

éCetverikov, Denis N.; Sørensen, Jesper R-V

Working Paper

Analytic and bootstrap-after-cross-validation methods for selecting penalty parameters of highdimensional M-estimators

cemmap working paper, No. CWP20/21

Provided in Cooperation with:

Institute for Fiscal Studies (IFS), London

Suggested Citation: éCetverikov, Denis N.; Sørensen, Jesper R-V (2021) : Analytic and bootstrap-after-cross-validation methods for selecting penalty parameters of highdimensional M-estimators, cemmap working paper, No. CWP20/21, Centre for Microdata Methods and Practice (cemmap), London,
<https://doi.org/10.47004/wp.cem.2021.2021>

This Version is available at:

<https://hdl.handle.net/10419/241956>

Standard-Nutzungsbedingungen:

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

Terms of use:

Documents in EconStor may be saved and copied for your personal and scholarly purposes.

You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.

If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.

Analytic and bootstrap-after-cross-validation methods for selecting penalty parameters of high-dimensional M-estimators

Denis Chetverikov
Jesper R-V Sørensen

The Institute for Fiscal Studies
Department of Economics, UCL

cemmap working paper CWP20/21



Analytic and Bootstrap-after-Cross-Validation Methods for Selecting Penalty Parameters of High-Dimensional M-Estimators*

Denis Chetverikov[†] Jesper R.-V. Sørensen[‡]

April 20, 2021

Abstract

We develop two new methods for selecting the penalty parameter for the ℓ^1 -penalized high-dimensional M-estimator, which we refer to as the analytic and bootstrap-after-cross-validation methods. For both methods, we derive nonasymptotic error bounds for the corresponding ℓ^1 -penalized M-estimator and show that the bounds converge to zero under mild conditions, thus providing a theoretical justification for these methods. We demonstrate via simulations that the finite-sample performance of our methods is much better than that of previously available and theoretically justified methods.

Keywords: Penalty parameter selection, penalized M-estimation, high-dimensional models, sparsity, cross-validation, bootstrap.

1 Introduction

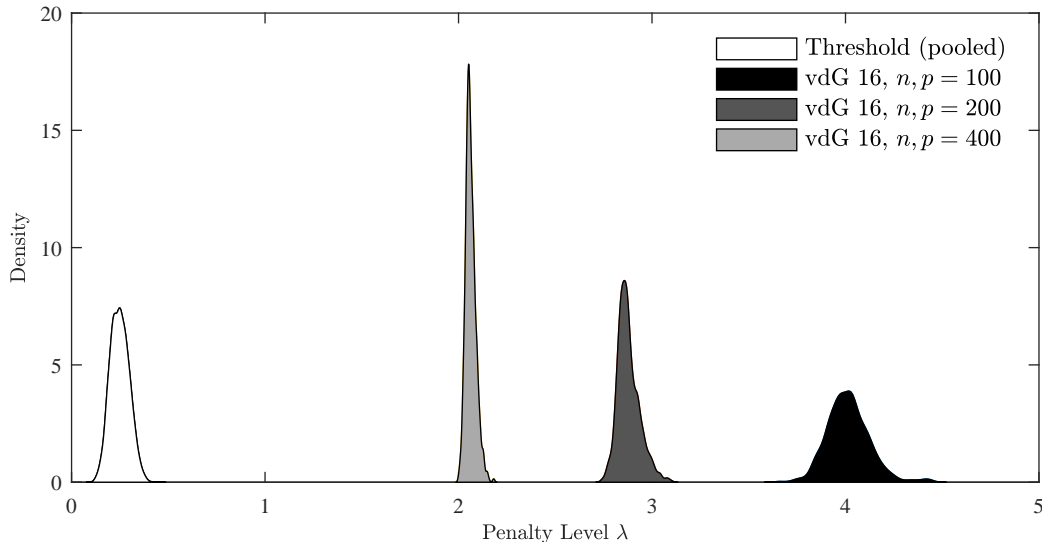
High-dimensional models have attracted substantial attention both in the econometrics and in the statistics/machine learning literature, e.g. see [Belloni et al. \(2018a\)](#) and [Hastie et al. \(2015\)](#), and ℓ^1 -penalized estimators have emerged among the most useful methods for learning parameters of such models. However, implementing these estimators requires a choice of the penalty parameter and with few notable exceptions, e.g. ℓ^1 -penalized linear mean and

*Date: April 2021. We thank Richard Blundell, Victor Chernozhukov, Bo Honoré, Whitney Newey, Joris Pinkse, Azeem Shaikh and Sara van de Geer for their insightful comments and discussions. Chetverikov's work was supported by NSF Grant SES - 1628889.

[†]Department of Economics, UCLA; e-mail: chetverikov@econ.ucla.edu.

[‡]Department of Economics, University of Copenhagen; e-mail: jrvs@econ.ku.dk.

Figure 1.1: Probability density functions of the smallest value (threshold) of the penalty parameter leading to all-zero estimated parameters and of the value of the penalty parameter obtained from the van de Geer (2016) method (vdG 16) in the setting of the ℓ^1 -penalized logit estimator. The figure demonstrates that the van de Geer penalty parameter value substantially exceeds the threshold value for the samples considered and thus yields the trivial, all-zero, estimates; see Section 6 for details.



quantile regression estimators, the choice of this penalty parameter in practice often remains unclear. Some methods, such as cross-validation and related sample splitting methods, tend to perform well in simulations but, as we discuss below, generally lack a sufficient theoretical justification. Other methods, such as those discussed in [van de Geer \(2016\)](#), are supported by a sound asymptotic theory but tend to perform poorly in moderate samples of practical relevance, often leading to trivial estimates, with all estimated parameters being exactly zero; see [Figure 1.1](#) for a demonstration in the case of the ℓ^1 -penalized logit estimator. In this paper, we deal with these problems and (i) propose two new methods for choosing penalty parameters in the context of ℓ^1 -penalized M-estimation, (ii) derive the supporting asymptotic theory, and (iii) demonstrate that our methods perform well in moderate samples.

We consider a model where the true value θ_0 of some parameter θ is given by the solution to an optimization problem

$$\theta_0 = \operatorname{argmin}_{\theta \in \Theta} \mathbb{E}[m(X^\top \theta, Y)], \quad (1.1)$$

where $m : \mathbf{R} \times \mathcal{Y} \rightarrow \mathbf{R}$ is a known (potentially nonsmooth) loss function that is convex in its first argument, $X = (X_1, \dots, X_p)^\top \in \mathcal{X} \subseteq \mathbf{R}^p$ a vector of candidate regressors, $Y \in \mathcal{Y}$ one or more outcome variables, and $\Theta \subseteq \mathbf{R}^p$ a convex parameter space. Prototypical loss functions are square-error loss and negative log-likelihood but the framework [\(1.1\)](#) also covers many other cross-sectional models and associated modern as well as classical estimation approaches

including logit and probit models, logistic calibration (Tan, 2017) covariate balancing (Imai and Ratkovic, 2014), and expectile regression (Newey and Powell, 1987). It also subsumes approaches to estimation of panel-data models such as the fixed-effects/conditional logit (Rasch, 1960) and trimmed least-absolute-deviations and least-squares methods for censored regression (Honoré, 1992), and partial likelihood estimation of heterogeneous panel models for duration (Chamberlain, 1985). We provide details on these examples in Section 2.¹

For the purpose of estimation, we assume access to a random sample $(X_1, Y_1), \dots, (X_n, Y_n)$ from the distribution of the pair (X, Y) , where the number p of candidate regressors in each $X_i = (X_{i1}, \dots, X_{ip})^\top$ may be (potentially much) larger than the sample size n , meaning that we cover high-dimensional models. Following the literature on high-dimensional models, we also assume that the vector $\theta_0 = (\theta_{01}, \dots, \theta_{0p})^\top$ is sparse in the sense that the number $s := \sum_{j=1}^p \mathbf{1}(\theta_{0j} \neq 0)$ of relevant regressors is much smaller than n .² With this sparsity assumption in mind, we study the ℓ^1 -penalized M-estimator

$$\widehat{\theta}(\lambda) \in \operatorname{argmin}_{\theta \in \Theta} \left\{ \frac{1}{n} \sum_{i=1}^n m(X_i^\top \theta, Y_i) + \lambda \|\theta\|_1 \right\}, \quad (1.2)$$

where $\|\theta\|_1 = \sum_{j=1}^p |\theta_j|$ denotes the ℓ^1 -norm of θ , and $\lambda \geq 0$ is a penalty parameter.

Implementing the estimator $\widehat{\theta}(\lambda)$ requires us to choose λ . To do so, we first extend the deterministic bound from Belloni and Chernozhukov (2011b) obtained for ℓ^1 -penalized quantile regression to the general setting of ℓ^1 -penalized M-estimators (1.2). In particular, we show that if $c_0 > 1$ is some constant, then there exists a constant C , depending on the distribution of the pair (X, Y) and c_0 , such that under mild regularity conditions, with probability approaching one, the event

$$\lambda \geq c_0 \max_{1 \leq j \leq p} \left| \frac{1}{n} \sum_{i=1}^n m'_1(X_i^\top \theta_0, Y_i) X_{ij} \right| \quad (1.3)$$

¹We consider the single-index setup solely for notational convenience. A more general setup entails L indices $\{X_{(\ell)}^\top \theta_{0(\ell)}\}_{\ell=1}^L$ and a loss function $m : \mathbf{R}^L \times \mathcal{Y} \rightarrow \mathbf{R}$. Multiple indices occur naturally in, e.g., multinomial models such as the multinomial and conditional logit models. Our treatment readily extends to L greater than one provided this constant does not depend on n .

²We take $s \geq 1$ throughout. This assumption is innocuous as we may always redefine s as $\max\{1, s\}$. Also, our notion of sparsity is *exact*: the number of nonzero coefficients is small. One may entertain weaker notions such as *approximate sparsity*, which allows for many nonzero but small coefficients. Simulation evidence suggests that our methods are relevant also in settings where only approximate sparsity is satisfied (see Section 6).

implies both

$$\|\widehat{\theta}(\lambda) - \theta_0\|_2 \leq C\sqrt{s} \left(\lambda + \sqrt{\frac{\ln(pn)}{n}} \right) \quad \text{and} \quad \|\widehat{\theta}(\lambda) - \theta_0\|_1 \leq Cs \left(\lambda + \sqrt{\frac{\ln(pn)}{n}} \right), \quad (1.4)$$

where $m'_1(t, y) := (\partial/\partial t) m(t, y)$ denotes the derivative of the loss function m with respect to its first argument (or a subgradient, if nondifferentiable). These bounds suggest the following principle: choose λ as small as possible subject to the event (1.3) occurring with high probability. We therefore wish to set $\lambda = c_0q(1 - \alpha)$, where

$$q(1 - \alpha) := (1 - \alpha)\text{-quantile of } \max_{1 \leq j \leq p} \left| \frac{1}{n} \sum_{i=1}^n m'_1(X_i^\top \theta_0, Y_i) X_{ij} \right|, \quad (1.5)$$

for some small user-specified probability tolerance level $\alpha \in (0, 1)$, e.g. $\alpha = .1$. This choice, however, is typically infeasible since the random variable in (1.5) depends on the unknown θ_0 . We thus have a vicious circle: to choose λ , we need an estimator of θ_0 , but to estimate θ_0 , we need to choose λ . In this paper, we offer two solutions to this problem, which constitute our key contributions.

To obtain our first solution, we show that whenever the loss function m is Lipschitz continuous with respect to its first argument, we can apply results from high-dimensional probability theory to derive an upper bound, say $\bar{q}(1 - \alpha)$, on $q(1 - \alpha)$ that does not depend on θ_0 and can be computed analytically from the available dataset. We can then set $\lambda = c_0\bar{q}(1 - \alpha)$, which we refer to as the *analytic method*. This method is computationally straightforward and, as we demonstrate by means of example, has several applications. On the other hand, it is not universally applicable as the loss function may or may not be Lipschitz continuous. For example, it works for the logit model but not for the probit model. Moreover, this method is somewhat conservative, in the sense that it yields a penalty satisfying $\lambda > c_0q(1 - \alpha)$.

To obtain our second solution, we show that even though the estimator $\widehat{\theta}(\lambda)$ based on λ chosen by cross-validation or its variants is generally difficult to analyze, it can be used to construct provably good (in a sense to be made clear later) estimators of the random vectors $m'_1(X_i^\top \theta_0, Y_i)X_i$. We are then able to derive an estimator, say $\widehat{q}(1 - \alpha)$, of $q(1 - \alpha)$ via bootstrapping, as discussed in Belloni et al. (2018a), and to set $\lambda = c_0\widehat{q}(1 - \alpha)$, which we refer to as the *bootstrap-after-cross-validation method*. This method is computationally somewhat more demanding than the analytic method, but it is generally much more widely applicable and nonconservative in the sense that it gives λ such that $\lambda \approx c_0q(1 - \alpha)$.³

³Both analytic and bootstrap-after-cross-validation methods require specifying the constant c_0 . While our

Drawing on simulations from a simple logit model, we illustrate the potential of our analytic and bootstrap-after-cross-validation methods. Our simulations indicate that, while both methods lead to useful estimates of θ_0 in the model (1.1) even in moderate samples, there may be significant gains from using the bootstrap-after-cross-validation method, even if the analytic method is also available. Moreover, both methods substantially outperform the choices of λ discussed in van de Geer (2016).

A key feature of our methods is that they yield bounds on both ℓ^1 and ℓ^2 estimation errors. In contrast, sample splitting methods typically yield bounds only for the excess risk $E_{X,Y}[m(X^\top \hat{\theta}(\lambda), Y) - m(X^\top \theta_0, Y)]$, e.g. see Lecue and Mitchell (2012). These bounds can be translated into the ℓ^2 estimation error $\|\hat{\theta}(\lambda) - \theta_0\|_2$, but it is not clear how to convert them into bounds on the ℓ^1 estimation error $\|\hat{\theta}(\lambda) - \theta_0\|_1$. A bound of the ℓ^1 type is crucial when we are interested in estimating dense functionals $a'\theta_0$ of θ_0 with $a \in \mathbf{R}^p$ being a vector of loadings with many nonzero components; see Belloni et al. (2018a) for details. Moreover, the bounds on the ℓ^1 estimation error are needed to perform inference on components of θ_0 via double machine learning, as in Belloni et al. (2018b). In the same fashion, when λ is selected using cross-validation, neither ℓ^1 nor ℓ^2 estimation error bounds are typically known. The only exception we are aware of is the linear mean regression model. The bounds for this model have been derived in Chetverikov et al. (2016) and Miolane and Montanari (2018), but even in this special case the bounds derived are not as sharp as those provided here.

The literature on learning parameters of high-dimensional models via ℓ^1 -penalized M-estimation is large. Instead of listing all existing papers, we therefore refer the interested reader to the excellent textbook treatment in Wainwright (2019) and focus here on only a few key references. van de Geer (2008, 2016) derives bounds on the estimation errors of general ℓ^1 -penalized M-estimators (1.2) and provides some choices of the penalty parameter λ . As discussed above, however, her recommendations give values of λ that are so large that the resulting estimators are typically trivial in moderate samples, with all coefficients being exactly zero (cf. Figure 1.1). Because of this issue, van de Geer (2008) remarks that her results should only be seen as an indication that her theory has something to say about finite sample sizes, and that other methods to choose λ should be used in practice. Negahban et al. (2012) develop results in a very general setting, and when specialized to our setting (1.2) their results become quite similar to our result that the bounds (1.4) hold under the event (1.3). The same authors also note that a challenge to using these results in practice is that the random variable in (1.3) is usually impossible to compute because it depends on

theory only requires that $c_0 > 1$, simulations suggest that both ℓ^1 and ℓ^2 estimation errors of the ℓ^1 -penalized M-estimator are increasing in c_0 for $c_0 > 1$. We therefore recommend setting $c_0 = 1.1$, which reflects one of the standard recommendations in the LASSO literature (see, e.g., Belloni and Chernozhukov, 2011a).

the unknown vector θ_0 . It is exactly this challenge that we overcome in this paper. [Belloni and Chernozhukov \(2011b\)](#) study high-dimensional quantile regression and note that the distribution of the random variable in (1.3) is in this case pivotal, making the choice of the penalty parameter simple. However, quantile regression is the only setting we are aware of in which the distribution of the random variable in (1.3) is pivotal.⁴ Finally, [Ninomiya and Kawano \(2016\)](#) consider information criteria for the choice of the penalty parameter λ but focus on fixed- p asymptotics, thus excluding high-dimensional models.

The rest of the paper is organized as follows. In Section 2 we provide a portfolio of examples that constitute possible applications of our methods. We refer to several of these examples in later sections. In Section 3 we develop bounds on the estimation error of the ℓ^1 -penalized M-estimator, which motivate our methods to choose the penalty parameter. We discuss the analytic method in Section 4 and the bootstrap-after-cross-validation method in Section 5. In Section 6 we illustrate our methods via a simulation study and compare them with existing methods. We give all the proofs in the Appendix, which also provides low-level conditions sufficient for some of the assumptions made in the main text.

Notation

Throughout $W_i := (X_i, Y_i), i \in \{1, \dots, n\}$, denotes n independent copies of a random vector $W := (X, Y) \in \mathcal{W}$. The distribution P of W , as well as the dimension p of the vector X and the number of nonzero components of the vector θ_0 may change with the sample size n , but we suppress this potential dependence. $\mathbb{E}[f(W)]$ denotes the expectation of a function f of W computed with respect to P , and $\mathbb{E}_n[f(W_i)] := n^{-1} \sum_{i=1}^n f(W_i)$ abbreviates the sample average. When only a nonempty subset $I \subsetneq \{1, \dots, n\}$ is in use, we write $\mathbb{E}_I[f(W_i)] := |I|^{-1} \sum_{i \in I} f(W_i)$ for the subsample average. For a set of indices $I \subseteq \{1, \dots, n\}$, I^c denotes the elements of $\{1, \dots, n\}$ not in I . Given a vector $\delta \in \mathbf{R}^p$ and a nonempty set of indices $J \subseteq \{1, \dots, p\}$, we let δ_J denote the vector in \mathbf{R}^p with coordinates given by $\delta_{Jj} = \delta_j$ if $j \in J$ and zero otherwise. We denote its ℓ^1 , ℓ^2 and ℓ^∞ norms by $\|\delta\|_1$, $\|\delta\|_2$, and $\|\delta\|_\infty$, respectively. We abbreviate $a \vee b := \max\{a, b\}$ and $a \wedge b := \min\{a, b\}$, and take $n \geq 3$ and $p \geq 2$ throughout. We introduce more notation as needed in the appendices.

⁴With a known censoring propensity, the linear programming estimator of [Buchinsky and Hahn \(1998\)](#) for censored quantile regression boils down to a variant of quantile regression and, therefore, leads to pivotality of the right-hand side of (1.3). However, known censoring propensity seems like a very special case.

2 Examples

In this section we discuss a variety of models that fit into the M-estimation framework (1.1) with the loss function $m(t, y)$ being convex in its first argument. We include models for cross-sectional data (Examples 1–5), panel data (Examples 6 and 7) and panel data for duration (Example 8). The examples cover both discrete and continuous outcomes in likelihood and nonlikelihood settings with smooth as well as kinked loss functions.

Example 1 (Binary Response Model). A relatively simple model fitting our framework is the *binary response model*, i.e. a model for an outcome $Y \in \{0, 1\}$ with

$$P(Y = 1|X) = F(X^\top \theta_0),$$

for a known cumulative distribution function (CDF) $F : \mathbf{R} \rightarrow [0, 1]$. The log-likelihood of this model yields the following loss function:

$$m(t, y) = -y \ln F(t) - (1 - y) \ln(1 - F(t)). \quad (2.1)$$

The *logit* model arises here by setting $F(t) = 1/(1 + e^{-t}) =: \Lambda(t)$, the standard logistic CDF, and the loss function reduces in this case to

$$m(t, y) = \ln(1 + e^t) - yt. \quad (2.2)$$

The *probit* model arises by setting $F(t) = \int_{-\infty}^t (2\pi)^{-1/2} e^{-u^2/2} du =: \Phi(t)$, the standard normal CDF, and the loss function in this case becomes

$$m(t, y) = -y \ln \Phi(t) - (1 - y) \ln(1 - \Phi(t)). \quad (2.3)$$

Note that the loss functions in both (2.2) and (2.3) are convex in t .

More generally, any binary response model with both F and $1 - F$ being log-concave leads to a loss (2.1) that is convex in t . For these log-concavities it suffices that F admits a probability density function (PDF) $f = F'$, which is itself log-concave (Pratt, 1981, Section 5). Both the standard logistic and standard normal PDFs are log-concave. Also, $\ln f$ is concave whenever $f(t) \propto e^{-|t|^\alpha}$ for some $\alpha \geq 1$ or $f(t) \propto t^{a-1}e^{-t}$ for $t \geq 0$ and some $a \geq 1$, the extreme cases being the Laplace and exponential distributions, respectively. Other examples of distributions for which f is log-concave can be found in the Gumbel, Weibull, Pareto and beta families (Pratt, 1981, Section 6). A t -distribution with $0 < \nu < \infty$ degrees of freedom (the standard Cauchy arising from $\nu = 1$) does not have a log-concave

density. However, both its CDF and complementary CDF are log-concave (*ibid.*). \square

Example 2 (Ordered Response Model). Consider the *ordered response model*, i.e. a model for an outcome $Y \in \{0, 1, \dots, J\}$ with

$$P(Y = j|X) = F(\alpha_{j+1} - X^\top \theta_0) - F(\alpha_j - X^\top \theta_0), \quad j \in \{0, 1, \dots, J\},$$

for a known CDF $F : \mathbf{R} \rightarrow [0, 1]$ and known cut-off points $-\infty = \alpha_0 < \alpha_1 < \dots < \alpha_J < \alpha_{J+1} = +\infty$. (We here interpret $F(-\infty)$ as zero and $F(+\infty)$ as one.) The log-likelihood of this model yields the loss function

$$m(t, y) = - \sum_{j=0}^J \mathbf{1}(y = j) \ln(F(\alpha_{j+1} - t) - F(\alpha_j - t)), \quad (2.4)$$

that is convex in t for any distribution F admitting a log-concave PDF $f = F'$ (Pratt, 1981, Section 3). See Example 1 for specific distributions satisfying this criterion. \square

Example 3 (Logistic Calibration). In the setting of average treatment effect estimation under a conditional independence assumption with high-dimensional vector of controls, consider the *logit propensity score model*

$$P(Y = 1|X) = \Lambda(X^\top \theta_0), \quad (2.5)$$

where $Y \in \{0, 1\}$ is a treatment indicator, X a vector of controls, and Λ the logistic CDF. Using (1.1), θ_0 can be identified with the logistic loss function in (2.2). However, as shown by Tan (2017), θ_0 can also be identified using (1.1) with the logistic *calibration* loss

$$m(t, y) = ye^{-t} + (1 - y)t, \quad (2.6)$$

which is convex in t as well. As demonstrated by Tan (2017), using this alternative loss function gives substantial advantages. In particular, it leads to average treatment effect estimators that enjoy particularly nice robustness properties. Specifically, under some conditions, these estimators remain root- n consistent and asymptotically normal even if the model (2.5) is misspecified. \square

Example 4 (Logistic Balancing). In the same setting as that of the previous example, the *covariate balancing approach* (Imai and Ratkovic, 2014) amounts to specifying a parametric model for the treatment indicator $Y \in \{0, 1\}$,

$$P(Y = 1|X) = F(X^\top \theta_0)$$

and ensuring covariate balance in the sense that

$$\mathbb{E} \left[\left\{ \frac{Y}{F(X^\top \theta_0)} - \frac{1 - Y}{1 - F(X^\top \theta_0)} \right\} X \right] = \mathbf{0}.$$

Balancing here amounts to enforcing a collection of moment conditions and is therefore naturally studied in a generalized method of moments (GMM) framework. However, specifying F to be the logistic CDF Λ , covariate balancing can be achieved via M-estimation of θ_0 based on the loss function

$$m(t, y) = (1 - y) e^t + y e^{-t} + (1 - 2y) t,$$

which is also convex in t . [See [Tan \(2017\)](#) for details.] □

Example 5 (Expectile Model). [Newey and Powell \(1987\)](#) study the conditional τ -*expectile model* $\mu_\tau(Y|X) = X^\top \theta_0$, where $\tau \in (0, 1)$, and propose the *asymmetric least squares* (ALS) estimator of θ_0 in this model. This estimator can be understood as an M-estimator with the loss function

$$m(t, y) = \rho_\tau(y - t), \tag{2.7}$$

where $\rho_\tau : \mathbf{R} \rightarrow \mathbf{R}$ is the piecewise quadratic and continuously differentiable function defined by

$$\rho_\tau(u) = |\tau - \mathbf{1}(u < 0)| u^2 = \begin{cases} (1 - \tau) u^2, & u < 0, \\ \tau u^2, & u \geq 0, \end{cases}$$

a smooth analogue of the ‘check’ function known from the quantile regression literature. This estimator can also be interpreted as a maximum likelihood estimator when model disturbances arise from a normal distribution with unequal weights placed on positive and negative disturbances ([Aigner et al., 1976](#)). Note that $m(t, y)$ in (2.7) is convex but not twice differentiable unless $\tau = 1/2$. □

