

Supplemental Appendix: Sparse Estimation and Uncertainty with Application to Subgroup Analysis

A Proof of Relative Efficiency of Oracle Estimator and OLS.

Proof: Denote as X_S the submatrix of X for which $\beta_k \neq 0$ and the Gram matrix for X as

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N X_i^\top X_i = \Sigma_X \quad (1)$$

and in block-partition form

$$\Sigma_X = \begin{pmatrix} \Sigma_{SS} & \Sigma_{SS^c} \\ \Sigma_{SS^c}^\top & \Sigma_{S^cS^c} \end{pmatrix} \quad (2)$$

Σ_X is invertible, since the least squares estimate exists and is unique. Since Σ_X is invertible, every square submatrix of Σ_X is also invertible.

The asymptotic relative efficiency of the least squares estimate and Oracle estimate is then

$$\lim_{N \rightarrow \infty} \frac{\frac{\sigma^2}{N} \text{Tr} \{ \Sigma^{-1} \}}{\frac{\sigma^2}{N} \text{Tr} \{ \Sigma_{SS}^{-1} \}} = \frac{\text{Tr} \{ \Sigma^{-1} \}}{\text{Tr} \{ \Sigma_{SS}^{-1} \}} \quad (3)$$

By the block inverse partition formula,

$$\text{Tr} (\Sigma_X^{-1}) = \text{Tr} \left\{ \begin{pmatrix} \Sigma_{SS} & \Sigma_{SS^c} \\ \Sigma_{SS^c}^\top & \Sigma_{S^cS^c} \end{pmatrix}^{-1} \right\} \quad (4)$$

$$= \text{Tr} \left\{ \left(\Sigma_{SS} - \Sigma_{SS^c} \Sigma_{S^cS^c}^{-1} \Sigma_{SS^c}^\top \right)^{-1} \right\} + \text{Tr} \left\{ \left(\Sigma_{S^cS^c} - \Sigma_{SS^c}^\top \Sigma_{SS}^{-1} \Sigma_{SS^c} \right)^{-1} \right\} \quad (5)$$

Consider the first summand inside the parentheses on the r.h.s. and apply Morrison-Woodbury-Sherman

$$\left(\Sigma_{SS} - \Sigma_{SS^c} \Sigma_{S^cS^c}^{-1} \Sigma_{SS^c}^\top \right)^{-1} = \Sigma_{SS}^{-1} + \Sigma_{SS}^{-1} \Sigma_{SS^c} \left(\Sigma_{S^cS^c} - \Sigma_{SS^c}^\top \Sigma_{SS}^{-1} \Sigma_{SS^c} \right)^{-1} \Sigma_{SS^c}^\top \Sigma_{SS}^{-1} \quad (6)$$

By Cauchy-Schwarz, the term $\Sigma_{S^cS^c} - \Sigma_{SS^c}^\top \Sigma_{SS}^{-1} \Sigma_{SS^c}$ is positive semi-definite, see e.g. Tripathi (1999, esp. the last line of the proof of Theorem 1.1.). By symmetry, we get an analogous result for the second summand in side the trace operator,

This gives

$$\text{Tr}(\Sigma_X^{-1}) = \text{Tr}(\Sigma_{SS}^{-1}) + \text{Tr}(\Sigma_{S^c S^c}^{-1}) + \quad (7)$$

$$\begin{aligned} & \text{Tr}\left(\Sigma_{SS}^{-1}\Sigma_{SS^c}\left(\Sigma_{S^c S^c} - \Sigma_{SS^c}^\top \Sigma_{SS}^{-1}\Sigma_{SS^c}\right)^{-1}\Sigma_{SS^c}^\top \Sigma_{SS}^{-1}\right) + \\ & \text{Tr}\left(\Sigma_{S^c S^c}^{-1}\Sigma_{SS^c}^\top \left(\Sigma_{SS} - \Sigma_{SS^c}\Sigma_{S^c S^c}^{-1}\Sigma_{SS^c}^\top\right)^{-1}\Sigma_{SS^c}\Sigma_{S^c S^c}^{-1}\right) \\ & \geq \text{Tr}(\Sigma_{SS}^{-1}) \end{aligned} \quad (8)$$

and therefore an estimator with the Oracle Property is asymptotically more efficient than least squares.

To establish when equality holds, if $X = X_S$, then clearly the asymptotic relative efficiency is

1. For only if, the inequality above is an equality only when $\text{Tr}(\Sigma_{S^c S^c}^{-1}) = 0$, which is not possible unless $X = X_S$.

B Preliminaries

We offer three sets of preliminary results. First, we show that the weights, \hat{w}_k , and magnitude of $|\hat{\beta}_k|$ are inversely related. Second, we formally differentiate between “large” and “small” estimates. This will help us derive bounds on \hat{w}_k . Third, we provide a bound on $\hat{\lambda}$. Note that we refer to the k^{th} order statistic of vector a as $a_{(k)}$, where $a_{(1)}$ is the smallest element of a .

B.1 Inverse relationship between weights and effect size.

