

Uniformly Valid Inference Based on the Lasso in Linear Mixed Models

joint work with U. Schneider and T. Krivobokova

Peter Kramlinger

Lima, SAE 2024

Model Selection and Inference

General Linear Model

$$\mathbf{y} = \mathbf{X}\beta_0 + \epsilon, \quad \epsilon \sim \mathcal{N}_n\{\mathbf{0}_n, \mathbf{V}(\theta_0)\}$$

for $\beta_0 \in \mathbb{R}^p$, $\theta_0 \in \Theta \subset \mathbb{R}^r$ unknown, $n, p, r \in \mathbb{N}$, $n > p$

Goal: Inference on β_0 after model selection.

- ▶ Model selection is data driven
- ▶ Selection and estimation are two sources of uncertainty
- ▶ Naive inference after model selection is not valid

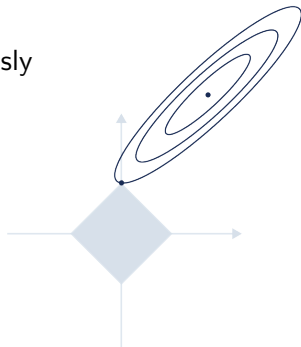
Least Absolute Shrinkage and Selection Operator (Lasso)

for $\lambda_j > 0, j = 1, \dots, p$ and REML estimator $\hat{\theta}$ for θ_0 ,

$$\hat{\beta}_L = \operatorname{argmin}_{\beta \in \mathbb{R}^p} \left\{ \ell(\beta, \hat{\theta}) + 2 \sum_{j=1}^p \lambda_j |\beta_j| \right\}$$

- Selection and estimation simultaneously
- Lasso not uniformly consistent

Pötscher and Leeb [2009]



Related Work

In mixed models:

- ▶ Penalization of θ_0 : Bondell et al. [2010], Ibrahim et al. [2011]
- ▶ Penalization of the random effects: Peng and Lu [2012]
- ▶ Many instances of penalized β_0 in mixed models.

In linear models:

- ▶ Uniformly valid inference based on the Lasso if $n > p$
Ewald and Schneider [2018]

$$\inf_{\beta_0} P_{\beta_0}(\beta_0 \in S) = \min_{\mathbf{d} \in \{-1,1\}^p} P(\beta_0 \in S_{\mathbf{d}}) = 1 - \alpha$$

Lasso

for $\lambda_j > 0, j = 1, \dots, p$ and REML estimator $\hat{\theta}$ for θ_0 ,

$$\hat{\beta}_L = \operatorname{argmin}_{\beta \in \mathbb{R}^p} \left\{ \ell(\beta, \hat{\theta}) + 2 \sum_{j=1}^p \lambda_j |\beta_j| \right\}$$

Goal: Uniformly valid inference on β_0 based on the Lasso.

- Find $S \subset \mathbb{R}^p$ such that

$$\inf_{\beta_0, \theta_0} P_{\beta_0, \theta_0}(\hat{\beta}_L \in S) \approx 1 - \alpha$$

- Idea: Separate estimation for θ_0 and β_0
- Requires that $\hat{\theta}$ is uniformly bounded

Restricted Maximum Likelihood (REML) Methodology

Estimate θ_0 from $\mathbf{A}^t \mathbf{y}$ instead of \mathbf{y} , for $\mathbf{A} \in \mathbb{R}^{n \times (n-p)}$ such that $\mathbf{A}^t \mathbf{X} = \mathbf{0}_{n-p}$. Then,

$$\hat{\theta} = \operatorname{argmax}_{\theta \in \Theta} \ell_R(\theta),$$

where $\mathbf{P}(\theta) = \mathbf{A} \{ \mathbf{A}^t \mathbf{V}(\theta) \mathbf{A} \}^{-1} \mathbf{A}^t$ and

$$\ell_R(\theta) = \operatorname{argmax}_{\theta \in \Theta} -\frac{1}{2} \ln |\mathbf{V}(\theta)| - \frac{1}{2} \ln |\mathbf{X}^t \mathbf{V}(\theta)^{-1} \mathbf{X}| - \frac{1}{2} \mathbf{y}^t \mathbf{P}(\theta) \mathbf{y}$$

- ▶ Requires low-dimensional setting $n > p$
- ▶ $\hat{\theta}$ is independent of β_0
- ▶ *Uniform* consistency has to be proved
- ▶ Maximum not unique, hence consider *Cramér* consistency

Theorem 1: Cramér Consistency of REML Estimator

Under some regularity conditions, set

$$\nu_k(\boldsymbol{\theta}_0) = \frac{\text{tr} \left\{ \mathbf{V}(\boldsymbol{\theta}_0)^{-1} \frac{\partial \mathbf{V}}{\partial \theta_k} \right\}}{\sqrt{\text{rk} \left(\frac{\partial \mathbf{V}}{\partial \theta_k} \right)}}, \quad k = 1, \dots, r.$$

There exists a sequence $\hat{\boldsymbol{\theta}}$ of local maximizers of ℓ_R , such that

$$\nu_k(\boldsymbol{\theta}_0) \left| \hat{\theta}_k - \theta_{0,k} \right| = O_P(1)$$

uniformly over $\boldsymbol{\theta}_0 \in \Theta$.

