# Approximate Leave-One-Out with `glmnet`

Linyun He        Wanchao Qin        Peng Xu        Yuze Zhou

August 29, 2018

## 1 ALO for Linear Regression

Recall the objective function for the elastic net problem:

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

Let $E = \{i : \beta_i \notin K, i = 1, \ldots, p\}$ be the active set, the ALO formula is

$$x_i^\top \tilde{\beta}^{\backslash j} \approx x_i^\top \hat{\beta} + \frac{H_{ii} \left( x_j^\top \hat{\beta} - y_j \right)}{1 - H_{ii}}, \qquad H = X_{\cdot,E} \left[ X_{\cdot,E}^\top X_{\cdot,E} + (1-\alpha)\lambda I_{E,E} \right]^{-1} X_{\cdot,E}^\top.$$

## 2 ALO for Logistic Regression

For binomial logistic regression, the primal problem is:

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

Let $E$ be the active set, the ALO formula is

$$x_i^\top \tilde{\beta}^{\backslash j} \approx x_i^\top \hat{\beta} + \frac{H_{ii} \left( 1 + e^{x_j^\top \hat{\beta}} \right) \left[ e^{x_j^\top \hat{\beta}} - y_j \left( 1 + e^{x_j^\top \hat{\beta}} \right) \right]}{\left( 1 + e^{x_j^\top \hat{\beta}} \right)^2 - H_{ii} e^{x_j^\top \hat{\beta}}},$$

where

$$H = X_{\cdot,E} \left[ X_{\cdot,E}^\top \operatorname{diag} \left( \frac{e^{x_j^\top \hat{\beta}}}{1 + 2e^{x_j^\top \hat{\beta}} + e^{2x_j^\top \hat{\beta}}} \right) X_{\cdot,E} + (1-\alpha)\lambda I_{A,A} \right]^{-1} X_{\cdot,E}^\top.$$

## 3 ALO for Poisson Regression

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 $E$ be the active set, the ALO formula is

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

where

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

## 4 ALO for Multinomial Regression

Assume that the response variable comes in as an $n \times K$ matrix indicator matrix, where $K$ is the number of classes. We re-parametrize by considering $\mathcal{B} = \operatorname{vec}(B)$:

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_K \end{bmatrix} = \begin{bmatrix} \begin{bmatrix} y_{11} \\ y_{12} \\ \vdots \\ y_{1K} \end{bmatrix} \\ \begin{bmatrix} y_{21} \\ y_{22} \\ \vdots \\ y_{2K} \end{bmatrix} \\ \vdots \\ \begin{bmatrix} y_{n1} \\ y_{n2} \\ \vdots \\ y_{nK} \end{bmatrix} \end{bmatrix}, \qquad \mathcal{X} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_K \end{bmatrix} = \begin{bmatrix} \begin{bmatrix} x_1^\top & 0 & \cdots & 0 \\ 0 & x_1^\top & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & x_1^\top \end{bmatrix} \\ \begin{bmatrix} x_2^\top & 0 & \cdots & 0 \\ 0 & x_2^\top & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & x_2^\top \end{bmatrix} \\ \vdots \\ \begin{bmatrix} x_n^\top & 0 & \cdots & 0 \\ 0 & x_n^\top & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & x_n^\top \end{bmatrix} \end{bmatrix}, \qquad \mathcal{B} = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_K \end{bmatrix}.$$

Again let $E$ denote the active set. Further, let

$$
\mathcal{A}(\mathcal{B}) := \begin{bmatrix} A_1(\mathcal{B}) \\ A_2(\mathcal{B}) \\ \vdots \\ A_n(\mathcal{B}) \end{bmatrix} = \begin{bmatrix} \begin{bmatrix} \frac{\exp(X_1^\top \beta_1)}{\sum_{k=1}^{K} \exp(X_1^\top \beta_k)} \\ \vdots \\ \frac{\exp(X_1^\top \beta_K)}{\sum_{k=1}^{K} \exp(X_1^\top \beta_k)} \end{bmatrix} \\ \begin{bmatrix} \frac{\exp(X_2^\top \beta_1)}{\sum_{k=1}^{K} \exp(X_2^\top \beta_k)} \\ \vdots \\ \frac{\exp(X_2^\top \beta_K)}{\sum_{k=1}^{K} \exp(X_2^\top \beta_k)} \end{bmatrix} \\ \vdots \\ \begin{bmatrix} \frac{\exp(X_n^\top \beta_1)}{\sum_{k=1}^{K} \exp(X_n^\top \beta_k)} \\ \vdots \\ \frac{\exp(X_n^\top \beta_K)}{\sum_{k=1}^{K} \exp(X_n^\top \beta_k)} \end{bmatrix} \end{bmatrix},
$$

and

$$
\mathcal{D}(\mathcal{B}) := \begin{bmatrix} \left[ \mathrm{diag}\left(A_1(\mathcal{B})\right) - A_1(\mathcal{B})A_1(\mathcal{B})^\top \right] & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \left[ \mathrm{diag}\left(A_n(\mathcal{B})\right) - A_n(\mathcal{B})A_n(\mathcal{B})^\top \right] \end{bmatrix}.
$$

