

Notes on Approximate Leave-One-Out for Elastic Net

Linyun He

Wanchao Qin

Peng Xu

Yuze Zhou

August 28, 2018

1 ALO for Elastic Net, Approximation in the Primal Domain

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

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

2 ALO for Elastic Net, Approximation in the Dual Domain

The original problem for elastic net is to solve for $\hat{\beta}$ such that:

$$\hat{\beta} = \arg \min_{\beta} \left(\frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda_1 \|\beta\|_1 + \lambda_2 \|\beta\|_2^2 \right) \quad (5)$$

By adding the Lagrangian, we get the formulation of L :

$$L = \frac{1}{2} \|\mathbf{y} - \mathbf{z}\|_2^2 + \lambda_1 \|\beta\|_1 + \lambda_2 \|\beta\|_2^2 + \mathbf{u}^\top (\mathbf{z} - \mathbf{X}\beta). \quad (6)$$

The original problem is solving the primal of the Lagrangian such that $p^* = \min_{\beta, \mathbf{z}} \max_{\mathbf{u}} L$ and the dual formulation $d^* = \max_{\mathbf{u}} \min_{\beta, \mathbf{z}} L$, to minimize over \mathbf{z} :

$$\frac{\partial L}{\partial \mathbf{z}} = \mathbf{z} - \mathbf{y} + \mathbf{u} = \mathbf{0} \implies \mathbf{y} = \mathbf{u} + \mathbf{z}.$$

Since β is penalized element-wisely, we can minimize over β by minimizing over each β_i , that is, we have to minimize $\lambda_1 |\beta_i| + \lambda_2 \beta_i^2 - \mathbf{u}^\top \mathbf{X}_i \beta$ for each dimension of β , where \mathbf{X}_i denotes the i th column of \mathbf{X} , therefore:

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

By taking all the above to the Lagrangian, we obtain the dual problem d^* as:

$$d^* = \min_{\mathbf{u}} \frac{1}{2} \|\mathbf{y} - \mathbf{u}\|_2^2 + \sum_{j: |\mathbf{X}_j^\top \mathbf{u}| > \lambda_1} \frac{(\lambda_1 - |\mathbf{u}^\top \mathbf{X}_j|)^2}{4\lambda_2}. \quad (7)$$

The minimizer $\hat{\mathbf{u}}$ could also be obtained from the dual problem through a proximal approach:

$$\hat{\mathbf{u}} = \mathbf{prox}_R(\mathbf{y}), \quad R(\mathbf{u}) = \sum_{j: |\mathbf{X}_j^\top \mathbf{u}| > \lambda_1} \frac{(\lambda_1 - |\mathbf{u}^\top \mathbf{X}_j|)^2}{4\lambda_2}.$$

By replacing the full data problem \mathbf{y} with $\mathbf{y}_\alpha = \mathbf{y} + (y_i^{\setminus i} - y_i)e_i$, where $y_i^{\setminus i}$ is the true LOO estimator and e_i is the i -th standard vector, and let $\mathbf{u}^{\setminus i} = \mathbf{prox}_R(\mathbf{y}_\alpha)$, we have:

$$\begin{aligned} 0 &= e_i^\top \mathbf{u}^{\setminus i} \\ &= e_i^\top \mathbf{prox}_R(\mathbf{y}_\alpha) \\ &\approx e_i^\top [\mathbf{prox}_R(\mathbf{y}) + \mathbf{J}_R(\mathbf{y})(\mathbf{y}_\alpha - \mathbf{y})] \\ &\approx \hat{u}_i + \mathbf{J}_{ii}(y_i^{\setminus i} - y_i). \end{aligned}$$

Here $\mathbf{J}_R(\mathbf{y})$ denotes the Jacobian matrix of the proximal operator at \mathbf{y} , thus the ALO estimator \tilde{y}_i is obtained as

$$\tilde{y}_i = y_i - \frac{\hat{u}_i}{\mathbf{J}_{ii}}. \quad (8)$$

The Jacobian could locally be obtained as:

$$\mathbf{J}_R(\mathbf{y}) = (\mathbf{I} + \nabla^2 R(\mathbf{prox}_R(\mathbf{y})))^{-1} = (\mathbf{I} + \nabla^2 R(\hat{\mathbf{u}}))^{-1} = \left(\mathbf{I} + \frac{1}{2\lambda_2} \mathbf{X}_E \mathbf{X}_E^\top \right)^{-1} \quad (9)$$

for $E = \{j : |\mathbf{X}_j^\top \mathbf{u}| > \lambda_1\}$.