PROPOSITION 1

$$\frac{\partial \hat{w}_k}{\partial |\hat{\beta}_k|} = -\hat{\lambda} \sqrt{\frac{1}{\sigma^2}} \text{Var}(w_k|\cdot) < 0. \quad (9)$$

Derivation: *The weights are calculated as*

$$\hat{w}_k = \mathbb{E}(w_k|\cdot) = \frac{\int_{w=0}^{\infty} w e^{-w\hat{\gamma} - \hat{\lambda}w\sqrt{\frac{1}{\sigma^2}}|\hat{\beta}_k|} dw}{\int_{w=0}^{\infty} e^{-w\hat{\gamma} - \hat{\lambda}w\sqrt{\frac{1}{\sigma^2}}|\hat{\beta}_k|} dw}. \quad (10)$$

Denote as $A = e^{-w^{\hat{\gamma}} - \hat{\lambda}w\sqrt{\frac{1}{\sigma^2}|\hat{\beta}_k|}}$. Then,

$$\frac{\partial \hat{w}_k}{\partial |\hat{\beta}_k|} = \frac{-\int_{w=0}^{\infty} A dw \times \int_{w=0}^{\infty} w^2 \hat{\lambda} \sqrt{\frac{1}{\sigma^2}} A dw + \int_{w=0}^{\infty} w A dw \int_{w=0}^{\infty} w \hat{\lambda} \sqrt{\frac{1}{\sigma^2}} A dw}{\left(\int_{w=0}^{\infty} A dw\right)^2} \quad (11)$$

$$= -\hat{\lambda} \sqrt{\frac{1}{\sigma^2}} \left\{ \frac{\int_{w=0}^{\infty} w^2 A dw}{\int_{w=0}^{\infty} A dw} - \left(\frac{\int_{w=0}^{\infty} w A dw}{\int_{w=0}^{\infty} A dw} \right)^2 \right\} \quad (12)$$

$$= -\hat{\lambda} \sqrt{\frac{1}{\sigma^2}} \text{Var}(w_k | \cdot) \quad (13)$$

where moving the derivative under the integral in the first line is allowed by the monotone convergence theorem.

This result allows us to associate the largest weight, $\hat{w}_{(K)}$ with the smallest estimate, $\hat{\beta}_{(1)}$, the second largest weight with the second smallest estimate, and so on. In general, weight $\hat{w}_{(k)}$ is associated with $|\hat{\beta}|_{(K-k+1)}$

B.2 Separating large and small weights and effect estimates.

We next distinguish between weights near zero from weights close to the maximal value $\hat{\gamma}$. This is our equivalent of either assuming the estimates are “well-separated” (Belloni and Chernozhukov, 2013), or separating “relevant” from “irrelevant” effects (Buhlmann and van de Geer, 2013). The key difference is that these authors separate large and small “true” effects, whereas we separate large and small estimated effects. As is common in the literature, our bounds will be more informative the better we can distinguish between zero- and non-zero effect estimates.

We separate the weights into two groups. In the kernel for $\Pr(w_k | \cdot)$, the numerator in Equation 10, is approximately exponential for large $|\beta_k|$, small w_k , and is approximately constant for $|\beta_k| \approx 0$, w_k large. Define as

$$p_k(C_1, C_2) = \max \left\{ \Pr \left(\hat{w}_k > \frac{C_1 \log(\hat{S})}{\lambda \hat{\sigma} |\hat{\beta}_k|} \right), \Pr \left(\hat{w}_k < C_2 \bar{\gamma} \right) \right\}; \quad C_1 > 0, 0 < C_2 < 1 \quad (14)$$

where the first inequality allows us to bound with some high probability small weights from above and the second, larger weights from below. We use this distinction to differentiate between weights tending to zero (the lefthand set) and those tending to the maximum (the righthand set).

$$\hat{S} = \left\{ k : \Pr \left(\hat{w}_k > \frac{C_1 \log(|\hat{S}|)}{\lambda \hat{\sigma} |\hat{\beta}_k|} \right) < \Pr \left(\hat{w}_k < C_2 \bar{\gamma} \right) \right\}. \quad (15)$$

The $\log(|\hat{S}|)$ term on the left comes from using the union bound applied to $\{p_k\}_{k=1}^K$ and a subexponential (rather than subgaussian) bound applied to each value p_k , as the kernel is approximately exponential in this range. Define

$$\Pr(\max(p_k) > C_3 \log(K)) = p_w(C_1, C_2, C_3). \quad (16)$$

such that, with probability at least $p_w(C_1, C_2, C_3)$, the weights can be bounded by one of the bounds above, i.e. is either “small” or “large.”

Lastly, denote as \underline{C}_1 the value that satisfies

$$\Pr\left(\hat{w}_k > \frac{C_1 \log(|\hat{S}|)}{\hat{\lambda} \hat{\sigma} |\hat{\beta}_k|}\right) = \Pr\left(\hat{w}_k \leq \frac{\underline{C}_1}{\hat{\lambda} \hat{\sigma} |\hat{\beta}_k| \log(|\hat{S}|)}\right) \quad (17)$$

which will give us a lower bound on all \hat{w}_k with probability at least $p_w(C_1, C_2, C_3)$.

