

Notes on Approximate Leave-One-Out for Elastic Net

Linyun He

Wanchao Qin

Peng Xu

Yuze Zhou

August 29, 2018

1 Approximation in the Primal Domain

1.1 ALO for Linear Regression

Recall the objective function for the elastic net problem:

$$\min_{\beta} \frac{1}{2} \sum_{j=1}^n (\mathbf{x}_j^\top \beta - y_j)^2 + \lambda \left(\alpha \|\beta\|_1 + \frac{1-\alpha}{2} \|\beta\|_2^2 \right). \quad (1)$$

Let $A = \{i : \beta_i \neq 0, i = 1, \dots, p\}$ be the active set, we have

$$\dot{\ell}(\mathbf{x}_j^\top \beta; y_j) = \mathbf{x}_j^\top \beta - y_j, \quad \ddot{\ell}(\mathbf{x}_j^\top \beta; y_j) = 1, \quad \nabla^2 R(\hat{\beta}_A) = (1-\alpha)\lambda \mathbf{I}_{A,A}.$$

Thus, Eqn. 31 reduces to

$$\mathbf{H} = \mathbf{X}_{\cdot,A} \left[\mathbf{X}_{\cdot,A}^\top \mathbf{X}_{\cdot,A} + (1-\alpha)\lambda \mathbf{I}_{A,A} \right]^{-1} \mathbf{X}_{\cdot,A}^\top. \quad (2)$$

By augmenting \mathbf{X} with an extra column of 1s, adding the intercept back to the model is straightforward, as Eqn. 31 now becomes

$$\mathbf{H} = [\mathbf{1}_n, \mathbf{X}_{\cdot,A}] \left\{ [\mathbf{1}_n, \mathbf{X}_{\cdot,A}]^\top \mathbf{D} [\mathbf{1}_n, \mathbf{X}_{\cdot,A}] + \nabla^2 R(\hat{\beta}_0, \hat{\beta}_A) \right\}^{-1} [\mathbf{1}_n, \mathbf{X}_{\cdot,A}]^\top, \quad (3)$$

where

$$\mathbf{D} = \text{diag} \left[\ddot{\ell}(\hat{\beta}_0 + \mathbf{x}_j^\top \hat{\beta}; y_j) \right]_{j \in A} = \mathbf{I}_{A,A}, \quad \nabla^2 R(\hat{\beta}_0, \hat{\beta}_A) = \begin{bmatrix} 0 & 0 & \dots & 0 \\ 0 & (1-\alpha)\lambda & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & (1-\alpha)\lambda \end{bmatrix}.$$

Then, the ALO can be computed as

$$\begin{bmatrix} 1 & \mathbf{x}_i^\top \end{bmatrix} \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_i \end{bmatrix} = (\hat{\beta}_0 + \mathbf{x}_i^\top \hat{\beta}) + \frac{\mathbf{H}_{ii}}{1 - \mathbf{H}_{ii} \ddot{\ell}(\hat{\beta}_0 + \mathbf{x}_i^\top \hat{\beta}; y_i)} \dot{\ell}(\hat{\beta}_0 + \mathbf{x}_i^\top \hat{\beta}; y_i) \quad (4)$$

1.2 ALO for Logistic Regression, Approximation in the Primal Domain

For binomial logistic regression, the primal problem is:

$$\min_{\beta} \sum_{j=1}^n \left[\ln \left(1 + e^{\mathbf{x}_j^\top \beta} \right) - y_j \mathbf{x}_j^\top \beta \right] + \lambda \left(\alpha \|\beta\|_1 + \frac{1-\alpha}{2} \|\beta\|_2^2 \right).$$

Let A be the active set, we have

$$\dot{\ell}(\mathbf{x}_j^\top \beta; y_j) = \frac{e^{\mathbf{x}_j^\top \beta}}{1 + e^{\mathbf{x}_j^\top \beta}} - y_j, \quad \ddot{\ell}(\mathbf{x}_j^\top \beta; y_j) = \frac{e^{\mathbf{x}_j^\top \beta}}{\left(1 + e^{\mathbf{x}_j^\top \beta} \right)^2}, \quad \nabla^2 R(\hat{\beta}_A) = (1 - \alpha) \lambda \mathbf{I}_{A,A}.$$

Therefore, the ALO formula is

$$\mathbf{x}_i^\top \tilde{\beta}^{\setminus j} \approx \mathbf{x}_i^\top \hat{\beta} + \frac{\mathbf{H}_{ii} \left(1 + e^{\mathbf{x}_j^\top \hat{\beta}} \right) \left[e^{\mathbf{x}_j^\top \hat{\beta}} - y_j \left(1 + e^{\mathbf{x}_j^\top \hat{\beta}} \right) \right]}{\left(1 + e^{\mathbf{x}_j^\top \hat{\beta}} \right)^2 - \mathbf{H}_{ii} e^{\mathbf{x}_j^\top \hat{\beta}}},$$

where

$$\mathbf{H} = \mathbf{X}_{\cdot,A} \left[\mathbf{X}_{\cdot,A}^\top \text{diag} \left(\frac{e^{\mathbf{x}_j^\top \hat{\beta}}}{1 + 2e^{\mathbf{x}_j^\top \hat{\beta}} + e^{2\mathbf{x}_j^\top \hat{\beta}}} \right) \mathbf{X}_{\cdot,A} + (1 - \alpha) \lambda \mathbf{I}_{A,A} \right]^{-1} \mathbf{X}_{\cdot,A}^\top.$$

