reweighted Least Squares (IRLS) [1]

2. Adding Lasso Penalty to Logistic Regression.

Note that the objective function of logistic regression with lasso penalty will be

$$Q(\beta) = -\log L(\beta \mid \mathbf{x}) + \lambda \sum_{j=1}^p |\beta_j| \quad (7)$$

However, we show that minimizing $Q(\beta)$ is equal to minimizing

$$Q(\beta)^N = \sum_{i=1}^n b_i (z_i - \beta' \mathbf{x}_i)^2 + \lambda \sum_{j=1}^p |\beta_j| \quad (8)$$

[3], [5] show that for weighted updates, the coordinate descent algorithm updates $\tilde{\beta}_j$ as

$$\tilde{\beta}_j \leftarrow \frac{S\left(\frac{1}{n} \sum_{i=1}^n w_i x_{ij} (y_i - \tilde{y}_i^{(j)}), \lambda \alpha\right)}{\frac{1}{n} \sum_{i=1}^n w_i x_{ij}^2 + \lambda(1 - \alpha)} \quad (9)$$

When $\alpha = 1$, it is usual lasso penalty. With $0 < \alpha < 1$, this is elastic penalty, which combines ℓ_1, ℓ_2 norm.

Reference

- [1] The elements of statistical learning: data mining, inference, and prediction, T.Hastie et al., 2009.
- [2] Statistical learning with sparsity: the lasso and generalizations, T.hastie et al., 2015
- [3] Regularization Paths for Generalized Linear Models via Coordinate Descent, J.Friedman et al., 2009
- [4] Coordinate Descent Algorithms For Lasso Penalized Regression, Wu et al., 2008
- [5] A coordinate majorization descent algorithm for L1 penalized learning, Yang et al., 2012
- [6] Package 'glmnet', J.Friedman et al., 2020