Example 6 (Panel Logit Model). Consider the *panel logit model*

$$P(Y_t = 1 | X, \alpha, Y_0, \dots, Y_{t-1}) = \Lambda(\alpha + X_t^\top \theta_0), \quad t = 1, 2,$$

where $Y = (Y_1, Y_2)^\top$ is a pair of outcome variables, $X = (X_1^\top, X_2^\top)^\top$ is a vector of regressors, and α is a unit-specific unobserved fixed effect. [Rasch \(1960\)](#) shows⁵ that θ_0 in this model can be identified by $\theta_0 = \operatorname{argmin}_{\theta \in \mathbf{R}^p} \mathbb{E}[m((X_1 - X_2)^\top \theta, Y)]$, where

$$m(t, y) = \mathbf{1}(y_1 \neq y_2) [\ln(1 + e^t) - y_1 t], \tag{2.8}$$

⁵See also [Chamberlain \(1984, Section 3.2\)](#) and [Wooldridge \(2010, Section 15.8.3\)](#).

which is convex in t . □

Example 7 (Panel Censored Model). Consider the *panel censored model*

$$Y_t = \max(0, \alpha + X_t^\top \theta_0 + \varepsilon_t), \quad t = 1, 2,$$

where $Y = (Y_1, Y_2)^\top \in \mathbf{R}_+^2$ is a pair of outcome variables, $X = (X_1^\top, X_2^\top)^\top$ is a vector of regressors, α is a unit-specific unobserved fixed effect, and ε_1 and ε_2 are unobserved error terms. [Honoré \(1992\)](#) shows that under certain conditions—including exchangeability of ε_1 and ε_2 conditional on X_1, X_2 , $\alpha - \theta_0$ in this model can be identified by $\theta_0 = \operatorname{argmin}_{\theta \in \mathbf{R}^p} \mathbb{E}[m((X_1 - X_2)^\top \theta, Y)]$, with m being the *trimmed loss* function

$$m(t, y) = \begin{cases} \Xi(y_1) - (y_2 + t) \xi(y_1), & t \leq -y_2, \\ \Xi(y_1 - y_2 - t), & -y_2 < t < y_1, \\ \Xi(-y_2) - (t - y_1) \xi(-y_2), & y_1 \leq t, \end{cases} \quad (2.9)$$

and either $\Xi = |\cdot|$ or $\Xi = (\cdot)^2$ and ξ its derivative (when defined).⁶ These choices lead to *trimmed least absolute deviations* (LAD) and *trimmed least squares* (LS) estimators, respectively, both of which have loss functions convex in t . Here, $\Xi = |\cdot|$ leads to a nondifferentiable loss. □

Example 8 (Panel Duration Model). Consider the *panel duration model* with a log-linear specification:

$$\ln h_t(y) = X_t^\top \theta_0 + h_0(y), \quad t = 1, 2,$$

where h_t denotes the hazard for spell t and both h_0 and h_t are allowed to be unit-specific. This model is a special case of the duration models studied in [Chamberlain \(1985, Section 3.1\)](#). Chamberlain presumes that the spells Y_1 and Y_2 are (conditionally) independent of each other and shows that the partial log-likelihood contribution is⁷

$$\theta \mapsto \mathbf{1}(Y_1 < Y_2) \ln \Lambda((X_1 - X_2)^\top \theta) + \mathbf{1}(Y_1 \geq Y_2) \ln (1 - \Lambda((X_1 - X_2)^\top \theta)).$$

The implied loss function

$$m(t, y) = \ln(1 + e^t) - \mathbf{1}(y_1 < y_2) t \quad (2.10)$$

is of the logit form (see [Example 1](#)), hence convex in t . With more than two completed

⁶When $\Xi = |\cdot|$, we set $\xi(0) := 0$ to make (2.9) consistent with formulas in [Honoré \(1992\)](#).

⁷See also [Lancaster \(1992, Chapter 9, Section 2.10.2\)](#).

spells, the partial log-likelihood takes a conditional-logit form (*ibid.*), and the resulting loss is therefore still a convex function (albeit involving multiple indices). \square

3 Nonasymptotic Bounds on Estimation Error

In this section, we derive bounds on the error of the ℓ^1 -penalized M-estimator (1.2) in the ℓ^1 and ℓ^2 norms. The argument reveals which quantities one needs to control in order to ensure good behavior of the estimator, motivating the choice of the penalty parameter λ in the following sections. We split the section into two subsections. In Section 3.1, we derive bounds via an empirical error function. In Section 3.2, we derive a bound on the empirical error function itself.

3.1 Bounds via Empirical Error Function

Denote

$$M(\theta) := \mathbb{E}[m(X^\top \theta, Y)] \quad \text{and} \quad \widehat{M}(\theta) := \mathbb{E}_n[m(X_i^\top \theta, Y_i)], \quad \theta \in \Theta,$$

Also, let

$$T := \{j \in \{1, \dots, p\}; \theta_{0j} \neq 0\}$$

and for any $c > 1$, let $\mathcal{R}(c)$ denote the *restricted set*

$$\mathcal{R}(c) := \{\delta \in \mathbf{R}^p; \|\delta_{T^c}\|_1 \leq c\|\delta_T\|_1\}.$$

In addition, fix $c_0 > 1$, and define the (random) *empirical error function* $\epsilon : \mathbf{R}_+ \rightarrow \mathbf{R}_+$ by

$$\epsilon(u) := \sup_{\substack{\delta \in \mathcal{R}(\bar{c}_0), \\ \|\delta\|_2 \leq u}} |(\mathbb{E}_n - \mathbb{E}) [m(X_i^\top (\theta_0 + \delta), Y_i) - m(X_i^\top \theta_0, Y_i)]|, \quad u \in \mathbf{R}_+,$$

where $\bar{c}_0 := (c_0 + 1) / (c_0 - 1)$. Moreover, define the *excess risk function* $\mathcal{E} : \Theta \rightarrow \mathbf{R}_+$ by

$$\mathcal{E}(\theta) := M(\theta) - M(\theta_0) = \mathbb{E} [m(X^\top \theta, Y) - m(X^\top \theta_0, Y)], \quad \theta \in \Theta.$$

In this subsection, we derive bounds on $\|\widehat{\theta}(\lambda) - \theta_0\|_1$ and $\|\widehat{\theta}(\lambda) - \theta_0\|_2$ via the empirical error function. Our bounds will be based on the following four assumptions.

Assumption 1 (Parameter Space). *The parameter space Θ is a convex subset of \mathbf{R}^p for which θ_0 is interior.*

Assumption 2 (Convexity). *The loss function $t \mapsto m(t, y)$ is convex for all $y \in \mathcal{Y}$.*

Assumption 3 (Differentiability and Integrability). *The derivative $m'_1(X^\top\theta, Y)$ exists almost surely and $E[|m(X^\top\theta, Y)|] < \infty$ for all $\theta \in \Theta$.*

Assumption 4 (Margin). *There exist finite constants $c_M, c'_M > 0$ such that*

$$\theta \in \Theta \text{ and } \|\theta - \theta_0\|_1 \leq c'_M \quad \text{imply} \quad \mathcal{E}(\theta) \geq c_M \|\theta - \theta_0\|_2^2.$$

Assumption 1 is a minor regularity condition. Assumption 2 is satisfied in all examples from the previous section. This assumption implies that the (random) function \widehat{M} is convex, hence subdifferentiable at the interior point θ_0 (Rockafellar, 1970, Theorem 23.4). The first part of Assumption 3 is satisfied in the sure sense in all examples from the previous section except for Example 7 with the trimmed LAD loss function, where it is satisfied if the conditional distribution of $(\varepsilon_1, \varepsilon_2)$ given (α, X_1, X_2) is continuous. In fact, for all our results except for those in Section 5.2, it would be sufficient to assume that the derivative $m'_1(X^\top\theta, Y)$ exists almost surely for $\theta = \theta_0$ only. The second part of Assumption 3 is a minor regularity condition. Assumption 4 is expected to be satisfied in most applications as well. We provide a set of conditions that are sufficient for Assumption 4 at the end of this subsection, which are then straightforward to verify in all of the examples from the previous section.

Next, for some score selection S from the subdifferential $\partial\widehat{M}(\theta_0)$, which always exist by Assumption 2 and is almost surely singleton by Assumption 3, and for some constants $\lambda_\epsilon, \bar{\lambda} > 0$, define the events

$$\begin{aligned} \mathcal{S} &:= \{\lambda \geq c_0 \|S\|_\infty\}, && \text{(score domination)} \\ \mathcal{L} &:= \{\lambda \leq \bar{\lambda}\}, && \text{(penalty majorization)} \\ \mathcal{E} &:= \{\epsilon(u_0) \leq \lambda_\epsilon u_0\}, && \text{(empirical error control)} \end{aligned}$$

where

$$u_0 := \frac{2}{c_M} (\lambda_\epsilon + (1 + \bar{c}_0) \bar{\lambda} \sqrt{s}). \quad (3.1)$$

Here, the event \mathcal{S} ensures that the penalty is large enough to provide a sufficient level of regularization and the event \mathcal{L} ensures that the penalty is not too large. For the purpose of the deterministic calculation of this section, the event \mathcal{L} plays little to no role, and one may enforce it by simply setting $\bar{\lambda} = \lambda$. However, in later sections, the penalty level will be a random quantity, and \mathcal{L} facilitates easy reference. The constant λ_ϵ appearing in the event \mathcal{E} represents a deterministic modulus of continuity of the empirical error function ϵ in a neighborhood of zero of size u_0 .

We now have the following result on the error bounds for the ℓ^1 -penalized M-estimator:

Theorem 1 (Nonasymptotic Bounds). *Let Assumptions 1, 2, 3, and 4 hold and suppose that $(1 + \bar{c}_0)u_0\sqrt{s} \leq c'_M$. Then on the event $\mathcal{S} \cap \mathcal{L} \cap \mathcal{E}$, we have both*

$$\|\widehat{\theta}(\lambda) - \theta_0\|_2 \leq \frac{2}{c_M} (\lambda_\epsilon + (1 + \bar{c}_0) \bar{\lambda} \sqrt{s}) \quad \text{and} \quad (3.2)$$

$$\|\widehat{\theta}(\lambda) - \theta_0\|_1 \leq \frac{2(1 + \bar{c}_0)}{c_M} (\lambda_\epsilon \sqrt{s} + (1 + \bar{c}_0) \bar{\lambda} s). \quad (3.3)$$

This theorem motivates our choices of the penalty parameter λ . In particular, we will show in the next subsection that empirical error control (\mathcal{E}) holds with probability approaching one if $\lambda_\epsilon = C_\epsilon \sqrt{s \log(pn)/n}$ for a sufficiently large constant $C_\epsilon > 0$. Therefore, setting $\bar{\lambda} = \lambda$, such that penalty majorization (\mathcal{L}) holds trivially, Theorem 1 gives the same bounds as those appearing in the Introduction and we arrive at the following principle: choose λ as small as possible subject to the constraint that the score domination event $\mathcal{S} = \{\lambda \geq c_0 \|S\|_\infty\}$ occurs with high probability. It is exactly this principle that guides our choices of λ in the following sections.

Remark 1 (On Uniqueness). Theorem 1 actually concerns the set of optimizers to the convex minimization problem (1.2) for a fixed value of λ . While objective function \widehat{M} is convex, it need not be strictly convex, such that the global minimum may be attained at more than one point $\widehat{\theta}(\lambda)$. The bounds stated here (and in what follows) hold for any of these optimizers. \square

Remark 2 (On Loss Structure). The proof of Theorem 1 requires neither the index structure placed on the loss function nor the separation of a datum W into regressors X and outcome(s) Y . The deterministic bounds in Theorem 1 continue to hold if $(w, \theta) \mapsto m(x^\top \theta, y)$ is replaced by a general loss $(w, \theta) \mapsto m_\theta(w)$, which is convex in θ and P -integrable in w . \square

Remark 3 (On Quadratic Margin). Our convexity, interiority, and differentiability assumptions suffice to show that the risk M is differentiable at θ_0 (Bertsekas, 1973, Proposition 2.3). Consequently, our estimand θ_0 must satisfy the population first-order condition $\nabla M(\theta_0) = \mathbf{0}$. Assumption 4 therefore amounts to assuming that the population criterion M admits a quadratic *margin* at θ_0 . The proof of Theorem 1 readily extends to more-than-quadratic margin behavior (thus leading to a flatter population criterion at θ_0 and, hence, worse identification). The name *margin condition* appears to originate from Tsybakov (2004, Assumption A1), who invokes a similar assumption in a classification context. van de Geer (2008, Assumption B) contains a more general formulation of margin behavior for estimation purposes. We concentrate on the (focal) quadratic case for the sake of simplicity. \square

We end this subsection with a proposition providing somewhat more primitive conditions

for ensuring a quadratic margin, i.e. Assumption 4.⁸

Proposition 1 (Quadratic Margin). *Let Assumptions 1 and 2 hold. In addition, suppose that $\|X\|_\infty \leq C_X$ with probability one, the smallest eigenvalue of $E[XX^\top]$ is at least $c_e \in \mathbf{R}_{++}$, and $t \mapsto E[m(t, Y)|X]$ is twice differentiable in a neighborhood $\{t \in \mathbf{R}; |t - X^\top \theta_0| < c_1\}$ of $X^\top \theta_0$ with radius $c_1 \in \mathbf{R}_{++}$ and with second derivative at least $c_2 \in \mathbf{R}_{++}$ with probability one. Then Assumption 4 holds with $c'_M \in (0, c_1/C_X)$ and $c_M \in (0, c_2 c_e/2]$.*

3.2 Empirical Error Function Control

In this subsection, we consider the problem of gaining control over the empirical error event $\mathcal{E} = \{\epsilon(u_0) \leq \lambda_\epsilon u_0\}$. More precisely, we present conditions under which one may ensure a linear modulus of continuity of the function ϵ in a neighborhood of zero with high probability. To do so, let \mathcal{S}_{p-1} denote the unit sphere in \mathbf{R}^p ,

$$\mathcal{S}_{p-1} := \{\delta \in \mathbf{R}^p; \|\delta\|_2 = 1\}.$$

We establish empirical error control using the following three assumptions.

Assumption 5 (Boundedness). *There exists a finite constant $C_X > 0$ such that $\|X\|_\infty \leq C_X$ almost surely.*

Assumption 6 (Locally Lipschitz Loss). *There exist finite constants $c_L, C_L > 0$ and a function $L : \mathcal{W} \rightarrow \mathbf{R}_+$ such that*

1. *for all $w = (x, y) \in \mathcal{W}$ and all $(t_1, t_2) \in \mathbf{R}^2$ satisfying $|t_1| \vee |t_2| \leq c_L$,*

$$|m(x^\top \theta_0 + t_1, y) - m(x^\top \theta_0 + t_2, y)| \leq L(w) |t_1 - t_2|;$$

2. $E[L(W)^4] \leq (C_L/2)^4$.

Assumption 7 (Weighted Population Design Matrix). *There exists a finite constant $C_{L,e} > 0$ such that all eigenvalues of the matrix $E[L(W)^2 XX^\top]$ are bounded from above by $C_{L,e}^2$.*

Assumption 5 is a regularity condition that some researchers may find rather strong but we emphasize that it can be relaxed. We have chosen to impose it in order to abstract from unnecessary technicalities. Note also that the same assumption was used, for example, in van de Geer (2008). Assumption 6.1 requires that the function $t \mapsto m(x^\top \theta_0 + t, y)$ is locally

⁸We state the result as a proposition, since we do not provide the regularity conditions necessary for interchanging the order of differentiation and integration.

Lipschitz continuous for all $w = (x, y) \in \mathcal{W}$ with Lipschitz constant $L(w)$, and Assumption 6.2 requires that the fourth moment of $L(W)$ is finite. Given that every convex function $f : C \rightarrow \mathbf{R}$ is Lipschitz relative to any compact subset S of the interior of its domain C (Rockafellar, 1970, Theorem 10.4), it follows that Assumption 6.1 is actually implied by Assumption 2, and so Assumption 6 should be regarded as a mild regularity condition restricting the moments of the random variable $L(W)$. At the end of this subsection, we illustrate the calculation of this random variable and implied restrictions on the data-generating process via the examples on binary and ordered choice models from Section 2. Finally, Assumption 7 restricts eigenvalues of the weighted population design matrix. This assumption is similar to conditions often imposed in the literature on high-dimensional models, where it is assumed that the eigenvalues of the matrix $\mathbb{E}[XX^\top]$ are bounded from above.

We now present a result showing that one may take $C_\epsilon \sqrt{s \ln(pn)}/n$ as the high-probability local modulus of continuity λ_ϵ appearing in the empirical error event $\mathcal{E} = \{\epsilon(u_0) \leq \lambda_\epsilon u_0\}$:

Lemma 1 (Empirical Error Bound). *Let Assumptions 5, 6, and 7 hold, and define the finite constant $C_\epsilon := 16\sqrt{2}(1 + \bar{c}_0)C_L C_X > 0$. Then provided $s \ln(pn) \geq 16C_{L,\epsilon}^2/C_\epsilon^2$ and $0 < u \leq c_L/[(1 + \bar{c}_0)C_X\sqrt{s}]$, we have*

$$\epsilon(u) \leq C_\epsilon u \sqrt{s \ln(pn)}/n$$

with probability at least $1 - 5n^{-1}$.

Remark 4 (Alternative Nonasymptotic Bounds). If the loss function m is globally Lipschitz in its first argument, then Assumption 6 holds with the function L being a constant. In this case, symmetrization, contraction, and concentration arguments may be used to bound the modified empirical error

$$\tilde{\epsilon}(u) := \sup_{\|\delta\|_1 \leq u} |(\mathbb{E}_n - \mathbb{E}) [m(X_i^\top(\theta_0 + \delta), Y_i) - m(X_i^\top\theta_0, Y_i)]|, \quad u \in \mathbf{R}_+,$$

now defined with respect to the ℓ^1 norm and without the restricted set $\mathcal{R}(\bar{c}_0)$. This is the approach taken by van de Geer (2008), who shows that there exists a constant \tilde{C} such that with probability approaching one,

$$\tilde{\epsilon}(u)/u \leq \tilde{C} \left(\sqrt{\frac{\ln p}{n}} + \frac{\ln p}{n} \right), \quad u \in \mathbf{R}_{++}.$$

She then demonstrates that useful bounds on the estimation error of $\hat{\theta}(\lambda)$ can be derived if λ is chosen to exceed the right-hand side of this inequality, which motivates alternative

methods to choose λ . Unfortunately, \tilde{C} typically relies on design constants unknown to the researcher. Moreover, even if these constants were known, the resulting values of \tilde{C} would typically be prohibitively large, yielding choices of λ leading to trivial estimates of the vector θ in moderate samples; see Section 6 for simulation results based on the choices in van de Geer (2008). Our bounds therefore seem more suitable for devising methods to choose λ . \square

We conclude this section by illustrating Assumption 6 via the binary and ordered response model examples from Section 2.

Example 1 (Binary Response Model, Continued). The logit loss function (2.2) is differentiable in t with $m'_1(t, y) = \Lambda(t) - y$. The logit loss function is therefore 1-Lipschitz in t , and so in this case one can take $L(w) = 1$ for all $w \in \mathcal{W}$, making Assumption 6 trivial. The probit loss function (2.3) is differentiable in t with

$$m'_1(t, y) = \frac{\varphi(t)}{\Phi(t)[1 - \Phi(t)]} [\Phi(t) - y],$$

and $\varphi(t) = (2\pi)^{-1/2} e^{-t^2/2}$ being the standard normal PDF. One can show that $m'_1(t, 1) = -\varphi(t)/\Phi(t) \sim t$ as $t \rightarrow -\infty$ and $m'_1(t, 0) = \varphi(t)/[1 - \Phi(t)] \sim t$ as $t \rightarrow \infty$, so the probit loss function is *not* globally Lipschitz in t . However, a mean-value argument shows that for any finite $c > 0$, $|t_1| \vee |t_2| \leq c$, and $w = (x, y) \in \mathcal{W}$,

$$|m(x^\top \theta_0 + t_1, y) - m(x^\top \theta_0 + t_2, y)| \leq \sup_{|t| \leq c} \frac{\varphi(x^\top \theta_0 + t)}{\Phi(x^\top \theta_0 + t)[1 - \Phi(x^\top \theta_0 + t)]} |t_1 - t_2|.$$

The function $\varphi/[\Phi(1 - \Phi)]$ is convex and even, so the right-hand side supremum is attained at both boundary points $\pm c$. It follows that the probit loss function is locally Lipschitz with

$$L(w) = \sup_{|t| \leq c} \frac{\varphi(x^\top \theta_0 + t)}{\Phi(x^\top \theta_0 + t)[1 - \Phi(x^\top \theta_0 + t)]} \leq 2(|x^\top \theta_0| + c),$$

see (1.2.2) in Adler and Taylor (2007). In this case, Assumption 6 therefore reduces to the requirement that $E[|X^\top \theta_0|^4]$ is bounded from above, which is a very mild regularity condition.

More generally, let F admit an everywhere positive log-concave PDF $f = F'$. Then the binary-response loss function (2.1) is differentiable with partial derivative

$$m'_1(t, y) = \frac{f(t)}{F(t)[1 - F(t)]} [F(t) - y]. \quad (3.4)$$

Given the binary outcome, we have $|F(t) - y| \leq 1$, and so, for any distribution such that $f/[F(1 - F)]$ is also bounded from above, the binary-response loss function is L -Lipschitz

in t with

$$L = \sup_{t \in \mathbf{R}} \frac{f(t)}{F(t)[1 - F(t)]}.$$

For example, one can show that the t -distribution with $0 < \nu < \infty$ degrees of freedom satisfies $f_\nu(t)/[1 - F_\nu(t)] \sim t/(1 + t^2/\nu)$ as $t \rightarrow \infty$ and $-f_\nu(t)/F_\nu(t) \sim t/(1 + t^2/\nu)$ as $t \rightarrow -\infty$, demonstrating that the resulting loss function is globally Lipschitz and that Assumption 6 holds trivially in this case. \square

Example 2 (Ordered Response Model, Continued). With F admitting an everywhere positive log-concave PDF $f = F'$, the loss (2.4) is differentiable with partial derivative

$$m'_1(t, y) = \sum_{j=0}^J \mathbf{1}(y = j) \frac{f(\alpha_{j+1} - t) - f(\alpha_j - t)}{F(\alpha_{j+1} - t) - F(\alpha_j - t)}. \quad (3.5)$$

(We here interpret $f(\pm\infty)$ and $F(-\infty)$ as zero and $F(+\infty)$ as one.) It follows from the mean-value theorem that the ordered-response loss is L -Lipschitz in t for any distribution F and cut-off points $-\infty = \alpha_0 < \alpha_1 < \dots < \alpha_J < \alpha_{J+1} = +\infty$ such that

$$L = \max_{0 \leq j \leq J} \sup_{t \in \mathbf{R}} \left| \frac{f(\alpha_{j+1} - t) - f(\alpha_j - t)}{F(\alpha_{j+1} - t) - F(\alpha_j - t)} \right|$$

is finite. For the logistic distribution $F = \Lambda$ and $f = \Lambda(1 - \Lambda)$, such that L simplifies to

$$L = \max_{0 \leq j \leq J} \sup_{t \in \mathbf{R}} |1 - \Lambda(\alpha_{j+1} - t) - \Lambda(\alpha_j - t)| = 1.$$

The ordered-logit loss is therefore globally Lipschitz, and Assumption 6 holds trivially. \square

4 Analytic Method

In this section we develop our analytic method to choose the penalty parameter λ . To do so, recall that we would like to choose the penalty parameter as small as possible while making score domination, $\mathcal{S} = \{\lambda \geq c_0 \|S\|_\infty\}$, a high-probability event. Recall also that S denotes a selection from the subdifferential $\partial \widehat{M}(\theta_0)$, which is nonempty under Assumption 2 and almost-surely singleton under Assumption 3. Therefore,

$$S = \mathbb{E}_n[(\partial/\partial\theta)m(X_i^\top\theta, Y_i)|_{\theta=\theta_0}] = \mathbb{E}_n[m'_1(X_i^\top\theta_0, Y_i)X_i] \quad \text{a.s.}$$

By analogy with the linear mean regression, we refer to $m'_1(X^\top\theta_0, Y)$ as the residual. Our analytic method makes use of the following two assumptions.