B.3 Bounding the tuning parameter $\hat{\lambda}$.

Given the results above, we can bound $\hat{\lambda}$. For the Oracle results below, we need to bound $\hat{\lambda}$ from below, though we note that a similar bound of the same order of N, K can be found using the strategy below.

As $\lambda^2 |\cdot| \sim \Gamma(\sqrt{N}K, \frac{1}{2} \sum_{k=1}^K \hat{\tau}_k^2 + \rho)$, this gives

$$\hat{\lambda}^2 = \frac{\sqrt{N}K}{\frac{1}{2} \sum_{k=1}^K \hat{\tau}_k^2 + \rho}. \quad (18)$$

Change of variables gives $\lambda |\cdot| \sim \text{generalizedGamma}\left(2 \times (\frac{1}{2} \sum_{k=1}^K \hat{\tau}_k^2 + \rho)^{-1/2}, 2\sqrt{N}K, 2\right)$, which gives the estimate

$$\hat{\lambda} = \frac{\tilde{\Gamma}(\sqrt{N}K + 1/2)/\tilde{\Gamma}(\sqrt{N}K)}{\sqrt{\frac{1}{2} \sum_{k=1}^K \hat{\tau}_k^2 + \rho}} \quad (19)$$

with $\tilde{\Gamma}()$ the Gamma function. Note $\hat{\lambda}^2 \geq (\hat{\lambda})^2$ and if $\sqrt{N}K > 1$, then $\Gamma(3/2)^2 (\hat{\lambda})^2 = \frac{4}{\pi} (\hat{\lambda})^2 > \hat{\lambda}^2$. Lastly,

$$1/\tau_k^2 |\cdot| \sim \text{InvGaussian}(\lambda w_k \sigma / |\beta_k|, w_k^2 \lambda^2) \Rightarrow \quad (20)$$

$$\sum_{k=1}^K \hat{\tau}_k^2 = \sum_{k=1}^K \frac{|\hat{\beta}_k|}{\hat{\lambda} \hat{w}_k \hat{\sigma}} + \frac{1}{\hat{\lambda}^2 \hat{w}_k^2} \quad (21)$$

and we use the bound

$$\sum_{k=1}^K \hat{\tau}_k^2 \leq \frac{|\hat{S}| \times |\hat{\beta}_{(K)}|}{\hat{\lambda} \hat{w}_{(1)} \hat{\sigma}} + \frac{|\hat{S}|}{(\hat{\lambda})^2 \hat{w}_{(1)}^2} + \frac{(K - |\hat{S}|) |\hat{\beta}|_{(K-|\hat{S}|-1)}}{\hat{\lambda} \hat{w}_{(K-|\hat{S}|-1)} \hat{\sigma}} + \frac{(K - |\hat{S}|)}{(\hat{\lambda})^2 \hat{w}_{(K-|\hat{S}|-1)}^2} \quad (22)$$

$$\leq \frac{|\hat{S}| \hat{\beta}_{(K)}^2}{\underline{C}_1 \log(\hat{S})} + \frac{|\hat{S}| \hat{\sigma}^2 \hat{\beta}_k^2}{\underline{C}_1^2 \log(\hat{S})^2} + \frac{(K - |\hat{S}|) |\hat{\beta}|_{(K-|\hat{S}|-1)}}{\hat{\lambda} C_2 \hat{\gamma} \hat{\sigma}} + \frac{(K - |\hat{S}|)}{(\hat{\lambda})^2 C_2^2 \hat{\gamma}^2} \quad (23)$$

$$= \frac{|\hat{S}| \hat{\beta}_{(K)}^2 (\underline{C}_1 \log(\hat{S}) + \hat{\sigma}^2)}{\underline{C}_1^2 \log(\hat{S})^2} + \frac{(K - |\hat{S}|) |\hat{\beta}|_{(K-|\hat{S}|-1)}}{\hat{\lambda} C_2 \hat{\gamma} \hat{\sigma}} + \frac{(K - |\hat{S}|)}{(\hat{\lambda})^2 C_2^2 \hat{\gamma}^2} \quad (24)$$

The first line follows from the inverse relationship between $|\hat{\beta}_k|$ and \hat{w}_k ; the second comes from the lower bounds on \hat{w}_k in \hat{S} and \hat{S}^c . The third line is just simplifying.

