

Proximal Estimation and Inference

Alberto Quaini^a **Fabio Trojani**^b

^aErasmus University

^bUniversity of Geneva, University of Turin and SFI

ESEM 2024, Rotterdam

Proximal Estimation

Let $\hat{\beta}_n^s$ be an initial estimator for a target parameter of interest $\beta_0 \in \mathbb{R}^p$

Definition (Proximal Estimator (PE))

Given symm. pos. def. matrix W_n , convex lsc f_n and $\lambda_n > 0$:

$$\hat{\beta}_n := \text{prox}_{\lambda_n f_n}^{W_n}(\hat{\beta}_n^s) := \arg \min_{\beta} \left\{ \frac{1}{2} \left\| \hat{\beta}_n^s - \beta \right\|_{W_n}^2 + \lambda_n f_n(\beta) \right\}$$

is called a proximal estimator of β_0

- PEs are penalized **minimum distance corrections** defined via a (differentiable) **proximal operator** $\text{prox}_{\lambda_n f_n}^{W_n}$

Well-known convex lsc penalties of class $\Gamma(\mathbb{R}^p)$

- Ridge: $f_n(\beta) = \frac{1}{2} \|\beta\|_2^2$
- Lasso: $f_n(\beta) = \|\beta\|_1$
- Elastic Net: $f_n(\beta) = \frac{\alpha}{2} \|\beta\|_2^2 + (1 - \alpha) \|\beta\|_1$ for some $\alpha \in (0, 1)$
- Adaptive Lasso: $f_n(\beta) = \sum_{k=1}^K \frac{|\beta_i|}{|\tilde{\beta}_{ni}|}$ for some consistent auxiliary estimator $\tilde{\beta}_n$
- Convex constraints: $f_n(\beta) = \begin{cases} 0 & \beta \in C \\ \infty & \beta \notin C \end{cases}$ for some convex set C

A convenient framework

- A wide class of PEs from different choices of $\hat{\beta}_n^s$, W_n and f_n

- Embeds naturally **Penalized Least Squares Estimators** (PLSEs):

$$\text{prox}_{\lambda_n f_n}^{\mathbf{X}'\mathbf{X}/n}(\hat{\beta}_n^{ls}) = \arg \min_{\beta} \left\{ \frac{1}{2} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 + \lambda_n f_n(\beta) \right\}$$

- PEs' properties derived within simple/unifying **convex analysis** framework:
 - Asymptotic distribution, Oracle properties,...
 - Weak/transparent assumptions on initial estimator $\hat{\beta}_n^{ls}$ and convex penalties f_n
 - Mainly determined by properties of penalties' **subgradient**
- PEs can be developed to handle **irregular designs** where PLSEs may be ill-behaved

Asymptotic properties of PEs

Main high-level assumptions

- Properties of PEs determined by triplet $(\hat{\beta}_n^s, \mathbf{W}_n, \lambda_n f_n)$:

$$\hat{\beta}_n := \text{prox}_{\lambda_n f_n}^{\mathbf{W}_n}(\hat{\beta}_n^s) := \arg \min_{\beta} \left\{ \frac{1}{2} \left\| \hat{\beta}_n^s - \beta \right\|_{\mathbf{W}_n}^2 + \lambda_n f_n(\beta) \right\}$$

- Main Assumptions (Consistency, Asymptotic Distribution)

A1 $\mathbf{W}_n \rightarrow_{\text{Pr}} \mathbf{W}_0$ for some positive definite matrix \mathbf{W}_0

A2 $r_n(\hat{\beta}_n^s - \beta_0) \rightarrow_d \eta$, for some rate $r_n \rightarrow \infty$ and some random variable η

A3 $r_n \lambda_n f_n \rightarrow_{\text{Pr}} g_0$ in epigraph, for some limit penalty g_0

- Applicable to general class of convex lsc penalties and to irregular designs

Asymptotic distribution

Proposition (Asymptotic distribution of PEs)