Assumption 8 (Conditional Mean Zero). *The residual $m'_1(X^\top \theta_0, Y)$ is such that with probability one, $\mathbb{E}[m'_1(X^\top \theta_0, Y) | X] = 0$.*

Assumption 9 (Residual: Analytic Method). *There exist functions $a, b : \mathcal{X} \rightarrow \mathbf{R}$ and a known finite constant $d > 0$ such that both $m'_1(X^\top \theta_0, Y) \subseteq [a(X), b(X)]$ and $b(X) - a(X) \leq d$ almost surely.*

Assumptions 8 and 9 presume that the residual $m'_1(X^\top \theta_0, Y)$ is centered conditional on the regressors and resides in a bounded interval of known width (*diameter*), respectively. The former assumption is satisfied in all of the examples from Section 2. As we explain at the end of this section, the latter assumption is satisfied in several, but not all, of the same examples.

Using these assumptions and appealing to Hoeffding's inequality (Vershynin, 2018, Theorem 2.2.6) conditional on the X_i 's, we see that for any coordinate j and any $t > 0$,

$$\mathbb{P}(|S_j| > t | \{X\}_{i=1}^n) \leq 2 \exp\left(-\frac{2nt^2}{d^2 \mathbb{E}_n[X_{ij}^2]}\right) \text{ a.s.}$$

A union bound then implies that for any $t > 0$,

$$\mathbb{P}(\|S\|_\infty > t | \{X_i\}_{i=1}^n) \leq 2p \exp\left(-\frac{2nt^2}{d^2 \max_{1 \leq j \leq p} \mathbb{E}_n[X_{ij}^2]}\right) \text{ a.s.}$$

Equating the right-hand side with $\alpha \in (0, 1)$ and solving for the resulting t , we arrive at the data-dependent penalty level

$$\widehat{\lambda}_\alpha^{\text{am}} := c_0 d \sqrt{\frac{\ln(2p/\alpha)}{2n} \max_{1 \leq j \leq p} \mathbb{E}_n[X_{ij}^2]}. \quad (4.1)$$

By construction, $\widehat{\lambda}_\alpha^{\text{am}} \geq c_0 \|S\|_\infty$ with conditional probability at least $1 - \alpha$ for almost every realization of the X_i 's, and so $\widehat{\lambda}_\alpha^{\text{am}} \geq c_0 \|S\|_\infty$ with probability at least $1 - \alpha$ also unconditionally. Given that this penalty level is available in closed form, we refer to this method for obtaining a penalty level as the *analytic method* (AM). Note that, under Assumption 5, the analytic penalty level admits the almost-sure bound

$$\widehat{\lambda}_\alpha^{\text{am}} \leq c_0 C_X d \sqrt{\frac{\ln(p/\alpha)}{n}} =: \bar{\lambda}_\alpha^{\text{am}}, \quad (4.2)$$

as long as $p \geq 2$. Use of the analytic method leads to the following result:

Theorem 2 (Nonasymptotic High-Probability Bounds: Analytic Method). *Let Assumptions 1–9 hold and let $\widehat{\theta} := \widehat{\theta}(\widehat{\lambda}_\alpha^{\text{am}})$ be a solution to the ℓ^1 -penalized M -estimation problem (1.2) with*

penalty level $\lambda = \widehat{\lambda}_\alpha^{\text{am}}$ given in (4.1). Define the finite constants $C_\epsilon := 16\sqrt{2}(1+\bar{c}_0)C_L C_X > 0$, $C_\lambda^{\text{am}} := c_0 C_X d > 0$, and

$$u_0 := \frac{2}{c_M} \left(C_\epsilon \sqrt{\frac{s \ln(pn)}{n}} + (1 + \bar{c}_0) C_\lambda^{\text{am}} \sqrt{\frac{s \ln(p/\alpha)}{n}} \right) > 0.$$

In addition, suppose that

$$s \ln(pn) \geq \frac{16C_{L,e}^2}{C_\epsilon^2} \text{ and } (1 + \bar{c}_0) u_0 \sqrt{s} \leq \frac{c_L}{C_X} \wedge c'_M. \quad (4.3)$$

Then both

$$\begin{aligned} \|\widehat{\theta} - \theta_0\|_2 &\leq \frac{2}{c_M} \left(C_\epsilon \sqrt{\frac{s \ln(pn)}{n}} + (1 + \bar{c}_0) C_\lambda^{\text{am}} \sqrt{\frac{s \ln(p/\alpha)}{n}} \right) \text{ and} \\ \|\widehat{\theta} - \theta_0\|_1 &\leq \frac{2(1 + \bar{c}_0)}{c_M} \left(C_\epsilon \sqrt{\frac{s^2 \ln(pn)}{n}} + (1 + \bar{c}_0) C_\lambda^{\text{am}} \sqrt{\frac{s^2 \ln(p/\alpha)}{n}} \right) \end{aligned}$$

with probability at least $1 - \alpha - 5n^{-1}$.

Theorem 2 gives nonasymptotic bounds on the estimation error of the ℓ^1 -penalized M-estimator based on the penalty parameter λ chosen according to the analytic method. From this theorem, we immediately obtain the corresponding convergence rates:

Corollary 1 (Convergence Rate Based on Analytic Method). *Let Assumptions 1–9 hold and let $\widehat{\theta} := \widehat{\theta}(\widehat{\lambda}_\alpha^{\text{am}})$ be a solution to the ℓ^1 -penalized M-estimation problem (1.2) with penalty level $\lambda = \widehat{\lambda}_\alpha^{\text{am}}$ given in (4.1). In addition, suppose that*

$$\frac{s^2 \ln(pn/\alpha)}{n} \rightarrow 0. \quad (4.4)$$

Then there exists a constant C depending only on the constants appearing in the aforementioned assumptions such that both

$$\|\widehat{\theta} - \theta_0\|_2 \leq C \sqrt{\frac{s \ln(pn/\alpha)}{n}} \text{ and } \|\widehat{\theta} - \theta_0\|_1 \leq C \sqrt{\frac{s^2 \ln(pn/\alpha)}{n}}$$

with probability $1 - \alpha - o(1)$.

We conclude this section by pointing out examples from Section 2 where Assumption 9 is satisfied, and so our analytic method can be applied.

Example 1 (Binary Response Model, Continued). The logit loss function (2.2) is differentiable in t with $m'_1(t, y) = \Lambda(t) - y$. The logit residual $m'_1(X^\top \theta_0, Y)$ thus resides in the

interval $[\Lambda(X^\top \theta_0) - 1, \Lambda(X^\top \theta_0)]$, and so satisfies Assumption 9 with $d = 1$.

More generally, let F admit an everywhere positive log-concave PDF $f = F'$. Then the binary-response loss (2.1) is differentiable with partial derivative (3.4). The binary nature of the outcome implies that

$$\min_{y \in \{0,1\}} m'_1(X^\top \theta_0, y) \leq m'_1(X^\top \theta_0, Y) \leq \max_{y \in \{0,1\}} m'_1(X^\top \theta_0, y).$$

From (3.4) we may deduce $m'_1(t, 1) = -f(t)/F(t) < 0 < f(t)/[1 - F(t)] = m'_1(t, 0)$. Inserting and simplifying, we therefore arrive at

$$\max_{y \in \{0,1\}} m'_1(t, y) - \min_{y \in \{0,1\}} m'_1(t, y) = \frac{f(t)}{F(t)[1 - F(t)]}.$$

Hence, for any distribution such that $f/[F(1 - F)]$ is also bounded from above, Assumption 9 is satisfied with

$$d = \sup_{t \in \mathbf{R}} \frac{f(t)}{F(t)[1 - F(t)]},$$

which only requires solving an unconstrained, univariate maximization problem. For example, as discussed earlier, $f/[F(1 - F)]$ is bounded from above if F is a t -distribution F_ν with $0 < \nu < \infty$ degrees of freedom. For $0 < \nu \leq 5$, the (unique) mode of $f_\nu/[F_\nu(1 - F_\nu)]$ is zero, such that $d = d_\nu = 4f_\nu(0) = 4\Gamma((\nu + 1)/2)/[\sqrt{\nu\pi}\Gamma(\nu/2)]$. For example, $d_1 = 4/\pi \approx 1.41$ for the standard Cauchy distribution. For higher degrees of freedom, the solution is more complicated, the exact d_ν being somewhat larger than the value $\frac{1}{2}\sqrt{\nu}$ of the (asymptotic) program $\sup_{t \in \mathbf{R}} |t|/(1 + t^2/\nu)$. For example, $\nu = 9$ produces $d_9 \approx 1.68 > \frac{3}{2}$.

As a side note, observe also that in contrast to the logit loss function, the probit loss function (2.3) does not satisfy Assumption 9. Indeed, this loss function is differentiable in t with

$$m'_1(t, y) = \frac{\varphi(t)}{\Phi(t)[1 - \Phi(t)]} [\Phi(t) - y]$$

but here $m'_1(t, 0) - m'_1(t, 1) = \varphi(t)/[1 - \Phi(t)] + \varphi(t)/\Phi(t) \sim t$ as $t \rightarrow \infty$. The probit residual is thus not confined to any bounded interval, violating Assumption 9. We could in principle reconcile the probit loss function with Assumption 9 by assuming that we know a constant $C_d > 0$ such that $\|\theta_0\|_1 \leq C_d$ and setting

$$d = \sup_{t \in [-\bar{X}C_d, \bar{X}C_d]} \frac{\varphi(t)}{\Phi(t)[1 - \Phi(t)]},$$

where $\bar{X} = \max_{1 \leq i \leq n} \|X_i\|_\infty$. While the resulting d is a known function of the X_i 's, this procedure would likely lead to very large values of the penalty parameter λ , thus making the

analytic method impractical. \square

Example 2 (Ordered Response Model, Continued). Provided the distribution F admits an everywhere positive log-concave PDF $f = F'$, the ordered-response loss (2.4) is differentiable in t with partial derivative (3.5). The discrete nature of the outcome and (3.5) imply that

$$\min_{0 \leq j \leq J} \frac{f(\alpha_{j+1} - t) - f(\alpha_j - t)}{F(\alpha_{j+1} - t) - F(\alpha_j - t)} \leq m'_1(t, y) \leq \max_{0 \leq j \leq J} \frac{f(\alpha_{j+1} - t) - f(\alpha_j - t)}{F(\alpha_{j+1} - t) - F(\alpha_j - t)},$$

where we interpret $f(\pm\infty)$ and $F(-\infty)$ as zero and $F(+\infty)$ as one. Hence, for a distribution F and cut-off points $\{\alpha_j\}$ such that the difference between the upper and lower bounds is bounded from above in t , Assumption 9 is satisfied with

$$d = \sup_{t \in \mathbf{R}} \left\{ \max_{0 \leq j \leq J-1} \frac{f(\alpha_{j+1} - t) - f(\alpha_j - t)}{F(\alpha_{j+1} - t) - F(\alpha_j - t)} - \min_{1 \leq j \leq J} \frac{f(\alpha_{j+1} - t) - f(\alpha_j - t)}{F(\alpha_{j+1} - t) - F(\alpha_j - t)} \right\},$$

where we have used our knowledge of the signs of the first and last elements to reduce the candidates for a minimum and maximum, respectively. With knowledge of F and the α_j 's, this quantity may at least in principle be computed. For the logistic distribution $F = \Lambda$ we have $f = \Lambda(1 - \Lambda)$, such that d simplifies to

$$\begin{aligned} d &= \sup_{t \in \mathbf{R}} \left\{ \max_{0 \leq j \leq J-1} \{1 - \Lambda(\alpha_{j+1} - t) - \Lambda(\alpha_j - t)\} - \min_{1 \leq j \leq J} \{1 - \Lambda(\alpha_{j+1} - t) - \Lambda(\alpha_j - t)\} \right\} \\ &= \sup_{t \in \mathbf{R}} \left\{ \max_{1 \leq j \leq J} \{\Lambda(\alpha_{j+1} - t) + \Lambda(\alpha_j - t)\} - \min_{0 \leq j \leq J-1} \{\Lambda(\alpha_{j+1} - t) + \Lambda(\alpha_j - t)\} \right\}. \end{aligned}$$

The second line shows that for *any* ordered logit we may use the possibly crude upper bound of $d = 2$. In the case of *trinary* ordered logit, the relevant pointwise maximum and minimum equal $1 + \Lambda(\alpha_2 - t)$ and $\Lambda(\alpha_1 - t)$, respectively. The resulting d is

$$d = \sup_{t \in \mathbf{R}} \{1 + \Lambda(\alpha_2 - t) - \Lambda(\alpha_1 - t)\} = 1 + \Lambda\left(\frac{\alpha_2 - \alpha_1}{2}\right) - \Lambda\left(\frac{\alpha_1 - \alpha_2}{2}\right),$$

the supremum being attained at $t = (\alpha_1 + \alpha_2)/2$. Further specializing to symmetric cut-offs $\alpha_2 = \alpha > 0$ and $\alpha_1 = -\alpha$ yields $d = 2\Lambda(\alpha) = 2/(1 + e^{-\alpha})$, which is in (1, 2) for all $\alpha > 0$ and asymptotes to two as $\alpha \rightarrow \infty$. For example, a value of $\alpha = 1$ produces $d = 2e/(1 + e) \approx 1.4621$. In the limiting case of $\alpha = 0$, we recover the value $d = 1$ for the binary logit, as expected. \square

Example 6 (Panel Logit Model, Continued). The panel logit loss function (2.8) is differentiable in t with $m'_1(t, y) = \mathbf{1}(y_1 \neq y_2)[\Lambda(t) - y_1]$. Thus, the residual $m'_1((X_1 - X_2)^\top \theta_0, Y)$ resides in the interval $[\Lambda((X_1 - X_2)^\top \theta_0) - 1, \Lambda((X_1 - X_2)^\top \theta_0)]$, and so satisfies Assumption

9 with $d = 1$. □

Example 7 (Panel Censored Model, Continued). The trimmed LAD loss function (2.9) is differentiable in t and satisfies $|m'_1(t, y)| \leq 1$ if $t \neq y_1 - y_2$ or $y_1 = y_2 = 0$. Thus, as long as the conditional distribution of $(\varepsilon_1, \varepsilon_2)$ given (α, X_1, X_2) is continuous (as implied by Honoré (1992, Assumption E.1)), this loss function satisfies Assumption 9 with $d = 2$. Note, however, that the trimmed LS loss function does not satisfy Assumption 9. □

Example 8 (Panel Duration Model, Continued). Since the loss function (2.10) here is of the logit form, it satisfies Assumption 9 with $d = 1$. □

5 Bootstrap-after-Cross-Validation Method

The analytic method of the previous section relies on Assumption 9. As explained there, this assumption is satisfied in quite a few applications. However, there are also many other applications where this assumption is not satisfied. Examples include the probit model, the logit model with estimation based on the logistic calibration loss function, and the panel censored model with estimation based on the trimmed LS loss function. Moreover, even if Assumption 9 is satisfied, the analytic penalty level $\hat{\lambda}_\alpha^{\text{am}}$ in (4.1) follows from a union-bound argument and may thus be quite conservative. In this section we therefore seek to provide a method to choose the penalty parameter which is not conservative and broadly available, yet amenable to theoretical analysis. We split the section into two subsections. In Section 5.1, we develop a generic bootstrap method that allows for choosing the penalty parameter λ assuming availability of some generic estimators \hat{U}_i of the residuals $U_i = m'_1(X_i^\top \theta_0, Y_i)$. In Section 5.2, we explain how to obtain suitable estimators \hat{U}_i via cross-validation.

5.1 Bootstrapping Penalty Level

To develop some intuition, suppose for the moment that residuals $U_i = m'_1(X_i^\top \theta_0, Y_i)$ are observable. In this case, we can estimate the $(1 - \alpha)$ quantile of the score $S = \mathbb{E}_n[U_i X_i]$,

$$q(1 - \alpha) := (1 - \alpha)\text{-quantile of } \max_{1 \leq j \leq p} |\mathbb{E}_n[U_i X_{ij}]|,$$

via the Gaussian multiplier bootstrap. To this end, let e_1, \dots, e_n be independent standard normal random variables that are independent of the data W_1, \dots, W_n . We then estimate $q(1 - \alpha)$ by

$$\tilde{q}(1 - \alpha) := (1 - \alpha)\text{-quantile of } \max_{1 \leq j \leq p} |\mathbb{E}_n[e_i U_i X_{ij}]| \text{ given } \{W_i\}_{i=1}^n.$$

It is rather standard to show that, under certain regularity conditions, $\tilde{q}(1 - \alpha)$ delivers a good approximation to $q(1 - \alpha)$, even if the dimension p of the X_i 's is much larger than the sample size n . To see why this is the case, let Z be any centered random vector in \mathbf{R}^p and let Z_1, \dots, Z_n be independent copies of Z . As established in [Chernozhukov et al. \(2013, 2017\)](#), the random vectors Z_1, \dots, Z_n satisfy the following high-dimensional versions of the central limit and Gaussian multiplier bootstrap theorems: If for some finite constants $b, B > 0$, one has

$$\min_{1 \leq j \leq p} \mathbb{E}[Z_{ij}^2] \geq b \text{ and } \mathbb{E} \left[\max_{1 \leq j \leq p} Z_{ij}^4 \right] \leq B,$$

then

$$\sup_{A \in \mathcal{A}_p} \left| \mathbb{P} \left(\frac{1}{\sqrt{n}} \sum_{i=1}^n Z_i \in A \right) - \mathbb{P}(\mathbf{N}(\mathbf{0}, \mathbb{E}[ZZ^\top]) \in A) \right| \leq C_{b,B} \left(\frac{\ln^7(pn)}{n} \right)^{1/6}, \quad (5.1)$$

and, with probability approaching one,

$$\sup_{A \in \mathcal{A}_p} \left| \mathbb{P} \left(\frac{1}{\sqrt{n}} \sum_{i=1}^n e_i Z_i \in A \mid \{Z_i\}_{i=1}^n \right) - \mathbb{P}(\mathbf{N}(\mathbf{0}, \mathbb{E}[ZZ^\top]) \in A) \right| \leq C'_{b,B} \left(\frac{\ln^6(pn)}{n} \right)^{1/6}, \quad (5.2)$$

where \mathcal{A}_p denotes the collection of all (hyper)rectangles in \mathbf{R}^p , and the constants $C_{b,B}$ and $C'_{b,B}$ depend only on b and B . Provided $\ln^7(pn)/n \rightarrow 0$, combination of these two results suggests that the Gaussian multiplier bootstrap yields a good approximation to the law of the potentially high-dimensional vector $n^{-1/2} \sum_{i=1}^n Z_i$ when restricted to (hyper)rectangles.

Consider now the family of rectangles defined by

$$A_t := \left\{ u \in \mathbf{R}^p; \max_{1 \leq j \leq p} |u_j| \leq t \right\}, \quad t \geq 0.$$

We can then write

$$\mathbb{P} \left(\max_{1 \leq j \leq p} |\mathbb{E}_n[U_i X_{ij}]| \leq t \right) = \mathbb{P} \left(\frac{1}{\sqrt{n}} \sum_{i=1}^n U_i X_i \in A_{t\sqrt{n}} \right).$$

The $U_i X_i$'s are centered under Assumption 8, and so the aforementioned results can be applied in our context of ℓ^1 -penalized M-estimation.

Of course, we typically do not observe the residuals $U_i = m'_1(X_i^\top \theta_0, Y_i)$, and so the method described above is infeasible. Fortunately, the result (5.2) continues to hold upon replacing $\{Z_i\}_{i=1}^n$ with estimators $\{\hat{Z}_i\}_{i=1}^n$, provided these estimators are ‘‘sufficiently good.’’ Suppose

therefore that residual estimators $\{\widehat{U}_i\}_{i=1}^n$ are available. We then compute

$$\widehat{q}(1 - \alpha) := (1 - \alpha)\text{-quantile of } \max_{1 \leq j \leq p} |\mathbb{E}_n[e_i \widehat{U}_i X_{ij}]| \text{ given } \{(W_i, \widehat{U}_i)\}_{i=1}^n. \quad (5.3)$$

A feasible penalty level is then given by

$$\widehat{\lambda}_\alpha^{\text{bm}} := c_0 \widehat{q}(1 - \alpha). \quad (5.4)$$

We refer to this method for obtaining a penalty level as the *bootstrap method* and to $\widehat{\lambda}_\alpha^{\text{bm}}$ itself as the *bootstrapped penalty level*. In Section 5.2, we show how to obtain residual estimators $\{\widehat{U}_i\}_{i=1}^n$ via cross-validation, thus obtaining the bootstrap-after-cross-validation method.

To ensure that $\widehat{q}(1 - \alpha)$ delivers a good approximation to $q(1 - \alpha)$, we invoke the following assumptions, where we denote $U = m'_1(X^\top \theta_0, Y)$ and $Z = (Z_1, \dots, Z_p)^\top = UX$.

Assumption 10 (Residual: Bootstrap Method). *There exist finite constants $c_U, C_U > 0$ such that (1) $\mathbb{E}[U^4] \leq C_U^4$ and (2) $\mathbb{E}[Z_j^2] \geq c_U$ for all $j \in \{1, \dots, p\}$.*

Assumption 11 (Residual Estimation). *There exist sequences β_n and δ_n of constants in \mathbf{R}_{++} converging to zero such that $\mathbb{E}_n[(\widehat{U}_i - U_i)^2] \leq \delta_n^2 / \ln^2(pn)$ with probability at least $1 - \beta_n$.*

Use of the bootstrap method leads to following result:

Theorem 3 (Nonasymptotic High-Probability Bounds: Bootstrap Method). *Let Assumptions 1–8, 10, and 11 hold and let $\widehat{\theta} := \widehat{\theta}(\widehat{\lambda}_\alpha^{\text{bm}})$ be a solution to the ℓ^1 -penalized M -estimation problem (1.2) with penalty level $\lambda = \widehat{\lambda}_\alpha^{\text{bm}}$ given in (5.4). Define the finite constants $C_\epsilon := 16\sqrt{2}(1 + \bar{c}_0)C_L C_X > 0$, $C_\lambda^{\text{bm}} := 2(2 + \sqrt{2})c_0 C_X (C_U + \delta_n / \ln(pn)) > 0$, and*

$$u_0 := \frac{2}{c_M} \left(C_\epsilon \sqrt{\frac{s \ln(pn)}{n}} + (1 + \bar{c}_0) C_\lambda^{\text{bm}} \sqrt{\frac{s \ln(p/\alpha)}{n}} \right) > 0.$$

In addition, suppose that

$$s \ln(pn) \geq \frac{16C_{L,e}^2}{C_\epsilon^2} \text{ and } (1 + \bar{c}_0) u_0 \sqrt{s} \leq \frac{c_L}{C_X} \wedge c'_M. \quad (5.5)$$

Then there exists a finite constant C , depending only on c_U , such that for

$$\rho_n := C \max \left\{ \beta_n, C_X \delta_n, \left(\frac{B^4 \ln^7(pn)}{n} \right)^{1/6}, \frac{1}{\ln^2(pn)} \right\} \text{ and } B := (C_X C_U)^3 \vee 1,$$

we have both

$$\begin{aligned}\|\widehat{\theta} - \theta_0\|_2 &\leq \frac{2}{c_M} \left(C_\epsilon \sqrt{\frac{s \ln(pn)}{n}} + (1 + \bar{c}_0) C_\lambda^{\text{bm}} \sqrt{\frac{s \ln(p/\alpha)}{n}} \right) \quad \text{and} \\ \|\widehat{\theta} - \theta_0\|_1 &\leq \frac{2(1 + \bar{c}_0)}{c_M} \left(C_\epsilon \sqrt{\frac{s^2 \ln(pn)}{n}} + (1 + \bar{c}_0) C_\lambda^{\text{bm}} \sqrt{\frac{s^2 \ln(p/\alpha)}{n}} \right)\end{aligned}$$

with probability at least $1 - \alpha - \rho_n - 6n^{-1}$.

5.2 Cross-Validating Residuals

Assumption 11 is ‘high-level’ in the sense that it does not specify how one performs residual estimation in practice. In this subsection, we explain how residual estimation can be performed via cross-validation (CV).

To describe our CV residual estimator, fix any integer $K \geq 2$, and let I_1, \dots, I_K partition the sample indices $\{1, \dots, n\}$. Provided n is divisible by K , the even partition

$$I_k = \{(k-1)n/K + 1, \dots, kn/K\}, \quad k \in \{1, \dots, K\}, \quad (5.6)$$

is natural, but not necessary. For the formal results below, we only require that each I_k specifies a ‘substantial’ subsample (see Assumption 12 below).