Combining inequalities gives

$$\frac{4}{\pi} (\hat{\lambda})^2 \geq \hat{\lambda}^2 = \frac{\sqrt{N} K}{\frac{1}{2} \sum_{k=1}^K \hat{\tau}_k^2 + \rho} \quad (25)$$

$$\Rightarrow (\hat{\lambda})^2 \geq \frac{\pi}{4} \times \frac{\sqrt{N} K}{\frac{|\hat{S}| \hat{\beta}_{(K)}^2 (\underline{C}_1 \log(\hat{S}) + \hat{\sigma}^2)}{2 \underline{C}_1^2 \log(\hat{S})^2} + \frac{(K - |\hat{S}|) |\hat{\beta}|_{(K-|\hat{S}|-1)}}{2 \hat{\lambda} C_2 \hat{\gamma} \hat{\sigma}} + \frac{(K - |\hat{S}|)}{2 (\hat{\lambda})^2 C_2^2 \hat{\gamma}^2} + \rho} \quad (26)$$

$$\Rightarrow \hat{\lambda} \geq \frac{\pi}{4} \times \frac{\sqrt{N} K}{\frac{\hat{\lambda} |\hat{S}| \hat{\beta}_{(K)}^2 (\underline{C}_1 \log(\hat{S}) + \hat{\sigma}^2)}{2 \underline{C}_1^2 \log(\hat{S})^2} + \frac{(K - |\hat{S}|) |\hat{\beta}|_{(K-|\hat{S}|-1)}}{2 C_2 \hat{\gamma} \hat{\sigma}} + \frac{(K - |\hat{S}|)}{2 \hat{\lambda} C_2^2 \hat{\gamma}^2} + \rho \hat{\lambda}} \quad (27)$$

where the second line comes from substituting from Inequality 24 and the third from multiplying both sides by $1/\hat{\lambda}$. Cross-multiplying gives a quadratic equation in $\hat{\lambda}$ of the form $\tilde{a}(\hat{\lambda})^2 + \tilde{b}\hat{\lambda} + \tilde{c} > 0$ where¹

$$\tilde{a} = \frac{|\hat{S}| \hat{\beta}_{(K)}^2 (\underline{C}_1 \log(\hat{S}) + \hat{\sigma}^2)}{2 \underline{C}_1^2 \log(\hat{S})^2} + \rho \quad (28)$$

$$\tilde{b} = \frac{(K - |\hat{S}|) |\hat{\beta}|_{(K-|\hat{S}|-1)}}{2 C_2 \hat{\gamma} \hat{\sigma}} \quad (29)$$

$$\tilde{c} = - \left(\frac{\pi}{4} \sqrt{N} K - \frac{(K - |\hat{S}|)}{2 C_2^2 \hat{\gamma}^2} \right). \quad (30)$$

¹We use the convention $0 \log 0 = 0$

The quadratic equation gives

$$\hat{\lambda} \geq \frac{-\frac{(K-|\hat{S}|)|\hat{\beta}|_{(K-|\hat{S}|-1)}}{2C_2\hat{\gamma}\hat{\sigma}} + \sqrt{\left\{\frac{(K-|\hat{S}|)|\hat{\beta}|_{(K-|\hat{S}|-1)}}{2C_2\hat{\gamma}\hat{\sigma}}\right\}^2 + 4\left\{\frac{\hat{\lambda}|\hat{S}|\hat{\beta}_{(K)}^2(C_1\log(\hat{S})+\hat{\sigma}^2)}{2C_1^2\log(\hat{S})^2} + \rho\right\} \times \left\{\frac{\pi}{4}\sqrt{N}K - \frac{(K-|\hat{S}|)}{2C_2^2\hat{\gamma}^2}\right\}}}{2\frac{|\hat{S}|\hat{\beta}_{(K)}^2(C_1\log(\hat{S})+\hat{\sigma}^2)}{2C_1^2\log(\hat{S})^2} + 2\rho} \quad (31)$$

which, for growing N and K , is of order $N^{1/4}K^{1/2}$ by the bound in 31.

C Variance Estimation

We sample from the approximate sampling distribution of the the LASSOplus estimator at each Gibbs update:

$$\beta_k \mathbf{1} \left(|\hat{\beta}_k^{sp}| \geq \frac{\lambda w_k \sigma_{sp}}{N-1} \right) \quad (32)$$

$$\approx \beta_k \Phi \left\{ \left| \left| \hat{\beta}_k^{sp} / \hat{\sigma}_{ls} \right| - \frac{\lambda w_k \sigma_{sp}}{\hat{\sigma}_{ls} \times (N-1)} \right| \right\} \quad (33)$$

$$= \beta_k \Phi \left\{ \sqrt{N-1} \left| \left| \hat{\beta}_k^{sp} / \sigma \right| - \frac{\lambda w_k \sigma_{sp}}{\sigma \times (N-1)} \right| \right\} \quad (34)$$

$$= g \left(\beta_k, \hat{\beta}_k^{sp}, \sigma, \sigma_{sp}, \lambda, w_k \right) \quad (35)$$

where $\Phi(a)$ is the cumulative distribution for a standard normal random variable and we approximate the standard error of the least squares coefficient as $\hat{\sigma}_{ls} \approx \sigma / \sqrt{N-1}$. Define

$$z_k = \sqrt{N-1} \left| \left| \hat{\beta}_k^{sp} / \sigma \right| - \frac{\lambda w_k \sigma_{sp}}{\sigma \times (N-1)} \right| \quad (36)$$