1.3 ALO for Poisson Regression, Approximation in the Primal Domain

For Poisson regression, the primal problem is:

$$\min_{\beta} \sum_{j=1}^n \left(e^{\mathbf{x}_j^\top \beta} - y_j \mathbf{x}_j^\top \beta \right) + \lambda \left(\alpha \|\beta\|_1 + \frac{1-\alpha}{2} \|\beta\|_2^2 \right).$$

Let A be the active set, we have

$$\dot{\ell}(\mathbf{x}_j^\top \beta; y_j) = e^{\mathbf{x}_j^\top \beta} - y_j, \quad \ddot{\ell}(\mathbf{x}_j^\top \beta; y_j) = e^{\mathbf{x}_j^\top \beta}, \quad \nabla^2 R(\hat{\beta}_A) = (1 - \alpha) \lambda \mathbf{I}_{A,A}.$$

Therefore, the ALO formula is

$$\mathbf{x}_i^\top \tilde{\beta}^{\setminus j} \approx \mathbf{x}_i^\top \hat{\beta} + \frac{\mathbf{H}_{ii} \left(e^{\mathbf{x}_j^\top \hat{\beta}} - y_j \right)}{1 - \mathbf{H}_{ii} e^{\mathbf{x}_j^\top \hat{\beta}}},$$

where

$$\mathbf{H} = \mathbf{X}_{\cdot,A} \left[\mathbf{X}_{\cdot,A}^\top \text{diag} \left(e^{\mathbf{x}_j^\top \hat{\beta}} \right) \mathbf{X}_{\cdot,A} + (1 - \alpha) \lambda \mathbf{I}_{A,A} \right]^{-1} \mathbf{X}_{\cdot,A}^\top.$$

2 Approximation in the Dual Domain

2.1 ALO for Linear Regression

First, to write the optimization problem of elastic net in a matrix form, and denote $D = [0, I]$, the optimization problem becomes:

$$\hat{\beta} = \arg \min_{\beta} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda_1 \|\mathbf{D}\beta\|_1 + \lambda_2 \|\mathbf{D}\beta\|_2^2$$

The augmented Lagrangian for the problem is:

$$L = \frac{1}{2} \|\mathbf{y} - \mathbf{z}\|_2^2 + \lambda_1 \|\omega\|_1 + \lambda_2 \|\omega\|_2^2 + \mathbf{u}^\top (\mathbf{z} - \mathbf{X}\beta) + \mathbf{v}^\top (\omega - \mathbf{D}\beta)$$

By taking the derivatives with respect to \mathbf{z} and β , we could obtain:

$$0 = \frac{\partial L}{\partial \mathbf{z}} = \mathbf{z} - \mathbf{y} + \mathbf{u}, \quad 0 = \frac{\partial L}{\partial \beta} = -\mathbf{X}^\top \mathbf{u} - \mathbf{D}^\top \mathbf{v}.$$

Since the first column of \mathbf{X} is filled with ones and the first column of \mathbf{D} is filled with zeros, the first row of $-\mathbf{X}^\top \mathbf{u} - \mathbf{D}^\top \mathbf{v} = 0$ gives $\mathbf{1}^\top \mathbf{u} = 0$ and due to $\mathbf{D} = [0, I]$, the rest rows give that $-\mathbf{X}_j^\top \mathbf{u} = \mathbf{v}_j$. Moreover, since the rest dimensions of ω is penalized element-wisely in the augmented Lagrangian, we can minimize over ω by minimizing over each ω_i , $i \geq 2$, that is, we have to minimize $\lambda_1 |\omega_i| + \lambda_2 \omega_i^2 - \mathbf{u}^\top \mathbf{X}_i \omega_i$ for each dimension of ω , where \mathbf{X}_i denotes the i th column of \mathbf{X} , therefore:

$$\min_{\omega_i} (\lambda_1 |\omega_i| + \lambda_2 \omega_i^2 - \mathbf{u}^\top \mathbf{X}_i \omega_i) = \begin{cases} 0 & |\mathbf{u}^\top \mathbf{X}_i| \leq \lambda_1, \\ -\frac{(\lambda_1 - |\mathbf{u}^\top \mathbf{X}_i|)^2}{4\lambda_2} & |\mathbf{u}^\top \mathbf{X}_i| > \lambda_1. \end{cases}$$

By taking all the above back to the Lagrangian, we obtain the dual problem as:

$$d^* = \min_{\mathbf{u}} \frac{1}{2} \|\mathbf{y} - \mathbf{u}\|_2^2 + \sum_{j: |\mathbf{X}_j^\top \mathbf{u}| > \lambda_1} \frac{(\lambda_1 - |\mathbf{u}^\top \mathbf{X}_i|)^2}{4\lambda_2}, \quad \text{subject to } \mathbf{1}^\top \mathbf{u} = 0.$$