- ▶ If $\boldsymbol{\theta}_0$ were fixed, $\nu_k(\boldsymbol{\theta}_0) \equiv \sqrt{n}$
- ▶ Note that $-\text{E}_{\boldsymbol{\theta}_0} \{ \partial^2 \ell_R(\boldsymbol{\theta}) / \partial \theta_k^2 |_{\boldsymbol{\theta}_0} \} = O\{\nu_k(\boldsymbol{\theta}_0)^2\}$

Theorem 2: Uniformly Valid Confidence Sets

There exists a sequence $\hat{\boldsymbol{\theta}}$ of local maximizers of ℓ_R such that for $E(\hat{\mathbf{C}}, \hat{\tau}) = \{\mathbf{z} \in \mathbb{R}^p \mid \mathbf{z}^t \hat{\mathbf{C}} \mathbf{z} \leq \hat{\tau}\}$, with $\hat{\mathbf{C}} = n^{-1} \mathbf{X}^t \mathbf{V}(\hat{\boldsymbol{\theta}})^{-1} \mathbf{X}$ and

$$\hat{\tau} = \max_{\mathbf{d} \in \{-1, 1\}^p} \chi_{p, 1-\alpha}^2 \left\{ n^{-1} \left\| \hat{\mathbf{C}}^{-1/2} \text{diag}(\lambda_1, \dots, \lambda_p) \mathbf{d} \right\|^2 \right\},$$

a quantile of a non-central χ_p^2 -distribution, it holds that

$$\inf_{\substack{\beta_0 \in \mathbb{R}^p \\ \theta_0 \in \Theta}} P_{\beta_0, \theta_0} \left\{ \sqrt{n} \left(\hat{\beta}_L - \beta_0 \right) \in E \left(\hat{\mathbf{C}}, \hat{\tau} \right) \right\} = 1 - \alpha + O \left(\frac{1}{\sqrt{n}} \right).$$

► Idea:

$$\inf_{\beta_0, \theta_0} P_{\beta_0, \theta_0} (\beta_0 \in S) = \inf_{\theta_0} \min_{\mathbf{d} \in \{-1, 1\}^p} P_{\beta_0, \theta_0} (\beta_0 \in S_{\mathbf{d}})$$

► Confidence set: $M_L = \left\{ \beta \in \mathbb{R}^p : n \left\| \hat{\mathbf{C}}^{1/2} \left(\hat{\beta}_L - \beta \right) \right\|^2 \leq \hat{\tau} \right\}$

Simulation Design

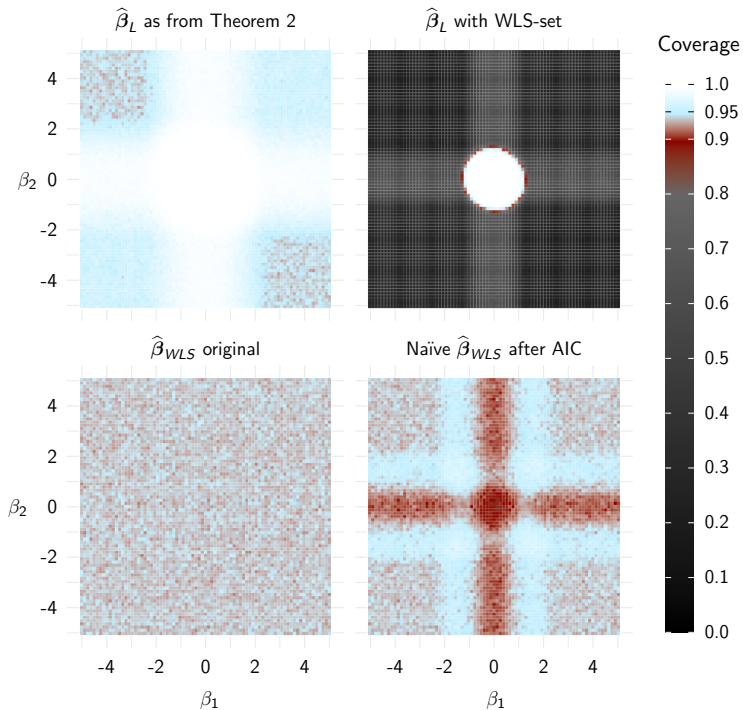
Random Intercept Model

For $i = 1, \dots, m$ and $j = 1, \dots, n_i$, let

$$y_{ij} = \mathbf{x}_{ij}^t \boldsymbol{\beta}_0 + v_i + u_{ij}, \quad u_{ij} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma_u^2), \quad v_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma_v^2)$$

and $m = 20$, $n = 400$, $n_i = 20$, $\sigma_u = \sigma_v = 4$, $\lambda_i = n^{1/2}/2$, $p = 2$.