Under Assumptions A1–A3:

$$r_n \left(\text{prox}_{\lambda_n f_n}^{\mathbf{W}_n}(\hat{\beta}_n^s) - \beta_0 \right) \rightarrow_d \text{prox}_{g_0'(\cdot; \beta_0)}^{\mathbf{W}_0}(\boldsymbol{\eta}) = \left(\text{Id} - P_{\partial g_0(\beta_0)}^{\mathbf{W}_0} \right) (\boldsymbol{\eta})$$

with *directional derivative* [*subgradient*] $g_0'(\cdot; \beta_0)$ [$\partial g_0(\beta_0)$] of limit penalty g_0 at β_0 , and *projection operator*:

$$P_{\partial g_0(\beta_0)}^{\mathbf{W}_0}(\boldsymbol{\eta}) := \arg \min_{\boldsymbol{\theta} \in \partial g_0(\beta_0)} \|\boldsymbol{\eta} - \boldsymbol{\theta}\|_{\mathbf{W}_0}$$

- **Functional characterization** of PEs' asy. distribution via limit penalty subgradient $\partial g_0(\beta_0)$
- **Closed-form** asymptotic distributions for established penalties in the literature

Example: Adaptive Lasso

- Adaptive Lasso penalty $f_n(\beta) = \sum_{k=1}^K \frac{|\beta_i|}{|\hat{\beta}_{ni}^s|}$
- If $\lambda_n r_n \rightarrow 0$ and $\lambda_n r_n^2 \rightarrow \infty$, then $\partial g_0(\beta_0) = \text{span}\{e_j : j \in \mathcal{A}\}^\perp$:

$$r_n \left(\text{prox}_{\lambda_n f_n}^{W_n}(\hat{\beta}_n^s) - \beta_0 \right) \rightarrow_d P_{\text{span}\{e_j; j \in \mathcal{A}\}}^{W_0}(\eta)$$

- Regular linear regression model with spherical errors:

$$- \hat{\beta}_n^s = \hat{\beta}_n^{ls}$$

$$- W_n = \mathbf{X}'\mathbf{X}/n \rightarrow \mathbf{Q}_0 := \mathbb{E}[\mathbf{X}_1\mathbf{X}_1'] \text{ and } \mathbf{X}'\epsilon/\sqrt{n} \rightarrow_d \mathbf{Z} \sim \mathcal{N}(\mathbf{0}, \sigma_\epsilon \mathbf{Q}_0)$$

- Adaptive Lasso PLSE's (efficient) asy. distribution:

$$\sqrt{n} \left(\text{prox}_{\lambda_n f_n}^{Q_n}(\hat{\beta}_n^{ls}) - \beta_0 \right) \rightarrow_d P_{\text{span}\{e_j; j \in \mathcal{A}\}}^{Q_0}(\mathbf{Q}_0^{-1}\mathbf{Z}) = \begin{cases} [(\mathbf{Q}_0)_{\mathcal{A}}]^{-1}(\mathbf{Z})_{\mathcal{A}} & \text{in } \mathcal{A} \\ 0 & \text{in } \mathcal{A}^c \end{cases}$$

Variable selection

Definition (Variable Selection)

Given a sequence of PE's estimated active sets:

$$\hat{\mathcal{A}}_n = \left\{ j : \left(\text{prox}_{\lambda_n f_n}^{\mathbf{W}_n}(\hat{\beta}_n^s) \right)_j \neq 0 \right\},$$

$\text{prox}_{\lambda_n f_n}^{\mathbf{W}_n}(\hat{\beta}_n^s)$ is said to perform consistent variable selection (VS) if $\mathbb{P}(\hat{\mathcal{A}}_n = \mathcal{A}) \rightarrow 1$.