Denote

$$f(\mathbf{u}) = \sum_{j: |\mathbf{X}_j^\top \mathbf{u}| > \lambda_1} \frac{(\lambda_1 - |\mathbf{u}^\top \mathbf{X}_i|)^2}{4\lambda_2},$$

clearly $f(\mathbf{u})$ is of quadratic form:

$$f(\mathbf{u}) = \frac{1}{2} \mathbf{u}^\top \mathbf{A} \mathbf{u} + \mathbf{a}^\top \mathbf{u} + b,$$

where b is a constant and does not matter in the optimization of the dual problem, \mathbf{A} and \mathbf{a} are:

$$\mathbf{A} = \frac{1}{2\lambda_2} \mathbf{X}_{\cdot, E} \mathbf{X}_{\cdot, E}^\top, \quad \mathbf{a} = \frac{\lambda_1}{2\lambda_2} \left(\sum_{i: \mathbf{X}_i^\top \leq \lambda_1} \mathbf{X}_i - \sum_{i: \mathbf{X}_i^\top > \lambda_1} \mathbf{X}_i \right)$$

where $E := \{i : |\mathbf{X}_i^\top \mathbf{u}| > \lambda\}$.

The dual problem could also be written in a proximal form:

$$\hat{\mathbf{u}} = \text{prox}_{\tilde{f}}(\mathbf{y}), \quad \tilde{f} = \mathbf{I}(\mathbf{1}^\top \mathbf{u} = 0)f(\mathbf{u}) + \mathbf{I}(\mathbf{1}^\top \mathbf{u} \neq 0)\infty.$$

After transforming $f(\mathbf{u})$ into a quadratic form, we could write the Lagrangian for the dual problem back again:

$$L = \frac{1}{2}\|\mathbf{y} - \mathbf{u}\|_2^2 + \frac{1}{2}\mathbf{u}^\top \mathbf{A}\mathbf{u} + \mathbf{a}^\top \mathbf{u} + b + \lambda \mathbf{1}^\top \mathbf{u}.$$

By taking the derivative with respect to \mathbf{u} , we could obtain:

$$\frac{\partial L}{\partial \mathbf{u}} = \mathbf{u} - \mathbf{y} + \mathbf{A}\mathbf{u} + \mathbf{a} + \lambda \mathbf{1} = 0$$

By shifting the terms, \mathbf{u} could be written as a formula of \mathbf{y} and λ : $\mathbf{u} = (\mathbf{I} + \mathbf{A})^{-1}(\mathbf{y} - \mathbf{a} - \lambda \mathbf{1})$, by taking the derivative with respect to \mathbf{y} at both sides, we could obtain the Jacobian matrix \mathbf{J} of the proximal operator $\text{prox}(\tilde{f})$ at \mathbf{y} as:

$$\mathbf{J} = (\mathbf{I} + \mathbf{A})^{-1} - (\mathbf{I} + \mathbf{A})^{-1} \mathbf{1} \nabla(\hat{\lambda})^\top$$

where $\nabla(\hat{\lambda})^\top$ denotes the gradient of λ as a function of \mathbf{y} . By taking $\mathbf{u} = (\mathbf{I} + \mathbf{A})^{-1}(\mathbf{y} - \mathbf{a} - \lambda \mathbf{1})$ back to the Lagrangian, the dual problem will become a second-order equation of λ :

$$\begin{aligned} d^* = \max_{\lambda} & \frac{1}{2}\|\mathbf{y} - (\mathbf{I} + \mathbf{A})^{-1}(\mathbf{y} - \mathbf{a} - \lambda \mathbf{1})\|_2^2 \\ & + \frac{1}{2}(\mathbf{y} - \mathbf{a} - \lambda \mathbf{1})^\top (\mathbf{I} + \mathbf{A})^{-1} \mathbf{A} (\mathbf{I} + \mathbf{A})^{-1} (\mathbf{y} - \mathbf{a} - \lambda \mathbf{1}) + \mathbf{a}^\top (\mathbf{I} + \mathbf{A})^{-1} (\mathbf{y} - \mathbf{a} - \lambda \mathbf{1}) \\ & + \lambda \mathbf{1}^\top (\mathbf{I} + \mathbf{A})^{-1} (\mathbf{y} - \mathbf{a} - \lambda \mathbf{1}) \end{aligned}$$

More specifically, the second-order term is:

$$\frac{1}{2} \mathbf{1}^\top (\mathbf{I} + \mathbf{A})^{-2} \mathbf{1} + \frac{1}{2} \mathbf{1}^\top (\mathbf{I} + \mathbf{A})^{-1} \mathbf{A} (\mathbf{I} + \mathbf{A})^{-1} \mathbf{1} - \mathbf{1}^\top (\mathbf{I} + \mathbf{A})^{-1} \mathbf{1}$$

and the first-order term is:

$$2\mathbf{1}^\top (\mathbf{I} + \mathbf{A})^{-1} (\mathbf{y} - \mathbf{a}) - \mathbf{1}^\top (\mathbf{I} + \mathbf{A})^{-2} (\mathbf{y} - \mathbf{a}) - \mathbf{1}^\top (\mathbf{I} + \mathbf{A})^{-1} \mathbf{A} (\mathbf{I} + \mathbf{A})^{-1} (\mathbf{y} - \mathbf{a})$$

