

1 Smooth Loss and Smooth Regularizer

1.1 Objective Function

Given smooth loss function, $l(\beta_0 + \mathbf{x}_j^T \boldsymbol{\beta}; y_j)$, and smooth regularization, $R(\beta_0, \boldsymbol{\beta})$, we can get the objective function under the full data like

$$Q^*(\beta_0, \boldsymbol{\beta}) = \sum_{j=1}^n l(\beta_0 + \mathbf{x}_j^T \boldsymbol{\beta}; y_j) + R(\beta_0, \boldsymbol{\beta}) \quad (1)$$

The leave- i -out objection function as

$$Q(\beta_0, \boldsymbol{\beta}) = \sum_{j \neq i} l(\beta_0 + \mathbf{x}_j^T \boldsymbol{\beta}; y_j) + R(\beta_0, \boldsymbol{\beta}) \quad (2)$$

1.2 Approximate Leave- i -Out Prediction

Assuming $\hat{\beta}_0, \hat{\boldsymbol{\beta}}$ are the estimated parameters under the full dataset.

Based on Newton's method, we have

$$\begin{bmatrix} \tilde{\beta}_0^{/i} \\ \tilde{\boldsymbol{\beta}}^{/i} \end{bmatrix} = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\boldsymbol{\beta}} \end{bmatrix} - \left[\nabla^2 Q(\hat{\beta}_0, \hat{\boldsymbol{\beta}}) \right]^{-1} \nabla Q(\hat{\beta}_0, \hat{\boldsymbol{\beta}}) \quad (3)$$

where,

$$\begin{cases} \nabla^2 Q(\hat{\beta}_0, \hat{\boldsymbol{\beta}}) = \sum_{j \neq i} \begin{bmatrix} 1 \\ \mathbf{x}_j \end{bmatrix} \begin{bmatrix} 1 & \mathbf{x}_j^T \end{bmatrix} \ddot{l}(\hat{\beta}_0 + \mathbf{x}_j^T \hat{\boldsymbol{\beta}}; y_j) + \nabla^2 R(\hat{\beta}_0, \hat{\boldsymbol{\beta}}) \\ \nabla Q(\hat{\beta}_0, \hat{\boldsymbol{\beta}}) = - \begin{bmatrix} 1 \\ \mathbf{x}_i \end{bmatrix} \dot{l}(\hat{\beta}_0 + \mathbf{x}_i^T \hat{\boldsymbol{\beta}}; y_i) \end{cases}$$

So, now we have

$$\begin{bmatrix} 1 & \mathbf{x}_i^T \end{bmatrix} \begin{bmatrix} \tilde{\beta}_0^{/i} \\ \tilde{\boldsymbol{\beta}}^{/i} \end{bmatrix} = \begin{bmatrix} 1 & \mathbf{x}_i^T \end{bmatrix} \begin{bmatrix} \hat{\beta}_0 \\ \hat{\boldsymbol{\beta}} \end{bmatrix} + \begin{bmatrix} 1 & \mathbf{x}_i^T \end{bmatrix} \left\{ \sum_{j \neq i} \begin{bmatrix} 1 \\ \mathbf{x}_j \end{bmatrix} \begin{bmatrix} 1 & \mathbf{x}_j^T \end{bmatrix} \ddot{l}(\hat{\beta}_0 + \mathbf{x}_j^T \hat{\boldsymbol{\beta}}; y_j) + \nabla^2 R(\hat{\beta}_0, \hat{\boldsymbol{\beta}}) \right\}^{-1} \begin{bmatrix} 1 \\ \mathbf{x}_i \end{bmatrix} \dot{l}(\hat{\beta}_0 + \mathbf{x}_i^T \hat{\boldsymbol{\beta}}; y_i)$$

Using the matrix inversion lemma, we get

$$\begin{bmatrix} 1 & \mathbf{x}_i^T \end{bmatrix} \begin{bmatrix} \tilde{\beta}_0^{/i} \\ \tilde{\boldsymbol{\beta}}^{/i} \end{bmatrix} = \hat{\beta}_0 + \mathbf{x}_i^T \hat{\boldsymbol{\beta}} + \frac{H_{ii}}{1 - H_{ii} \ddot{l}(\hat{\beta}_0 + \mathbf{x}_i^T \hat{\boldsymbol{\beta}}; y_i)} \dot{l}(\hat{\beta}_0 + \mathbf{x}_i^T \hat{\boldsymbol{\beta}}; y_i) \quad (4)$$

where,

$$\begin{aligned} \mathbf{H} &= [\mathbf{1}_n, \mathbf{X}] \left\{ [\mathbf{1}_n, \mathbf{X}]^T \mathbf{D} [\mathbf{1}_n, \mathbf{X}] + \nabla^2 R(\hat{\beta}_0, \hat{\boldsymbol{\beta}}) \right\}^{-1} [\mathbf{1}_n, \mathbf{X}]^T \\ \mathbf{D} &= \text{diag} \left(\ddot{l}(\hat{\beta}_0 + \mathbf{x}_1^T \hat{\boldsymbol{\beta}}; y_1), \dots, \ddot{l}(\hat{\beta}_0 + \mathbf{x}_n^T \hat{\boldsymbol{\beta}}; y_n) \right) \end{aligned}$$

2 Elastic Net

2.1 Objective Function

Given the loss function,

$$l(\beta_0 + \mathbf{x}_j^T \boldsymbol{\beta}; y_j) = \frac{1}{2} (y_j - \beta_0 - \mathbf{x}_j^T \boldsymbol{\beta})^2$$