- ▶ For each $\boldsymbol{\beta}_0 \in [-4, 4]^2$, perform 3,000 replications
- ▶ For each replications, check if $\boldsymbol{\beta}_0 \in M_L$



Real Data Example

Acidity in US Lakes Opsomer et al. [2009]

$$\mathbf{y} = \mathbf{1}_n \beta_1 + \mathbf{x} \beta_2 + \mathbf{Z} \mathbf{v} + \mathbf{D} \mathbf{u} + \boldsymbol{\epsilon},$$

$$\mathbf{v} \sim \mathcal{N}(\mathbf{0}_m, \sigma_v^2 \mathbf{I}_m), \quad \mathbf{u} \sim \mathcal{N}(\mathbf{0}_K, \sigma_u^2 \mathbf{I}_K) \quad \text{and} \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}_n, \sigma_e^2 \mathbf{I}_n)$$

- ▶ WLS: $\hat{\boldsymbol{\beta}}_{WLS} = (483, -1.20)^t$; $sd(\hat{\boldsymbol{\beta}}_{WLS}) = (1718, 0.14)^t$
- ▶ Lasso: $\hat{\boldsymbol{\beta}}_L = (0, -1.12)^t$
- ▶ Lasso-based confidence set ~ 200 times larger, but

$$M_L(\beta_2) = [-1.49, -0.75] \quad M_{WLS}(\beta_2) = [-1.57, -0.83]$$

- H. Bondell, A. Krishna, and S. Ghosh. Joint Variable Selection for Fixed and Random Effects in Linear Mixed-Effects Models. *Biometrics*, 66:1069–1077, 2010.
- K. Ewald and U. Schneider. Uniformly Valid Confidence Sets Based on the Lasso. *Electronic Journal of Statistics*, 12:1358–1387, 2018.
- J. Ibrahim, H. Zhu, R. Garcia, and R. Guo. Fixed and Random Effects Selection in Mixed Effects Models. *Biometrics*, 67(2):1358–1387, 2011.
- P. Kramlinger, U. Schneider, and T. Krivobokova. Uniformly valid inference based on the lasso in linear mixed models. *Journal of Multivariate Analysis*, 198:105230, 2023.
- J. D. Opsomer, Claeskens, G., M. G. Ranalli, G. Kauermann, and F. J. Breidt. Non-parametric small area estimation using penalized spline regression. *JRSS B*, 70:265286, 2009.
- H. Peng and Y. Lu. Model Selection in Linear Mixed Models. *Journal of Multivariate Analysis*, 109:109–129, 2012.
- B. Pötscher and H. Leeb. On the distribution of penalized maximum likelihood estimators: The LASSO, SCAD, and thresholding. *Journal of Multivariate Analysis*, 100:2065–2082, 2009.
- L. Weiss. Asymptotic Properties of Maximum Likelihood Estimators in Some Nonstandard Cases. *Journal of the American Statistical Association*, 66: 345–350, 1971.

Conclusions

- ▶ REML estimator for θ_0 is uniformly Cramér consistent
- ▶ Inference on β_0 based on $\hat{\beta}_L$ uniformly valid over $\mathbb{R}^p \times \Theta$
- ▶ Resulting coverage probability depends on β_0
- ▶ If interest in inference only, do no selection
- ▶ Lasso-based confidence sets are large
- ▶ Valid inference (for $n > p$) if nuisance and target parameters can be separated

Regularity Conditions

(A) $\boldsymbol{\theta}_0 \in \Theta = \{\boldsymbol{\theta} \in \mathbb{R}_{>0}^r \mid \max(\boldsymbol{\theta})/\min(\boldsymbol{\theta}) \leq c\}; \infty > c \text{ const.}$

(B) $\mathbf{V}(\boldsymbol{\theta}_0) = \sum_{k=1}^r \theta_{0,k} \mathbf{H}_k; \mathbf{H}_k \geq 0, k = 1 \dots, r-1 \text{ and } \mathbf{H}_r > 0$

(C) $\text{rk}(\mathbf{X}) = p < n$

(D) For constants $0 < \underline{\omega} \leq \bar{\omega} < \infty$ it holds

$$\underline{\omega} \leq \frac{\text{tr} \{ \mathbf{V}(\mathbf{1}_r)^{-1} \mathbf{H}_k \}}{\text{rk}(\mathbf{H}_k)} \leq \bar{\omega},$$

where $\mathbf{1}_r = (1, \dots, 1)' \in \mathbb{R}^r$

(E) $m = \min\{\text{rk}(\mathbf{H}_1), \dots, \text{rk}(\mathbf{H}_r)\} \rightarrow \infty$