Thus by solving the second-order equation, we could obtain

$$\hat{\lambda} = \frac{2\mathbf{1}^\top (\mathbf{I} + \mathbf{A})^{-1} (\mathbf{y} - \mathbf{a}) - \mathbf{1}^\top (\mathbf{I} + \mathbf{A})^{-2} (\mathbf{y} - \mathbf{a}) - \mathbf{1}^\top (\mathbf{I} + \mathbf{A})^{-1} \mathbf{A} (\mathbf{I} + \mathbf{A})^{-1} (\mathbf{y} - \mathbf{a})}{\mathbf{1}^\top (\mathbf{I} + \mathbf{A})^{-2} \mathbf{1} + \mathbf{1}^\top (\mathbf{I} + \mathbf{A})^{-1} \mathbf{A} (\mathbf{I} + \mathbf{A})^{-1} \mathbf{1} - 2\mathbf{1}^\top (\mathbf{I} + \mathbf{A})^{-1} \mathbf{1}}$$

and the gradient

$$\nabla(\hat{\lambda}) = \frac{2(\mathbf{I} + \mathbf{A})^{-1} \mathbf{1} - (\mathbf{I} + \mathbf{A})^{-2} \mathbf{1} - (\mathbf{I} + \mathbf{A})^{-1} \mathbf{A} (\mathbf{I} + \mathbf{A})^{-1} \mathbf{1}}{\mathbf{1}^\top (\mathbf{I} + \mathbf{A})^{-2} \mathbf{1} + \mathbf{1}^\top (\mathbf{I} + \mathbf{A})^{-1} \mathbf{A} (\mathbf{I} + \mathbf{A})^{-1} \mathbf{1} - 2\mathbf{1}^\top (\mathbf{I} + \mathbf{A})^{-1} \mathbf{1}}$$

By taking the gradient back to

$$\mathbf{J} = (\mathbf{I} + \mathbf{A})^{-1} - (\mathbf{I} + \mathbf{A})^{-1} \mathbf{1} \nabla(\hat{\lambda})^\top = (\mathbf{I} + \mathbf{A})^{-1} - \frac{(\mathbf{I} + \mathbf{A})^{-1} \mathbf{1} \mathbf{1}^\top (\mathbf{I} + \mathbf{A})^{-1}}{\mathbf{1}^\top (\mathbf{I} + \mathbf{A})^{-1} \mathbf{1}},$$

we could obtain the Jacobian.

2.2 Proof of the Equivalence between the primal and the dual solutions

First recall the ALO formula from the dual approach

$$\mathbf{y}^{/i} = \mathbf{y}_i - \frac{\mathbf{u}_i}{J_{ii}} = \frac{J_{ii} - 1}{J_{ii}} \mathbf{y}_i + \frac{1}{J_{ii}} \mathbf{x}_i \hat{\beta}$$

and the primal formula from the primal approach

$$\mathbf{y}^{/i} = \mathbf{x}_i \hat{\beta} + \frac{H_{ii}}{1 - H_{ii}} (\mathbf{x}_i \hat{\beta} - \mathbf{y}_i) = -\frac{H_{ii}}{1 - H_{ii}} \mathbf{y}_i + \frac{1}{1 - H_{ii}} \mathbf{x}_i \hat{\beta}.$$

In this section, we're going to show that $\mathbf{H} + \mathbf{J} = \mathbf{I}$, thus giving $H_{ii} + J_{ii} = 1$, and that the solutions given by both the primal and the dual approach are equivalent.

First, using matrix inverse lemma, we could calculate the inverse of $(\mathbf{I} + \mathbf{A}) = (\mathbf{I} + \frac{1}{2\lambda_2} \mathbf{X}_{\cdot,E} \mathbf{X}_{\cdot,E}^\top)$ as:

$$(\mathbf{I} + \mathbf{A})^{-1} = (\mathbf{I} + \frac{1}{2\lambda_2} \mathbf{X}_{\cdot,E} \mathbf{X}_{\cdot,E}^\top)^{-1} = \mathbf{I} - \mathbf{X}_{\cdot,E} (2\lambda_2 \mathbf{I} + \mathbf{X}_{\cdot,E}^\top \mathbf{X}_{\cdot,E})^{-1} \mathbf{X}_{\cdot,E}^\top,$$

therefore the matrix \mathbf{J} is:

$$\begin{aligned} \mathbf{J} &= (\mathbf{I} + \mathbf{A})^{-1} - \frac{(\mathbf{I} + \mathbf{A})^{-1} \mathbf{1} \mathbf{1}^\top (\mathbf{I} + \mathbf{A})^{-1}}{\mathbf{1}^\top (\mathbf{I} + \mathbf{A})^{-1} \mathbf{1}} \\ &= \mathbf{I} - \mathbf{X}_{\cdot,E} (2\lambda_2 \mathbf{I} + \mathbf{X}_{\cdot,E}^\top \mathbf{X}_{\cdot,E})^{-1} \mathbf{X}_{\cdot,E}^\top \\ &\quad - \frac{(\mathbf{I} - \mathbf{X}_{\cdot,E} (2\lambda_2 \mathbf{I} + \mathbf{X}_{\cdot,E}^\top \mathbf{X}_{\cdot,E})^{-1} \mathbf{X}_{\cdot,E}^\top \mathbf{1})(\mathbf{1}^\top - \mathbf{1}^\top \mathbf{X}_{\cdot,E} (2\lambda_2 \mathbf{I} + \mathbf{X}_{\cdot,E}^\top \mathbf{X}_{\cdot,E})^{-1} \mathbf{X}_{\cdot,E}^\top)}{\mathbf{1}^\top (\mathbf{I} - \mathbf{X}_{\cdot,E} (2\lambda_2 \mathbf{I} + \mathbf{X}_{\cdot,E}^\top \mathbf{X}_{\cdot,E})^{-1} \mathbf{X}_{\cdot,E}^\top) \mathbf{1}} \end{aligned}$$