$$p_k = \Phi \{z_k\} \quad (37)$$

Define the vector of partial derivatives

$$\nabla g \left(\beta_k, \hat{\beta}_k^{sp}, \sigma, \sigma_{sp}, \lambda, w_k \right) = \left[\frac{\partial g(\cdot)}{\partial \beta_k}, \frac{\partial g(\cdot)}{\partial \hat{\beta}_k^{sp}}, \frac{\partial g(\cdot)}{\partial \sigma}, \frac{\partial g(\cdot)}{\partial \sigma_{sp}}, \frac{\partial g(\cdot)}{\partial \lambda}, \frac{\partial g(\cdot)}{\partial w_k} \right]^\top \quad (38)$$

$$= \begin{bmatrix} p_k \\ \beta_k \times \phi(z_k) \times \sqrt{N-1} / \sigma \operatorname{sgn}(\hat{\beta}_k^{sp}) \\ \beta_k \times \phi(z_k) \times \sqrt{N-1} \left(-\frac{|\hat{\beta}_k^{sp}|}{\sigma^2} + \frac{\lambda w_k \sigma_{sp}}{\sigma^2 \times (N-1)} \right) \\ \beta_k \times \phi(z_k) \times \sqrt{N-1} \times \frac{\lambda w_k}{\sigma \times (N-1)} \\ \beta_k \times \phi(z_k) \times \sqrt{N-1} \times \frac{w_k \sigma_{sp}}{\sigma \times (N-1)} \\ \beta_k \times \phi(z_k) \times \sqrt{N-1} \times \frac{\lambda \sigma_{sp}}{\sigma \times (N-1)} \end{bmatrix} \quad (39)$$

and the 6×6 matrix

$$V = \text{diag} \left[\text{Var}(\beta_k), \text{Var}(\hat{\beta}_k^{sp}), \text{Var}(\sigma), \text{Var}(\sigma_{sp}), \text{Var}(\lambda), \text{Var}(w_k) \right] \quad (40)$$

where we are assuming zero covariance between elements. All elements of V are calculated analytically from the variance of the conditional pseudoposterior densities except for $\text{Var}(w_k)$ which is calculated from the approximate density used in the griddy Gibbs sampler. Our approximate variance is then

$$\hat{\sigma}_j^2 = \nabla g^\top(\cdot) V \nabla g(\cdot) \quad (41)$$

D EM Updates for LASSOplus-EM

For our EM implementation, we treat in $\beta^{plus-EM}$ and σ^2 as parameters and the remaining parameters as “missing,” i.e. to be estimated. As we have already calculated the conditional posterior densities for all parameters, the EM updates is straightforward.

Standardize Y and all columns of X to be mean-zero, sample variance one. Initialize $\forall k : \hat{\beta}_k \leftarrow u_k$ with $u_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 0.01)$; $\hat{\sigma}^2 \leftarrow \|Y - X\hat{\beta}\|_2^2/N$; $\hat{\lambda} \leftarrow 1$; $\hat{w}_k \leftarrow 1$.

At each given step, the most current updates from the previous steps are used. To convergence,

- E-steps

1. $\forall k : (\widehat{1/\tau_k^2}) \leftarrow \hat{\lambda} \hat{w}_k \hat{\sigma} / |\hat{\beta}_k|$; $\hat{\tau}_k^2 \leftarrow |\hat{\beta}_k| / (\hat{\lambda} \hat{w}_k \hat{\sigma}) + 1 / (\hat{\lambda}^2 \hat{w}_k^2)$
2. $\hat{\lambda} \leftarrow \frac{\tilde{\Gamma}(\sqrt{NK}+1/2)/\tilde{\Gamma}(\sqrt{NK})}{\sqrt{\frac{1}{2} \sum_{k=1}^K \hat{\tau}_k^2 + \rho}}$; $\hat{\lambda}^2 \leftarrow \frac{\sqrt{NK}}{\frac{1}{2} \sum_{k=1}^K \hat{\tau}_k^2 + \rho}$ with $\tilde{\Gamma}()$ the gamma function.
3. $\forall k$: update \hat{w}_k via numerical integration using kernel $\Pr(w_k|\cdot) \propto e^{-w\hat{\gamma} - \hat{\lambda}w\sqrt{\frac{1}{\sigma^2}}|\hat{\beta}_k|}$
4. Update $\hat{\gamma}$ via numerical integration using kernel $\Pr(\gamma|\cdot) \propto \gamma e^{-\sum_{k=1}^K \hat{w}_k \gamma - \gamma}$

- M-Steps

- $\hat{\sigma}^2 \leftarrow \frac{\sum_{i=1}^N (Y_i - X_i^\top \hat{\beta})^2 + \sum_{k=1}^K (\hat{\beta}_k)^2 \times \frac{1}{\hat{\tau}_k^2}}{N+K}$; $\frac{1}{\hat{\sigma}^2} \leftarrow \frac{N+K-2}{\sum_{i=1}^N (Y_i - X_i^\top \hat{\beta})^2 + \sum_{k=1}^K (\hat{\beta}_k)^2 \times \frac{1}{\hat{\tau}_k^2}}$
- Conditional M -steps: $\forall k : \hat{\beta}^k \leftarrow \frac{\sum_{i=1}^N X_{ik}(Y_i - \sum_{k' \neq k} X_{ik'} \hat{\beta}_{k'})}{(N-1) + \frac{1}{\hat{\tau}_k^2}}$ where it is understood that at update \tilde{k} , updated estimates of $\{\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_{\tilde{k}-1}\}$ are used.