Finally, define

$$
\mathcal{K}(X, \mathcal{B}) := X^\top \mathcal{D}(\mathcal{B}) X + \nabla^2 R(\mathcal{B}), \qquad \mathcal{G}_{i,E}(X, \mathcal{B}) := X_{i,E} \mathcal{K}(X_{\cdot,E}, \hat{\mathcal{B}})^+ X_{i,E}^\top.
$$

Then, with Newton's method, we can approximate the leave-$i$-out prediction as

$$
X_i \tilde{\mathcal{B}}^{\backslash i} = X_i \hat{\mathcal{B}} + \mathcal{G}_{i,E}(X, \mathcal{B}) \left( A_i(\hat{\mathcal{B}}) - y_i \right)
$$

$$
- \mathcal{G}_{i,E}(X, \hat{\mathcal{B}}) \left\{ \mathcal{G}_{i,E}(X, \hat{\mathcal{B}}) - \left[ \mathrm{diag}\left(A_i(\hat{\mathcal{B}})\right) - A_i(\hat{\mathcal{B}})A_i(\hat{\mathcal{B}})^\top \right]^+ \right\}^+ \mathcal{G}_{i,E}(X, \hat{\mathcal{B}}) \left( A_i(\hat{\mathcal{B}}) - y_i \right)
$$

# 5   ALO with Intercept

Including the intercept is straightforward. As we can augment $X$ with an extra column of 1s, i.e. $X^* = [1_n, X]$. Since the intercept is not reugularized, we need to change the corresponding second partial derivatives to 0, e.g.

$$
\nabla^2 R\left(\hat{\beta}_0, \hat{\beta}_A\right) = \begin{bmatrix} 0 & 0 & \cdots & 0 \\ 0 & (1-\alpha)\lambda & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & (1-\alpha)\lambda \end{bmatrix}.
$$

For multinomial model it can be a bit more complicated since there are now $K$ intercepts. For programming convenience we augment $\boldsymbol{X}$ by block of $\boldsymbol{I}_K$ and stack the intercepts on tops of $\mathcal{B}$, i.e.

$$
\boldsymbol{X}^* = \begin{bmatrix} \begin{bmatrix} \boldsymbol{I}_K & \boldsymbol{X}_1 \\ \boldsymbol{I}_K & \boldsymbol{X}_2 \end{bmatrix} \\ \vdots \\ \begin{bmatrix} \boldsymbol{I}_K & \boldsymbol{X}_K \end{bmatrix} \end{bmatrix}, \qquad \mathcal{B}^* = \begin{bmatrix} \beta_{01} \\ \vdots \\ \beta_{0K} \\ \beta_1 \\ \vdots \\ \beta_K \end{bmatrix}.
$$

Accordingly, the first $K$ diagonal elements of $\nabla^2 R$ will be set to 0.

# 6 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 $\lambda$ sequence (and then unstandardizes the resulting coefficients)". So it is necessary to rescale $y$ by the MLE $\hat{\sigma}_y$ before fitting the model. For instance, `glmnet` is in fact optimizing the following problem for linear regression (assuming $\boldsymbol{X}$ is already standardized):

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

We thus have

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

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

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

Further complications present when option `standardization = T` is given, in which case `glmnet` first standardize the data $\boldsymbol{X}$ using $\hat{\sigma}_{\boldsymbol{X}}$:

- If `intercept = F`, compute $\boldsymbol{X}^* = \text{diag}[\hat{\sigma}_y \hat{\sigma}_{\boldsymbol{X}}]^{-1} \boldsymbol{X}$.

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

Afterwards, the the coefficients are returned unstandardized, i.e. let $(\beta_0, \beta)$ denotes the original intercept and coefficients, `glmnet` reports

$$
\beta^* = \hat{\sigma}_y \text{diag}[\hat{\sigma}_{\boldsymbol{X}}]^{-1} \beta, \qquad \beta_0^* = \beta_0 - \bar{\boldsymbol{X}} \beta^*.
$$

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 $y$ is now either categorical or count data.

# 7  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: Averaged (over 10 runs) elapsed time (in seconds) comparison, 25 $\lambda$s, $\alpha = 0.5$.

# References

[1]  Trevor Hastie & Junyang Qian, *Glmnet Vignette*.