Proposition (Variable Selection)

Let Assumptions A1–A3 be satisfied. IF VS holds then, as $n \rightarrow \infty$:

$$\mathbb{P} \left((\boldsymbol{\eta})_{\mathcal{A}^c} = \left(P_{\partial \mathcal{G}_0(\beta_0)}^{\mathbf{W}_0}(\boldsymbol{\eta}) \right)_{\mathcal{A}^c} \right) = 1.$$

Conversely, VS holds if optimal subgradient vectors $\mathbf{v}_n^{\text{opt}} = \mathbf{W}_n(\hat{\beta}_n^s - \hat{\beta}_n) \in \lambda_n \partial f_n(\hat{\beta}_n)$ are such that:

$$r_n \left\| (\mathbf{v}_n^{\text{opt}})_{\mathcal{A}^c} \right\|_1 \xrightarrow{\text{Pr}} \infty \quad \text{as } n \rightarrow \infty$$

Example: Adaptive Lasso

- Adaptive Lasso **optimal subgradient vector** implied by penalty $f_n(\beta) = \sum_{k=1}^K \frac{|\beta_i|}{|\tilde{\beta}_{ni}|}$ yields:

$$r_n \left\| (\mathbf{v}_n^{opt})_{\mathcal{A}^c} \right\|_1 = \lambda_n r_n \sum_{j \in \mathcal{A}^c} (1/|\tilde{\beta}_{nj}|)$$

- Whenever $\lambda_n r_n^2 \rightarrow \infty$ and $(\tilde{\beta}_n)_{\mathcal{A}^c} = O_{Pr}(1/r_n)$:

$$r_n \left\| (\mathbf{v}_n^{opt})_{\mathcal{A}^c} \right\|_1 \rightarrow_{Pr} \infty$$

i.e., **VS holds**.

PEs for irregular designs

Linear regression model with irregular design

- Linear model $\mathbf{Y} = \mathbf{X}\beta_0 + \epsilon$
- Sample design matrix $\mathbf{Q}_n := \mathbf{X}'\mathbf{X}/n$
- Population design matrix $\mathbf{Q}_{0n} := \mathbb{E}[\mathbf{Q}_n]$

Definition (Irregular design)

(i) **Singular design.** There exists singular matrix \mathbf{Q}_0 such that:

$$\mathbf{Q}_{0n} = \mathbf{Q}_0, \quad \text{for all } n$$

(ii) **Nearly-singular design.** \mathbf{Q}_{0n} is regular for all n and there exists singular matrix \mathbf{Q}_0 such that:

$$\mathbf{Q}_{0n} \rightarrow \mathbf{Q}_0 \text{ as } n \rightarrow \infty$$

Proximal estimation approach to irregular designs

- Under an irregular design, the set of limit population LS solutions is **not a singleton**
- Introduce a convenient **identifiable** parameter $\beta_0 \in \mathbb{R}^p$
- Build initial estimator $\hat{\beta}_n^s$ of β_0 , which is **well-behaved** under both regular and irregular designs
- Build suitable proximal estimator $\text{prox}_{\lambda_n f_n}^{W_n}(\hat{\beta}_n^s)$ of β_0 , which ideally satisfies the **Oracle property**

Ridgeless (limit) population parameter

- Let $\delta_0 := \lim_{n \rightarrow \infty} \mathbb{E}[\mathbf{X}'\mathbf{Y}/n]$

Definition (Ridgeless target parameter)

Given Moore-Penrose inverse \mathbf{Q}_0^+ , the **Ridgeless population parameter** is given by:

$$\beta_0^+ := \arg \min_{\beta} \{\|\beta\|_2 : \mathbf{Q}_0\beta = \delta_0\} = \mathbf{Q}_0^+ \delta_0$$

- β_0^+ is **identified** under both a regular and an irregular design

Ridgeless estimator

Definition (Ridgeless estimator)

$$\hat{\beta}_n^+ := \arg \min_{\beta} \{ \|\beta\|_2 : \mathbf{Q}_n \beta = \mathbf{X}' \mathbf{Y} / n \} = \mathbf{Q}_n^+ \mathbf{X}' \mathbf{Y} / n$$

- $\hat{\beta}_n^+ = \hat{\beta}_n^{ls}$ if \mathbf{Q}_n is regular
- Using standard assumptions, $\hat{\beta}_n^+$ is \sqrt{n} -consistent for β_0^+ and asymptotically normal, both under a regular and a singular design
- $\hat{\beta}_n^+$ is **not consistent** under a nearly-singular design, because \mathbf{Q}_n is not a **rank-consistent** estimator of \mathbf{Q}_0