Define the *subsample criterion* \widehat{M}_I to be the sample criterion

$$\widehat{M}_I(\theta) := \mathbb{E}_I [m(X_i^\top \theta, Y_i)], \quad \theta \in \Theta, \quad \emptyset \neq I \subsetneq \{1, \dots, n\}, \quad (5.7)$$

based only on observations $i \in I$, and let Λ_n denote a finite subset of \mathbf{R}_{++} composed by candidate penalty levels. We require Λ_n to be ‘sufficiently rich’ (see Assumption 13 below). Our CV procedure then goes as follows. First, estimate parameters θ_0 by

$$\widehat{\theta}_{I_k^c}(\lambda) \in \operatorname{argmin}_{\theta \in \Theta} \left\{ \widehat{M}_{I_k^c}(\theta) + \lambda \|\theta\|_1 \right\},$$

for each candidate penalty level $\lambda \in \Lambda_n$ and holding out each subsample $k \in \{1, \dots, K\}$ in turn. Second, determine the penalty level

$$\widehat{\lambda}^{\text{cv}} \in \operatorname{argmin}_{\lambda \in \Lambda_n} \sum_{k=1}^K \sum_{i \in I_k} m(X_i^\top \widehat{\theta}_{I_k^c}(\lambda), Y_i) \quad (5.8)$$

by minimizing the out-of-sample loss over the set of candidate penalties. Third, estimate residuals $U_i = m_1(X_i^\top \theta_0, Y_i)$, $i \in \{1, \dots, n\}$, by predicting out of each estimation sample,

i.e.,

$$\widehat{U}_i^{\text{cv}} := m'_1(X_i^\top \widehat{\theta}_{I_k^c}(\widehat{\lambda}^{\text{cv}}), Y_i), \quad i \in I_k, \quad k \in \{1, \dots, K\}. \quad (5.9)$$

Combining the bootstrap estimate $\widehat{\lambda}_\alpha^{\text{bm}} = c_0 \widehat{q}(1 - \alpha)$ from the previous subsection with the CV residual estimates $\widehat{U}_i = \widehat{U}_i^{\text{cv}}$ from this subsection, we obtain the *bootstrap-after-cross-validation* (BCV) method for choosing the penalty parameter $\lambda = \widehat{\lambda}_\alpha^{\text{bcv}}$.

To ensure good performance of the estimator resulting from initiating the bootstrap method with the CV residual estimates (5.9), we invoke the following assumptions.

Assumption 12 (Data Partition). *The number $K \in \{2, 3, \dots\}$ of subsamples is constant and does not depend on n . There exists a constant $c_D \in (0, 1)$ such that $\min_{1 \leq k \leq K} |I_k| \geq c_D n$.*

Assumption 13 (Candidate Penalties). *There exists finite constants $c_\Lambda, C_\Lambda > 0$ and $a \in (0, 1)$ such that*

$$\Lambda_n = \{C_\Lambda a^\ell; a^\ell \geq c_\Lambda/n, \ell \in \{0, 1, 2, \dots\}\}.$$

Assumption 14 (Residual). *There exist finite constants $\sigma, C_{ms1} > 0$ such that:*

1. For all $t \in \mathbf{R}$,

$$\ln \mathbb{E} \left[\exp \left(t m'_1(X^\top \theta_0, Y) \right) \mid X \right] \leq \frac{\sigma^2 t^2}{2} \text{ a.s.}$$

2. For all $\theta \in \Theta$,

$$\mathbb{E} \left[\left\{ m'_1(X^\top \theta, Y) - m'_1(X^\top \theta_0, Y) \right\}^2 \right] \leq C_{ms1}^2 \left(\sqrt{\mathcal{E}(\theta)} \vee \mathcal{E}(\theta) \right).$$

Assumption 15 (Mean-Square Loss Continuity). *There exists a finite constant $C_{ms} > 0$ such that for all $\theta \in \Theta$,*

$$\mathbb{E} \left[\left\{ m(X^\top \theta, Y) - m(X^\top \theta_0, Y) \right\}^2 \right] \leq C_{ms}^2 \left(\mathcal{E}(\theta) \vee \mathcal{E}(\theta)^2 \right).$$

Assumption 12 means that we rely upon the classical K -fold CV with fixed K . This assumption does rule out leave-one-out CV, since $K = n$ and $I_k = \{k\}$ implies $|I_k|/n \rightarrow 0$. Assumption 13 allows for a rather large candidate set Λ_n of penalty values. Note that the largest penalty value (C_Λ) can be set arbitrarily large and the smallest value (c_Λ/n) converges rapidly to zero. In Lemma B.1 we show that these properties ensure that the set Λ_n eventually contains a “good” penalty candidate, say λ_* , in the sense of leading to a uniform bound on the excess risk of subsample estimators $\widehat{\theta}_{I_k^c}(\lambda_*), k \in \{1, \dots, K\}$. Other candidate penalty sets leading to a bound on the subsample estimator excess risk are certainly possible. Assumptions 14 and 15 are high-level but rather mild. We provide a set of low-level conditions suitable for each of the examples from Section 2 in Appendix A.

Use of CV residual estimators leads to the following result:

Theorem 4 (High-Probability CV-Residual Error Bound). *Let Assumptions 1–8 and 12–15 hold. Define the finite constants $C_\epsilon := 16\sqrt{2}(1+\bar{c}_0)C_L C_X > 0$, $C_S := 2C_X\sigma/\sqrt{(K-1)c_D} > 0$,*

$$C_\mathcal{E} := \sqrt{\frac{2}{c_M}} \left(\frac{C_\epsilon}{(K-1)c_D} + \frac{(1+\bar{c}_0)c_0 C_S}{a} \right) > 0 \quad (5.10)$$

and

$$\tilde{u}_0 := \frac{2}{c_M} \left(\frac{C_\epsilon}{(K-1)c_D} + \frac{(1+\bar{c}_0)c_0 C_S}{a} \right) \sqrt{\frac{s \ln(pn)}{n}} > 0. \quad (5.11)$$

In addition, suppose that

$$s \ln(pn) \geq 16(K-1)c_D C_{L,e}^2 / C_\epsilon^2, \quad (1+\bar{c}_0)\tilde{u}_0\sqrt{s} \leq (c_L/C_X) \wedge c'_M, \quad (5.12)$$

$$\ln(pn)/n \leq (C_\Lambda a / c_0 C_S)^2, \quad n \ln(pn) \geq (c_\Lambda / c_0 C_S)^2, \quad n \geq 1/c_\Lambda. \quad (5.13)$$

Then for any $t \in \mathbf{R}_{++}$ satisfying

$$\frac{32C_{ms}^2}{c_D^2 \ln(1/a)} \frac{t \ln n}{n} + \frac{8C_\mathcal{E}^2 s \ln(pn)}{c_D n} \leq 1, \quad (5.14)$$

we have

$$\mathbb{E}_n [(\widehat{U}_i^{\text{cv}} - U_i)^2] \leq \frac{8C_{ms1}^2 t \ln n}{\ln(1/a)} \left(\frac{2C_{ms}^2}{c_D^2 \ln(1/a)} \frac{t \ln n}{n} + \frac{C_\mathcal{E}^2 s \ln(pn)}{2c_D n} \right)^{1/2} \quad (5.15)$$

with probability at least $1 - K(4n^{-1} + 2[(K-1)c_D n]^{-1} + 2t^{-1})$.

This theorem provides an avenue for verification of Assumption 11. Specifically, it implies that for any sequence t_n of constants in \mathbf{R}_{++} satisfying (5.14), we can take δ_n^2 and β_n in Assumption 11 to be the right-hand side of (5.15) multiplied by $\ln^2(pn)$ and $K(4n^{-1} + 2[(K-1)c_D n]^{-1} + 2t_n^{-1})$, respectively. Combining Theorems 3 and 4, we obtain convergence rates for the ℓ^1 -penalized M-estimator based on the penalty parameter λ chosen according to the BCV method:

Corollary 2 (Convergence Rate Based on Bootstrap after CV Method). *Let Assumptions 1–8, 10, and 12–15 hold and let $\widehat{\theta} := \widehat{\theta}(\widehat{\lambda}_\alpha^{\text{bcv}})$ be a solution to the ℓ^1 -penalized M-estimation problem (1.2) with penalty level $\lambda = \widehat{\lambda}_\alpha^{\text{bcv}}$. In addition, suppose that*

$$\frac{s^2 \ln(pn/\alpha)}{n} \rightarrow 0, \quad \frac{s \ln^5(pn)(\ln n)^2}{n} \rightarrow 0, \quad \text{and} \quad \frac{\ln^7(pn/\alpha)}{n} \rightarrow 0. \quad (5.16)$$

Then there exists a constant C depending only on the constants appearing in the aforementioned assumptions such that both

$$\|\hat{\theta} - \theta_0\|_2 \leq C \sqrt{\frac{s \ln(pn/\alpha)}{n}} \quad \text{and} \quad \|\hat{\theta} - \theta_0\|_1 \leq C \sqrt{\frac{s^2 \ln(pn/\alpha)}{n}}$$

with probability $1 - \alpha - o(1)$.

Corollaries 1 and 2 demonstrate that both analytic and bootstrap-after-cross-validation methods with α , for example, equal to $1/n$ yield ℓ^1 -penalized M-estimators whose convergence rates in the ℓ^2 and ℓ^1 norms are $\sqrt{s \ln(pn)/n}$ and $\sqrt{s^2 \ln(pn)/n}$, respectively. These are typical rates that one expects in the high-dimensional settings under sparsity. For example, it is well-known that these rates are minimax optimal in the case of the high-dimensional linear mean regression model; see [Rigollet and Tsybakov \(2011\)](#) and [Chetverikov et al. \(2016\)](#).

6 Simulations

In this section we investigate the finite-sample behavior of estimators based on the analytic and bootstrap-after-cross-validation methods for obtaining penalty levels proposed in Sections 4 and 5, respectively.

6.1 Simulation Design

For concreteness, we consider a data-generating process (DGP) of the form

$$Y_i = \mathbf{1}\left(\sum_{j=1}^p \theta_{0j} X_{ij} + \varepsilon_i > 0\right), \quad \varepsilon_i | X_{i1}, \dots, X_{ip} \sim \text{Logistic}(0, 1), \quad i \in \{1, \dots, n\},$$

thus leading to a binary logit model. The regressors $X = (X_1, \dots, X_p)$ are jointly centered Gaussian with a covariance matrix of the Toeplitz form

$$\text{cov}(X_{ij}, X_{ik}) = \text{E}[X_{ij} X_{ik}] = \rho^{|j-k|}, \quad j, k \in \{1, \dots, p\},$$

such that ρ determines the overall correlation level. We allow $\rho \in \{0, .1, \dots, .9\}$, thus running the gamut of (positive) correlation levels. Since the ε_i 's are standard Logistic, the ‘‘noise’’ in our DGP is fixed at $\text{var}(\varepsilon_i) = \pi^2/3 \approx 3.3$. However, the ‘‘signal’’ $\text{var}(\sum_{j=1}^p \theta_{0j} X_{ij}) = \theta_0^\top \text{E}[X_i X_i^\top] \theta_0$ depends on both the correlation level and coefficient pattern. We consider both sparse and dense coefficient patterns.

The *sparse coefficient pattern* has only nonzero coefficients for the first couple of regressors,

$$\text{Pattern 1 (Sparse): } \theta_0 = (1, 1, 0, \dots, 0)^\top,$$

thus yielding $s = 2$ relevant regressors among the p candidates. The implied *signals* are here given by

$$\text{var} \left(\sum_{j=1}^p \theta_{0j} X_{ij} \right) = 2(1 + \rho) \in \{2, 2.2, \dots, 3.8\},$$

further implying a *signal-to-noise ratio* (SNR) range of about .6 to about 1.2. Compared to existing simulations studies for the high-dimensional logit, the signals considered here are rather low.⁹

The *dense coefficient pattern* have all nonzero coefficients,

$$\text{Pattern 2 (Dense): } \theta_{0j} = (1/\sqrt{2})^{j-1}, \quad j \in \{1, \dots, p\},$$

thus implying $s = p$. The base $(1/\sqrt{2})$ was here chosen to (approximately) equate the signals arising from the dense and sparse coefficient patterns in the baseline case of uncorrelated regressors ($\rho = 0$), which, in turn, amounts to $\|\theta_0\|_2^2$. We attempt sample sizes $n \in \{100, 200, 400\}$ and limit attention to the high-dimensional regime “ $p \geq n$ ” by fixing $p = n$ throughout.

With a sparse coefficient pattern, the nonzero coefficients are well separated from zero and should be relatively easy to detect—at least with larger sample sizes. With a dense coefficient pattern, every regressor is in principle relevant, and our implicit assumption of exact sparsity fails ($s = p = n$). Note, however, that the relevance of the regressors, as measured by their coefficient, is rapidly decaying in the regressor index, such that the vast majority of the signal is still captured by a fraction of the regressors. For example, in the baseline case of uncorrelated regressors ($\rho = 0$), the first 10 regressors account for 99.9 pct. of the total signal, and the model may be interpreted as effectively sparse. One may therefore hope that our methods also apply (to some extent) under this *approximate sparsity* alone.

6.2 Estimators and Implementation

For estimation purposes, we mark up the score by $c_0 = 1.1$ and specify the tolerance as $\alpha = \alpha_n = 10/n$, thus leading to an α of 10, 5 and 2.5 percent for $n = 100, 200$ and 400, respectively. We let α decrease with n , such that the error bounds in Theorems 2 and 3

⁹For example, the design in [Friedman, Hastie, and Tibshirani \(2010, Section 5.2\)](#) implies a SNR of three. In [Ng \(2004, Section 5\)](#), the SNR is over 30.

may be interpreted as holding with probability approaching one. We consider three feasible estimators based on the analytic and bootstrap methods. With our binary-logit design, the *analytic method* (4.1) for specifying the penalty level is justified with $d = 1$, and $\widehat{\theta}(\widehat{\lambda}_\alpha^{\text{am}})$ constitutes our *first* estimator. Our *second* estimator is based on the *bootstrap method* (5.3) initiated with residual estimates

$$\widehat{U}_i^{\text{am}} := m'_1(X_i^\top \widehat{\theta}(\widehat{\lambda}_\alpha^{\text{am}}), Y_i), \quad i \in \{1, \dots, n\},$$

resulting from the analytic method. Our *third* estimator follows similarly, except that we initiate the bootstrap method with *cross-validation residual estimates* (5.9).¹⁰ To introduce a benchmark, we also consider the infeasible estimator arising from the bootstrap method using the true residuals. We refer to the latter three bootstrap-based estimators as *bootstrapping after the analytic method* (BAM), *bootstrapping after cross validation* (BCV), and the *oracle bootstrap* (Oracle). All simulations are carried out in Matlab[®] with optimization and cross validation done using the user-contributed `glmnet` package.¹¹ For each sample size $n(=p)$, each correlation level ρ , and each coefficient pattern (sparse or dense), we use 2,000 simulation draws and 1,000 standard Gaussian bootstrap draws per simulation draw (when applicable). In constructing the candidate penalty set Λ_n , we use the `glmnet` default setting, which constructs a log-scale equi-distant grid of a 100 candidate penalties from the threshold penalty level to essentially zero. The threshold is the (approximately) smallest level of penalization needed to set every coefficient to zero, thus resulting in a trivial (null) model.¹²

6.3 Simulation Results

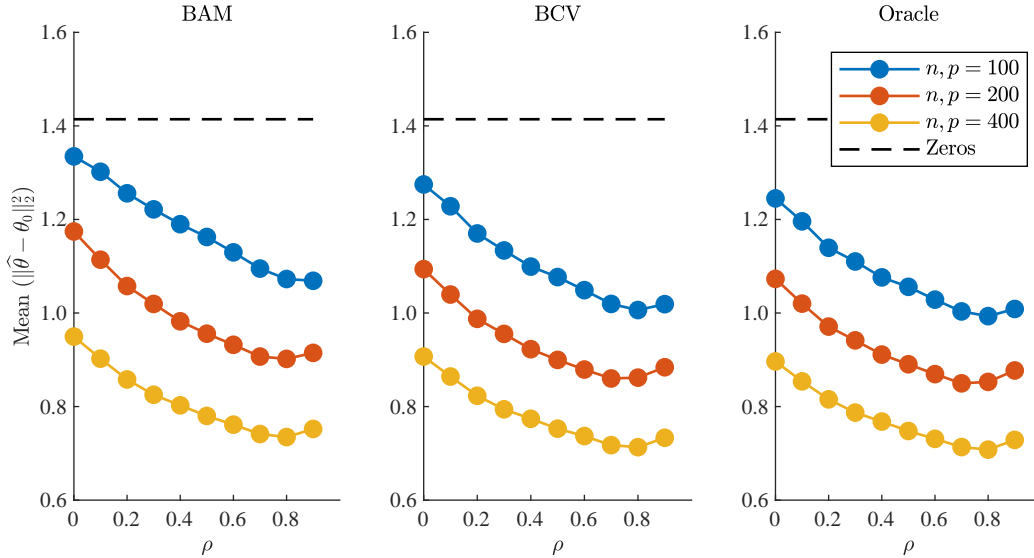
Figure 6.1 shows the mean-square ℓ^2 estimation error (over the 2,000 simulation draws) as a function of the correlation level ρ for each of the three bootstrap-based estimators and each sample size $n(=p)$, obtained with a sparse coefficient pattern. For each of these three estimators, we see that the mean estimation error decreases with sample size. Convergence appears to take place even though the number of candidate regressors matches the sample size and no matter the level of regressor correlation. This finding indicates that our bootstrap

¹⁰We use 10-fold cross validation, splitting the data evenly. As a result, $K = 10$ and $c_D = \frac{1}{10}$.

¹¹We use the August 30, 2013 version of `glmnet` for Matlab[®], available for download at https://web.stanford.edu/~hastie/glmnet_matlab/. Cross validation is conveniently done using `cvglmnet`, which automatically stores the out-of-fold predictions $X_i^\top \widehat{\theta}_{I_k^c}(\lambda)$, $i \in I_k$, for each candidate penalty.

¹²Log-scale equi-distance from a “large” candidate value to essentially zero fits well with the form of Λ_n in our Assumption 13 (interpreting $c_\Lambda/n \approx 0$). However, the threshold penalty is a function of the data and, thus, random. The resulting candidate penalty set used in our simulations is therefore also random, and thus, strictly speaking, not allowed by Assumption 13. We believe this deviation from our theory to be only a minor issue.

Figure 6.1: Consistency of Bootstrap-Based Estimators with Exact Sparsity



method is useful for high-dimensional estimation, not only in the best-case scenario where residuals are observed, but also when residuals are estimated by a pilot method—whether it be analytic or computational.

Figure 6.2 rearranges the plots in Figure 6.1 in order to facilitate comparison of the various estimators, now including the estimator based on the analytic method (AM). For each of the three sample sizes/number of candidate regressors, we see that the oracle performs better than the other two bootstrap-based estimators. Bootstrapping after cross validation appears to outperform bootstrapping after the analytic method, which, in turn, improves greatly upon the analytic method itself. While residual estimation comes at a price, bootstrapping after cross-validation achieves near-oracle performance even with our smallest sample size—and is essentially indistinguishable from the oracle at $n = 400$. Bootstrapping after the analytic method here comes in close second place among the feasible estimators, which indicates that BAM provides a computationally inexpensive way of obtaining quality results.

Figures 6.3 and 6.4 reproduce Figures 6.1 and 6.2, respectively, with results stemming from the dense coefficient pattern (approximate sparsity). The plots in Figure 6.3 are also indicative of consistency, although convergence is slowed down by the lack of exact sparsity (compare with Figure 6.1). The ranking of estimators in Figure 6.2 is preserved in Figure 6.4. These findings suggest that our methods remain relevant under a less stringent assumption than exact sparsity.

As a final exercise, we compare our analytic and bootstrap methods to existing penalty methods formally justifiable in our binary logit model. Specifically, we here compare with the analytic penalty levels provided in Bunea (2008b, Theorem 2.4), van de Geer (2008, Theorem

Figure 6.2: Comparing Estimators with Exact Sparsity

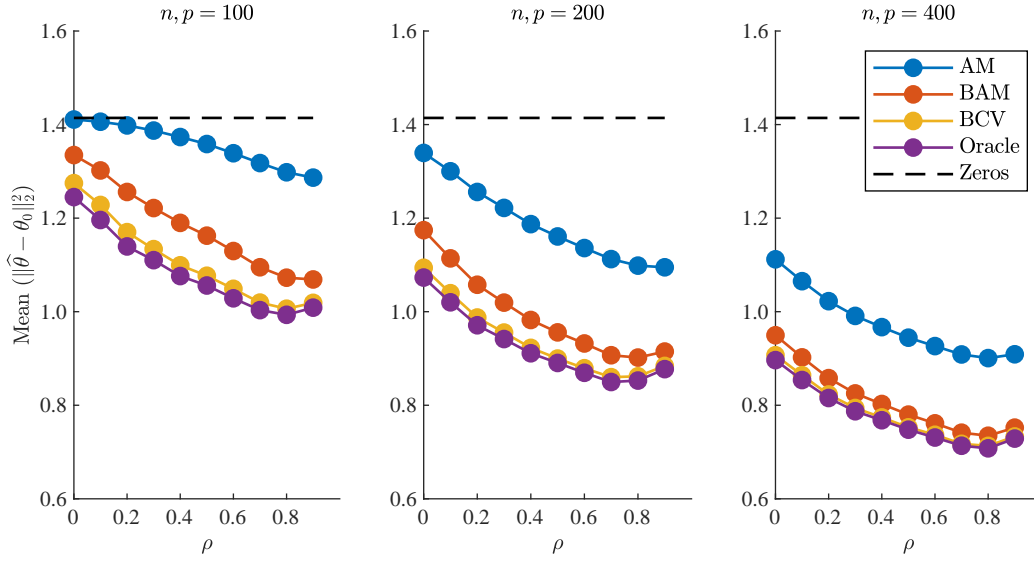


Figure 6.3: Consistency of Bootstrap-Based Estimators with Approximate Sparsity

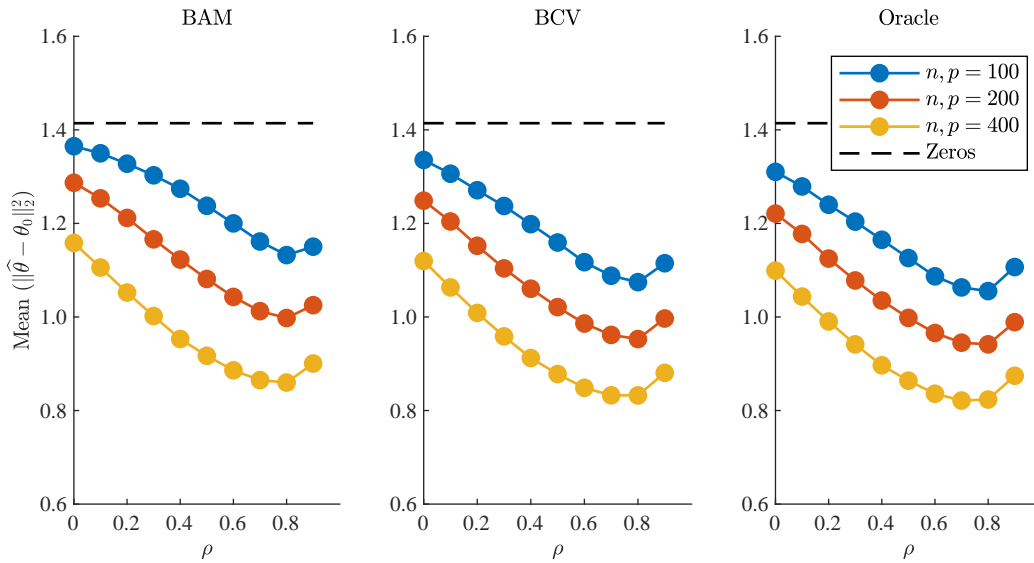
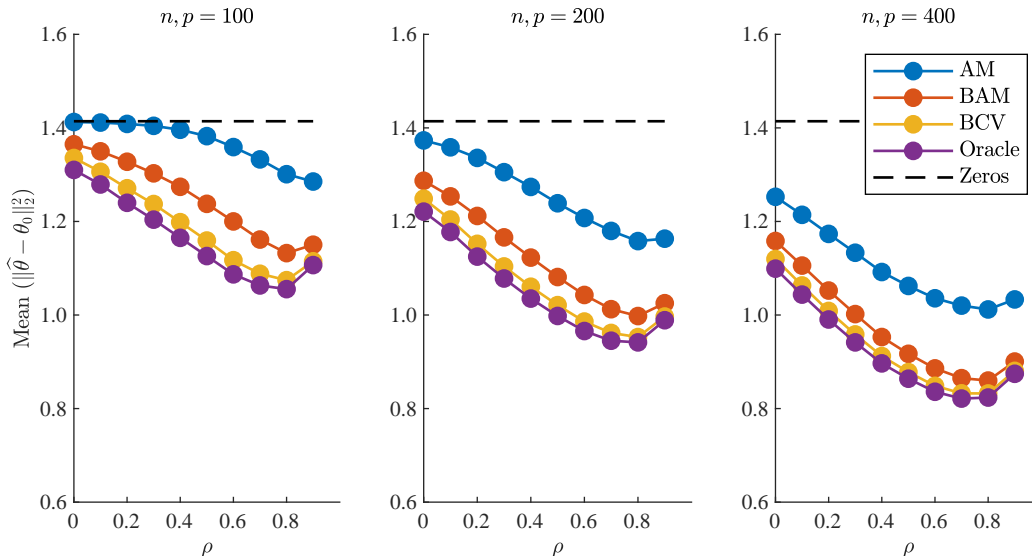