3 ALO for Elastic Net, Approximation with Proximal Formulation

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

$$\mathbf{prox}_R(\mathbf{z}) = \gamma \operatorname{sgn}(\mathbf{z}) \odot (|\mathbf{z}| - \lambda \mathbf{1}_p)_+, \quad \gamma = \frac{1}{1 + (1 - \alpha)\lambda}. \quad (10)$$

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

$$\mathbf{H} = \gamma \mathbf{X}_{\cdot, E} \left[\gamma \mathbf{X}_{\cdot, E}^\top \mathbf{X}_{\cdot, E} + (1 - \gamma) \mathbf{I}_{E, E} \right]^{-1} \mathbf{X}_{\cdot, E}^\top. \quad (11)$$

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

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

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

$$\mathbf{J}_{E, E} = [\mathbf{J}(\mathbf{u})]_{E, E} = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & (1 - \alpha)\lambda & 0 & \dots & 0 \\ 0 & 0 & (1 - \alpha)\lambda & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & (1 - \alpha)\lambda \end{bmatrix}^{-1}. \quad (12)$$

The ALO formula is then immediate by Thm. 5.1.

4 ALO for LASSO, with Intercept through Generalized LASSO

For the generalized LASSO:

$$\min_{\boldsymbol{\beta}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda \|\mathbf{D}\boldsymbol{\beta}\|_1, \quad (13)$$

the dual problem can be derived as:

$$\min_{\boldsymbol{\theta}} \frac{1}{2} \|\mathbf{y} - \boldsymbol{\theta}\|_2^2, \quad \boldsymbol{\theta} \in \{\mathbf{X}^\top \boldsymbol{\theta} = \mathbf{D}^\top \mathbf{u}, \|\mathbf{u}\|_\infty \leq \lambda\}. \quad (14)$$

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:

$$\mathbf{y}^{\setminus i} = \mathbf{y}_i - \frac{\hat{\mathbf{u}}_i}{\mathbf{J}_{ii}}. \quad (15)$$

For the case of LASSO with an intercept, we could expand the \mathbf{X} with a column of ones in the first column, expand $\boldsymbol{\beta}$ 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 $\mathbf{I} - \tilde{\mathbf{X}}_E(\tilde{\mathbf{X}}_E^\top \tilde{\mathbf{X}}_E)^{-1} \tilde{\mathbf{X}}_E^\top$.

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

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

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

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:

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

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

$$\begin{aligned} \begin{bmatrix} \hat{\beta}_0 \\ \hat{\boldsymbol{\beta}} \end{bmatrix} &= \arg \min \frac{1}{2} \left\| \mathbf{y}^* - \beta_0 \begin{bmatrix} \mathbf{1}_n \\ \mathbf{0}_p \end{bmatrix} - \mathbf{X}^* \boldsymbol{\beta} \right\|_2^2 + \lambda_1 \|\boldsymbol{\beta}\|_1 \\ &= \arg \min \frac{1}{2} \left\| \begin{bmatrix} \mathbf{y} \\ \mathbf{0}_p \end{bmatrix} - \begin{bmatrix} \mathbf{1}_n & \mathbf{X} \\ \mathbf{0}_p & \sqrt{\lambda_2} \mathbf{I}_p \end{bmatrix} \begin{bmatrix} \beta_0 \\ \boldsymbol{\beta} \end{bmatrix} \right\|_2^2 + \lambda_1 \|\boldsymbol{\beta}\|_1, \end{aligned} \quad (16)$$

which is a special case of the general LASSO.

6 Generalized Loss function

6.1 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

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

where

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

Adding in intercept terms is then straightforward.

6.2 ALO for Logistic Regression, Approximation in the Dual Domain

6.2.1 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:

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

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

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

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 $\mathbf{y}_u = \hat{\theta}_i \ddot{\ell}^*(-\hat{\theta}_i; y_i) + \hat{\theta}_i / \sqrt{\ddot{\ell}^*(-\hat{\theta}_i; y_i)}$ could locally be treated as the projection onto the orthogonal complement of the polyhedron $\{\|X^\top K u\|_\infty \leq \lambda_1\}$, thus $J = I - X_{u,E} (X_{u,E}^\top X_{u,E})^{-1} X_{u,E}$, where $X_{u,E}$ are the columns of $X_u = X^\top K$, such that the columns in the set $E = \{i : |X_i^\top \theta| = \lambda_1\}$ 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.

6.2.2 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 $J = (I + X_{u,E} X_{u,E}^\top / 2\lambda_2)$, where $X_{u,E}$ are the columns of $X_u = X^\top K$, such that the columns in the set $E = \{i : |X_i^\top \theta| = \lambda_1\}$ 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.