, and the regularization $R(\beta_0, \boldsymbol{\beta}) = \lambda_1 \|\boldsymbol{\beta}\|_1 + \lambda_2 \|\boldsymbol{\beta}\|_2^2$, we can get the objective function under the full data like

$$Q^*(\beta_0, \boldsymbol{\beta}) = \sum_{j=1}^n l(\beta_0 + \mathbf{x}_j^T \boldsymbol{\beta}; y_j) + R(\beta_0, \boldsymbol{\beta})$$

The leave- i -out objection function as

$$\begin{aligned}
Q(\beta_0, \boldsymbol{\beta}) &= \sum_{j \neq i} l(\beta_0 + \mathbf{x}_j^T \boldsymbol{\beta}; y_j) + R(\beta_0, \boldsymbol{\beta}) \\
&= \frac{1}{2} \sum_{j \neq i} (y_j - \beta_0 - \mathbf{x}_j^T \boldsymbol{\beta})^2 + R(\beta_0, \boldsymbol{\beta}) \\
&= \frac{1}{2} \sum_{j=1}^n (y_j - \beta_0 - \mathbf{x}_j^T \boldsymbol{\beta})^2 + R(\beta_0, \boldsymbol{\beta}) - \frac{1}{2} (y_i - \beta_0 - \mathbf{x}_i^T \boldsymbol{\beta})^2
\end{aligned}$$

2.2 Primal Domain

Under the primal domain, we can notice that the loss function is smooth and the regularizer is not smooth, which has one zero-order singularity $K = \{0\}$ for all of the parameters except β_0 .

Assuming $\hat{\beta}_0, \hat{\boldsymbol{\beta}}$ are the estimated parameters under the full dataset.

Define active set $A = \{i : \beta_i \notin K, i = 1, \dots, p\}$. Then based on the formula (31), we can compute \mathbf{H} as

$$\mathbf{H} = [\mathbf{1}_n, \mathbf{X}_{\cdot A}] \left\{ [\mathbf{1}_n, \mathbf{X}_{\cdot A}]^T \mathbf{D} [\mathbf{1}_n, \mathbf{X}_{\cdot A}] + \nabla^2 R(\hat{\beta}_0, \hat{\boldsymbol{\beta}}_A) \right\}^{-1} [\mathbf{1}_n, \mathbf{X}_{\cdot A}]^T$$

where,

$$\begin{aligned}
\mathbf{D} &= \text{diag} \left(\ddot{l}(\hat{\beta}_0 + \mathbf{x}_1^T \hat{\boldsymbol{\beta}}; y_1), \dots, \ddot{l}(\hat{\beta}_0 + \mathbf{x}_n^T \hat{\boldsymbol{\beta}}; y_n) \right) = \mathbf{I} \\
\nabla^2 R(\hat{\beta}_0, \hat{\boldsymbol{\beta}}_A) &= \begin{bmatrix} 0 & & & \\ & 2\lambda_2 & & \\ & & \ddots & \\ & & & 2\lambda_2 \end{bmatrix}
\end{aligned}$$

At last, we have

$$\begin{bmatrix} 1 & \mathbf{x}_i^T \end{bmatrix} \begin{bmatrix} \tilde{\beta}_0^{/i} \\ \tilde{\boldsymbol{\beta}}^{/i} \end{bmatrix} = \hat{\beta}_0 + \mathbf{x}_i^T \hat{\boldsymbol{\beta}} + \frac{H_{ii}}{1 - H_{ii} \ddot{l}(\hat{\beta}_0 + \mathbf{x}_i^T \hat{\boldsymbol{\beta}}; y_i)} \ddot{l}(\hat{\beta}_0 + \mathbf{x}_i^T \hat{\boldsymbol{\beta}}; y_i)$$

2.3 Proximal Operator

Similar as formula (33), we can define the proximal operator as

$$\begin{bmatrix} \hat{\beta}_0^{/i} \\ \hat{\boldsymbol{\beta}}^{/i} \end{bmatrix} = \mathbf{prox}_{R(\cdot)} \left(\begin{bmatrix} \hat{\beta}_0^{/i} \\ \hat{\boldsymbol{\beta}}^{/i} \end{bmatrix} - \sum_{j \neq i} \begin{bmatrix} 1 \\ \mathbf{x}_j \end{bmatrix} \ddot{l}(\hat{\beta}_0^{/i} + \mathbf{x}_j^T \hat{\boldsymbol{\beta}}^{/i}; y_j) \right)$$

where, based on the first order condition,

$$\sum_{j \neq i} \ddot{l}(\hat{\beta}_0^{/i} + \mathbf{x}_j^T \hat{\boldsymbol{\beta}}^{/i}; y_j) = 0$$

Define

$$\mathbf{u} = \begin{bmatrix} \hat{\beta}_0^{/i} \\ \hat{\boldsymbol{\beta}}^{/i} \end{bmatrix} - \sum_{j \neq i} \begin{bmatrix} 1 \\ \mathbf{x}_j \end{bmatrix} \ddot{l}(\hat{\beta}_0^{/i} + \mathbf{x}_j^T \hat{\boldsymbol{\beta}}^{/i}; y_j)$$

Hence, for all of the point in the active set, $A = \{i : \beta_i \notin K, i = 1, \dots, p\}$, the Jacobian of proximal operator equals to

$$\mathbf{J}_{E,E} = [\mathbf{J}(\mathbf{u})]_{E,E} = \begin{bmatrix} 1 & & & & \\ & 1 + 2\lambda_2 & & & \\ & & 1 + 2\lambda_2 & & \\ & & & \ddots & \\ & & & & 1 + 2\lambda_2 \end{bmatrix}^{-1}, \quad E = 1 \cup \{i+1 : i \in A\}$$