Consistent modified Ridgeless estimation of β_0^+

Definition (Modified Ridgeless estimator)

Modified Ridgeless estimator of parameter β_0^+ is:

$$\check{\beta}_n^+ := \arg \min_{\beta} \{ \|\beta\|_2 : \check{Q}_n \beta = \mathbf{X}' \mathbf{Y} / n \} = \check{Q}_n^+ \mathbf{X}' \mathbf{Y} / n ,$$

where \check{Q}_n is a consistent estimator of Q_0 that is rank consistent:

$$\mathbb{P}(\text{Range}(\check{Q}_n) = \text{Range}(Q_0)) \rightarrow 1, \quad \text{as } n \rightarrow \infty$$

- \check{Q}_n obtained by **truncating eigenvalues** of Q_n with hard threshold $\mu_n \propto n^{-\alpha}$ for suitable $\alpha > 0$
- $\check{\beta}_n^+$ is clearly a **consistent** estimator of β_0^+ under both a regular and an irregular design

Asymptotic distribution of modified Ridgeless estimator

- MR1 $\mathbf{X}'\epsilon/\sqrt{n} \rightarrow_d \mathbf{Z} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Omega}_0)$, with $\mathbf{\Omega}_0$ symm. pos. semi definite
- MR2 $\sqrt{n}(\mathbf{Q}_n - \mathbf{Q}_{0n}) \rightarrow_d \mathbf{\Theta}$, for some random matrix $\mathbf{\Theta}$
- MR3 $\mathbf{P}_0(\mathbf{Q}_n - \mathbf{Q}_{0n})\mathbf{P}_0^\perp = o_p(1/\sqrt{n})$

Proposition (Asymptotic distribution of modified Ridgeless estimator)

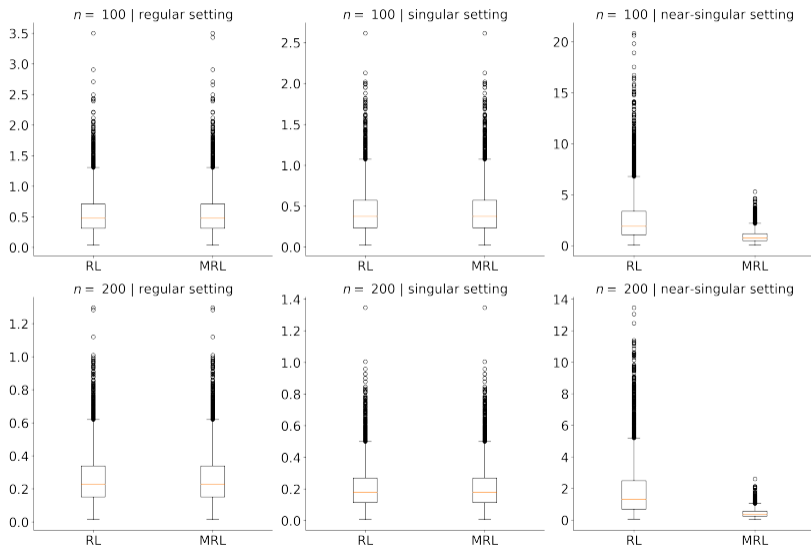
If $\mathbf{Q}_{0n} - \mathbf{Q}_0 = \frac{\Delta}{\tau_n} + o(1/\tau_n)$, where $\tau_n \rightarrow \infty$ as $n \rightarrow \infty$, and Assumptions MR1-MR2 hold, then:

$$\sqrt{n}(\check{\beta}_n^+ - \beta_0^+) \rightarrow_d \eta := \mathbf{P}_0(\mathbf{\Theta} + c\Delta)\mathbf{Q}_0^+\beta_0^+ + \mathbf{Q}_0^+\mathbf{Z},$$