Now recall that

$$\mathbf{H} = [\mathbf{1}, \mathbf{X}_{\cdot,E}] ([\mathbf{1}, \mathbf{X}_{\cdot,E}]^\top [\mathbf{1}, \mathbf{X}_{\cdot,E}] + \text{diag}[0, 2\lambda_2, \dots, 2\lambda_2]) [\mathbf{1}, \mathbf{X}_{\cdot,E}]^\top,$$

by adopting block inverse, we could derive \mathbf{H} as:

$$\begin{aligned} \mathbf{H} &= [\mathbf{1}, \mathbf{X}_{\cdot,E}] \begin{pmatrix} n & \mathbf{1}^\top \mathbf{X}_{\cdot,E} \\ \mathbf{X}_{\cdot,E}^\top \mathbf{1} & \mathbf{X}_{\cdot,E}^\top \mathbf{X}_{\cdot,E} + 2\lambda_2 \mathbf{I} \end{pmatrix} [\mathbf{1}, \mathbf{X}_{\cdot,E}]^\top \\ &= [\mathbf{1}, \mathbf{X}_{\cdot,E}] \begin{pmatrix} \frac{1}{\mathbf{1}^\top (\mathbf{I} - \mathbf{X}_{\cdot,E} (2\lambda_2 \mathbf{I} + \mathbf{X}_{\cdot,E}^\top \mathbf{X}_{\cdot,E})^{-1} \mathbf{X}_{\cdot,E}^\top) \mathbf{1}} & \frac{-\mathbf{1}^\top \mathbf{X}_{\cdot,E} (2\lambda_2 \mathbf{I} + \mathbf{X}_{\cdot,E}^\top \mathbf{X}_{\cdot,E})^{-1}}{\mathbf{1}^\top (\mathbf{I} - \mathbf{X}_{\cdot,E} (2\lambda_2 \mathbf{I} + \mathbf{X}_{\cdot,E}^\top \mathbf{X}_{\cdot,E})^{-1} \mathbf{X}_{\cdot,E}^\top) \mathbf{1}} \\ \frac{-(2\lambda_2 \mathbf{I} + \mathbf{X}_{\cdot,E}^\top \mathbf{X}_{\cdot,E})^{-1} \mathbf{X}_{\cdot,E}^\top \mathbf{1}}{\mathbf{1}^\top (\mathbf{I} - \mathbf{X}_{\cdot,E} (2\lambda_2 \mathbf{I} + \mathbf{X}_{\cdot,E}^\top \mathbf{X}_{\cdot,E})^{-1} \mathbf{X}_{\cdot,E}^\top) \mathbf{1}} & (2\lambda_2 \mathbf{I} + \mathbf{X}_{\cdot,E}^\top \mathbf{X}_{\cdot,E})^{-1} + \frac{(2\lambda_2 \mathbf{I} + \mathbf{X}_{\cdot,E}^\top \mathbf{X}_{\cdot,E})^{-1} \mathbf{X}_{\cdot,E}^\top \mathbf{1} \mathbf{1}^\top \mathbf{X}_{\cdot,E} (2\lambda_2 \mathbf{I} + \mathbf{X}_{\cdot,E}^\top \mathbf{X}_{\cdot,E})^{-1}}{\mathbf{1}^\top (\mathbf{I} - \mathbf{X}_{\cdot,E} (2\lambda_2 \mathbf{I} + \mathbf{X}_{\cdot,E}^\top \mathbf{X}_{\cdot,E})^{-1} \mathbf{X}_{\cdot,E}^\top) \mathbf{1}} \end{pmatrix} \\ &\quad \times [\mathbf{1}, \mathbf{X}_{\cdot,E}]^\top \\ &= \mathbf{X}_{\cdot,E} (2\lambda_2 \mathbf{I} + \mathbf{X}_{\cdot,E}^\top \mathbf{X}_{\cdot,E})^{-1} \mathbf{X}_{\cdot,E}^\top \\ &\quad + \frac{(\mathbf{I} - \mathbf{X}_{\cdot,E} (2\lambda_2 \mathbf{I} + \mathbf{X}_{\cdot,E}^\top \mathbf{X}_{\cdot,E})^{-1} \mathbf{X}_{\cdot,E}^\top \mathbf{1})(\mathbf{1}^\top - \mathbf{1}^\top \mathbf{X}_{\cdot,E} (2\lambda_2 \mathbf{I} + \mathbf{X}_{\cdot,E}^\top \mathbf{X}_{\cdot,E})^{-1} \mathbf{X}_{\cdot,E}^\top)}{\mathbf{1}^\top (\mathbf{I} - \mathbf{X}_{\cdot,E} (2\lambda_2 \mathbf{I} + \mathbf{X}_{\cdot,E}^\top \mathbf{X}_{\cdot,E})^{-1} \mathbf{X}_{\cdot,E}^\top) \mathbf{1}} \\ &= \mathbf{I} - \mathbf{J} \end{aligned}$$

Now that we have showed that $\mathbf{H} + \mathbf{J} = \mathbf{I}$, we could also conclude that $J_{ii} + H_{ii} = 1$ and that the ALO solutions given by both the primal and dual approaches are equivalent.