At last, using formula (45) and (45), we have

$$\begin{bmatrix} 1 & \mathbf{x}_i^T \end{bmatrix} \begin{bmatrix} \tilde{\beta}_0^{/i} \\ \tilde{\boldsymbol{\beta}}^{/i} \end{bmatrix} = \hat{\beta}_0 + \mathbf{x}_i^T \hat{\boldsymbol{\beta}} + \frac{H_{ii}}{1 - H_{ii}} i \left(\hat{\beta}_0 + \mathbf{x}_i^T \hat{\boldsymbol{\beta}}; y_i \right)$$

where, defining $\mathbf{X}^* = \begin{bmatrix} \mathbf{1}_n & \mathbf{X} \end{bmatrix}$,

$$\begin{aligned} \mathbf{H} &= \mathbf{X}_{:,E}^* \left[\mathbf{J}_{E,E} \mathbf{X}_{:,E}^{*T} \mathbf{D} \mathbf{X}_{:,E}^* + \mathbf{I}_{E,E} - \mathbf{J}_{E,E} \right]^{-1} \mathbf{J}_{E,E} \mathbf{X}_{:,E}^{*T} \\ \mathbf{D} &= \text{diag} \left(i \left(\hat{\beta}_0 + \mathbf{x}_1^T \hat{\boldsymbol{\beta}}; y_1 \right), \dots, i \left(\hat{\beta}_0 + \mathbf{x}_n^T \hat{\boldsymbol{\beta}}; y_n \right) \right) \end{aligned}$$

3 GLMNET Optimization

Considering the Elastic Net optimization problem,

$$\min_{\beta_0, \boldsymbol{\beta}} \frac{1}{2n} \sum_{j=1}^n (y_j - \beta_0 - \mathbf{x}_j^T \boldsymbol{\beta})^2 + \frac{1}{2} \lambda (1 - \alpha) \|\boldsymbol{\beta}\|_2^2 + \lambda \alpha \|\boldsymbol{\beta}\|_1$$

However, in the documentation of 'glmnet' package, it notes that for "gaussian", glmnet standardizes y to have unit variance (using 1/n rather than 1/(n-1) formula) before computing its lambda sequence (and then unstandardizes the resulting coefficients); if you wish to reproduce/compare results with other software, best to supply a standardized y. The coefficients for any predictor variables with zero variance are set to zero for all values of lambda.

It means that the function is actually optimizing the following problem.

$$\min_{\beta_0, \boldsymbol{\beta}} \frac{1}{2n} \sum_{j=1}^n (y_j^* - \beta_0^* - \mathbf{x}_j^T \boldsymbol{\beta}^*)^2 + \frac{1}{2} \lambda (1 - \alpha) \|\boldsymbol{\beta}^*\|_2^2 + \lambda \alpha \|\boldsymbol{\beta}^*\|_1$$

where,

$$\begin{cases} y_j^* = \frac{y_j}{\text{sd}(\mathbf{y})} \\ \beta_0^* = \frac{\beta_0}{\text{sd}(\mathbf{y})} \\ \boldsymbol{\beta}^* = \frac{\boldsymbol{\beta}}{\text{sd}(\mathbf{y})} \end{cases}$$

So, the 'glmnet' function is actually optimizing

$$\begin{aligned} & \min_{\beta_0, \boldsymbol{\beta}} \frac{1}{2n} \sum_{j=1}^n (y_j^* - \beta_0^* - \mathbf{x}_j^T \boldsymbol{\beta}^*)^2 + \frac{1}{2} \lambda (1 - \alpha) \|\boldsymbol{\beta}^*\|_2^2 + \lambda \alpha \|\boldsymbol{\beta}^*\|_1 \\ &= \min_{\beta_0, \boldsymbol{\beta}} \frac{1}{2n} \sum_{j=1}^n \left(\frac{y_j}{\text{sd}(\mathbf{y})} - \frac{\beta_0}{\text{sd}(\mathbf{y})} - \frac{\mathbf{x}_j^T \boldsymbol{\beta}}{\text{sd}(\mathbf{y})} \right)^2 + \frac{1}{2} \frac{\lambda}{\text{sd}(\mathbf{y})^2} (1 - \alpha) \|\boldsymbol{\beta}\|_2^2 + \frac{\lambda}{\text{sd}(\mathbf{y})} \alpha \|\boldsymbol{\beta}\|_1 \end{aligned}$$

So, in order to coincide with the target optimization problem, we can rescale \mathbf{y} and \mathbf{X} by $\text{sd}(\mathbf{y})$, λ by $\text{sd}(\mathbf{y})^2$, then we can get

$$\begin{aligned} & \min_{\beta_0, \boldsymbol{\beta}} \frac{1}{2n} \sum_{j=1}^n (y_j - \beta_0 - \mathbf{x}_j^T \boldsymbol{\beta})^2 + \frac{1}{2} \lambda (1 - \alpha) \|\boldsymbol{\beta}\|_2^2 + \lambda \alpha \|\boldsymbol{\beta}\|_1 \\ &= \min_{\beta_0, \boldsymbol{\beta}} \frac{1}{2n} \sum_{j=1}^n \left(\frac{y_j}{\text{sd}(\mathbf{y})} - \beta_0^* - \frac{\mathbf{x}_j^T \boldsymbol{\beta}}{\text{sd}(\mathbf{y})} \right)^2 + \frac{1}{2} \frac{\lambda}{\text{sd}(\mathbf{y})^2} (1 - \alpha) \|\boldsymbol{\beta}\|_2^2 + \frac{\lambda}{\text{sd}(\mathbf{y})^2} \alpha \|\boldsymbol{\beta}\|_1 \\ &= \min_{\beta_0, \boldsymbol{\beta}} \frac{1}{2n} \sum_{j=1}^n (y_j^* - \beta_0^* - \mathbf{x}_j^T \boldsymbol{\beta})^2 + \frac{1}{2} \lambda^* (1 - \alpha) \|\boldsymbol{\beta}\|_2^2 + \lambda^* \alpha \|\boldsymbol{\beta}\|_1 \end{aligned}$$

