

Supplementary Material for

**Conditional Selective Inference
for Robust Regression and Outlier Detection
using Piecewise-Linear Homotopy Continuation**

by Toshiaki Tsukurimichi, Yu Inatsu, Vo Nguyen Le Duy and Ichiro Takeuchi.

Appendix B: Breakpoints in Huber regression

In this section, we give the lemma for calculating breakpoints. Let \mathcal{A}_z , $S_{\mathcal{A}_z^c}$, $\hat{\mathbf{r}}(z)$, $\hat{\mathbf{u}}_{\mathcal{A}_z}(z)$, $\mathbf{h}_{\mathcal{A}_z^c}(z)$, $\psi(z)$ and $\gamma(z)$ be the same definitions as in Appendix Appendix A:. Then, the following lemma holds:

Lemma 5 *Consider a real value z . Then, $\mathcal{A}_z = \mathcal{A}_{z'}$ for any real value z' in the interval $[z, z + t_z]$, where $z + t_z$ is the value of transition point,*

$$t_z = \min\{t_z^1, t_z^2\},$$

$$t_z^1 = \min_{j \in \mathcal{A}_z^c} \left(-\frac{(S_{\mathcal{A}_z^c} \hat{\mathbf{r}}(z) - \mathbf{h}_{\mathcal{A}_z^c}(z))_j}{(S_{\mathcal{A}_z^c} \psi(z))_j} \right)_{++} \quad \text{and} \quad t_z^2 = \min_{j \in \mathcal{A}_z} \left(-\frac{(\hat{\mathbf{u}}_{\mathcal{A}_z}(z))_j}{\gamma_j(z)} \right)_{++}.$$

Here, $(a)_{++} = a$ if $a \geq 0$, and otherwise $(a)_{++} = +\infty$.

Proof We first show how to derive t_z^1 . From (29), we have

$$\hat{\mathbf{r}}(z') = \hat{\mathbf{r}}(z) + \psi(z) \times (z' - z).$$

Then, we need to guarantee

$$\begin{aligned} S_{\mathcal{A}_z^c} \hat{\mathbf{r}}(z) - \mathbf{h}_{\mathcal{A}_z^c}(z) &\leq 0, \\ S_{\mathcal{A}_z^c} (\hat{\mathbf{r}}(z') + \psi(z) \times (z' - z)) - \mathbf{h}_{\mathcal{A}_z^c}(z) &\leq 0, \\ S_{\mathcal{A}_z^c} \psi(z) \times (z' - z) &\leq -(S_{\mathcal{A}_z^c} \hat{\mathbf{r}}(z) - \mathbf{h}_{\mathcal{A}_z^c}(z)). \end{aligned} \quad (36)$$

The right hand side of (36) is positive since $S_{\mathcal{A}_z^c} \hat{\mathbf{r}}(z) - \mathbf{h}_{\mathcal{A}_z^c}(z) \leq 0$. Therefore, satisfying equation (36) implies that

$$z' - z \leq \min_{j \in \mathcal{A}_z^c} \left(-\frac{(S_{\mathcal{A}_z^c} \hat{\mathbf{r}}(z) - \mathbf{h}_{\mathcal{A}_z^c}(z))_j}{(S_{\mathcal{A}_z^c} \psi(z))_j} \right)_{++} = t_z^1.$$

Next, we show how to derive t_z^2 . From (30), we have

$$\hat{\mathbf{u}}_{\mathcal{A}_z}(z') = \hat{\mathbf{u}}_{\mathcal{A}_z}(z) + \gamma(z) \times (z' - z).$$

Thus, noting that $\hat{\mathbf{u}}_{\mathcal{A}_z}(z') \geq 0$ we need guarantee

$$\hat{\mathbf{u}}_{\mathcal{A}_z}(z') = \hat{\mathbf{u}}_{\mathcal{A}_z}(z) + \gamma(z) \times (z' - z) \geq 0. \quad (37)$$

Hence, satisfying equation (37) means that

$$z' - z \leq \min_{j \in \mathcal{A}_z} \left(-\frac{(\hat{\mathbf{u}}_{\mathcal{A}_z}(z))_j}{\gamma_j(z)} \right)_{++} = t_z^2.$$

Appendix C: Huberized Lasso

According to Equation (8) in Chen and Bien (2020), we consider the following optimization problem

$$(\hat{\beta}, \hat{\mathbf{u}}) = \arg \min_{\beta \in \mathbb{R}^p, \mathbf{u} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{y} - X\beta - \mathbf{u}\|_2^2 + \lambda \|\mathbf{u}\|_1, \quad (38)$$

given $X = (x_1, \dots, x_n)^\top$ and $\mathbf{y} = (y_1, \dots, y_n)^\top$. Following Sec 3.2 in the supplementary of Chen and Bien (2020), the optimization in (38) can be transformed to

$$\hat{\mathbf{u}} = \arg \min_{\mathbf{u} \in \mathbb{R}^n} \frac{1}{2} \|\tilde{\mathbf{y}} - \tilde{X}\mathbf{u}\|_2^2 + \lambda \|\mathbf{u}\|_1,$$

where $\tilde{\mathbf{y}} = P_X^\perp \mathbf{y}$ and $\tilde{X} = P_X^\perp X$. We can obtain $\hat{\mathbf{u}}$ in (38) by using Lasso algorithm \mathcal{A} . Then, the set of the observed outliers is defined as

$$\mathcal{A}(\mathbf{y}) = \{j : \hat{u}_j \neq 0\}.$$

Finally, the inference for a selected outlier is defined as follows

$$\boldsymbol{\eta}^\top \mathbf{Y} \mid \{\mathcal{A}(\mathbf{Y}) = \mathcal{A}(\mathbf{y}), \mathbf{q}(\mathbf{Y}) = \mathbf{q}(\mathbf{y})\}.$$

Unfortunately, as pointed out in Lee et al. (2016), characterizing $\mathcal{A}(\mathbf{Y}) = \mathcal{A}(\mathbf{y})$ in (39) is computationally intractable because we have to consider $2^{|\mathcal{A}(\mathbf{y})|}$ possible sign vectors. As suggested in Lee et al. (2016), we need to consider inference conditional not only on the selected features but also on their signs to overcome the aforementioned issue. Specifically, let $\mathbf{s}(\mathbf{y})$ denote the sign vector of the selected features when applying Lasso on \mathbf{y} , the conditional inference we need to focus is

$$\boldsymbol{\eta}^\top \mathbf{Y} \mid \{\mathcal{A}(\mathbf{Y}) = \mathcal{A}(\mathbf{y}), \mathbf{s}(\mathbf{Y}) = \mathbf{s}(\mathbf{y}), \mathbf{q}(\mathbf{Y}) = \mathbf{q}(\mathbf{y})\}. \quad (39)$$

However, additionally considering the signs leads to low statistical power because of *over-conditioning*. This is well-known as the major drawback in SI literature.