2.3 ALO for Logistic Regression with Lasso penalty

First, let's rewrite the optimization problem with the loss functions separated for each observation, therefore the loss function goes:

$$-\sum_{i=1}^n \left[y_i x_i^\top \beta + \ln(1 + e^{x_i^\top \beta}) \right] + \lambda_1 \|\beta\|_1,$$

where the individual loss function is $\ell(x_i^\top \beta; y_i) = y_i x_i^\top \beta + \ln(1 + e^{x_i^\top \beta})$ and the regularizer is $R(\beta) = \lambda_1 \|\beta\|_1$, from which we could derive the dual optimal and the conjugate functions:

$$\hat{\theta} = y - \frac{e^{X\hat{\beta}}}{1 + e^{X\hat{\beta}}}, \quad \ell^*(-\theta_i; y_i) = (y_i - \theta_i) \ln \frac{y_i - \theta_i}{1 - (y_i - \theta_i)} - \ln \frac{1}{1 - (y_i - \theta_i)}, \quad R^*(\beta) = \begin{cases} 0 & \|\beta\|_\infty \leq \lambda_1, \\ \infty & \text{otherwise.} \end{cases}$$

From the results of the conjugate functions above, we could also obtain the derivatives of the loss functions and the Jacobian of the regularizer:

$$\dot{\ell}^*(-\theta_i; y_i) = \ln \frac{y_i - \theta_i}{1 - (y_i - \theta_i)}, \quad \ddot{\ell}^*(-\theta_i; y_i) = \frac{1}{(y_i - \theta_i)[1 - (y_i - \theta_i)]}$$

Recall Eqn. 20 from the main paper, the quadratic surrogate of the dual problem is

$$\min_u \frac{1}{2} \sum_{i=1}^n \left(u_i - \frac{\hat{\theta}_i \ddot{\ell}^*(-\hat{\theta}_i; y_i) + \hat{y}_i}{\sqrt{\ddot{\ell}^*(-\hat{\theta}_i; y_i)}} \right)^2 + R^*(X^\top K u),$$

where $K = \text{diag} \sqrt{\ddot{\ell}^*(-\hat{\theta}_i; y_i)}$. Therefore the Jacobian at $\mathbf{y}_u = \hat{\theta}_i \ddot{\ell}^*(-\hat{\theta}_i; y_i) + \hat{y}_i / \sqrt{\ddot{\ell}^*(-\hat{\theta}_i; y_i)}$ could locally be treated as the projection onto the orthogonal complement of the polyhedron $\{\|X^\top K u\|_\infty \leq \lambda_1\}$, thus $\mathbf{J} = \mathbf{I} - \mathbf{X}_{u,E} (\mathbf{X}_{u,E}^\top \mathbf{X}_{u,E})^{-1} \mathbf{X}_{u,E}$, where $\mathbf{X}_{u,E}$ are the columns of $\mathbf{X}_u = X^\top K$, such that the columns in the set $E = \{i : |\mathbf{X}_i^\top \theta| = \lambda_1\}$ are selected. Take everything to Eqn. 22, $y^{/i} = K_{ii}(y_{u,i} - K_{ii} \hat{\theta}_i / J_{ii})$, we could obtain the ALO for the i -th observation.

2.4 ALO for Logistic Regression with Elastic Net Penalty

The optimization problem for logistic regression with elastic net penalty is:

$$-\sum_{i=1}^n \left[y_i x_i^\top \beta + \ln(1 + e^{x_i^\top \beta}) \right] + \lambda_1 \|\beta\|_1 + \lambda_2 \|\beta\|_2^2.$$

The optimization problem is the same except the regularizer is changed, therefore the only thing different is the conjugate function of the regularizer, R^* and the corresponding Jacobian, here $R(\beta) = \lambda_1 \|\beta\|_1 + \lambda_2 \|\beta\|_2^2$:

$$R^*(\beta) = \sum_{|u_i| > \lambda_1} \frac{(\lambda_1 - |u_i|)^2}{4\lambda_2}.$$

The corresponding Jacobian is $\mathbf{J} = (\mathbf{I} + \mathbf{X}_{u,E} \mathbf{X}_{u,E}^\top / 2\lambda_2)$, where $\mathbf{X}_{u,E}$ are the columns of $\mathbf{X}_u = X^\top K$, such that the columns in the set $E = \{i : |\mathbf{X}_i^\top \theta| = \lambda_1\}$ are selected. Take everything to Eqn. 22, $y^{/i} = K_{ii}(y_{u,i} - K_{ii} \hat{\theta}_i / J_{ii})$, we could obtain the ALO for the i -th observation.