where,

$$\begin{cases} \beta_0^* = \frac{\beta_0}{\text{sd}(\mathbf{y})} \\ \lambda^* = \frac{\lambda}{\text{sd}(\mathbf{y})^2} \end{cases}$$

4 Multinomial

4.1 Loss function

Assume we have K classes, n observations \mathbf{X} , and \mathbf{R}^p parameters β_k for each class.

Define leave- i -out variables as

$$\mathbf{y}_{(n-1)K \times 1}^{/i} = \{y_{jk}\}_{j \neq i}^{k=1, \dots, K} = \begin{bmatrix} y_{11} \\ y_{12} \\ \vdots \\ y_{1K} \\ y_{21} \\ y_{22} \\ \vdots \\ y_{2K} \\ \vdots \\ y_{n1} \\ y_{n2} \\ \vdots \\ y_{nK} \end{bmatrix}, \quad \mathbf{X}_{(n-1)K \times pK}^{/i} = \begin{bmatrix} \mathbf{x}_1^T & 0 & \cdots & 0 \\ 0 & \mathbf{x}_1^T & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \mathbf{x}_1^T \\ \mathbf{x}_2^T & 0 & \cdots & 0 \\ 0 & \mathbf{x}_2^T & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \mathbf{x}_2^T \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{x}_n^T & 0 & \cdots & 0 \\ 0 & \mathbf{x}_n^T & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \mathbf{x}_n^T \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_K \end{bmatrix}$$

The loss function would be

$$\begin{aligned} l(\mathbf{B}) &= - \left\{ \sum_{j \neq i} \left[\sum_{k=1}^K y_{jk} \mathbf{x}_j^T \beta_k - \log \left(\sum_{k=1}^K e^{\mathbf{x}_j^T \beta_k} \right) \right] \right\} \\ &= \sum_{j \neq i} \log \left(\sum_{k=1}^K e^{\mathbf{x}_j^T \beta_k} \right) - \sum_{j \neq i} \sum_{k=1}^K y_{jk} \mathbf{x}_j^T \beta_k \end{aligned}$$

By taking the first order derivative, we get

$$\begin{aligned} \frac{\partial l(\mathbf{B})}{\partial \mathbf{B}} &= \begin{bmatrix} \frac{\partial l(\mathbf{B})}{\partial \beta_1} \\ \vdots \\ \frac{\partial l(\mathbf{B})}{\partial \beta_K} \end{bmatrix} = \begin{bmatrix} \sum_{j \neq i} \frac{\exp(\mathbf{x}_j^T \beta_1)}{\sum_{k=1}^K \exp(\mathbf{x}_j^T \beta_k)} \mathbf{x}_j - \sum_{j \neq i} y_{j1} \mathbf{x}_j \\ \vdots \\ \sum_{j \neq i} \frac{\exp(\mathbf{x}_j^T \beta_K)}{\sum_{k=1}^K \exp(\mathbf{x}_j^T \beta_k)} \mathbf{x}_j - \sum_{j \neq i} y_{jK} \mathbf{x}_j \end{bmatrix} = \mathbf{X}^{/iT} \begin{bmatrix} \frac{\exp(\mathbf{x}_1^T \beta_1)}{\sum_{k=1}^K \exp(\mathbf{x}_1^T \beta_k)} \\ \vdots \\ \frac{\exp(\mathbf{x}_1^T \beta_K)}{\sum_{k=1}^K \exp(\mathbf{x}_1^T \beta_k)} \\ \frac{\exp(\mathbf{x}_2^T \beta_1)}{\sum_{k=1}^K \exp(\mathbf{x}_2^T \beta_k)} \\ \vdots \\ \frac{\exp(\mathbf{x}_2^T \beta_K)}{\sum_{k=1}^K \exp(\mathbf{x}_2^T \beta_k)} \\ \vdots \\ \frac{\exp(\mathbf{x}_n^T \beta_1)}{\sum_{k=1}^K \exp(\mathbf{x}_n^T \beta_k)} \\ \vdots \\ \frac{\exp(\mathbf{x}_n^T \beta_K)}{\sum_{k=1}^K \exp(\mathbf{x}_n^T \beta_k)} \end{bmatrix}_{(n-1)K \times 1} - \mathbf{X}^{/iT} \mathbf{y}^{/i} \\ &= \mathbf{X}^{/iT} \left[\mathbf{A}^{/i}(\beta) - \mathbf{y}^{/i} \right] = \mathbf{X}^{/iT} \left(\begin{bmatrix} \mathbf{A}_1(\beta) \\ \mathbf{A}_2(\beta) \\ \vdots \\ \mathbf{A}_n(\beta) \end{bmatrix} - \mathbf{y}^{/i} \right) \end{aligned}$$