Figure 6.4: Comparing Estimators with Approximate Sparsity



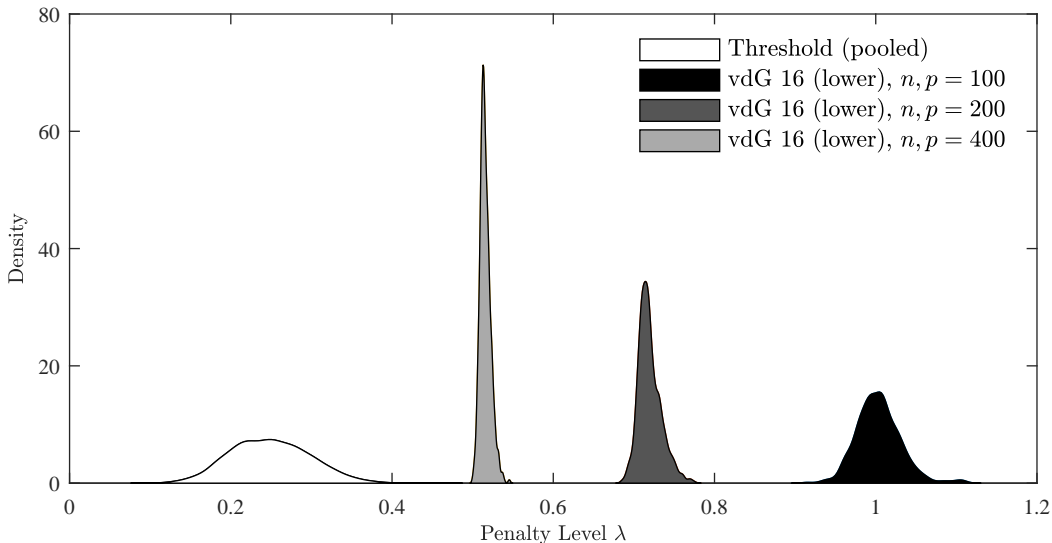
2.1) and van de Geer (2016, Theorem 12.1). Across all of our designs and simulation draws, the *smallest* Bunea (2008a) penalty is larger than the *largest* van de Geer (2008) penalty, which, in turn, is similar in size to her (2016) penalty level. We therefore restrict attention to the latter. In our notation, the van de Geer (2016, Theorem 12.1) penalty takes the form

$$\widehat{\lambda}_\alpha^{\text{vdG16}} = 8c_0 \sqrt{\frac{2 \ln(2p/\alpha)}{n} \max_{1 \leq j \leq p} \mathbb{E}_n [X_{ij}^2]}, \quad (6.1)$$

which is nothing more than 16 times our analytic penalty level (4.1). (Recall that $d = 1$.) Figure 1.1 in the Introduction displays the distribution of the van de Geer penalty level as a function of the sample size $n(= p)$, pooling over both correlation levels and coefficient patterns. For comparison, we include the distribution of the threshold penalty pooled over all designs.¹³ The latter threshold is the (approximately) smallest level of penalization needed to set every coefficient to zero, thus resulting in a trivial (null) model. The figure shows that the distribution of the threshold penalty is an order of magnitude closer to the origin than the van de Geer penalty. As a consequence, the latter penalty results in a trivial model estimate across *all* of our designs and simulation draws. The estimators resulting from the Bunea and van de Geer penalties are therefore all represented by the “Zeros” lines in Figures 6.1–6.4. Inspection of the proof underlying van de Geer (2016, Theorem 12.1) suggests that the factor of 8 in (6.1) may be reduced to a 2, when restricting attention to our framework. However,

¹³All density estimates in Figures 1.1 and 6.5 are created using the Matlab[®] package `ksdensity` with default settings.

Figure 6.5: Kernel Density Estimates of Penalty Distributions



even with this lower bound on the multiplier, the supports of these penalty distributions remain separated (Figure 6.5).

Our findings should not be interpreted as a critique of these authors, whose work were intended as primarily of theoretical interest. For example, [van de Geer \(2008, p. 621\)](#) explicitly states that other penalty choices should be used in practice. It is, however, not immediately clear how one should modify the penalty choices of these authors without disconnecting theory from practice. In contrast, the simulation results of this section demonstrate that our analytic and bootstrap methods are not only theoretically justifiable, but also practically useful.

References

- ADLER, R. J. AND J. E. TAYLOR (2007): *Random fields and geometry*, Springer Science & Business Media.
- AIGNER, D. J., T. AMEMIYA, AND D. J. POIRIER (1976): “On the estimation of production frontiers: maximum likelihood estimation of the parameters of a discontinuous density function,” *International Economic Review*, 377–396.
- BELLONI, A. AND V. CHERNOZHUKOV (2011a): *High dimensional sparse econometric models: An introduction*, Springer.

- (2011b): “l1-penalized quantile regression in high-dimensional sparse models,” *The Annals of Statistics*, 39, 82–130.
- BELLONI, A., V. CHERNOZHUKOV, D. CHETVERIKOV, C. HANSEN, AND K. KATO (2018a): “High-dimensional econometrics and regularized GMM,” *Working Paper*.
- BELLONI, A., V. CHERNOZHUKOV, D. CHETVERIKOV, AND Y. WEI (2018b): “Uniformly Valid Post-Regularization Confidence Regions for Many Functional Parameters in Z-Estimation Framework,” *The Annals of Statistics*, 3643–3675.
- BERTSEKAS, D. P. (1973): “Stochastic optimization problems with nondifferentiable cost functionals,” *Journal of Optimization Theory and Applications*, 12, 218–231.
- BOUCHERON, S., G. LUGOSI, AND P. MASSART (2012): *Concentration inequalities: A nonasymptotic theory of independence*, Clarendon press, Oxford.
- BUCHINSKY, M. AND J. HAHN (1998): “An alternative estimator for the censored quantile regression model,” *Econometrica*, 653–671.
- BUNEA, F. (2008a): “Consistent selection via the Lasso for high dimensional approximating regression models,” in *Pushing the limits of contemporary statistics: contributions in honor of Jayanta K. Ghosh*, Institute of Mathematical Statistics, 122–137.
- (2008b): “Honest variable selection in linear and logistic regression models via l1 and l1+ l2 penalization,” *Electronic Journal of Statistics*, 2, 1153–1194, publisher: The Institute of Mathematical Statistics and the Bernoulli Society.
- CHAMBERLAIN, G. (1984): “Panel data,” *Handbook of econometrics*, 2, 1247–1318.
- (1985): “Heterogeneity, omitted variable bias, and duration dependence,” in *Longitudinal Analysis of Labor Market Data*, ed. by J. J. Heckman and B. Singer, Cambridge University Press, 3–38.
- CHERNOZHUKOV, V., D. CHETVERIKOV, AND K. KATO (2013): “Gaussian approximations and multiplier bootstrap for maxima of sums of high-dimensional random vectors,” *The Annals of Statistics*, 41, 2786–2819.
- (2017): “Central limit theorems and bootstrap in high dimensions,” *The Annals of Probability*, 45, 2309–2352.
- CHETVERIKOV, D., Z. LIAO, AND V. CHERNOZHUKOV (2016): “On cross-validated Lasso,” *arXiv preprint arXiv:1605.02214*.

- FRIEDMAN, J., T. HASTIE, AND R. TIBSHIRANI (2010): “Regularization paths for generalized linear models via coordinate descent,” *Journal of statistical software*, 33, 1.
- HASTIE, T., R. TIBSHIRANI, AND M. WAINWRIGHT (2015): *Statistical Learning with Sparsity: The Lasso and Generalizations*, CRC Press.
- HONORÉ, B. E. (1992): “Trimmed LAD and least squares estimation of truncated and censored regression models with fixed effects,” *Econometrica: journal of the Econometric Society*, 533–565.
- IMAI, K. AND M. RATKOVIC (2014): “Covariate balancing propensity score,” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 76, 243–263.
- LANCASTER, T. (1992): *The econometric analysis of transition data*, 17, Cambridge university press.
- LECUE, G. AND G. MITCHELL (2012): “Oracle inequalities for cross-validation type procedures,” *Electronic Journal of Statistics*, 1803–1837.
- LEDOUX, M. AND M. TALAGRAND (1991): *Probability in Banach Spaces: Isoperimetry and Processes*, vol. 23, Springer Science & Business Media.
- MIOLANE, L. AND A. MONTANARI (2018): “The distribution of the Lasso: uniform control over sparse balls and adaptive parameter tuning,” *arXiv:1811.01212*.
- NEGAHBAN, S., P. RAVIKUMAR, M. WAINWRIGHT, AND B. YU (2012): “A unified framework for high-dimensional analysis of M-estimators with decomposable regularizers,” *Statistical Science*.
- NEWBY, W. K. AND J. L. POWELL (1987): “Asymmetric least squares estimation and testing,” *Econometrica: Journal of the Econometric Society*, 819–847.
- NG, A. Y. (2004): “Feature selection, L 1 vs. L 2 regularization, and rotational invariance,” in *Proceedings of the twenty-first international conference on Machine learning*, ACM, 78.
- NINOMIYA, Y. AND S. KAWANO (2016): “AIC for the Lasso in generalized linear models,” *Electronic Journal of Statistics*.
- PRATT, J. W. (1981): “Concavity of the log likelihood,” *Journal of the American Statistical Association*, 76, 103–106, publisher: Taylor & Francis.
- RASCH, G. (1960): “Studies in mathematical psychology: I. Probabilistic models for some intelligence and attainment tests.” .

- RIGOLLET, P. AND A. TSYBAKOV (2011): “Exponential screening and optimal rates of sparse estimation,” *Annals of Statistics*, 731–771.
- ROCKAFELLAR, R. T. (1970): “Convex Analysis (Princeton Mathematical Series),” *Princeton University Press*, 46, 49.
- TALAGRAND, M. (2010): *Mean field models for spin glasses: Volume I: Basic examples*, vol. 54, Springer Science & Business Media.
- TAN, Z. (2017): “Regularized calibrated estimation of propensity scores with model misspecification and high-dimensional data,” *arXiv:1710.08074v1*.
- TSYBAKOV, A. B. (2004): “Optimal aggregation of classifiers in statistical learning,” *The Annals of Statistics*, 32, 135–166, publisher: Institute of Mathematical Statistics.
- VAN DE GEER, S. (2016): “Estimation and testing under sparsity,” *Lecture notes in mathematics*, 2159, publisher: Springer.
- VAN DE GEER, S. A. (2008): “High-Dimensional Generalized Linear Models and the Lasso,” *The Annals of Statistics*, 36, 614–645.
- VERSHYNIN, R. (2018): *High-dimensional probability: An introduction with applications in data science*, Cambridge university press.
- WAINWRIGHT, M. (2019): *High-dimensional statistics: a non-asymptotic viewpoint*, Cambridge university press.
- WOOLDRIDGE, J. M. (2010): *Econometric analysis of cross section and panel data*, MIT press.

Online Appendix

We split the online appendix into four parts. In Appendix [A](#), we provide low-level conditions that are sufficient for Assumptions [14](#) and [15](#) in each of the examples considered in the main text. In Appendix [B](#), we provide proofs of the results stated in the main text. In Appendix [C](#), we provide auxiliary proofs. In Appendix [D](#), we provide a collection of technical tools used to prove the main results.

A Verification of Assumptions [14](#) and [15](#)

In this section, we verify Assumptions [14](#) and [15](#) in each of the examples from Section [2](#). Throughout this section, we use $m''_{11}(t, y)$ to denote the second derivative of the function $t \mapsto m(t, y)$, whenever it exists. Also, we suppose that Assumptions [1](#), [2](#), [4](#), [5](#) and [8](#) hold and that there exists a finite constant $C_d > 0$ such that

$$\|\theta\|_1 \leq C_d \quad \text{for all } \theta \in \Theta. \quad (\text{A.1})$$

Under these assumptions, there exists a finite constant $c_e > 0$ such that

$$\mathcal{E}(\theta) \geq c_e \|\theta - \theta_0\|_2^2 \quad \text{for all } \theta \in \Theta, \quad (\text{A.2})$$

which we use extensively throughout this section. To see why this bound holds, fix any $\theta \in \Theta$ and observe that if $\|\theta - \theta_0\|_1 \leq c'_M$, then $\mathcal{E}(\theta) \geq c_M \|\theta - \theta_0\|_2^2$ by Assumption [4](#). Therefore, we only need to consider the case $\|\theta - \theta_0\|_1 > c'_M$. In this case, for $t := c'_M / \|\theta - \theta_0\|_1$, we have by Assumption [2](#) that

$$t\mathcal{E}(\theta) + (1-t)\mathcal{E}(\theta_0) \geq \mathcal{E}(t\theta + (1-t)\theta_0),$$

and so, given that $\mathcal{E}(\theta_0) = 0$, we have

$$\mathcal{E}(\theta) \geq \frac{\mathcal{E}(\theta_0 + t(\theta - \theta_0))}{t} \geq \frac{c_M t^2 \|\theta - \theta_0\|_2^2}{t} = \frac{c_M c'_M \|\theta - \theta_0\|_2^2}{\|\theta - \theta_0\|_1} \geq \frac{c_M c'_M}{2C_d} \|\theta - \theta_0\|_2^2.$$

Hence, [\(A.2\)](#) holds with $c_e = c_M \wedge (c_M c'_M / 2C_d)$.

Finally, again throughout the section, we assume that there exists a finite constant $C_{ev} > 0$ such that

$$\lambda_{\max}(\mathbb{E}[XX^\top]) \leq C_{ev}, \quad (\text{A.3})$$

where $\lambda_{\max}(\mathbb{E}[XX^\top])$ denotes the largest eigenvalue of the matrix $\mathbb{E}[XX^\top]$.

Example 1 (Binary Response Model, Continued). In the case of the logit loss function (2.2), we have $m'_1(t, y) = \Lambda(t) - y$, so that $|m'_1(t, y)| \leq 1$ for all $t \in \mathbf{R}$ and $y \in \mathcal{Y}$. Hence, Assumption 14.1 holds by Hoeffding's Lemma (Boucheron et al., 2012, Lemma 2.2) and Assumption 15 holds by noting that for all $\theta \in \Theta$, we have

$$\begin{aligned} \mathbb{E} \left[\{m(X^\top \theta, Y) - m(X^\top \theta_0, Y)\}^2 \right] &\leq \mathbb{E} \left[m'_1(X^\top \tilde{\theta}, Y)^2 |X^\top (\theta - \theta_0)|^2 \right] \\ &\leq \mathbb{E} [|X^\top (\theta - \theta_0)|^2] \leq C_{ev} \|\theta - \theta_0\|_2^2 \leq (C_{ev}/c_e) \mathcal{E}(\theta), \end{aligned}$$

where the first inequality follows from the mean-value theorem with $\tilde{\theta}$ being a value on the line connecting θ_0 and θ , the third from (A.3), and the fourth from (A.2). In addition, $m''_{11}(t, y) = \Lambda'(t) = e^t/(1 + e^t)^2$, so that $|m''_{11}(t, y)| \leq 1$ for all $t \in \mathbf{R}$ and $y \in \mathcal{Y}$. Hence, Assumption 14.2 holds since for all $\theta \in \Theta$, we have

$$\begin{aligned} \mathbb{E} \left[\{m'_1(X^\top \theta, Y) - m'_1(X^\top \theta_0, Y)\}^2 \right] &\leq \mathbb{E} \left[m''_{11}(X^\top \tilde{\theta}, Y)^2 |X^\top (\theta - \theta_0)|^2 \right] \\ &\leq \mathbb{E} [|X^\top (\theta - \theta_0)|^2] \leq C_{ev} \|\theta - \theta_0\|_2^2 \leq (C_{ev}/c_e) \mathcal{E}(\theta), \end{aligned}$$

where the first inequality follows from the mean-value theorem with $\tilde{\theta}$ being a value on the line connecting θ_0 and θ , the third from (A.3), and the fourth from (A.2).

In the case of the probit loss function (2.3), we have

$$m'_1(t, y) = -\frac{y\varphi(t)}{\Phi(t)} + \frac{(1-y)\varphi(t)}{1-\Phi(t)}$$

and

$$m''_{11}(t, y) = \frac{yt\varphi(t)}{\Phi(t)} + \frac{y\varphi(t)^2}{\Phi(t)^2} - \frac{(1-y)t\varphi(t)}{1-\Phi(t)} + \frac{(1-y)\varphi(t)^2}{[1-\Phi(t)]^2}.$$

Therefore, given that the functions

$$t \mapsto \frac{\varphi(t)}{\Phi(t)} + \frac{\varphi(t)}{1-\Phi(t)}$$

and

$$t \mapsto \frac{|t|\varphi(t)}{\Phi(t)} + \frac{\varphi(t)^2}{\Phi(t)^2} + \frac{|t|\varphi(t)}{1-\Phi(t)} + \frac{\varphi(t)^2}{[1-\Phi(t)]^2}$$

are continuous, it follows that both $|m'_1(t, y)|$ and $m''_{11}(t, y)$ are bounded from above uniformly over $y \in \mathcal{Y}$ and $t \in [-C_X C_d, C_X C_d]$. Hence, Assumptions 14 and 15 hold by the same argument as that in the logit case since $|X^\top \theta_0| \leq C_X C_d$ almost surely by Assumption 5 and (A.1). \square

Example 2 (Ordered Response Model, Continued). In the ordered logit model

$$m'_1(t, y) = \sum_{j=0}^J \mathbf{1}(y = j) [1 - \Lambda(\alpha_{j+1} - t) - \Lambda(\alpha_j - t)],$$

thus implying

$$m''_{11}(t, y) = \sum_{j=0}^J \mathbf{1}(y = j) [\Lambda'(\alpha_{j+1} - t) + \Lambda'(\alpha_j - t)],$$

where $\Lambda'(\pm\infty)$ and $\Lambda(-\infty)$ are interpreted as zero and $\Lambda(+\infty)$ as one. Both Λ and Λ' are bounded (by one). Hence, verification of Assumptions 14 and 15 follows from an argument parallel to the one given for the binary logit. The argument for the ordered probit similarly runs in parallel to the one given in the binary case (see Example 1). \square

Example 3 (Logistic Calibration, Continued). Since $m(t, y) = ye^{-t} + (1-y)t$ in this example, we have

$$m'_1(t, y) = -ye^{-t} + 1 - y \quad \text{and} \quad m''_{11}(t, y) = ye^{-t}.$$

Therefore, given that the function $t \mapsto e^{-t}$ is continuous, it follows that both $|m'_1(t, y)|$ and $m''_{11}(t, y)$ are bounded from above uniformly over $y \in \mathcal{Y}$ and $t \in [-C_X C_d, C_X C_d]$. Hence, Assumptions 14 and 15 hold by noting that $|X^\top \theta_0| \leq C_X C_d$ almost surely by Assumption 5 and (A.1) like in Example 1. \square

Example 4 (Logistic Balancing, Continued). Since $m(t, y) = (1-y)e^t + ye^{-t} + (1-2y)t$ in this example, we have

$$m'_1(t, y) = (1-y)e^t - ye^{-t} + 1 - 2y \quad \text{and} \quad m''_{11}(t, y) = (1-y)e^t + ye^{-t}.$$

Therefore, given that the functions $t \mapsto e^t$ and $t \mapsto e^{-t}$ are continuous, it follows that both $|m'_1(t, y)|$ and $|m''_{11}(t, y)|$ are bounded from above uniformly over $y \in \mathcal{Y}$ and $t \in [-C_X C_d, C_X C_d]$. Hence, Assumptions 14 and 15 hold by noting that $|X^\top \theta_0| \leq C_X C_d$ almost surely by Assumption 5 and (A.1) like in Example 1. \square

Example 5 (Expectile Model, Continued). Since $m(t, y) = |\tau - \mathbf{1}(y - t < 0)|(y - t)^2$ in this example, we have

$$m'_1(t, y) = 2|\tau - \mathbf{1}(y - t < 0)|(t - y)$$

and

$$m''_{11}(t, y) = 2|\tau - \mathbf{1}(y - t < 0)| \quad \text{if } y \neq t.$$

Therefore, given that the function $t \mapsto |\tau - \mathbf{1}(y - t < 0)|$ is bounded (by one), Assumption

14.1 holds by noting that $|X^\top \theta_0| \leq C_X C_d$ almost surely by Assumption 5 and (A.1) and assuming that there exist a finite constant $C > 0$ such that

$$P(|Y| \geq t|X) \leq 2 \exp(-t^2/C) \quad \text{for all } t > 0 \text{ a.s.}; \quad (\text{A.4})$$

see Proposition 2.5.2 in Vershynin (2018). If we further assume that Y is continuously distributed conditional on X , then $m''_{11}(X^\top \theta, Y)$ exists with probability one for each $\theta \in \Theta$, and the mean-value theorem yields

$$\begin{aligned} \mathbb{E} \left[\{m'_1(X^\top \theta, Y) - m'_1(X^\top \theta_0, Y)\}^2 \right] &\leq \mathbb{E} \left[m''_{11}(X^\top \tilde{\theta}, Y)^2 |X^\top (\theta - \theta_0)|^2 \right] \\ &\leq 4\mathbb{E} [|X^\top (\theta - \theta_0)|^2] \leq (4C_{ev}/c_e)\mathcal{E}(\theta), \end{aligned}$$

where $\tilde{\theta}$ is some convex combination of θ and θ_0 and the rest is similar to the logit case in the binary response model above (Example 1). Assumption 14.2 follows.

To verify Assumption 15, note that there exists a finite constant $\tilde{C} > 0$ such that for all $\theta \in \Theta$, we have

$$\begin{aligned} \mathbb{E} \left[\{m(X^\top \theta, Y) - m(X^\top \theta_0, Y)\}^2 \right] &\leq \mathbb{E} \left[m'_1(X^\top \tilde{\theta}, Y)^2 |X^\top (\theta - \theta_0)|^2 \right] \\ &\leq 4\mathbb{E} [(Y - X^\top \theta_0)^2 |X^\top (\theta - \theta_0)|^2] \\ &\leq 8\mathbb{E} [(Y^2 + (X^\top \theta_0)^2) |X^\top (\theta - \theta_0)|^2] \\ &\leq \tilde{C}\mathbb{E} [|X^\top (\theta - \theta_0)|^2] \leq (\tilde{C}C_{ev}/c_e)\mathcal{E}(\theta), \end{aligned}$$

where the fourth inequality follows from (A.4) and (A.1) and Assumption 5 and the rest is similar to the logit case in the binary response model above (Example 1). \square

Example 6 (Panel Logit Model, Continued). Since $m(t, y) = \mathbf{1}(y_1 \neq y_2)[\ln(1 + e^t) - y_1 t]$ in this example, the verification of Assumptions 14 and 15 is analogous to that in the logit case of the binary response model above (Example 1). \square

Example 7 (Panel Censored Model, Continued). In the case $\Xi = (\cdot)^2$, we have

$$m'_1(t, y) = \begin{cases} -2y_1, & t \leq -y_2, \\ 2(t + y_2 - y_1), & -y_2 < t < y_1, \\ 2y_2, & y_1 \leq t, \end{cases}$$

and $m'_1(\cdot, y)$ is 2-Lipschitz for all $y \in \mathbf{R}_+^2$. Therefore, Assumptions 14 and 15 hold by noting that $|X^\top \theta_0| \leq C_X C_d$ almost surely by Assumption 5 and (A.1) and assuming that there

exist a finite constant $C > 0$ such that $P(|Y_1| \geq t|X) \leq 2 \exp(-t^2/C)$ and $P(|Y_2| \geq t|X) \leq 2 \exp(-t^2/C)$ for all $t > 0$ almost surely; see Proposition 2.5.2 in [Vershynin \(2018\)](#).