2.5 ALO for Poisson Regression with Lasso penalty

The optimization function for Poisson regression with lasso penalty is:

$$\sum_{i=1}^n \left[-y_i \mathbf{x}_i^\top \boldsymbol{\beta} + e^{\mathbf{x}_i^\top \boldsymbol{\beta}} + \ln(y_i!) \right] + \lambda_1 \|\boldsymbol{\beta}\|_1$$

The regularizer is the same as the logistic regression with the lasso penalty case, thus the Jacobian will also be the same, therefore we only have to focus on the loss function $\ell(\mathbf{x}_i^\top \boldsymbol{\beta}; y_i) = -y_i \mathbf{x}_i^\top \boldsymbol{\beta} + e^{\mathbf{x}_i^\top \boldsymbol{\beta}} + \ln(y_i!)$. The optimal solution for the dual problem $\hat{\theta} = y - e^{\mathbf{X}\hat{\boldsymbol{\beta}}}$ and the conjugate of the loss function is $\ell^*(-\theta_i; y_i) = (y_i - \theta_i) \ln(y_i - \theta_i) - (y_i - \theta_i)$, the corresponding derivatives are therefore:

$$\dot{\ell}^*(-\theta_i; y_i) = \ln(y_i - \theta_i), \quad \ddot{\ell}^*(-\theta_i; y_i) = \frac{1}{y_i - \theta_i}.$$

By plugging everything into Eqn. 22, we obtain the ALO for Poisson regression with the lasso penalty.

2.6 ALO for Poisson Regression with Elastic Net Penalty

The loss function for Poisson regression with elastic net penalty is the same as that of Poisson regression with the lasso penalty and the regularizer of it is the same as that of logistic regression with elastic net penalty. Thus by plugging everything into Eqn. 22, we could obtain the ALO for Poisson regression with elastic net penalty.

3 Approximation with Proximal Formulation

3.1 ALO for Linear Regression

For the elastic net problem, the proximal mapping is known to be

$$\mathbf{prox}_R(\mathbf{z}) = \gamma \operatorname{sgn}(\mathbf{z}) \odot (|\mathbf{z}| - \lambda \mathbf{1}_p)_+, \quad \gamma = \frac{1}{1 + (1 - \alpha)\lambda}. \quad (5)$$

Let E be the active set, if $z_i \in E$, then

$$\frac{\partial}{\partial z_i} \gamma \operatorname{sgn}(z_i)(|z_i| - \lambda)_+ = \gamma.$$

Plug in $\mathbf{z} = \hat{\boldsymbol{\beta}} - \sum_{j=1}^n \dot{\ell}(\mathbf{x}_j^\top \hat{\boldsymbol{\beta}}; y_j) \mathbf{x}_j$, Eqn. 46 thus reduce to

$$\mathbf{H} = \gamma \mathbf{X}_{:,E} \left[\gamma \mathbf{X}_{:,E}^\top \mathbf{X}_{:,E} + (1 - \gamma) \mathbf{I}_{E,E} \right]^{-1} \mathbf{X}_{:,E}^\top. \quad (6)$$

Bringing back the intercept term is straightforward as well. Noted that

$$\begin{bmatrix} \hat{\beta}_0^{\setminus i} \\ \hat{\boldsymbol{\beta}}^{\setminus i} \end{bmatrix} = \mathbf{prox}_R(\mathbf{z}), \quad \mathbf{z} = \begin{bmatrix} \hat{\beta}_0^{\setminus i} \\ \hat{\boldsymbol{\beta}}^{\setminus i} \end{bmatrix} - \sum_{j \neq i} \begin{bmatrix} 1 \\ \mathbf{x}_j \end{bmatrix} \dot{\ell} \left(\hat{\beta}_0^{\setminus i} + \mathbf{x}_j^\top \hat{\boldsymbol{\beta}}^{\setminus i}; y_j \right).$$

Hence, from the first-order condition $\sum_{j \neq i} \dot{\ell}(\hat{\beta}_0^{\setminus i} + \mathbf{x}_j^\top \hat{\beta}^{\setminus i}; y_j) = 0$, we can derive that

$$\mathbf{J}_{E,E} = [\mathbf{J}(\mathbf{u})]_{E,E} = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & (1-\alpha)\lambda & 0 & \dots & 0 \\ 0 & 0 & (1-\alpha)\lambda & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & (1-\alpha)\lambda \end{bmatrix}^{-1}. \quad (7)$$

The ALO formula is then immediate by Thm. 5.1.

3.2 ALO for LASSO, with Intercept through Generalized LASSO

For the generalized LASSO:

$$\min_{\beta} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\mathbf{D}\beta\|_1, \quad (8)$$

the dual problem can be derived as:

$$\min_{\mathbf{u}} \frac{1}{2} \|\mathbf{y} - \boldsymbol{\theta}\|_2^2, \quad \boldsymbol{\theta} \in \{\mathbf{X}^\top \boldsymbol{\theta} = \mathbf{D}^\top \mathbf{u}, \|\mathbf{u}\|_\infty \leq \lambda\}. \quad (9)$$

The dual problem could be written in a proximal approach, such that:

$$\hat{\mathbf{u}} = \mathbf{prox}_R(\mathbf{y}), \quad R(\mathbf{u}) = \begin{cases} 0 & \boldsymbol{\theta} \in \{\mathbf{X}^\top \boldsymbol{\theta} = \mathbf{D}^\top \mathbf{u}, \|\mathbf{u}\|_\infty \leq \lambda\}, \\ \infty & \text{otherwise.} \end{cases}$$