Similarly, we can get

$$\frac{\partial^2 l(\mathcal{B})}{\partial \mathcal{B} \partial \mathcal{B}^T} = \mathcal{X}^{/iT} \frac{\partial \mathcal{A}^{/i}(\mathcal{B})}{\partial \mathcal{B}^T} = \mathcal{X}^{/iT} \begin{bmatrix} \frac{\partial \mathbf{A}_1(\mathcal{B})}{\partial \mathcal{B}^T} \\ \frac{\partial \mathbf{A}_2(\mathcal{B})}{\partial \mathcal{B}^T} \\ \vdots \\ \frac{\partial \mathbf{A}_n(\mathcal{B})}{\partial \mathcal{B}^T} \end{bmatrix}$$

where,

$$\begin{aligned} \frac{\partial \mathbf{A}_j(\mathcal{B})}{\partial \mathcal{B}^T} &= \begin{bmatrix} \frac{\exp(\mathbf{x}_j^T \beta_1)}{\sum_{k=1}^K \exp(\mathbf{x}_j^T \beta_k)} & 0 & \cdots & 0 \\ 0 & \frac{\exp(\mathbf{x}_j^T \beta_2)}{\sum_{k=1}^K \exp(\mathbf{x}_j^T \beta_k)} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \frac{\exp(\mathbf{x}_j^T \beta_K)}{\sum_{k=1}^K \exp(\mathbf{x}_j^T \beta_k)} \end{bmatrix}_{K \times K} \begin{bmatrix} \mathbf{x}_j^T & 0 & \cdots & 0 \\ 0 & \mathbf{x}_j^T & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \mathbf{x}_j^T \end{bmatrix}_{K \times pK} \\ &- \begin{bmatrix} \frac{\exp(\mathbf{x}_j^T \beta_1)}{\sum_{k=1}^K \exp(\mathbf{x}_j^T \beta_k)} \\ \frac{\exp(\mathbf{x}_j^T \beta_2)}{\sum_{k=1}^K \exp(\mathbf{x}_j^T \beta_k)} \\ \vdots \\ \frac{\exp(\mathbf{x}_j^T \beta_K)}{\sum_{k=1}^K \exp(\mathbf{x}_j^T \beta_k)} \end{bmatrix}_{K \times 1} \begin{bmatrix} \frac{\exp(\mathbf{x}_j^T \beta_1)}{\sum_{k=1}^K \exp(\mathbf{x}_j^T \beta_k)} & \frac{\exp(\mathbf{x}_j^T \beta_2)}{\sum_{k=1}^K \exp(\mathbf{x}_j^T \beta_k)} & \cdots & \frac{\exp(\mathbf{x}_j^T \beta_K)}{\sum_{k=1}^K \exp(\mathbf{x}_j^T \beta_k)} \end{bmatrix}_{1 \times K} \begin{bmatrix} \mathbf{x}_j^T & 0 & \cdots & 0 \\ 0 & \mathbf{x}_j^T & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \mathbf{x}_j^T \end{bmatrix}_{K \times pK} \\ &= [\text{diag}(\mathbf{A}_j(\mathcal{B})) - \mathbf{A}_j(\mathcal{B}) \mathbf{A}_j(\mathcal{B})^T]_{K \times K} \begin{bmatrix} \mathbf{x}_j^T & 0 & \cdots & 0 \\ 0 & \mathbf{x}_j^T & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \mathbf{x}_j^T \end{bmatrix}_{K \times pK} \end{aligned}$$

So, we have

$$\begin{aligned} \frac{\partial^2 l(\mathcal{B})}{\partial \mathcal{B} \partial \mathcal{B}^T} &= \mathcal{X}^{/iT} \frac{\partial \mathcal{A}^{/i}(\mathcal{B})}{\partial \mathcal{B}^T} = \mathcal{X}^{/iT} \begin{bmatrix} \frac{\partial \mathbf{A}_1(\mathcal{B})}{\partial \mathcal{B}^T} \\ \frac{\partial \mathbf{A}_2(\mathcal{B})}{\partial \mathcal{B}^T} \\ \vdots \\ \frac{\partial \mathbf{A}_n(\mathcal{B})}{\partial \mathcal{B}^T} \end{bmatrix} \\ &= \mathcal{X}^{/iT} \begin{bmatrix} [\text{diag}(\mathbf{A}_1(\mathcal{B})) - \mathbf{A}_1(\mathcal{B}) \mathbf{A}_1(\mathcal{B})^T]_{K \times K} \begin{bmatrix} \mathbf{x}_1^T & 0 & \cdots & 0 \\ 0 & \mathbf{x}_1^T & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \mathbf{x}_1^T \end{bmatrix}_{K \times pK} \\ [\text{diag}(\mathbf{A}_2(\mathcal{B})) - \mathbf{A}_2(\mathcal{B}) \mathbf{A}_2(\mathcal{B})^T]_{K \times K} \begin{bmatrix} \mathbf{x}_2^T & 0 & \cdots & 0 \\ 0 & \mathbf{x}_2^T & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \mathbf{x}_2^T \end{bmatrix}_{K \times pK} \\ \vdots \\ [\text{diag}(\mathbf{A}_n(\mathcal{B})) - \mathbf{A}_n(\mathcal{B}) \mathbf{A}_n(\mathcal{B})^T]_{K \times K} \begin{bmatrix} \mathbf{x}_n^T & 0 & \cdots & 0 \\ 0 & \mathbf{x}_n^T & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \mathbf{x}_n^T \end{bmatrix}_{K \times pK} \end{bmatrix}_{(n-1)K \times pK} \\ &= \mathcal{X}^{/iT} \mathcal{D}_{(n-1)K \times (n-1)K}^{/i}(\mathcal{B}) \mathcal{X}^{/i} \end{aligned}$$

where,

$$\mathcal{D}^{/i}(\mathcal{B}) = \begin{bmatrix} [\text{diag}(\mathbf{A}_1(\mathcal{B})) - \mathbf{A}_1(\mathcal{B})\mathbf{A}_1(\mathcal{B})^T] & 0 & \cdots & 0 \\ 0 & [\text{diag}(\mathbf{A}_2(\mathcal{B})) - \mathbf{A}_2(\mathcal{B})\mathbf{A}_2(\mathcal{B})^T] & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & [\text{diag}(\mathbf{A}_n(\mathcal{B})) - \mathbf{A}_n(\mathcal{B})\mathbf{A}_n(\mathcal{B})^T] \end{bmatrix}$$