LASSOplus updates:

- $\hat{\sigma}_{sp}^2 \leftarrow \leftarrow \frac{\sum_{i=1}^N (Y_i - X_i^\top \hat{\beta})^2 + \sum_{k=1}^K (\hat{\beta}_k)^2 \times \frac{1}{\tau_k^2}}{\sqrt{N+K}}$
- $\hat{\beta}_k^{plus} \leftarrow \hat{\beta}_k \mathbf{1} \left(\left| \sum_{i=1}^N X_{ik} (Y_i - \sum_{k' \neq k} X_{ik'} \hat{\beta}_{k'}) \right| > \hat{\lambda} \hat{w}_k \sqrt{\hat{\sigma}_{sp}^2} \right)$

E Independence between Adjusted Higher-Order Terms and Lower-Order Terms

We prove first that, under the residualized construction, the least squares coefficient on the a higher-order interaction term is uncorrelated with the coefficients on lower-order terms. By this means, the effect of the higher-order term does not vary with its lower-order components, and hence can be interpreted on its own. We then extend the result to the conditional pseudoposterior density of the estimates.

Denote the $N \times 1$ vector of outcomes Y , $N \times L$ matrix of lower-order terms \mathbf{X}_{lower} and vector of mean-zero, equivariant errors ϵ . Define as $X_{inter} = [X_{inter}]_i = \prod_{1 \leq l' \leq L} x_{il'}$, the elementwise product of the lower-order terms. Assume $[\mathbf{X}_{lower} | X_{inter}]$ is full rank. Using parameters $\{\beta_0, \vec{\beta}_l, \beta_{inter}\}$, define the model, with $\vec{\beta}_l$ an $L \times 1$ vector and the others scalars, as

$$Y = \mathbf{X}_{lower} \vec{\beta}_l + X_{inter} \beta_{inter} + \epsilon. \quad (42)$$

Define the matrices

$$\mathbf{X} = [\mathbf{X}_{lower} | X_{inter}] \quad (43)$$

$$\mathbf{M}_{lower} = \mathbf{I}_L - (\mathbf{X}_{lower})(\mathbf{X}_{lower}^\top \mathbf{X}_{lower})^{-1} \mathbf{X}_{lower}^\top \quad (44)$$

$$\tilde{X}_{inter} = \mathbf{M}_{lower} X_{inter} \quad (45)$$

$$\mathbf{X}_{adjust} = [\mathbf{X}_{lower} | \tilde{X}_{inter}] \quad (46)$$

The vector \tilde{X}_{inter} is the residualized interaction term described in the text, giving parameterization

$$Y = \mathbf{X}_{lower} \vec{\beta}_l + \tilde{X}_{inter} \tilde{\beta}_{inter} + \epsilon \quad (47)$$

where the error vector ϵ , stays unchanged since the two parameterizations differ only by a linear transformation of the covariates.

The covariance of the least squares estimates in the first parameterization is proportional to the inverse of the cross product of the design matrix. Using the block-partition matrix formula gives

$$(\mathbf{X}^\top \mathbf{X})^{-1} = \begin{bmatrix} (\mathbf{X}_{lower}^\top \mathbf{X}_{lower}) & \mathbf{X}_{lower}^\top X_{inter} \\ X_{inter}^\top \mathbf{X}_{lower} & X_{inter}^\top X_{inter} \end{bmatrix}^{-1} \quad (48)$$

$$= \begin{bmatrix} \left(\mathbf{X}_{lower}^\top \mathbf{X}_{lower} - \frac{1}{c_0} \mathbf{X}_{lower}^\top X_{inter} X_{inter}^\top \mathbf{X}_{lower} \right)^{-1} & -\frac{1}{c_0} (\mathbf{X}_{lower}^\top \mathbf{X}_{lower})^{-1} X_{lower}^\top X_{inter} \\ -\frac{1}{c_0} X_{inter}^\top \mathbf{X}_{lower} (\mathbf{X}_{lower}^\top \mathbf{X}_{lower})^{-1} & \frac{1}{c_0} \end{bmatrix} \quad (49)$$

with the constant $c_0 = X_{inter}^\top X_{inter} - X_{inter}^\top \mathbf{X}_{lower} (\mathbf{X}_{lower}^\top \mathbf{X}_{lower})^{-1} \mathbf{X}_{lower}^\top X_{inter}$. This implies

$$\text{Cov}(\hat{\beta}_{inter}, \hat{\beta}_k) \propto - \left[\frac{1}{c_0} (\mathbf{X}_{lower}^\top \mathbf{X}_{lower})^{-1} \mathbf{X}_{lower}^\top X_{inter} \right]_j \quad \text{for } j \in \{1, 2, \dots, L\} \quad (50)$$