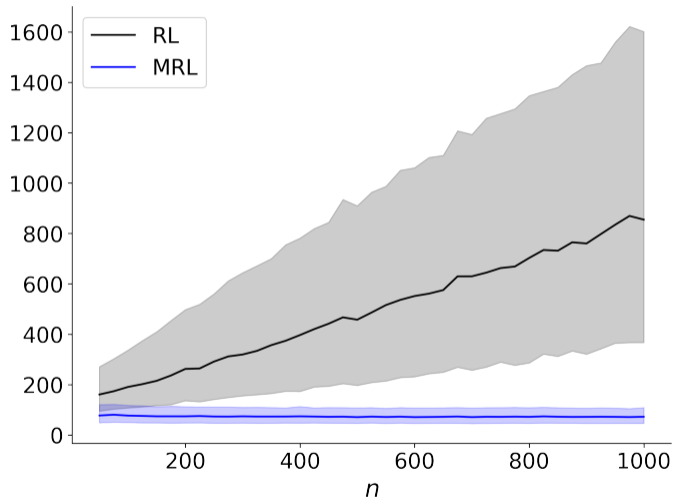
where $\sqrt{n}/\tau_n \rightarrow c \in \{0, 1\}$. Thus, if Assumption MR3 also holds and $c = 0$, then $\eta = \mathbf{Q}_0^+\mathbf{Z}$.

- $\check{\beta}_n^+$ is always asymptotically normally distributed whenever $\text{vec}(\mathbf{\Theta})$ is Gaussian

Squared errors $\left\| \hat{\beta}_n^s - \beta_0^+ \right\|_2^2$ of Ridgeless/modified Ridgeless



Nearly singular: normalized squared errors $n \left\| \hat{\beta}_n^s - \beta_0^+ \right\|_2^2$ of Ridgeless/modified Ridgeless



PEs for irregular designs

- A weighting matrix satisfying Assumption A1 under Assumptions MR1-MR2:

$$\mathbf{W}_n = \bar{\mathbf{Q}}_n := \check{\mathbf{Q}}_n + \mathbf{I} - \check{\mathbf{Q}}_n \check{\mathbf{Q}}_n^+ \xrightarrow{\text{Pr}} \bar{\mathbf{Q}}_0 := \mathbf{Q}_0 + \mathbf{I} - \mathbf{Q}_0 \mathbf{Q}_0^+$$

Proposition (Asymptotic distribution of PEs for irregular designs)

Consider following PE:

$$\text{prox}_{\lambda_n f_n}^{\bar{\mathbf{Q}}_n}(\check{\beta}_n^+) = \arg \min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{2} \|\check{\beta}_n^+ - \beta\|_{\bar{\mathbf{Q}}_n}^2 + \lambda_n f_n(\beta) \right\}$$

If Assumptions MR1, MR2 and A3 hold, then:

$$\sqrt{n} \left(\text{prox}_{\lambda_n f_n}^{\bar{\mathbf{Q}}_n}(\check{\beta}_n^+) - \beta_0^+ \right) \rightarrow_d \text{prox}_{g_0'(\beta_0^+; \cdot)}^{\bar{\mathbf{Q}}_0}(\eta) = \left(\text{Id} - P_{\partial g_0(\beta_0^+)}^{\bar{\mathbf{Q}}_0} \right) (\eta)$$

Oracle PEs for irregular designs

Proposition (Oracle PE for irregular designs)

Given Adaptive Lasso penalty $f_n(\beta) = \sum_{i=1}^p |\beta_i|/|\check{\beta}_{in}^+|$, define PE:

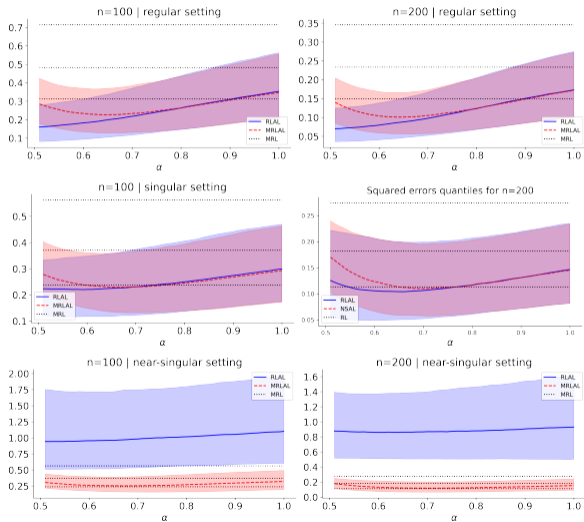
$$\text{prox}_{\lambda_n f_n}^{\bar{\mathbf{Q}}_n}(\check{\beta}_n^+) = \arg \min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{2} \|\check{\beta}_n^+ - \beta\|_{\bar{\mathbf{Q}}_n}^2 + \lambda_n f_n(\beta) \right\}$$

Let Assumptions **A3** and **MR1-MR3** hold with $c = 0$. If $\lambda_n \sqrt{n} \rightarrow 0$ and $\lambda_n n \rightarrow +\infty$ then:

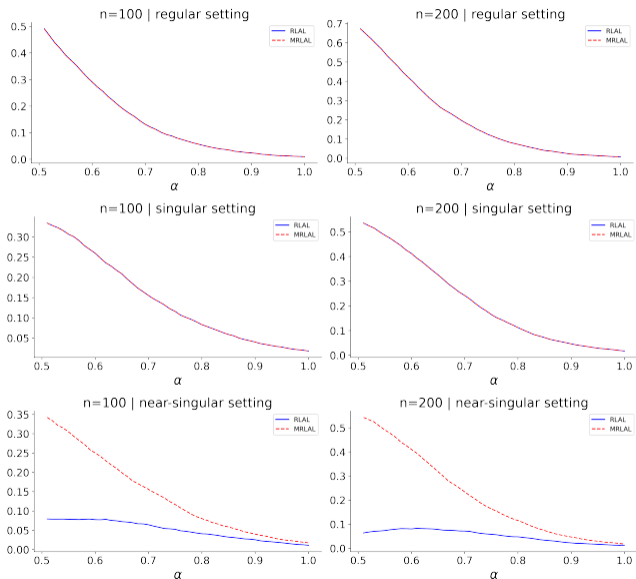
$$\sqrt{n} \begin{pmatrix} (\text{prox}_{\lambda_n f_n}^{\bar{\mathbf{Q}}_n}(\check{\beta}_n^+) - \beta_0^+)_{\mathcal{A}} \\ (\text{prox}_{\lambda_n f_n}^{\bar{\mathbf{Q}}_n}(\check{\beta}_n^+) - \beta_0^+)_{\mathcal{A}^c} \end{pmatrix} \rightarrow_d \begin{pmatrix} N((\mathbf{0})_{\mathcal{A}}, \sigma_0^2 [(\mathbf{Q}_0)_{\mathcal{A}}]^+) \\ (\mathbf{0})_{\mathcal{A}^c} \end{pmatrix}.$$

Moreover, the **consistent variable selection** property holds.

Squared errors $\left\| \hat{\beta}_n^s - \beta_0^+ \right\|_2^2$ for RLAL and MRLAL vs. $\lambda_n = n^{-\alpha}$ [$\alpha \in (0.5, 1)$]



Detection probabilities $\mathbb{P}(\hat{\mathcal{A}}_n = \mathcal{A})$ of RLAL and MRLAL vs. $\lambda_n = n^{-\alpha}$ [$\alpha \in (0.5, 1)$]



Conclusions

- A convenient class of PEs built with **minimum distance corrections** defined through **smooth proximal operators**
- A **unifying convex analysis framework** characterizing PEs' asymptotic properties:
 - Asymptotic distribution, Oracle property,...
- **Oracle** PE of minimum norm parameter in linear regression models with an **irregular** design
- Extensions:
 - Instrumental variables proximal estimation and inference under **weak instruments**
 - Estimation of stochastic discount factors in economies with **nearly redundant** payoffs

Thank you!