In the case $\Xi = |\cdot|$, the subdifferential $\partial_1 m(t, y)$ of the loss with respect to its first argument is contained in $[-1, 1]$ for all $t \in \mathbf{R}$ and all $y \in \mathbf{R}_+^2$. Since the subgradients of $m(\cdot, y)$ are bounded in absolute value by one, Assumptions [14.1](#) and [15](#) hold as in the logit case of the binary response model above (Example [1](#)). To verify Assumption [14.2](#), note that for all $\theta \in \Theta$, $m((X_1 - X_2)^\top \theta, Y)$ is differentiable at θ except when $Y_1 - Y_2 = (X_1 - X_2)^\top \theta$ and at least one of the Y_t 's is positive. Moreover, the piecewise linearity of the loss implies that the subgradients $g(\theta) \in \partial_1 m((X_1 - X_2)^\top \theta, Y)$ and $g(\theta_0) \in \partial_1 m((X_1 - X_2)^\top \theta_0, Y)$ may differ only if $(X_1 - X_2)^\top \theta$ and $(X_1 - X_2)^\top \theta_0$ are located at or on opposite sides of $Y_1 - Y_2$ and at least one of the Y_t 's is nonzero. Since these subgradients can be at most 2 apart, assuming that $(\varepsilon_1, \varepsilon_2)$ is continuously distributed conditional on (X_1, X_2, α) and that both the conditional PDF of ε_2 given $(\varepsilon_1, X_1, X_2, \alpha)$ and the conditional PDF of ε_1 given $(\varepsilon_2, X_1, X_2, \alpha)$ are bounded from above by some finite constant $C > 0$, we get

$$\begin{aligned} & \mathbb{E} \left[\left\{ m'_1 \left((X_1 - X_2)^\top \theta, Y \right) - m'_1 \left((X_1 - X_2)^\top \theta_0, Y \right) \right\}^2 \right] \\ & \leq 4 \mathbb{E} \left[\mathbf{1}(Y_1 > 0 \text{ or } Y_2 > 0) \mathbf{1} \left(\text{sgn} \left(Y_1 - Y_2 - (X_1 - X_2)^\top \theta \right) \neq \text{sgn} \left(Y_1 - Y_2 - (X_1 - X_2)^\top \theta_0 \right) \right) \right] \\ & \leq 4CE \left[|(X_1 - X_2)^\top (\theta - \theta_0)| \right] \\ & \leq 4C \left(\mathbb{E} \left[|(X_1 - X_2)^\top (\theta - \theta_0)|^2 \right] \right)^{1/2} \\ & \leq 4C \sqrt{C_{ev}} \|\theta - \theta_0\|_2 \leq (4C \sqrt{C_{ev}/c_e}) \sqrt{\mathcal{E}(\theta)}, \end{aligned}$$

where the second inequality follows from Jensen's inequality, the third from [\(A.3\)](#) (with C_{ev} now denoting the upper bound on the eigenvalues of $E[(X_1 - X_2)(X_1 - X_2)^\top]$), and the fourth from [\(A.2\)](#). This gives Assumption [14.2](#), as desired. \square

Example 8 (Panel Duration Model, Continued). Since $m(t, y) = \ln(1 + e^t) - \mathbf{1}(y_1 < y_2)t$ in this example, the verification of Assumptions [14](#) and [15](#) is analogous to that in the logit case of the binary response model above (Example [1](#)). \square

B Proofs for Statements in Main Text

In this section, we provide proofs of all results stated in the main text.

B.1 Proofs for Section [3](#)

PROOF OF THEOREM [1](#). We proceed in two steps.

Step 1: Abbreviate $\widehat{\theta} := \widehat{\theta}(\lambda)$. By minimization and the triangle inequality,

$$\mathbb{E}_n[m(X_i^\top \widehat{\theta}, Y_i) - m(X_i^\top \theta_0, Y_i)] \leq \lambda(\|\theta_0\|_1 - \|\widehat{\theta}\|_1) \leq \lambda(\|\widehat{\delta}_T\|_1 - \|\widehat{\delta}_{T^c}\|_1),$$

where $\widehat{\delta} := \widehat{\theta} - \theta_0$. By convexity followed by Hölder's inequality and score domination (\mathcal{S}),

$$\mathbb{E}_n[m(X_i^\top \widehat{\theta}, Y_i) - m(X_i^\top \theta_0, Y_i)] \geq S^\top(\widehat{\theta} - \theta_0) \geq -\|S\|_\infty \|\widehat{\delta}\|_1 \geq -\frac{\lambda}{c_0}(\|\widehat{\delta}_T\|_1 + \|\widehat{\delta}_{T^c}\|_1).$$

Combining the two previous displays, we get

$$\|\widehat{\delta}_{T^c}\|_1 \leq \frac{c_0 + 1}{c_0 - 1} \|\widehat{\delta}_T\|_1 = \bar{c}_0 \|\widehat{\delta}_T\|_1.$$

Therefore, on the event \mathcal{S} , we have $\widehat{\delta} \in \mathcal{R}(\bar{c}_0)$.

Step 2: Seeking a contradiction, suppose that we are on the event $\mathcal{S} \cap \mathcal{L} \cap \mathcal{E}$ but $\|\widehat{\delta}\|_2 > u_0$. Since Step 1 implies that $\widehat{\delta} \in \mathcal{R}(\bar{c}_0)$, it then follows by minimization that

$$\begin{aligned} 0 &\geq \inf_{\substack{\delta \in \mathcal{R}(\bar{c}_0), \\ \|\delta\|_2 > u_0}} \left\{ \mathbb{E}_n[m(X_i^\top(\theta_0 + \delta), Y_i) - m(X_i^\top \theta_0, Y_i)] + \lambda(\|\theta_0 + \delta\|_1 - \|\theta_0\|_1) \right\} \\ &\geq \inf_{\substack{\delta \in \mathcal{R}(\bar{c}_0), \\ \|\delta\|_2 \geq u_0}} \left\{ \mathbb{E}_n[m(X_i^\top(\theta_0 + \delta), Y_i) - m(X_i^\top \theta_0, Y_i)] + \lambda(\|\theta_0 + \delta\|_1 - \|\theta_0\|_1) \right\}. \end{aligned}$$

Now, since $\delta \mapsto \{\mathbb{E}_n[m(X_i^\top(\theta_0 + \delta), Y_i) - m(X_i^\top \theta_0, Y_i)] + \lambda(\|\theta_0 + \delta\|_1 - \|\theta_0\|_1)\}$ is a (random) convex function taking the value 0 when $\delta = \mathbf{0} \in \mathbf{R}^p$ and $\mathcal{R}(\bar{c}_0)$ a cone (i.e., $\delta \in \mathcal{R}(\bar{c}_0)$ implies $t\delta \in \mathcal{R}(\bar{c}_0)$ for any $t \in \mathbf{R}_{++}$), the previous display implies that

$$0 \geq \inf_{\substack{\delta \in \mathcal{R}(\bar{c}_0), \\ \|\delta\|_2 = u_0}} \left\{ \mathbb{E}_n[m(X_i^\top(\theta_0 + \delta), Y_i) - m(X_i^\top \theta_0, Y_i)] + \lambda(\|\theta_0 + \delta\|_1 - \|\theta_0\|_1) \right\}.$$

By superadditivity of infima and definition of the empirical error function, on the event \mathcal{L} , the right-hand side here is bounded from below by

$$\begin{aligned} &\inf_{\substack{\delta \in \mathcal{R}(\bar{c}_0), \\ \|\delta\|_2 = u_0}} \mathbb{E}[m(X^\top(\theta_0 + \delta), Y_i) - m(X^\top \theta_0, Y_i)] \\ &+ \inf_{\substack{\delta \in \mathcal{R}(\bar{c}_0), \\ \|\delta\|_2 = u_0}} (\mathbb{E}_n - \mathbb{E})[m(X_i^\top(\theta_0 + \delta), Y_i) - m(X_i^\top \theta_0, Y_i)] + \lambda \inf_{\substack{\delta \in \mathcal{R}(\bar{c}_0), \\ \|\delta\|_2 = u_0}} \{\|\theta_0 + \delta\|_1 - \|\theta_0\|_1\} \\ &\geq \inf_{\substack{\delta \in \mathcal{R}(\bar{c}_0), \\ \|\delta\|_2 = u_0}} \mathbb{E}[m(X^\top(\theta_0 + \delta), Y_i) - m(X^\top \theta_0, Y_i)] - \epsilon(u_0) - \bar{\lambda} \sup_{\substack{\delta \in \mathcal{R}(\bar{c}_0), \\ \|\delta\|_2 = u_0}} \|\theta_0 + \delta\|_1 - \|\theta_0\|_1. \end{aligned}$$

Next, since we assume that $(1 + \bar{c}_0)u_0\sqrt{s} \leq c'_M$, any $\delta \in \mathcal{R}(\bar{c}_0)$ such that $\|\delta\|_2 = u_0$ must satisfy

$$\|\delta\|_1 \leq (1 + \bar{c}_0)\|\delta_T\|_1 \leq (1 + \bar{c}_0)\sqrt{s}\|\delta_T\|_2 \leq (1 + \bar{c}_0)\sqrt{s}\|\delta\|_2 \leq c'_M. \quad (\text{B.1})$$

Therefore, by Assumption 4,

$$\inf_{\substack{\delta \in \mathcal{R}(\bar{c}_0), \\ \|\delta\|_2 = u_0}} \mathbb{E}[m(X_i^\top(\theta_0 + \delta), Y_i) - m(X_i^\top\theta_0, Y_i)] \geq c_M \inf_{\substack{\delta \in \mathcal{R}(\bar{c}_0), \\ \|\delta\|_2 = u_0}} \|\delta\|_2^2 = c_M u_0^2.$$

Also, by the triangle inequality,

$$\sup_{\substack{\delta \in \mathcal{R}(\bar{c}_0), \\ \|\delta\|_2 = u_0}} \left| \|\theta_0 + \delta\|_1 - \|\theta_0\|_1 \right| \leq \sup_{\substack{\delta \in \mathcal{R}(\bar{c}_0), \\ \|\delta\|_2 = u_0}} \|\delta\|_1 \leq (1 + \bar{c}_0)u_0\sqrt{s}.$$

In addition, $\epsilon(u_0) \leq \lambda_\epsilon u_0$ on the event \mathcal{E} . Therefore, it follows that

$$\begin{aligned} 0 &\geq \inf_{\substack{\delta \in \mathcal{R}(\bar{c}_0), \\ \|\delta\|_2 = u_0}} \mathbb{E}[m(X^\top(\theta_0 + \delta), Y) - m(X^\top\theta_0, Y)] - \epsilon(u_0) - \bar{\lambda} \sup_{\substack{\delta \in \mathcal{R}(\bar{c}_0), \\ \|\delta\|_2 = u_0}} \{ \|\theta_0\|_1 - \|\theta_0 + \delta\|_1 \}. \\ &\geq c_M u_0^2 - \lambda_\epsilon u_0 - (1 + \bar{c}_0)\bar{\lambda}u_0\sqrt{s}. \end{aligned}$$

However, by definition of u_0 ,

$$c_M u_0^2 - \lambda_\epsilon u_0 - (1 + \bar{c}_0)\bar{\lambda}u_0\sqrt{s} = c_M u_0^2 - (\lambda_\epsilon + (1 + \bar{c}_0)\bar{\lambda}\sqrt{s})u_0 = \frac{c_M}{2}u_0^2 > 0,$$

yielding the desired contradiction. We therefore conclude that on the event $\mathcal{S} \cap \mathcal{L} \cap \mathcal{E}$ we have $\|\widehat{\delta}\|_2 \leq u_0$, which establishes the ℓ^2 bound. The ℓ^1 bound then follows from the ℓ^2 bound and (B.1). \square

PROOF OF PROPOSITION 1. The interiority, convexity and differentiability assumptions imply that $\nabla M(\theta_0) = \mathbf{0}$. Let $\theta \in \Theta$ satisfy $\|\theta - \theta_0\|_1 < c_1/C_X$. Then for any convex combination $\bar{\theta}$ of θ and θ_0

$$\|X^\top\bar{\theta} - X^\top\theta_0\| \leq \|X\|_\infty \|\bar{\theta} - \theta_0\|_1 \leq \|X\|_\infty \|\theta - \theta_0\|_1 < c_1$$

almost surely. Provided that we may interchange the order of differentiation and integration,

a second-order mean-value expansion yields

$$\begin{aligned}
M(\theta) - M(\theta_0) &= \frac{1}{2} (\theta - \theta_0)^\top \nabla^2 M(\bar{\theta}) (\theta - \theta_0) \\
&= \frac{1}{2} \mathbf{E} \left[\frac{\partial^2}{\partial^2 t} \mathbf{E} [m(t, Y) | X] \Big|_{t=X^\top \bar{\theta}} \{X^\top (\theta - \theta_0)\}^2 \right] \\
&\geq \frac{1}{2} c_2 \mathbf{E} [\{X^\top (\theta - \theta_0)\}^2] \\
&\geq \frac{1}{2} c_2 c_e \|\theta - \theta_0\|_2^2,
\end{aligned}$$

hence yielding the asserted claim and completing the proof of the proposition. \square

B.2 Proofs for Section 3.2

PROOF OF LEMMA 1. The claim will follow from an application of the maximal inequality in Lemma D.1. Setting up for such an application, fix $0 < u \leq c_L / [C_X (1 + \bar{c}) \sqrt{s}]$ and define $\Delta(u) := \mathcal{R}(\bar{c}_0) \cap \{\delta \in \mathbf{R}^p; \|\delta\|_2 \leq u\}$. The zero vector in \mathbf{R}^p belongs to both $\mathcal{R}(\bar{c}_0)$ and $\{\delta \in \mathbf{R}^p; \|\delta\|_2 \leq u\}$, so $\Delta(u)$ is a nonempty subset of \mathbf{R}^p . By definition of $\Delta(u)$, any $\delta \in \Delta(u)$ must satisfy

$$\|\delta\|_1 \leq (1 + \bar{c}_0) \|\delta_T\|_1 \leq (1 + \bar{c}_0) \sqrt{s} \|\delta_T\|_2 \leq (1 + \bar{c}_0) \sqrt{s} \|\delta\|_2 \leq (1 + \bar{c}_0) u \sqrt{s},$$

thus implying

$$\|\Delta(u)\|_1 := \sup_{\delta \in \Delta(u)} \|\delta\|_1 \leq (1 + \bar{c}_0) u \sqrt{s}. \quad (\text{B.2})$$

Next, define $h : \mathbf{R} \times \mathcal{W} \rightarrow \mathbf{R}$ by $h(t, w) := m(x^\top \theta_0 + t, y) - m(x^\top \theta_0, y)$ for all $t \in \mathbf{R}$ and $w = (x, y) \in \mathcal{W}$. By Assumption 6.1, $h : [-c_L, c_L] \times \mathcal{W} \rightarrow \mathbf{R}$ is Lipschitz in its first argument and satisfies $h(0, \cdot) \equiv 0$, thus verifying Condition 1 of Lemma D.1. Condition 2 of the same lemma follows from Hölder's inequality, Assumption 5, (B.2), the upper bound on u and the calculation

$$\sup_{\delta \in \Delta(u)} |X^\top \delta| \leq \|X\|_\infty \|\Delta(u)\|_1 \leq C_X (1 + \bar{c}_0) u \sqrt{s} \leq c_L \text{ a.s.}$$

Given that $h : [-c_L, c_L] \times \mathcal{W} \rightarrow \mathbf{R}$ is both Lipschitz in its first argument and satisfies $h(0, \cdot) \equiv 0$ as well as $\sup_{\delta \in \Delta(u)} |X^\top \delta| \leq c_L$ almost surely,

$$\begin{aligned} \sup_{\delta \in \Delta(u)} \mathbb{E}[h(X^\top \delta, W)^2] &\leq \sup_{\delta \in \Delta(u)} \mathbb{E}[L(W)^2 (X^\top \delta)^2] \\ &\leq u^2 \sup_{\delta \in \mathcal{S}_{p-1}} \mathbb{E}[L(W)^2 (X^\top \delta)^2] \leq C_{L,e}^2 u^2 \end{aligned}$$

by Assumption 7. Condition 3 of Lemma D.1 therefore holds for $B_{1n} = C_{L,e}u$. Given that $\|X\|_\infty \leq C_X$ almost surely and $\mathbb{E}_n[L(W_i)^2] \leq C_L^2$ with probability at least $1 - n^{-1}$ by Assumption 6.2 and Chebyshev's inequality, it follows that

$$\max_{1 \leq j \leq p} \mathbb{E}_n[L(W_i)^2 X_{ij}^2] \leq C_L^2 C_X^2$$

with the same probability. Condition 4 of the lemma therefore holds with $B_{2n} = C_L C_X$ and $\gamma_n = n^{-1}$. Lemma D.1 combined with the bounds on $\|\Delta(u)\|_1$ from (B.2) and $\ln(8pn) \leq 4 \ln(pn)$ (which follows from $p \geq 2$) therefore shows that for all $n \in \mathbf{N}$,

$$\begin{aligned} &\mathbb{P}\left(\epsilon(u) > (\{4C_{L,e}\} \vee \{C_\epsilon \sqrt{s \ln(pn)}\})u/\sqrt{n}\right) \\ &= \mathbb{P}\left(\sup_{\delta \in \Delta(u)} |\mathbb{G}_n[h(X_i^\top \delta, W_i)]| > (\{4C_{L,e}\} \vee \{C_\epsilon \sqrt{s \ln(pn)}\})u\right) \leq 5n^{-1}. \end{aligned}$$

The claim now follows since we assume that $s \ln(pn) \geq 16C_{L,e}^2/C_\epsilon^2$. \square

B.3 Proofs for Section 4

PROOF OF THEOREM 2. We set up for an application of Theorem 1. To this end, define $\lambda_\epsilon := C_\epsilon \sqrt{s \ln(pn)/n}$ and $\bar{\lambda} := \bar{\lambda}_\alpha^{\text{am}}$ [see (4.2)], which are positive and finite under our assumptions. Then it follows from Lemma 1, whose application is justified by the inequalities in (4.3), that $\epsilon(u_0) \leq \lambda_\epsilon u_0$ with probability at least $1 - 5n^{-1}$, meaning that $\mathbb{P}(\mathcal{E}) \geq 1 - 5n^{-1}$. Also, observe that by the choice of penalty (4.1) and Assumptions 8 and 9, the event $\hat{\lambda}_\alpha^{\text{am}} \geq c_0 \|S\|_\infty$ occurs with probability at least $1 - \alpha$, as discussed in the main text, meaning that $\mathbb{P}(\mathcal{S}) \geq 1 - \alpha$. In addition, $\mathbb{P}(\mathcal{L}) = 1$ by (4.2). Therefore, the asserted claims follow from Theorem 1, whose application is again justified by inequalities in (4.3). \square

PROOF OF COROLLARY 1. The assumption (4.4) ensures that (4.3) holds for all n large enough. Therefore, the asserted claim follows immediately from Theorem 2. \square

B.4 Proofs for Section 5.1

PROOF OF THEOREM 3. We set up for an application of Theorem 1. To this end, define $\lambda_\epsilon := C_\epsilon \sqrt{s \ln(pn)/n}$. Then it follows from Lemma 1, whose application is justified by inequalities in (5.5), that $\epsilon(u_0) \leq \lambda_\epsilon u_0$ with probability at least $1 - 5n^{-1}$, meaning that $P(\mathcal{E}) \geq 1 - 5n^{-1}$.

Further, observe that conditional on $\{(W_i, \widehat{U}_i)\}_{i=1}^n$, the random vector $\mathbb{E}_n[e_i \widehat{U}_i X_i]$ is centered Gaussian in \mathbf{R}^p with j th coordinate variance $n^{-1} \mathbb{E}_n[\widehat{U}_i^2 X_{ij}^2]$. Lemma D.2 therefore shows that

$$\widehat{q}(1 - \alpha) \leq (2 + \sqrt{2}) \sqrt{\frac{\ln(p/\alpha)}{n} \max_{1 \leq j \leq p} \mathbb{E}_n[\widehat{U}_i^2 X_{ij}^2]}.$$

In addition, with probability at least $1 - \beta_n - n^{-1}$,

$$\max_{1 \leq j \leq p} \mathbb{E}_n[\widehat{U}_i^2 X_{ij}^2] \leq C_X^2 \mathbb{E}_n[\widehat{U}_i^2] \leq 2C_X^2 \left(\mathbb{E}_n[U_i^2] + \mathbb{E}_n[(\widehat{U}_i - U_i)^2] \right) \leq 4C_X^2 (C_U^2 + \delta_n^2 / \ln^2(pn)),$$

where the first inequality follows from Assumption 5, the second from the elementary inequality $(a + b)^2 \leq 2a^2 + 2b^2$, and the third from Assumptions 10 and 11 and Chebyshev's inequality. Hence, with the same probability,

$$\widehat{\lambda}_\alpha^{\text{bm}} \leq \bar{\lambda} := 2(2 + \sqrt{2})c_0 C_X (C_U + \delta_n / \ln(pn)) \sqrt{\frac{\ln(p/\alpha)}{n}},$$

meaning that $P(\mathcal{L}) \geq 1 - \beta_n - n^{-1}$.

Next, Assumptions 5 and 10 imply that the moment conditions (D.1) for $Z_{ij} = U_i X_{ij}$ hold with b replaced by c_U and B_n replaced by $B = (C_X C_U)^3 \vee 1$. Further, the same assumptions also imply the estimation error condition (D.2) for $\widehat{Z}_{ij} = \widehat{U}_i X_{ij}$ hold with δ_n replaced by $C_X \delta_n$. Since the Z_i 's are centered (by Assumption 8), Theorem D.4 therefore shows that there exists a finite constant C depending only on c_U such that¹⁴

$$\begin{aligned} & \sup_{\alpha \in (0,1)} |P(\|S\|_\infty > \widehat{q}(1 - \alpha)) - \alpha| \\ & \leq C \max \left\{ \beta_n, C_X \delta_n, \left(\frac{B^4 \ln^7(pn)}{n} \right)^{1/6}, \frac{1}{\ln^2(pn)} \right\} = \rho_n. \end{aligned}$$

It thus follows by construction of the bootstrap penalty level $\widehat{\lambda}_\alpha^{\text{bm}} = c_0 \widehat{q}(1 - \alpha)$ that the event $\widehat{\lambda}_\alpha^{\text{bm}} \geq c_0 \|S\|_\infty$ occurs with probability at least $1 - \alpha - \rho_n$, meaning that $P(\mathcal{S}) \geq 1 - \alpha - \rho_n$.

Therefore, the asserted claims follow from Theorem 1, whose application is again justified

¹⁴We here invoke the scaling property that $q_{tV}(\alpha) = tq_V(\alpha)$ for $t > 0$ and $\alpha \in (0, 1)$ and $q_V(\alpha)$ denoting the α quantile of the random variable V .

by the inequalities in (5.5). □

B.5 Proofs for Section 5.2

For the arguments in this section, we introduce some additional notation. For any nonempty $I \subsetneq \{1, \dots, n\}$, define the *subsample score*

$$S_I := \mathbb{E}_I [m'_1 (X_i^\top \theta_0, Y_i) X_i]$$

and *subsample empirical error*

$$\epsilon_I(u) := \sup_{\substack{\delta \in \mathcal{R}(\bar{c}_0), \\ \|\delta\|_2 \leq u}} |(\mathbb{E}_I - \mathbb{E}) [m (X_i^\top (\theta_0 + \delta), Y_i) - m (X_i^\top \theta_0, Y_i)]|, \quad u \in \mathbf{R}_+.$$

In proving Theorem 4, we rely on the following lemmas.

Lemma B.1. *Let Assumption 13 hold. Then for any finite constant $C > 0$ satisfying $n^{-1} \ln(pn) \leq (C_\Lambda a/C)^2$ and $n \ln(pn) \geq (c_\Lambda/C)^2$, the candidate penalty set Λ_n and the interval $[C\sqrt{n^{-1} \ln(pn)}, (C/a)\sqrt{n^{-1} \ln(pn)}]$ have an element in common.*

Lemma B.2. *Let Assumptions 5, 6, 7, and 12 hold, and define the finite constant $C_\epsilon := 16\sqrt{2}(1 + \bar{c}_0)C_L C_X > 0$. Then provided*

$$s \ln(pn) \geq 16(K-1)c_D \left(\frac{C_{L,e}}{C_\epsilon} \right)^2 \quad \text{and} \quad 0 < u \leq \frac{c_L}{(1 + \bar{c}_0)C_X \sqrt{s}},$$

we have

$$\max_{1 \leq k \leq K} \epsilon_{I_k^c}(u) \leq \frac{C_\epsilon}{(K-1)c_D} u \sqrt{\frac{s \ln(pn)}{n}}$$

with probability at least $1 - K(4n^{-1} + [(K-1)c_D n]^{-1})$.

Lemma B.3. *Let Assumptions 5, 12 and 14 hold and define the finite constant $C_S := 2C_X \sigma / \sqrt{(K-1)c_D}$. Then*

$$\max_{1 \leq k \leq K} \|S_{I_k^c}\|_\infty \leq C_S \sqrt{\ln(pn)/n}$$

with probability at least $1 - K[(K-1)c_D n]^{-1}$.