4.2 Newton's Method

With Newton's method, we have the one step update as

$$\begin{aligned} \tilde{\mathcal{B}}^{/i} &= \hat{\mathcal{B}} - \left[\mathcal{X}^{/iT} \mathcal{D}^{/i}(\mathcal{B}) \mathcal{X}^{/i} + \nabla^2 R(\mathcal{B}) \right]^{-1} \left[\mathcal{X}^{/iT} (\mathcal{A}^{/iT}(\mathcal{B}) - \mathbf{y}^{/i}) + \nabla R(\mathcal{B}) \right] \\ &= \hat{\mathcal{B}} + \left[\mathcal{X}^T \mathcal{D}(\mathcal{B}) \mathcal{X} + \nabla^2 R(\mathcal{B}) - \mathbf{X}_i^T [\text{diag}(\mathbf{A}_i(\mathcal{B})) - \mathbf{A}_i(\mathcal{B})\mathbf{A}_i(\mathcal{B})^T] \mathbf{X}_i \right]^{-1} \mathbf{X}_i^T (\mathbf{A}_i(\mathcal{B}) - \mathbf{y}_i) \end{aligned}$$

where,

$$\mathbf{X}_i = \begin{bmatrix} \mathbf{x}_i^T & 0 & \cdots & 0 \\ 0 & \mathbf{x}_i^T & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \mathbf{x}_i^T \end{bmatrix}_{K \times pK}, \quad \mathbf{y}_i = \begin{bmatrix} y_{i1} \\ y_{i2} \\ \vdots \\ y_{iK} \end{bmatrix}$$

Defining $\mathcal{K}(\mathcal{B}) = \mathcal{X}^T \mathcal{D}(\mathcal{B}) \mathcal{X} + \nabla^2 R(\mathcal{B})$, with matrix inversion lemma, we can get

$$\begin{aligned} \tilde{\mathcal{B}}^{/i} &= \hat{\mathcal{B}} + [\mathcal{K}(\mathcal{B}) - \mathbf{X}_i^T [\text{diag}(\mathbf{A}_i(\mathcal{B})) - \mathbf{A}_i(\mathcal{B})\mathbf{A}_i(\mathcal{B})^T] \mathbf{X}_i]^{-1} \mathbf{X}_i^T (\mathbf{A}_i(\mathcal{B}) - \mathbf{y}_i) \\ &= \hat{\mathcal{B}} + \mathcal{K}(\mathcal{B})^{-1} \mathbf{X}_i^T (\mathbf{A}_i(\mathcal{B}) - \mathbf{y}_i) \\ &\quad - \mathcal{K}(\mathcal{B})^{-1} \mathbf{X}_i^T \left\{ -[\text{diag}(\mathbf{A}_i(\mathcal{B})) - \mathbf{A}_i(\mathcal{B})\mathbf{A}_i(\mathcal{B})^T]^{-1} + \mathbf{X}_i \mathcal{K}(\mathcal{B})^{-1} \mathbf{X}_i^T \right\}^{-1} \mathbf{X}_i \mathcal{K}(\mathcal{B})^{-1} \mathbf{X}_i^T (\mathbf{A}_i(\mathcal{B}) - \mathbf{y}_i) \end{aligned}$$

4.3 Approximate Leave- i -Out Prediction

Given the approximate leave- i -out estimation, we can do the approximate leave- i -out prediction as

$$\begin{aligned} \mathbf{y}_i^{/i} &= \begin{bmatrix} y_{i1}^{/i} \\ \vdots \\ y_{iK}^{/i} \end{bmatrix} = \mathbf{X}_i \tilde{\mathcal{B}}^{/i} \\ &= \mathbf{X}_i \hat{\mathcal{B}} + \mathbf{X}_i \mathcal{K}(\mathcal{B})^{-1} \mathbf{X}_i^T (\mathbf{A}_i(\mathcal{B}) - \mathbf{y}_i) \\ &\quad - \mathbf{X}_i \mathcal{K}(\mathcal{B})^{-1} \mathbf{X}_i^T \left\{ -[\text{diag}(\mathbf{A}_i(\mathcal{B})) - \mathbf{A}_i(\mathcal{B})\mathbf{A}_i(\mathcal{B})^T]^{-1} + \mathbf{X}_i \mathcal{K}(\mathcal{B})^{-1} \mathbf{X}_i^T \right\}^{-1} \mathbf{X}_i \mathcal{K}(\mathcal{B})^{-1} \mathbf{X}_i^T (\mathbf{A}_i(\mathcal{B}) - \mathbf{y}_i) \end{aligned}$$

5 Computational Complexity Analysis

5.1 Regular ALO

Here, we take the loss function

$$l(\mathbf{x}_j^T \boldsymbol{\beta}; y_j) = (y_j - \mathbf{x}_j^T \boldsymbol{\beta})^2$$

and regularizer

$$R(\boldsymbol{\beta}) = \lambda_1 \|\boldsymbol{\beta}\|_1 + \lambda_2 \|\boldsymbol{\beta}\|_2^2$$

Under the Regular ALO case, we need to do the following linear algebra computation.