6.3 ALO for Poisson Regression, Approximation in the Primal Domain

For Poisson regression, the primal problem is:

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

Let A be the active set, we have

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

Therefore, the ALO formula is

$$x_i^\top \tilde{\beta}^{/j} \approx x_i^\top \hat{\beta} + \frac{H_{ii} (e^{x_j^\top \hat{\beta}} - y_j)}{1 - H_{ii} e^{x_j^\top \hat{\beta}}},$$

where

$$H = X_{\cdot,A} \left[X_{\cdot,A}^\top \text{diag} \left(e^{x_j^\top \hat{\beta}} \right) X_{\cdot,A} + (1 - \alpha) \lambda I_{A,A} \right]^{-1} X_{\cdot,A}^\top.$$

6.4 ALO for Poisson Regression, Approximation in the Dual Domain

6.4.1 ALO for Poisson Regression with Lasso penalty

The optimization function for Poisson regression with lasso penalty is:

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

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 $\ell(x_i^\top \beta; y_i) = -y_i x_i^\top \beta + e^{x_i^\top \beta} + \ln(y_i!)$. The optimal solution for the dual problem $\hat{\theta} = y - e^{\mathbf{X}\hat{\beta}}$ 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.

6.4.2 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.

7 Usage of ALO formulae with glmnet package

The `glmnet` package scales the elastic net loss function by a factor of $1/n$. Furthermore, for linear problems `glmnet` implicitly “standardizes y to have unit variance before computing its λ sequence (and then unstandardizes the resulting coefficients)” (cf. [Glmnet Vignette]). So to get comparable results, it is necessary to rescale y by the MLE $\hat{\sigma}_y$ before fitting the model. More precisely, `glmnet` is in fact optimizing the following problem:

$$\min_{\beta} \frac{1}{2n} \sum_{j=1}^n \left(\frac{x_j^\top \beta}{\hat{\sigma}_y} - \frac{y_j}{\hat{\sigma}_y} \right)^2 + \frac{\lambda}{\hat{\sigma}_y} \alpha \|\beta\|_1 + \frac{\lambda}{\hat{\sigma}_y^2} \frac{1-\alpha}{2} \|\beta\|_2^2. \quad (17)$$

We thus have

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

Hence, for the linear elastic net problem, the primal ALO is:

$$\hat{y}_j^{\setminus i} = \hat{y}_j + \frac{\mathbf{H}_{ii}(\hat{y}_j - y_j)}{n \hat{\sigma}_y - \mathbf{H}_{ii}}, \quad \mathbf{H} = \mathbf{X}_{\cdot,A} \left[\frac{1}{n \hat{\sigma}_y} \mathbf{X}_{\cdot,A}^\top \mathbf{X}_{\cdot,A} + \frac{(1-\alpha)\lambda}{\hat{\sigma}_y^2} \mathbf{I}_{A,A} \right]^{-1} \mathbf{X}_{\cdot,A}^\top.$$

We get further complication when option `standardization = T` is given, in which case `glmnet` first standardize the data \mathbf{X} using $\hat{\sigma}_{\mathbf{X}}$ (assuming \mathbf{X} is standardized):

- If `intercept = F`, compute $\mathbf{X}^* = \text{diag}[\hat{\sigma}_y \hat{\sigma}_{\mathbf{X}}]^{-1} \mathbf{X}$.
- If `intercept = T`, compute $\mathbf{X}^* = \text{diag}[\hat{\sigma}_y \hat{\sigma}_{\mathbf{X}}]^{-1} (\mathbf{X} - \bar{\mathbf{X}} \mathbf{1} \mathbf{1}^\top)$.

Afterwards, the the coefficients are returned de-standardized:

- Let (β_0, β) be the original intercept and coefficients, compute $\beta^* = \hat{\sigma}_y \text{diag}[\hat{\sigma}_{\mathbf{X}}]^{-1} \beta$ then $\beta_0^* = \beta_0 - \bar{\mathbf{X}} \beta^*$.
- Return pair (β_0^*, β^*) .

For logistics and Poisson regression the standardization procedure is basically the same, except `glmnet` no longer standardize by $\hat{\sigma}_y$, which make sense since \mathbf{y} is now either categorical or count data.

8 Benchmark

n	p	k	Average ALO	Average 5-fold CV	Relative	n full fit
300	100	60	0.016	0.053	3.313	1.2
500	800	500	0.251	0.533	2.124	36.5
1000	1200	800	0.489	1.200	2.454	211.0
2500	2000	1200	2.267	3.623	1.598	1577.5
5000	2500	2000	5.017	8.097	1.614	7740.0
10000	10000	2500	27.236	36.520	1.341	62530.0

Table 1: Average elapsed time (in seconds) comparison, 10 runs, $\alpha = 0.5$.