Lemma B.4. *Let Assumptions 1, 2, 3, and 4 hold. Fix some finite constants $\lambda_\epsilon, \bar{\lambda} > 0$ and $k \in \{1, \dots, K\}$ and define $u_0 := (2/c_M)(\lambda_\epsilon + (1 + \bar{c}_0)\bar{\lambda}\sqrt{s})$. In addition, suppose that*

$(1 + \bar{c}_0)u_0\sqrt{s} \leq c'_M$. Then for any (possibly random) $\lambda \in \Lambda_n$, on the event $\{\lambda \geq c_0\|S_{I_k^c}\|_\infty\} \cap \{\lambda \leq \bar{\lambda}\} \cap \{\epsilon_{I_k^c}(u_0) \leq \lambda_\epsilon u_0\}$, we have

$$\mathcal{E}(\widehat{\theta}_{I_k^c}(\lambda)) \leq \frac{2}{c_M} (\lambda_\epsilon + (1 + \bar{c}_0)\bar{\lambda}\sqrt{s})^2.$$

Lemma B.5. Let Assumptions 1–7 and 12–14 hold and define the finite constants $C_\epsilon := 16\sqrt{2}(1 + \bar{c}_0)C_L C_X > 0$. $C_S := 2C_X\sigma/\sqrt{(K-1)c_D} > 0$, and

$$\tilde{u}_0 := \frac{2}{c_M} \left(\frac{C_\epsilon}{(K-1)c_D} + \frac{(1 + \bar{c}_0)c_0 C_S}{a} \right) \sqrt{\frac{s \ln(pn)}{n}} > 0. \quad (\text{B.3})$$

In addition, suppose that the following inequalities hold:

$$\left\{ \begin{array}{l} s \ln(pn) \geq 16(K-1)c_D C_{L,e}^2 / C_\epsilon^2, \\ (1 + \bar{c}_0)\tilde{u}_0\sqrt{s} \leq (c_L/C_X) \wedge c'_M, \\ n^{-1} \ln(pn) \leq (C_\Lambda a / c_0 C_S)^2, \\ \text{and } n \ln(pn) \geq (c_\Lambda / c_0 C_S)^2, \end{array} \right\}. \quad (\text{B.4})$$

Then there exists a candidate penalty level $\lambda_* \in \Lambda_n$ (possibly depending on n), such that

$$\max_{1 \leq k \leq K} \mathcal{E}(\widehat{\theta}_{I_k^c}(\lambda_*)) \leq \frac{2}{c_M} \left(\frac{C_\epsilon}{(K-1)c_D} + \frac{(1 + \bar{c}_0)c_0 C_S}{a} \right)^2 \frac{s \ln(pn)}{n}$$

with probability at least $1 - K(4n^{-1} + 2[(K-1)c_D n]^{-1})$,

Lemma B.6. Let Assumptions 1–7 and 12–15 hold and define the finite constants $C_\epsilon := 16\sqrt{2}(1 + \bar{c}_0)C_L C_X > 0$, $C_S := 2C_X\sigma/\sqrt{(K-1)c_D} > 0$ and

$$C_\mathcal{E} := \sqrt{\frac{2}{c_M}} \left(\frac{C_\epsilon}{(K-1)c_D} + \frac{(1 + \bar{c}_0)c_0 C_S}{a} \right) > 0.$$

In addition, suppose that the inequalities (B.4) hold with \tilde{u}_0 appearing in (B.3). Then for any $t \in \mathbf{R}_{++}$ such that

$$\left\{ n \geq \frac{1}{c_\Lambda}, \quad \frac{C_\mathcal{E}^2 s \ln(pn)}{n} \leq 1, \quad \text{and} \quad 2C_{ms} \sqrt{\frac{t \ln n}{c_D \ln(1/a) n}} \leq \frac{1}{2} \right\}, \quad (\text{B.5})$$

we have

$$\max_{1 \leq k \leq K} \mathcal{E}(\widehat{\theta}_{I_k^c}(\widehat{\lambda}^{\text{cv}})) \leq \frac{32C_{ms}^2}{c_D^2 \ln(1/a)} \frac{t \ln n}{n} + \frac{8C_\mathcal{E}^2}{c_D} \frac{s \ln(pn)}{n} \quad (\text{B.6})$$

with probability at least $1 - K(4n^{-1} + 2[(K-1)c_D n]^{-1} + t^{-1})$.

PROOF OF THEOREM 4. Fix any $t \in \mathbf{R}_{++}$ satisfying (5.14) and $\lambda \in \Lambda_n$. For all $k \in \{1, \dots, K\}$, by Assumption 14 and Markov's inequality applied conditional on $\{W_i\}_{i \in I_k^c}$, we have

$$\mathbb{P}\left(\mathbb{E}_{I_k} \left[\left\{ m'_1(X_i^\top \widehat{\theta}_{I_k^c}(\lambda), Y_i) - m'_1(X_i^\top \theta_0, Y_i) \right\}^2 \right] > C_{ms1}^2 t \left[\sqrt{\mathcal{E}(\widehat{\theta}_{I_k^c}(\lambda))} \vee \mathcal{E}(\widehat{\theta}_{I_k^c}(\lambda)) \right] \right) \leq t^{-1}.$$

In addition, since $n \geq 1/c_\Lambda$ by (5.13), Assumption 13 implies that $|\Lambda_n| \leq 2(\ln n) / \ln(1/a)$. Therefore, by the union bound, for all $k \in \{1, \dots, K\}$,

$$\begin{aligned} & \mathbb{P}\left(\exists \lambda \in \Lambda_n \text{ s.t. } \mathbb{E}_{I_k} \left[\left\{ m'_1(X_i^\top \widehat{\theta}_{I_k^c}(\lambda), Y_i) - m'_1(X_i^\top \theta_0, Y_i) \right\}^2 \right] \right. \\ & \quad \left. > \frac{2C_{ms1}^2 t \ln n}{\ln(1/a)} \left[\sqrt{\mathcal{E}(\widehat{\theta}_{I_k^c}(\lambda))} \vee \mathcal{E}(\widehat{\theta}_{I_k^c}(\lambda)) \right] \right) \leq t^{-1}. \end{aligned}$$

Next, introduce events $\mathcal{C} := \bigcap_{k=1}^K \mathcal{C}_k$, where

$$\begin{aligned} \mathcal{C}_k &:= \left\{ \mathbb{E}_{I_k} \left[\left\{ m'_1(X_i^\top \widehat{\theta}_{I_k^c}(\widehat{\lambda}^{\text{cv}}), Y_i) - m'_1(X_i^\top \theta_0, Y_i) \right\}^2 \right] \right. \\ & \quad \left. \leq \frac{2C_{ms1}^2 t \ln n}{\ln(1/a)} \left[\sqrt{\mathcal{E}(\widehat{\theta}_{I_k^c}(\lambda))} \vee \mathcal{E}(\widehat{\theta}_{I_k^c}(\lambda)) \right] \right\}, \end{aligned}$$

and

$$\mathcal{R} := \left\{ \max_{1 \leq k \leq K} \mathcal{E}(\widehat{\theta}_{I_k^c}(\widehat{\lambda}^{\text{cv}})) \leq \frac{32C_{ms}^2}{c_D^2 \ln(1/a)} \frac{t \ln n}{n} + \frac{8C_\mathcal{E}^2 s \ln(pn)}{c_D n} \right\}.$$

Given that the cross-validated penalty $\widehat{\lambda}^{\text{cv}}$ is a random element of Λ_n , it follows that $\max_{1 \leq k \leq K} \mathbb{P}(\mathcal{C}_k^c) \leq 1/t$, and so, by the union bound, $\mathbb{P}(\mathcal{C}^c) \leq K/t$. Moreover, by Lemma B.6, whose application is justified by the inequalities in (5.12), (5.13), and (5.14), we have $\mathbb{P}(\mathcal{R}^c) \leq K(4n^{-1} + 2[(K-1)c_D n]^{-1} + t^{-1})$. Therefore, again by the union bound, $\mathbb{P}(\mathcal{C} \cap \mathcal{R}) \geq 1 - K(4n^{-1} + 2[(K-1)c_D n]^{-1} + 2t^{-1})$. But on $\mathcal{C} \cap \mathcal{R}$, we have

$$\begin{aligned} \mathbb{E}_n [(\widehat{U}_i^{\text{cv}} - U_i)^2] &= \sum_{k=1}^K \frac{|I_k|}{n} \mathbb{E}_{I_k} \left[\left\{ m'_1(X_i^\top \widehat{\theta}_{I_k^c}(\widehat{\lambda}^{\text{cv}}), Y_i) - m'_1(X_i^\top \theta_0, Y_i) \right\}^2 \right] \\ &\leq \frac{2C_{ms1}^2 t \ln n}{\ln(1/a)} \sum_{k=1}^K \frac{|I_k|}{n} \left[\sqrt{\mathcal{E}(\widehat{\theta}_{I_k^c}(\lambda))} \vee \mathcal{E}(\widehat{\theta}_{I_k^c}(\lambda)) \right] \\ &\leq \frac{2C_{ms1}^2 t \ln n}{\ln(1/a)} \left(\frac{32C_{ms}^2}{c_D^2 \ln(1/a)} \frac{t \ln n}{n} + \frac{8C_\mathcal{E}^2 s \ln(pn)}{c_D n} \right)^{1/2}, \end{aligned}$$

where the first inequality follows from \mathcal{C} and the second from \mathcal{R} and (5.14). This gives the asserted claim and completes the proof. \square

PROOF OF COROLLARY 2. The assumption (5.16) ensures that there exists a sequence t_n of constants \mathbf{R}_+ such that both

$$t_n \rightarrow \infty \quad \text{and} \quad \frac{t_n^3 s \ln^5(pn)(\ln n)^2}{n} \rightarrow 0. \quad (\text{B.7})$$

Therefore, setting $\delta_n > 0$ such that

$$\delta_n^2 := \frac{8C_{ms1}^2 t_n \ln n}{\ln(1/a)} \left(\frac{2C_{ms}^2}{c_D^2 \ln(1/a)} \frac{t_n \ln n}{n} + \frac{C_{\mathcal{E}}^2 s \ln(pn)}{2c_D} \frac{1}{n} \right)^{1/2} \ln^2(pn)$$

and

$$\beta_n := K (4n^{-1} + 2[(K-1)c_D n]^{-1} + 2t_n^{-1}),$$

we have both $\delta_n \rightarrow 0$ and $\beta_n \rightarrow 0$. In addition, (B.7) implies that (5.14) with $t = t_n$ holds for all n large enough. Also, (5.16) ensures that (5.12) and (5.13) hold for all n large enough as well. Hence, Theorem 4 implies that Assumption 11 with δ_n and β_n thus chosen holds for all n large enough. Thus, given that (5.16) also ensures that (5.5) holds for all n large enough, the asserted claim follows from Theorem 3. \square

C Proofs for Supporting Lemmas

In this section, we prove Lemmas B.1–B.6 used in the proof of Theorem 4.

PROOF OF LEMMA B.1. Denote $b_n := C\sqrt{n^{-1} \ln(pn)}$. We will show that there exists an integer ℓ_0 such that

$$c_\Lambda/n \leq b_n \leq C_\Lambda a^{\ell_0} \leq b_n/a \leq C_\Lambda. \quad (\text{C.1})$$

By Assumption 13, this will imply that $C_\Lambda a^{\ell_0}$ belongs to both the candidate penalty set Λ_n and the interval $[C\sqrt{n^{-1} \ln(pn)}, (C/a)\sqrt{n^{-1} \ln(pn)}]$.

To prove (C.1), note that the condition $n^{-1} \ln(pn) \leq (C_\Lambda a/C)^2$ implies that

$$0 \leq \frac{\ln(b_n/C_\Lambda)}{\ln a} - 1. \quad (\text{C.2})$$

In addition, there exists an integer ℓ_0 such that

$$\frac{\ln(b_n/C_\Lambda)}{\ln a} - 1 \leq \ell_0 \leq \frac{\ln(b_n/C_\Lambda)}{\ln a}. \quad (\text{C.3})$$

Combining (C.2) and (C.3), we obtain $b_n \leq C_\Lambda a^{\ell_0} \leq b_n/a \leq C_\Lambda$. Moreover, the condition $n \ln(pn) \geq (c_\Lambda/C)^2$ implies that $c_\Lambda/n \leq b$. Combining these inequalities gives (C.1) and

completes the proof of the lemma. \square

PROOF OF LEMMA B.2. The claim will follow from an application of the maximal inequality in Lemma D.1 in a manner very similar to the proof of Lemma 1. Verification of Conditions 1–3 of Lemma D.1 are unrelated to sampling. It thus remains to verify Condition 4. To do so, fix a (hold-out) subsample $k \in \{1, \dots, K\}$. Given that $\|X\|_\infty \leq C_X$ almost surely (Assumption 5), $|I_k| \geq c_D n$ (Assumption 12) and $\mathbb{E}_n[L(W_i)^2] \leq C_L^2$ with probability at least $1 - n^{-1}$ (Assumption 6 and Chebyshev's inequality), it follows that

$$\max_{1 \leq j \leq p} \mathbb{E}_{I_k^c}[L(W_i)^2 X_{ij}^2] \leq \frac{1}{(K-1)c_D} \max_{1 \leq j \leq p} \mathbb{E}_n[L(W_i)^2 X_{ij}^2] \leq \frac{C_L^2 C_X^2}{(K-1)c_D}.$$

with the same probability. Condition 4 of the lemma thus holds with $\gamma_n = n^{-1}$ and the now (K, c_D) -dependent $B_{2n} = C_L C_X / \sqrt{(K-1)c_D}$. Lemma D.1 then shows that for any $0 < u \leq c_L / [(1 + \bar{c}_0) C_X \sqrt{s}]$,

$$\begin{aligned} \mathbb{P} \left(\sqrt{|I_k^c|} \epsilon_{I_k^c}(u) > \left(\{4C_{L,e}\} \vee \left\{ C_\epsilon \sqrt{\frac{s \ln(pn)}{(K-1)c_D}} \right\} \right) u \right) \\ \leq 4n^{-1} + |I_k^c|^{-1} \leq 4n^{-1} + [(K-1)c_D n]^{-1}, \end{aligned}$$

where the second inequality follows from $|I_k| \geq c_D n$. Now $s \ln(pn) \geq 16(K-1)c_D C_{L,e}^2 / C_\epsilon^2$ is equivalent to $4C_{L,e} \leq C_\epsilon \{s \ln(pn) / [(K-1)c_D]\}^{1/2}$, and so with probability at least $1 - (4n^{-1} + [(K-1)c_D n]^{-1})$,

$$\epsilon_{I_k^c}(u) \leq (C_\epsilon / [(K-1)c_D]) u \sqrt{s \ln(pn) / n},$$

where we used the bound $|I_k^c| \geq (K-1)c_D n$. The asserted claim now follows from combining this inequality and the union bound. \square

PROOF OF LEMMA B.3. Fix a (hold-out) subsample $k \in \{1, \dots, K\}$. Assumptions 5 and 14 imply that for each $t \in \mathbf{R}$ and each $j \in \{1, \dots, p\}$, the random variables $S_{I_k^c, j} = |I_k^c|^{-1} \sum_{i \in I_k^c} m_1'(X_i^\top \theta_0, Y_i) X_{ij}$ satisfy

$$\ln \mathbb{E} \left[e^{t S_{I_k^c, j}} \mid \{X_i\}_{i=1}^n \right] \leq \frac{C_X^2 \sigma^2 t^2}{2|I_k^c|} \text{ a.s.}$$

Hence, by Chernoff's inequality, for any $t > 0$,

$$\mathbb{P}(|S_{I_k^c, j}| > t) \leq 2 \exp \left(-\frac{|I_k^c| t^2}{2C_X^2 \sigma^2} \right).$$

The union bound then implies that

$$\mathbb{P} \left(\|S_{I_k^c}\|_\infty > t \right) \leq 2p \exp \left(-\frac{|I_k^c| t^2}{2C_X^2 \sigma^2} \right), \quad t > 0,$$

from which we obtain

$$\mathbb{P} \left(\|S_{I_k^c}\|_\infty > C_X \sigma \sqrt{\frac{2 \ln(2p|I_k^c|)}{|I_k^c|}} \right) \leq |I_k^c|^{-1} \leq [(K-1)c_D n]^{-1},$$

where the second inequality follows from Assumption 12. In addition,

$$\frac{\ln(2p|I_k^c|)}{|I_k^c|} \leq \frac{2 \ln(p|I_k^c|)}{|I_k^c|} \leq \frac{2 \ln(pn)}{(K-1)c_D n},$$

where we again used Assumption 12 (recall also that we take $p \geq 2$). Combining these inequalities and applying the union bound, we obtain the asserted claim. \square

PROOF OF LEMMA B.4. Denote $\hat{\theta} := \hat{\theta}_{I_k^c}(\lambda)$ and $\hat{\delta} := \hat{\theta} - \theta_0$. By Theorem 1, we then have $\|\hat{\delta}\|_1 \leq (1 + \bar{c}_0)u_0\sqrt{s}$ and $\|\hat{\delta}\|_2 \leq u_0$. An argument parallel to Step 1 of the proof of Theorem 1 also shows that the assumed $\lambda \geq c_0\|S_{I_k^c}\|_\infty$ implies $\hat{\delta} \in \mathcal{R}(\bar{c}_0)$. Therefore,

$$\begin{aligned} \mathcal{E}(\hat{\theta}) &= \widehat{M}_{I_k^c}(\hat{\theta}) - \widehat{M}_{I_k^c}(\theta_0) - [\widehat{M}_{I_k^c}(\theta_0 + \hat{\delta}) - \widehat{M}_{I_k^c}(\theta_0) - M(\theta_0 + \hat{\delta}) + M(\theta_0)] \\ &\leq \lambda(\|\theta_0\|_1 - \|\hat{\theta}\|_1) + |\widehat{M}_{I_k^c}(\theta_0 + \hat{\delta}) - \widehat{M}_{I_k^c}(\theta_0) - M(\theta_0 + \hat{\delta}) + M(\theta_0)| \\ &\leq \bar{\lambda}\|\hat{\delta}\|_1 + \epsilon_{I_k^c}(u_0) \leq (\lambda_\epsilon + (1 + \bar{c}_0)\bar{\lambda}\sqrt{s})u_0 = \frac{2}{c_M}(\lambda_\epsilon + (1 + \bar{c}_0)\bar{\lambda}\sqrt{s})^2, \end{aligned}$$

where the second line follows from the definition of $\hat{\theta}$ and the third from $\hat{\delta} \in \mathcal{R}(\bar{c}_0)$, the definition of $\epsilon_{I_k^c}(u_0)$, imposed conditions, and the triangle inequality. This gives the asserted claim. \square

PROOF OF LEMMA B.5. By (B.4) and Lemma B.1,

$$\left[c_0 C_S \sqrt{\frac{\ln(pn)}{n}}, \frac{c_0 C_S}{a} \sqrt{\frac{\ln(pn)}{n}} \right] \cap \Lambda_n \neq \emptyset,$$

and so we can fix a penalty $\lambda_* \in \Lambda_n$ satisfying

$$c_0 C_S \sqrt{\frac{\ln(pn)}{n}} \leq \lambda_* \leq \frac{c_0 C_S}{a} \sqrt{\frac{\ln(pn)}{n}} =: \bar{\lambda}.$$

Further, denote

$$\lambda_\epsilon := \frac{C_\epsilon}{(K-1)c_D} \sqrt{\frac{s \ln(pn)}{n}}$$

and for all $k \in \{1, \dots, K\}$, consider events

$$\mathcal{Z}_k := \left\{ \|S_{I_k^c}\|_\infty \leq C_S \sqrt{n^{-1} \ln(pn)} \right\} \quad \text{and} \quad \mathcal{E}_k := \left\{ \epsilon_{I_k^c}(\tilde{u}_0) \leq \lambda_\epsilon \tilde{u}_0 \right\}.$$

Also, note that using λ_ϵ and $\bar{\lambda}$, \tilde{u}_0 can be written as

$$\tilde{u}_0 = \frac{2}{c_M} (\lambda_\epsilon + (1 + \bar{c}_0) \bar{\lambda} \sqrt{s}).$$

Lemma B.4 and (B.4) therefore imply that on $\mathcal{Z}_k \cap \mathcal{E}_k$,

$$\mathcal{E}(\widehat{\theta}_{I_k^c}(\lambda_*)) \leq \frac{2}{c_M} \left(\frac{C_\epsilon}{(K-1)c_D} + \frac{(1 + \bar{c}_0) c_0 C_S}{a} \right)^2 \frac{s \ln(pn)}{n}. \quad (\text{C.4})$$

In turn, Lemma B.2 and (B.4) show that

$$\mathbb{P}((\cap_{k=1}^K \mathcal{E}_k)^c) \leq K (4n^{-1} + [(K-1)c_D n]^{-1}).$$

Also, Lemma B.3 shows that

$$\mathbb{P}((\cap_{k=1}^K \mathcal{Z}_k)^c) \leq K [(K-1)c_D n]^{-1}.$$

It thus follows from the union bound that (C.4) holds simultaneously for all $k \in \{1, \dots, K\}$ with probability at least $1 - K (4n^{-1} + 2[(K-1)c_D n]^{-1})$. \square

PROOF OF LEMMA B.6. For any $\theta_1, \theta_2 \in \Theta$ and $k \in \{1, \dots, K\}$, let

$$f_k(\theta_1, \theta_2) := (\mathbb{E}_{I_k} - \mathbb{E})[m(X_i^\top \theta_1, Y_i) - m(X_i^\top \theta_2, Y_i)].$$

Also, let $\lambda_* \in \Lambda_n$ be a value of λ satisfying the bound of Lemma B.5 and consider events

$$\mathcal{R} := \left\{ \max_{1 \leq k \leq K} \mathcal{E}(\widehat{\theta}_{I_k^c}(\lambda_*)) \leq C_\mathcal{E}^2 \frac{s \ln(pn)}{n} \right\} \quad \text{and} \quad \mathcal{C}(t) := \bigcap_{k=1}^K \mathcal{C}_k(t),$$

where for each $k \in \{1, \dots, K\}$,

$$\begin{aligned} \mathcal{C}_k(t) &:= \left\{ |f_k(\widehat{\theta}_{I_k^c}(\widehat{\lambda}^{\text{cv}}), \widehat{\theta}_{I_k^c}(\lambda_*))| \right. \\ &\leq \left. \sqrt{\frac{2t \ln n}{c_D \ln(1/a) n}} \sqrt{\mathbb{E}_{X,Y} \left[\left\{ m(X^\top \widehat{\theta}_{I_k^c}(\widehat{\lambda}^{\text{cv}}), Y) - m(X^\top \widehat{\theta}_{I_k^c}(\lambda_*), Y) \right\}^2 \right]} \right\}. \end{aligned}$$

Now, fix a subsample $k \in \{1, \dots, K\}$ and observe that for any $\lambda \in \Lambda_n$, the variance of the conditional distribution of

$$\mathbb{E}_{I_k} \left[m(X_i^\top \widehat{\theta}_{I_k^c}(\lambda), Y_i) - m(X_i^\top \widehat{\theta}_{I_k^c}(\lambda_*), Y_i) \right]$$

given $\{(X_i, Y_i)\}_{i \in I_k^c}$ is bounded from above by

$$|I_k|^{-1} \mathbb{E}_{X,Y} \left[\left\{ m(X^\top \widehat{\theta}_{I_k^c}(\lambda), Y) - m(X^\top \widehat{\theta}_{I_k^c}(\lambda_*), Y) \right\}^2 \right].$$

In addition, by (B.5) and Assumption 13, we have $|\Lambda_n| \leq 2(\ln n) / \ln(1/a)$. Therefore, the union bound followed by Chebyshev's inequality applied conditional on $\{(X_i, Y_i)\}_{i \in I_k^c}$ gives

$$\begin{aligned} &\mathbb{P} \left(\exists \lambda \in \Lambda_n \text{ s.t. } |f_k(\widehat{\theta}_{I_k^c}(\lambda), \widehat{\theta}_{I_k^c}(\lambda_*))| \right. \\ &> \left. \sqrt{\frac{2t \ln n}{c_D n \ln(1/a)}} \sqrt{\mathbb{E}_{X,Y} \left[\left\{ m(X^\top \widehat{\theta}_{I_k^c}(\lambda), Y) - m(X^\top \widehat{\theta}_{I_k^c}(\lambda_*), Y) \right\}^2 \right]} \right) \\ &\leq \sum_{\lambda \in \Lambda_n} \frac{c_D n \ln(1/a)}{2t |I_k| \ln n} \leq \frac{1}{t}, \end{aligned}$$

where the second inequality follows from Assumption 12. Hence, by the union bound and Lemma B.5,

$$\mathbb{P}((\mathcal{R} \cap \mathcal{C}(t))^c) \leq K(4n^{-1} + 2[(K-1)c_D n]^{-1} + t^{-1}).$$

We will now prove that (B.6) holds on $\mathcal{R} \cap \mathcal{C}(t)$. For the rest of proof, we therefore remain on this event.