Assuming the active set is E , then we have

$$(\mathbf{X}^T \mathbf{X})_{E,E} + \nabla^2 R(\boldsymbol{\beta}_E), \quad O(|E|^2)$$

$$\begin{aligned}
& \left[(\mathbf{X}^T \mathbf{X})_{E,E} + \nabla^2 R(\boldsymbol{\beta}_E) \right]^{-1}, \quad O(|E|^3) \\
\mathbf{H} &= \mathbf{X}_{\cdot E} \left[(\mathbf{X}^T \mathbf{X})_{E,E} + \nabla^2 R(\boldsymbol{\beta}_E) \right]^{-1} \mathbf{X}_{\cdot E}^T, \quad O(n|E|^2) + O(n|E|) \rightarrow O(n|E|^2) \\
& \frac{\text{diag}(\mathbf{H})}{1 - \text{diag}(\mathbf{H})} \cdot (\mathbf{X}\boldsymbol{\beta} - \mathbf{y}), \quad O(n) + O(np) + O(n) + O(n) \rightarrow O(np)
\end{aligned}$$

So, the final computational complexity for regular ALO method is

$$O(|E|^3) + O(n|E|^2) \rightarrow O(|E|^2 \max\{n, |E|\})$$

5.2 Block Inversion Lemma in ALO

Here, we take the loss function

$$l(\mathbf{x}_j^T \boldsymbol{\beta}; y_j) = (y_j - \mathbf{x}_j^T \boldsymbol{\beta})^2$$

and regularizer

$$R(\boldsymbol{\beta}) = \lambda_1 \|\boldsymbol{\beta}\|_1 + \lambda_2 \|\boldsymbol{\beta}\|_2^2$$

Under the Block Inversion ALO case, we need to do the following linear algebra computation.

Assume the last active set is E_t .

At this step, we need to do the following

$$E_t = E_{t1} \cup E_{t2} \xrightarrow{\text{drop}} E_{t1} \xrightarrow{\text{add}} E_{t+1} = E_{t1} \cup E_{t3}$$

At the drop step,

$$\begin{aligned}
\mathbf{A}_t^{-1} &= (\mathbf{X}^T \mathbf{X} + \nabla^2 R(\boldsymbol{\beta}))_{E_t E_t}^{-1} = \begin{bmatrix} (\mathbf{X}^T \mathbf{X} + \nabla^2 R(\boldsymbol{\beta}))_{E_{t1} E_{t1}} & (\mathbf{X}^T \mathbf{X} + \nabla^2 R(\boldsymbol{\beta}))_{E_{t1} E_{t2}} \\ (\mathbf{X}^T \mathbf{X} + \nabla^2 R(\boldsymbol{\beta}))_{E_{t2} E_{t1}} & (\mathbf{X}^T \mathbf{X} + \nabla^2 R(\boldsymbol{\beta}))_{E_{t2} E_{t2}} \end{bmatrix}^{-1} \\
&= \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \\
\Rightarrow (\mathbf{X}^T \mathbf{X} + \nabla^2 R(\boldsymbol{\beta}))_{E_{t1} E_{t1}}^{-1} &= \mathbf{A}_{11} - \mathbf{A}_{12} \mathbf{A}_{22}^{-1} \mathbf{A}_{21} \\
&\Rightarrow O(|E_{t2}|^3) + O(|E_{t1}| |E_{t2}|^2) + O(|E_{t1}|^2 |E_{t2}|) + O(|E_{t1}|^2) \rightarrow O(|E_{t2}| \max\{|E_{t1}|, |E_{t2}|\}^2)
\end{aligned}$$

At the addition step,

$$\begin{aligned}
\mathbf{A}_{t+1}^{-1} &= (\mathbf{X}^T \mathbf{X} + \nabla^2 R(\boldsymbol{\beta}))_{E_{t+1} E_{t+1}}^{-1} = \begin{bmatrix} (\mathbf{X}^T \mathbf{X} + \nabla^2 R(\boldsymbol{\beta}))_{E_{t1} E_{t1}} & (\mathbf{X}^T \mathbf{X} + \nabla^2 R(\boldsymbol{\beta}))_{E_{t1} E_{t3}} \\ (\mathbf{X}^T \mathbf{X} + \nabla^2 R(\boldsymbol{\beta}))_{E_{t3} E_{t1}} & (\mathbf{X}^T \mathbf{X} + \nabla^2 R(\boldsymbol{\beta}))_{E_{t3} E_{t3}} \end{bmatrix}^{-1} \\
&= \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{13} \\ \mathbf{A}_{31} & \mathbf{A}_{33} \end{bmatrix}^{-1} \\
&= \begin{bmatrix} \mathbf{A}_{11}^{-1} + \mathbf{A}_{11}^{-1} \mathbf{A}_{13} \mathbf{E} \mathbf{A}_{31} \mathbf{A}_{11}^{-1} & -\mathbf{A}_{11}^{-1} \mathbf{A}_{13} \mathbf{E} \\ -\mathbf{E} \mathbf{A}_{31} \mathbf{A}_{11}^{-1} & \mathbf{E} \end{bmatrix}, \quad \text{where } \mathbf{E} = (\mathbf{A}_{33} - \mathbf{A}_{31} \mathbf{A}_{11}^{-1} \mathbf{A}_{13})^{-1} \\
\Rightarrow \begin{cases} O(|E_{t1}|^2 |E_{t3}|) + O(|E_{t1}| |E_{t3}|^2) + O(|E_{t3}|^3) \rightarrow O(|E_{t3}| \max\{|E_{t1}|, |E_{t3}|\}^2) \\ O(|E_{t1}| |E_{t3}|^2) + O(|E_{t1}|^2 |E_{t3}|) \rightarrow O(|E_{t1}| |E_{t3}| \max\{|E_{t1}|, |E_{t3}|\}) \\ O(|E_{t1}| |E_{t3}|^2) \end{cases} \\
&\Rightarrow O(|E_{t3}| \max\{|E_{t1}|, |E_{t3}|\}^2)
\end{aligned}$$