In general, this covariance will not be zero, suggesting that under the normal parameterization the effect of the interaction term varies with movements in its lower order terms. Repeating the same exercise with a model parameterized in terms of \tilde{X}_{inter} gives

$$\text{Cov}(\hat{\tilde{\beta}}_{inter}, \hat{\tilde{\beta}}_k) \propto - \left[\frac{1}{c_0} (\mathbf{X}_{lower}^\top \mathbf{X}_{lower})^{-1} \mathbf{X}_{lower}^\top \tilde{X}_{inter} \right]_j \quad (51)$$

$$= - \left[\frac{1}{c_0} (\mathbf{X}_{lower}^\top \mathbf{X}_{lower})^{-1} \mathbf{X}_{lower}^\top \mathbf{M}_{lower} X_{inter} \right]_j \quad (52)$$

$$= 0 \quad \text{for } j \in \{1, 2, \dots, L\} \quad (53)$$

Therefore, under the parameterization with residualized interaction terms, the marginal effect of each interaction term is uncorrelated with that of its lower order terms. To extend to the multivariate case, assume the full design matrix of all effects is full-rank, and all other effects have been partialled out. The case of $K > N$ requires an assumption similar to the restricted eigenvalue assumption (Bickel, Ritov and Tsybakov, 2009), that all submatrices of size $L + 2$ are full rank and all components of the submatrices not in \mathbf{X} are linearly independent of \mathbf{X} . Partialing out with respect to the other covariates in either case leaves the results unchanged.

Next, we show the result holds for the conditional pseudoposterior density under a conditional independent normal prior, as with the augmented LASSOplus. Assume $[\vec{\beta}_l^\top, \beta_{inter}]^\top \sim \mathcal{N}(0_{L+1}, D)$ with D an $(L + 1) \times (L + 1)$ diagonal matrix with positive entries along the diagonal. In this case, the conditional posterior of $[\vec{\beta}_l^\top, \beta_{inter}]^\top$ under a normal likelihood takes the form

$$\Pr([\vec{\beta}_l^\top, \beta_{inter}]^\top | \cdot) \sim \mathcal{N}(A^{-1} \mathbf{X}^\top Y, \sigma^2 A^{-1}) \quad (54)$$

with $A = \mathbf{X}^\top \mathbf{X} + D$. Carrying through the same derivation as above gives the posterior covariance between β_{L+1} , the parameter on the interaction term, and β_k , $1 \leq j \leq L$, as

$$A_{j,L+1}^{-1} \propto - \left[\frac{1}{c'_0} (\mathbf{X}_{lower}^\top \mathbf{X}_{lower} + D_{1:L,1:L})^{-1} \mathbf{X}_{lower}^\top X_{inter} \right]_j \quad \text{for } j \in \{1, 2, \dots, L\} \quad (55)$$

which will not be 0, in general. In this case, the constant $c'_0 = (X_{inter} + D_{L+1,L+1})^\top (X_{inter} + D_{L+1,L+1}) - X_{inter}^\top \mathbf{X}_{lower} (\mathbf{X}_{lower}^\top \mathbf{X}_{lower} + D_{1:L,1:L})^{-1} \mathbf{X}_{lower}^\top X_{inter}$.

Considering the residualized interaction term instead of the standard term gives

$$A_{j,L+1}^{-1} \propto - \left[\frac{1}{c'_0} (\mathbf{X}_{lower}^\top \mathbf{X}_{lower} + D_{1:L,1:L})^{-1} \mathbf{X}_{lower}^\top \tilde{X}_{inter} \right]_j = 0 \quad \text{for } j \in \{1, 2, \dots, L\} \quad (56)$$

F Alternative Estimators

For the LASSO and adaptive LASSO, we found the BIC statistic of Wang and Leng (2007) performed poorly when $K > N$, sometimes including dozens of false positives. We instead use a standard BIC statistic where we take the degrees of freedom as the number of non-zero coefficients Zou, Hastie and Tibshirani (2007).

In terms of uncertainty estimates, we implement the approximate confidence intervals for the LASSOplus. We use the posterior intervals for the horseshoe model. For the frequentist LASSO and adaptive LASSO, we implement the perturbation method of Minnier, Tian and Cai (2011). For $p \in \{1, 2, \dots, P\}$ for some large P , the method requires fitting

$$\hat{\beta}^{alasso,p}(\lambda|w., g.) = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^N g_i^p (Y_i - X_i^\top \beta)^2 + \lambda \sum_{k=1}^K w_k |\beta_k|; \quad (57)$$

$$w_k = 1/|\beta_k^0| \quad (58)$$

where the weights are $g_i^p \stackrel{\text{i.i.d.}}{\sim} \exp(1)$. For the LASSO, we simply take $w_k = 1$ for all k . Minnier, Tian and Cai (2011) prove that the set $\{\hat{\beta}^{alasso,p}(\lambda|w., g.)\}_{p=1}^P$ will achieve nominal coverage asymptotically, though the result does not hold for the LASSO. We fit the perturbation method to both for comparison. We found the perturbation method performs better than the parametric bootstrap suggested by Chatterjee and Lahiri (2011, 2013), so we do not present the results.