Denote \mathbf{J} as the Jacobian of the proximal operator at the full data problem \mathbf{y} , then the ALO estimator could be obtained as:

$$\mathbf{y}^{\setminus i} = \mathbf{y}_i - \frac{\hat{\mathbf{u}}_i}{J_{ii}}. \quad (10)$$

For the case of LASSO with an intercept, we could expand the \mathbf{X} with a column of ones in the first column, expand β with another dimension and choose $\mathbf{D} = [\mathbf{0}, \mathbf{I}]$. Let $E := \{j : |\mathbf{X}_j^\top \boldsymbol{\theta}| = \lambda\}$ denote the active set. The Jacobian is locally given as the projection onto the orthogonal complement of the span of \mathbf{X}_E and the vector of ones. Further denote $\tilde{\mathbf{X}}_E = [\mathbf{1}, \mathbf{X}_E]$, then the Jacobian is given as $\mathbf{I} - \tilde{\mathbf{X}}_E(\tilde{\mathbf{X}}_E^\top \tilde{\mathbf{X}}_E)^{-1} \tilde{\mathbf{X}}_E^\top$.

3.3 ALO for Elastic Net, without Penalty on Intercept through Generalized LASSO

Without penalty on intercept, the elastic net problem can be written as:

$$\begin{aligned} \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta} \end{bmatrix} &= \arg \min \frac{1}{2} \|\mathbf{y} - \beta_0 - \mathbf{X}\beta\|_2^2 + \lambda_1 \|\beta\|_1 + \lambda_2 \|\beta\|_2^2 \\ &= \arg \min \frac{1}{2} \begin{bmatrix} \beta_0 \\ \beta \end{bmatrix}^\top \left(\begin{bmatrix} 1 & \mathbf{X} \end{bmatrix}^\top \begin{bmatrix} 1 & \mathbf{X} \end{bmatrix} + \lambda_2 \text{diag}(0; \mathbf{1}_p) \right) \begin{bmatrix} \beta_0 \\ \beta \end{bmatrix} - \mathbf{y}^\top \begin{bmatrix} 1 & \mathbf{X} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta \end{bmatrix} + \lambda_1 \|\beta\|_1 \end{aligned}$$

where we assume that the size of \mathbf{X} is $n \times p$. In the mean time, note the LASSO problem (also without penalty on intercept) is:

$$\begin{aligned} \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta} \end{bmatrix} &= \arg \min \frac{1}{2} \|\mathbf{y} - \beta_0 - \mathbf{X}\beta\|_2^2 + \lambda_1 \|\beta\|_1 \\ &= \arg \min \frac{1}{2} \begin{bmatrix} \beta_0 \\ \beta \end{bmatrix}^\top \begin{bmatrix} \mathbf{1} & \mathbf{X} \end{bmatrix}^\top \begin{bmatrix} \mathbf{1} & \mathbf{X} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta \end{bmatrix} - \mathbf{y}^\top \begin{bmatrix} \mathbf{1} & \mathbf{X} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta \end{bmatrix} + \lambda_1 \|\beta\|_1 \end{aligned}$$

Thus we can add some ‘‘observations’’ to the data and let

$$\mathbf{y}^* = \begin{bmatrix} \mathbf{y} \\ \mathbf{0}_p \end{bmatrix}, \quad \mathbf{X}^* = \begin{bmatrix} \mathbf{X} \\ \sqrt{\lambda_2} \mathbf{I}_p \end{bmatrix},$$

then the elastic net becomes

$$\begin{aligned} \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta} \end{bmatrix} &= \arg \min \frac{1}{2} \left\| \mathbf{y}^* - \beta_0 \begin{bmatrix} \mathbf{1}_n \\ \mathbf{0}_p \end{bmatrix} - \mathbf{X}^* \beta \right\|_2^2 + \lambda_1 \|\beta\|_1 \\ &= \arg \min \frac{1}{2} \left\| \begin{bmatrix} \mathbf{y} \\ \mathbf{0}_p \end{bmatrix} - \begin{bmatrix} \mathbf{1}_n & \mathbf{X} \\ \mathbf{0}_p & \sqrt{\lambda_2} \mathbf{I}_p \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta \end{bmatrix} \right\|_2^2 + \lambda_1 \|\beta\|_1, \end{aligned} \tag{11}$$

which is a special case of the general LASSO.

4 Derivation for Multinomial ALO

4.1 Loss Function

Assume we have K classes, n observations \mathbf{X} , and R^p parameters β_k for each class. Define leave- i -out variables as

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

The loss function would be

$$\ell(\mathcal{B}) = - \left\{ \sum_{j \neq i} \left[\sum_{k=1}^K y_{jk} \mathbf{X}_j^\top \beta_k - \log \left(\sum_{k=1}^K e^{\mathbf{X}_j^\top \beta_k} \right) \right] \right\} = \sum_{j \neq i} \log \left(\sum_{k=1}^K e^{\mathbf{X}_j^\top \beta_k} \right) - \sum_{j \neq i} \sum_{k=1}^K y_{jk} \mathbf{X}_j^\top \beta_k$$

By taking the first order derivative, we get

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

Similarly, we can get

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

where,

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

So, we have

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

where,

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

4.2 Newton's Method

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

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

where,

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

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

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

4.3 Approximate Leave- i -Out Prediction

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

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