Given that

$$\widehat{\lambda}^{\text{cv}} \in \operatorname{argmin}_{\lambda \in \Lambda_n} \sum_{k=1}^K \sum_{i \in I_k} m(X_i^\top \widehat{\theta}_{I_k^c}(\lambda), Y_i),$$

a problem for which λ_* is feasible, we must have

$$\sum_{k=1}^K |I_k| \mathbb{E}_{I_k} [m(X_i^\top \widehat{\theta}_{I_k^c}(\lambda_*), Y_i)] \geq \sum_{k=1}^K |I_k| \mathbb{E}_{I_k} [m(X_i^\top \widehat{\theta}_{I_k^c}(\widehat{\lambda}^{\text{cv}}), Y_i)]$$

It therefore follows from the triangle inequality and $\mathcal{C}(t)$ that

$$\begin{aligned} & \sum_{k=1}^K \frac{|I_k|}{n} \left[\mathcal{E}(\widehat{\theta}_{I_k^c}(\widehat{\lambda}^{\text{cv}})) - \mathcal{E}(\widehat{\theta}_{I_k^c}(\lambda_*)) \right] \\ &= \sum_{k=1}^K \frac{|I_k|}{n} \mathbb{E}_{X,Y} \left[m(X^\top \widehat{\theta}_{I_k^c}(\widehat{\lambda}^{\text{cv}}), Y) - m(X^\top \widehat{\theta}_{I_k^c}(\lambda_*), Y) \right] \\ &\leq \sum_{k=1}^K \frac{|I_k|}{n} |f_k(\widehat{\theta}_{I_k^c}(\widehat{\lambda}^{\text{cv}}), \widehat{\theta}_{I_k^c}(\lambda_*))| \\ &\leq \sum_{k=1}^K \frac{|I_k|}{n} \sqrt{\frac{2t \ln n}{c_D \ln(1/a) n}} \sqrt{\mathbb{E}_{X,Y} \left[\left\{ m(X^\top \widehat{\theta}_{I_k^c}(\widehat{\lambda}^{\text{cv}}), Y) - m(X^\top \widehat{\theta}_{I_k^c}(\lambda_*), Y) \right\}^2 \right]}. \end{aligned} \quad (\text{C.5})$$

In addition, on \mathcal{R} , we have

$$\max_{1 \leq k \leq K} \mathcal{E}(\widehat{\theta}_{I_k^c}(\lambda_*)) \leq \frac{C_{\mathcal{E}}^2 s \ln(pn)}{n} \leq 1, \quad (\text{C.6})$$

where the second inequality follows from (B.5). Assumption 15 therefore yields

$$\mathbb{E}_{X,Y} \left[\left\{ m(X^\top \widehat{\theta}_{I_k^c}(\widehat{\lambda}^{\text{cv}}), Y) - m(X^\top \theta_0, Y) \right\}^2 \right] \leq C_{ms}^2 \left[\mathcal{E}(\widehat{\theta}_{I_k^c}(\widehat{\lambda}^{\text{cv}})) \vee \mathcal{E}(\widehat{\theta}_{I_k^c}(\widehat{\lambda}^{\text{cv}}))^2 \right],$$

and

$$\mathbb{E}_{X,Y} \left[\left\{ m(X^\top \widehat{\theta}_{I_k^c}(\lambda_*), Y) - m(X^\top \theta_0, Y) \right\}^2 \right] \leq C_{ms}^2 \mathcal{E}(\widehat{\theta}_{I_k^c}(\lambda_*))$$

for all $k \in \{1, \dots, K\}$. Thus, using the well-known inequality $(a+b)^2 \leq 2a^2 + 2b^2$, we get

$$\begin{aligned} & \mathbb{E}_{X,Y} \left[\left\{ m(X^\top \widehat{\theta}_{I_k^c}(\widehat{\lambda}^{\text{cv}}), Y) - m(X^\top \widehat{\theta}_{I_k^c}(\lambda_*), Y) \right\}^2 \right] \\ &\leq 2\mathbb{E}_{X,Y} \left[\left\{ m(X^\top \widehat{\theta}_{I_k^c}(\widehat{\lambda}^{\text{cv}}), Y) - m(X^\top \theta_0, Y) \right\}^2 \right] \\ &\quad + 2\mathbb{E}_{X,Y} \left[\left\{ m(X^\top \widehat{\theta}_{I_k^c}(\lambda_*), Y) - m(X^\top \theta_0, Y) \right\}^2 \right] \\ &\leq 2C_{ms}^2 \left[\mathcal{E}(\widehat{\theta}_{I_k^c}(\widehat{\lambda}^{\text{cv}})) + \mathcal{E}(\widehat{\theta}_{I_k^c}(\widehat{\lambda}^{\text{cv}}))^2 + \mathcal{E}(\widehat{\theta}_{I_k^c}(\lambda_*)) \right]. \end{aligned}$$

Substituting this bound into (C.5), we obtain

$$\begin{aligned}
& \sqrt{\frac{c_D \ln(1/a)n}{2t \ln n}} \sum_{k=1}^K \frac{|I_k|}{n} \left[\mathcal{E}(\widehat{\theta}_{I_k^c}(\widehat{\lambda}^{\text{cv}})) - \mathcal{E}(\widehat{\theta}_{I_k^c}(\lambda_*)) \right] \\
& \leq \sum_{k=1}^K \frac{|I_k|}{n} \sqrt{2C_{ms}^2 [\mathcal{E}(\widehat{\theta}_{I_k^c}(\widehat{\lambda}^{\text{cv}})) + \mathcal{E}(\widehat{\theta}_{I_k^c}(\widehat{\lambda}^{\text{cv}}))^2 + \mathcal{E}(\widehat{\theta}_{I_k^c}(\lambda_*))]} \\
& \leq \sqrt{2}C_{ms} \sum_{k=1}^K \frac{|I_k|}{n} \left(\sqrt{\mathcal{E}(\widehat{\theta}_{I_k^c}(\widehat{\lambda}^{\text{cv}}))} + \mathcal{E}(\widehat{\theta}_{I_k^c}(\widehat{\lambda}^{\text{cv}})) + \sqrt{\mathcal{E}(\widehat{\theta}_{I_k^c}(\lambda_*))} \right) \\
& \leq \sqrt{2}C_{ms} \left(\sqrt{\sum_{k=1}^K \frac{|I_k|}{n} \mathcal{E}(\widehat{\theta}_{I_k^c}(\widehat{\lambda}^{\text{cv}}))} + \sum_{k=1}^K \frac{|I_k|}{n} \mathcal{E}(\widehat{\theta}_{I_k^c}(\widehat{\lambda}^{\text{cv}})) + \sqrt{\sum_{k=1}^K \frac{|I_k|}{n} \mathcal{E}(\widehat{\theta}_{I_k^c}(\lambda_*))} \right)
\end{aligned}$$

or rearranging and using the last inequality in (B.5),

$$\begin{aligned}
\sum_{k=1}^K \frac{|I_k|}{n} \mathcal{E}(\widehat{\theta}_{I_k^c}(\widehat{\lambda}^{\text{cv}})) & \leq 4C_{ms} \sqrt{\frac{t \ln n}{c_D \ln(1/a)n}} \left(\sqrt{\sum_{k=1}^K \frac{|I_k|}{n} \mathcal{E}(\widehat{\theta}_{I_k^c}(\widehat{\lambda}^{\text{cv}}))} + \sqrt{\sum_{k=1}^K \frac{|I_k|}{n} \mathcal{E}(\widehat{\theta}_{I_k^c}(\lambda_*))} \right) \\
& \quad + 2 \sum_{k=1}^K \frac{|I_k|}{n} \mathcal{E}(\widehat{\theta}_{I_k^c}(\lambda_*)).
\end{aligned}$$

Thus, given that the inequality $x \leq 2a(\sqrt{x} + \sqrt{y}) + 2y$ for $x, y \geq 0$ implies that $\sqrt{x} \leq a + [(a + \sqrt{y})^2 + y]^{1/2} \leq 2a + 2\sqrt{y}$, so that $x \leq 8a^2 + 8y$, it follows that

$$\sum_{k=1}^K \frac{|I_k|}{n} \mathcal{E}(\widehat{\theta}_{I_k^c}(\widehat{\lambda}^{\text{cv}})) \leq \frac{32C_{ms}^2}{c_D \ln(1/a)} \frac{t \ln n}{n} + 8 \sum_{k=1}^K \frac{|I_k|}{n} \mathcal{E}(\widehat{\theta}_{I_k^c}(\lambda_*)).$$

Combining this bound with Assumption 12 and using (C.6), we obtain

$$\max_{1 \leq k \leq K} \mathcal{E}(\widehat{\theta}_{I_k^c}(\widehat{\lambda}^{\text{cv}})) \leq \frac{32C_{ms}^2}{c_D^2 \ln(1/a)} \frac{t \ln n}{n} + \frac{8C_{\xi}^2}{c_D} \frac{s \ln(pn)}{n},$$

which completes the proof of the lemma. \square

D Fundamental Tools

D.1 Maximal Inequality

Let $\mathbb{G}_n[f(W_i)] := \sqrt{n} \{\mathbb{E}_n[f(W_i)] - \mathbb{E}[f(W)]\}$ abbreviate the centered and scaled empirical average.

Lemma D.1 (Maximal Inequality Based on Contraction Principle). *Let $\{W_i\}_{i=1}^n$ be independent copies of a random vector W , with support \mathcal{W} , of which X is a p -dimensional subvector, let Δ be a nonempty subset of \mathbf{R}^p , and let $h : \mathbf{R} \times \mathcal{W} \rightarrow \mathbf{R}$ be a measurable map satisfying $h(0, \cdot) \equiv 0$. Suppose that there exists constants $C_h, B_{1n}, B_{2n} \in \mathbf{R}_+, \gamma_n \in (0, 1)$ and a measurable function $L : \mathcal{W} \rightarrow \mathbf{R}_+$ such that*

1. *for all $w \in \mathcal{W}$ and all $t_1, t_2 \in \mathbf{R}$ satisfying $|t_1| \vee |t_2| \leq C_h$,*

$$|h(t_1, w) - h(t_2, w)| \leq L(w) |t_1 - t_2|;$$

2. *$\sup_{\delta \in \Delta} |X^\top \delta| \leq C_h$ a.s.;*

3. *$\sup_{\delta \in \Delta} \mathbb{E}[h(X^\top \delta, W)^2] \leq B_{1n}^2$; and,*

4. *$\max_{1 \leq j \leq p} \mathbb{E}_n[L(W_i)^2 X_{ij}^2] \leq B_{2n}^2$ with probability at least $1 - \gamma_n$.*

Then, denoting $\|\Delta\|_1 := \sup_{\delta \in \Delta} \|\delta\|_1$, we have

$$\mathbb{P} \left(\sup_{\delta \in \Delta} |\mathbb{G}_n[h(X_i^\top \delta, W_i)]| > u \right) \leq 4\gamma_n + n^{-1},$$

provided $u \geq \{4B_{1n}\} \vee \{8\sqrt{2}B_{2n} \|\Delta\|_1 \sqrt{\ln(8pn)}\}$.

Proof. The claim follows from [Belloni et al. \(2018a, Lemma D.3\)](#), which, in turn, follows from a variant of an argument given in [Ledoux and Talagrand \(1991\)](#). \square

D.2 Gaussian Inequality

Lemma D.2 (Gaussian Quantile Bound). *Let (Y_1, \dots, Y_p) be centered Gaussian in \mathbf{R}^p with $\sigma^2 := \max_{1 \leq j \leq p} \mathbb{E}[Y_j^2]$ and $p \geq 2$. Let $q^Y(1 - \alpha)$ denote the $(1 - \alpha)$ -quantile of $\max_{1 \leq j \leq p} |Y_j|$ for $\alpha \in (0, 1)$. Then $q^Y(1 - \alpha) \leq (2 + \sqrt{2})\sigma \sqrt{\ln(p/\alpha)}$.*

Proof. By the Borell-TIS (Tsirelson-Ibragimov-Sudakov) inequality ([Adler and Taylor, 2007, Theorem 2.1.1](#)), for any $t > 0$ we have

$$\mathbb{P} \left(\max_{1 \leq j \leq p} |Y_j| > \mathbb{E} \left[\max_{1 \leq j \leq p} |Y_j| \right] + \sigma t \right) \leq e^{-t^2/2}.$$

This inequality translates to the quantile bound

$$q^Y(1 - \alpha) \leq \mathbb{E} \left[\max_{1 \leq j \leq p} |Y_j| \right] + \sigma \sqrt{2 \ln(1/\alpha)}.$$

Talagrand (2010, Proposition A.3.1) shows that

$$\mathbb{E}\left[\max_{1 \leq j \leq p} |Y_j|\right] \leq \sigma \sqrt{2 \ln(2p)},$$

thus implying

$$q^Y (1 - \alpha) \leq \sigma \left(\sqrt{2 \ln(2p)} + \sqrt{2 \ln(1/\alpha)} \right).$$

The claim now follows from $p \geq 2$. □

D.3 CLT and Bootstrap in High Dimensions

Throughout this section we let Z_1, \dots, Z_n be independent centered \mathbf{R}^p -valued random variables and denote their scaled average and variance by

$$S_n := \frac{1}{\sqrt{n}} \sum_{i=1}^n Z_i \quad \text{and} \quad \Sigma := \frac{1}{n} \sum_{i=1}^n \mathbb{E}[Z_i Z_i^\top],$$

respectively. (The existence of Σ is guaranteed by our assumptions below.) For \mathbf{R}^p -valued random variables U and V , define the distributional measure of distance

$$\rho(U, V) := \sup_{A \in \mathcal{A}_p} |\mathbb{P}(U \in A) - \mathbb{P}(V \in A)|,$$

where \mathcal{A}_p denotes the collection of hyperrectangles in \mathbf{R}^p . Also, for $M \in \mathbf{R}^{p \times p}$ symmetric positive definite, write $N_M := \mathcal{N}(\mathbf{0}, M)$.

Theorem D.1 (High-Dimensional CLT). *If for some finite constants $b > 0$ and $B_n \geq 1$,*

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E}[Z_{ij}^2] \geq b, \quad \frac{1}{n} \sum_{i=1}^n \mathbb{E}[|Z_{ij}|^{2+k}] \leq B_n^k \quad \text{and} \quad \mathbb{E}\left[\max_{1 \leq j \leq p} Z_{ij}^4\right] \leq B_n^4, \quad (\text{D.1})$$

for all $i \in \{1, \dots, n\}$, $j \in \{1, \dots, p\}$ and $k \in \{1, 2\}$, then there exists a finite constant C_b , depending only on b , such that

$$\rho(S_n, N_\Sigma) \leq C_b \left(\frac{B_n^4 \ln^7(pn)}{n} \right)^{1/6}.$$

Proof. The claim follows from Chernozhukov et al. (2017, Proposition 2.1). □

Let \widehat{Z}_i be an estimator of Z_i , and let e_1, \dots, e_n be i.i.d. $\mathcal{N}(0, 1)$ and independent of both the Z_i 's and the \widehat{Z}_i 's. Define $\widehat{S}_n^e := n^{-1/2} \sum_{i=1}^n e_i \widehat{Z}_i$ and let \mathbb{P}_e denote the (conditional)

probability measure computed with respect to the e_i 's for fixed Z_i 's and \widehat{Z}_i 's. Also, abbreviate

$$\widetilde{\rho}(\widehat{S}_n^e, N_\Sigma) := \sup_{A \in \mathcal{A}_p} \left| \mathbb{P}_e(\widehat{S}_n^e \in A) - \mathbb{P}(N_\Sigma \in A) \right|,$$

with the tilde stressing that $\widetilde{\rho}(\widehat{S}_n^e, N_\Sigma)$ is a random quantity.

Theorem D.2 (Multiplier Bootstrap for Many Approximate Means). *Let (D.1) hold for some finite constants $b > 0$ and $B_n \geq 1$, and let $\{\beta_n\}_1^\infty$ and $\{\delta_n\}_1^\infty$ be sequences in \mathbf{R}_{++} both converging to zero such that*

$$\mathbb{P}\left(\max_{1 \leq j \leq p} \frac{1}{n} \sum_{i=1}^n (\widehat{Z}_{ij} - Z_{ij})^2 > \frac{\delta_n^2}{\ln^2(pn)}\right) \leq \beta_n. \quad (\text{D.2})$$

Then there exists a finite constant C_b , depending only on b , such that with probability at least $1 - \beta_n - 1/\ln^2(pn)$,

$$\widetilde{\rho}(\widehat{S}_n^e, N_\Sigma) \leq C_b \max \left\{ \delta_n, \left(\frac{B_n^4 \ln^6(pn)}{n} \right)^{1/6} \right\}.$$

Proof. The claim follows from [Belloni et al. \(2018a\)](#), Theorem 2.2), which is here restated in order to highlight the dependence on the sequences β_n and δ_n . [Note that their Theorem 2.2 does not actually require their Condition A(i).] \square

For any M symmetric positive definite, define $q_M^N : \mathbf{R} \rightarrow \mathbf{R} \cup \{\pm\infty\}$ as the (extended) quantile function of $\|N_M\|_\infty$,

$$q_M^N(\alpha) := \inf \{t \in \mathbf{R}; \mathbb{P}(\|N_M\|_\infty \leq t) \geq \alpha\}, \quad \alpha \in \mathbf{R}.$$

Here we interpret $q_M^N(\alpha)$ as $+\infty (= \inf \emptyset)$ if $\alpha \geq 1$, and $-\infty (= \inf \mathbf{R})$ if $\alpha \leq 0$, such that q_M^N is monotone increasing.

Lemma D.3. *Let $M \in \mathbf{R}^{p \times p}$ be symmetric positive definite, let U be an \mathbf{R}^p -valued random variable, and let q denote the quantile function of $\|U\|_\infty$. Then*

$$q_M^N(\alpha - 2\rho(U, N_M)) \leq q(\alpha) \leq q_M^N(\alpha + \rho(U, N_M)) \text{ for all } \alpha \in (0, 1).$$

Proof. Given positive definiteness of M , for any $t \in \mathbf{R}$,

$$\mathbb{P}(\|N_M\|_\infty = t) \leq \sum_{j=1}^p \mathbb{P}(|N(0, M_{jj})| = t) = 0.$$

It follows that for each $\alpha \in (0, 1)$, $q_M^N(\alpha)$ is uniquely defined by

$$\mathbb{P}(\|N_M\|_\infty \leq q_M^N(\alpha)) = \alpha.$$

In establishing the *lower* bound we may take $\rho(U, N_M) < \alpha$. (Otherwise $q_M^N(\alpha - \rho(U, N_M)) = -\infty$ and there is nothing to show.) Then $[-q_M^N(\alpha - \rho(U, N_M)), q_M^N(\alpha - \rho(U, N_M))]^p$ is a rectangle and

$$\mathbb{P}(\|U\|_\infty \leq q_M^N(\alpha - 2\rho(U, N_M))) \leq \mathbb{P}(\|N_M\|_\infty \leq q_M^N(\alpha - 2\rho(U, N_M))) + \rho(U, N_M) < \alpha,$$

which implies the lower bound. In establishing the *upper* bound we may assume $\rho(U, N_M) < 1 - \alpha$. (Otherwise $q_M^N(\alpha + \rho(U, N_M)) = +\infty$ and there is nothing to show.) Then from the rectangle $[-q_M^N(\alpha + \rho(U, N_M)), q_M^N(\alpha + \rho(U, N_M))]^p$, a parallel calculation shows

$$\mathbb{P}(\|U\|_\infty \leq q_M^N(\alpha + \rho(U, N_M))) \geq \alpha,$$

which by definition of quantiles implies the upper bound. \square

Now, define $q_n(\alpha)$ as the α -quantile of $\|S_n\|_\infty$

$$q_n(\alpha) := \inf \{t \in \mathbf{R}; \mathbb{P}(\|S_n\|_\infty \leq t) \geq \alpha\}, \quad \alpha \in (0, 1),$$

and let $\hat{q}_n(\alpha)$ be the α -quantile of $\|\hat{S}_n^e\|_\infty$ computed conditional on X_i 's and \hat{X}_i 's,

$$\hat{q}_n(\alpha) := \inf \left\{ t \in \mathbf{R}; \mathbb{P}_e(\|\hat{S}_n^e\|_\infty \leq t) \geq \alpha \right\}, \quad \alpha \in (0, 1).$$

Theorem D.3 (Quantile Comparison). *If (D.1) holds for some finite constants $b > 0$ and $B_n \geq 1$, and*

$$\rho_n := 2C_b \left(\frac{B_n^4 \ln^7(pn)}{n} \right)^{1/6}$$

denotes the upper bound in Theorem D.1 multiplied by two, then

$$q_\Sigma^N(1 - \alpha - \rho_n) \leq q_n(1 - \alpha) \leq q_\Sigma^N(1 - \alpha + \rho_n) \text{ for all } \alpha \in (0, 1).$$

(ii) If, in addition, (D.2) holds for some sequences $\{\delta_n\}_1^\infty$ and $\{\beta_n\}_1^\infty$ in \mathbf{R}_{++} both converging to zero, and

$$\rho'_n := 2C'_b \max \left\{ \delta_n, \left(\frac{B_n^4 \ln^6(pn)}{n} \right)^{1/6} \right\}$$

denotes the upper bound in Theorem D.2 multiplied by two, then with probability at least

$$1 - \beta_n - 1/\ln^2(pn),$$

$$q_\Sigma^N (1 - \alpha - \rho'_n) \leq \widehat{q}_n (1 - \alpha) \leq q_\Sigma^N (1 - \alpha + \rho'_n) \text{ for all } \alpha \in (0, 1).$$

Proof. Apply Lemma D.3 with $U = S_n$ to obtain

$$q_\Sigma^N (1 - \alpha - 2\rho(S_n, N_\Sigma)) \leq q_n (1 - \alpha) \leq q_\Sigma^N (1 - \alpha + \rho(S_n, N_\Sigma)) \text{ for all } \alpha \in (0, 1).$$

The first pair of inequalities then follows from $2\rho(S_n, N_\Sigma) \leq \rho_n$ (Theorem D.1). To establish the second claim, apply Lemma D.3 with $U = \widehat{S}_n^e$ and conditional on the X_i 's and \widehat{X}_i 's to obtain

$$q_\Sigma^N (1 - \alpha - 2\widetilde{\rho}(\widehat{S}_n^e, N_\Sigma)) \leq \widehat{q}_n (1 - \alpha) \leq q_\Sigma^N (1 - \alpha + \widetilde{\rho}(\widehat{S}_n^e, N_\Sigma)) \text{ for all } \alpha \in (0, 1).$$

The second pair of inequalities then follows on the event $2\widetilde{\rho}(\widehat{S}_n^e, N_\Sigma) \leq \rho'_n$, which by Theorem D.2 occurs with probability at least $1 - \beta_n - 1/\ln^2(pn)$. \square

Theorem D.4 (Multiplier Bootstrap Consistency). *Let (D.1) and (D.2) hold for some finite constants $b > 0$ and $B_n \geq 1$ and some sequences $\{\delta_n\}_1^\infty$ and $\{\beta_n\}_1^\infty$ in \mathbf{R}_{++} both converging to zero. Then there exists a finite constant C_b , depending only on b , such that*

$$\sup_{\alpha \in (0,1)} |\mathbb{P}(\|S_n\|_\infty > \widehat{q}_n (1 - \alpha)) - \alpha| \leq C_b \max \left\{ \beta_n, \delta_n, \left(\frac{B_n^4 \ln^7(pn)}{n} \right)^{1/6}, \frac{1}{\ln^2(pn)} \right\}.$$

Thus, if in addition $B_n^4 \ln^7(pn)/n \rightarrow 0$, then

$$\sup_{\alpha \in (0,1)} |\mathbb{P}(\|S_n\|_\infty > \widehat{q}_n (1 - \alpha)) - \alpha| \rightarrow 0.$$

Proof. By Theorems D.1 and D.3,

$$\begin{aligned} \mathbb{P}(\|S_n\|_\infty \leq \widehat{q}_n (1 - \alpha)) &\leq \mathbb{P}(\|S_n\|_\infty \leq q_\Sigma^N (1 - \alpha + \rho'_n)) + \beta_n + \frac{1}{\ln^2(pn)} \\ &\leq \mathbb{P}(\|N_\Sigma\|_\infty \leq q_\Sigma^N (1 - \alpha + \rho'_n)) + \rho_n + \beta_n + \frac{1}{\ln^2(pn)} \\ &\leq 1 - \alpha + \rho'_n + \rho_n + \beta_n + \frac{1}{\ln^2(pn)}. \end{aligned}$$

A parallel argument shows

$$\mathbb{P}(\|S_n\|_\infty \leq \widehat{q}_n (1 - \alpha)) \geq 1 - \alpha - \left(\rho'_n + \rho_n + \beta_n + \frac{1}{\ln^2(pn)} \right).$$

The claim now follows from combining and rearranging the previous two displays. □