We next move on to the LASSO+OLS method of Belloni and Chernozhukov (2013), hereafter BC. The empirical process approach selects the tuning parameter in order to bound $2 \max(\epsilon^\top X_{\cdot k})$ with some high probability. BC note that, up to a scale parameter σ , the tuning parameter value can

be simulated quite easily, and they define $\Lambda(1 - \alpha_{sig}|X)$ as the $1 - \alpha_{sig}$ quantile of $2 \max(\epsilon^\top X_{\cdot k} / \sigma_b)$ for $\mathbb{E}(\epsilon_i|X_i) = 0$; $\text{Var}(\epsilon_i|X_i) = \sigma_b^2$ as approximated through a simulation.

Second stage variable selection. Tuning λ in order to satisfy the Oracle Inequality will generally over-select effects. The reason is that the LASSO induces bias in the coefficient estimates, and that bias leaves a gap for irrelevant effects that are correlated with the relevant effects to be drawn in and selected. Several methods in the empirical process framework have used the Oracle Inequality-tuned LASSO to over-select covariates and then, in a second stage, select a subset of these.

One way to do so is simply thresholding the LASSO estimates, so

$$\hat{\beta}^{thresh} = \hat{\beta}^L \odot \mathbf{1} \left(|\hat{\beta}^L| > \tau \right) \quad (59)$$

where the inequality and multiplication \odot are taken elementwise. A second option is to take then re-run OLS on variables that survive the threshold. Define X_{thresh} as the submatrix of X corresponding with elements of $\hat{\beta}^{thresh}(\tau)$ that are not zero. Then,

$$\hat{\beta}^{thresh+OLS}(\tau) = (X_{thresh}^\top X_{thresh})^{-1} X_{thresh}^\top Y. \quad (60)$$

In the case X_{thresh} is rank-deficient, either ridge regression or partial least squares can be used (Liu and Yu, 2013). The post LASSO OLS estimator is then $\hat{\beta}^{thresh+OLS}(0)$, which is simply OLS used on all selected LASSO covariates.

Belloni and Chernozhukov (2013) propose a different means of selecting a subset of relevant effects and eliminating the first-stage LASSO bias. Denote $Q(\theta) = \|Y - X\theta\|_2^2$. The select τ such that

$$t_\gamma = \max_{t \geq 0} Q \left(\hat{\beta}^{thresh+OLS}(\tau) \right) - Q \left(\hat{\beta}^L \right) \leq \gamma \quad (61)$$

for $\gamma \leq 0$. Taking $\gamma = 0$ returns the sparsest OLS-reflated model that generates a lower residual sum of squares than the LASSO estimator. We follow the suggestion of Belloni and Chernozhukov (2013, expr 2.14) and take $\gamma = \left\{ Q \left(\hat{\beta}^{thresh+OLS}(0) \right) - Q \left(\hat{\beta}^L \right) \right\} / 2$ in the simulations.

References

Belloni, Alexandre and Victor Chernozhukov. 2013. “Least squares after model selection in high-dimensional sparse models.” *Bernoulli* 19(2):521–547.

- Bickel, Peter, Ya'acov Ritov and Alexandre Tsybakov. 2009. "Simultaneous Analysis of Lasso and Dantzig Selector." *Annals of Statistics* 37(4):1705–1732.
- Buhlmann, Peter and Sara van de Geer. 2013. *Statistics for High-Dimensional Data*. Springer.
- Chatterjee, A and SN Lahiri. 2011. "Bootstrapping Lasso Estimators." *Journal of the American Statistical Association* 106(494):608–625.
- Chatterjee, A and SN Lahiri. 2013. "Rates of convergence of the adaptive LASSO estimators to the oracle distribution and higher order refinements by the bootstrap." *The Annals of Statistics* 41(3):1232–1259.
- Liu, H. and B. Yu. 2013. "Asymptotic properties of Lasso+mLS and Lasso+Ridge in sparse high-dimensional linear regression." *Electronic Journal of Statistics* 7(3124–3169).
- Minnier, Jessica, Lu Tian and Tianxi Cai. 2011. "A perturbation method for inference on regularized regression estimates." *Journal of the American Statistical Association* 106(496).
- Tripathi, Gautam. 1999. "A matrix extension of the Cauchy-Schwarz inequality." *Economics Letters* 63:1–3.
- Wang, Hansheng and Chenlei Leng. 2007. "Unified LASSO Estimation by Least Squares Approximation." *Journal of the American Statistical Association* 102(479):1039–1048.
- Zou, Hui, Trevor Hastie and Robert Tibshirani. 2007. "On the degrees of freedom of the lasso." *The Annals of Statistics* 35(5):2173–2192.