At the Schulz Iteration step, for start point,

$$\mathbf{V}_0 = [\mathbf{I} - (\lambda_{t+1,2} - \lambda_{t,2}) \mathbf{A}_{t+1}^{-1}] \mathbf{A}_{t+1}^{-1}, \quad O(|E_{t+1}|^3)$$

for each iteration,

$$\mathbf{V}_{k+1} = \mathbf{V}_k(2\mathbf{I} - \mathbf{F}_{t+1}\mathbf{V}_k), \text{ where } \mathbf{F}_{t+1} = (\mathbf{X}^T\mathbf{X} + \nabla^2 R(\boldsymbol{\beta}_{t+1}))_{E_{t+1}E_{t+1}}, \quad O(|E_{t+1}|^3)$$

At the final step,

$$\begin{aligned} \mathbf{H} &= \mathbf{X}_{\cdot E_{t+1}} \mathbf{V}_{\inf} \mathbf{X}_{\cdot E_{t+1}}^T, \quad O(n|E_{t+1}|^2) + O(n|E_{t+1}|) \rightarrow O(n|E_{t+1}|^2) \\ \frac{\text{diag}(\mathbf{H})}{1 - \text{diag}(\mathbf{H})} &\cdot (\mathbf{X}\boldsymbol{\beta} - \mathbf{y}), \quad O(n) + O(np) + O(n) + O(n) \rightarrow O(np) \end{aligned}$$

So, the final computational complexity for Block Inversion ALO method is

$$O(|E_{t+1}|^3) + O(n|E_{t+1}|^2) \rightarrow O(|E_{t+1}|^2 \max\{n, |E_{t+1}|\})$$

6 Cholesky Decompostion in Elastic Net

6.1 Definition

Define the original design matrix as $\mathbf{X}_{n \times p}$, active set E_t corresponding to the parameter λ_t . The second order derivative of regularization function is $\nabla^2 R(\boldsymbol{\beta})$

6.2 Cholesky Decompostion

Here, assuming we already know the Cholesky decomposition of

$$\mathbf{X}_{\cdot E_{t-1}}^T \mathbf{X}_{\cdot E_{t-1}} = \mathbf{L}_{t-1} \mathbf{L}_{t-1}^T$$

and at this step, we want to find out the Cholesky decomposition of

$$\mathbf{X}_{\cdot E_t}^T \mathbf{X}_{\cdot E_t} + \nabla^2 R(\boldsymbol{\beta}_t)$$

Under the Elastic Net case, we have

$$\begin{aligned} \nabla^2 R(\boldsymbol{\beta}_t) &= \begin{bmatrix} 0 & 0 & 0 & \cdots & 0 \\ 0 & 2\lambda_2 & 0 & \cdots & 0 \\ 0 & 0 & 2\lambda_2 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 2\lambda_2 \end{bmatrix} \\ &= \begin{bmatrix} 0 \\ \sqrt{2\lambda_2} \\ 0 \\ \vdots \\ 0 \end{bmatrix} \begin{bmatrix} 0 & \sqrt{2\lambda_2} & 0 & \cdots & 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ \sqrt{2\lambda_2} \\ \vdots \\ 0 \end{bmatrix} \begin{bmatrix} 0 & 0 & \sqrt{2\lambda_2} & \cdots & 0 \end{bmatrix} + \cdots \\ &= \mathbf{r}_2 \mathbf{r}_2^T + \mathbf{r}_3 \mathbf{r}_3^T + \cdots + \mathbf{r}_{|E_t|} \mathbf{r}_{|E_t|}^T \end{aligned}$$

So, at each step, we can first update the Cholesky decomposition of

$$\mathbf{X}_{\cdot E_{t-1}}^T \mathbf{X}_{\cdot E_{t-1}} = \mathbf{L}_{t-1} \mathbf{L}_{t-1}^T \xrightarrow{\text{Cholesky Update}} \mathbf{X}_{\cdot E_t}^T \mathbf{X}_{\cdot E_t} = \mathbf{L}_t \mathbf{L}_t^T$$

Then, based on the rank-one update, we can find out

$$\mathbf{X}_{\cdot E_t}^T \mathbf{X}_{\cdot E_t} + \nabla^2 R(\boldsymbol{\beta}_t) = \mathbf{X}_{\cdot E_t}^T \mathbf{X}_{\cdot E_t} + \mathbf{r}_2 \mathbf{r}_2^T + \mathbf{r}_3 \mathbf{r}_3^T + \cdots + \mathbf{r}_{|E_t|} \mathbf{r}_{|E_t|}^T$$

Here we can use Cholesky rank-one update to find the Cholesky decomposition of $\mathbf{X}_{\cdot E_t}^T \mathbf{X}_{\cdot E_t} + \nabla^2 R(\boldsymbol{\beta}_t)$.