

# Linear programming approach to partially identified econometric models

Andrei Voronin\*

February 22, 2024

PRELIMINARY AND INCOMPLETE

[Link to latest version](#)

## Abstract

Sharp bounds on partially identified parameters are often given by the values of linear programs (LPs). This paper introduces a novel estimator of the LP value  $B(\theta) = \min_{Mx \geq c} p'x$ , where  $\theta = (p, M, c)$  is estimated. Unlike existing procedures, our estimator is  $\sqrt{n}$ -consistent whenever the true LP is feasible and finite, and remains valid under point-identification, over-identifying constraints, and solution multiplicity. Turning to uniform properties, we prove that the LP value cannot be uniformly consistently estimated over the unrestricted set of measures. We then show that our estimator achieves uniform consistency under a condition that appears minimal for the existence of any such estimator. We also obtain computationally efficient, asymptotically normal inference with exact asymptotic coverage. To complement our estimation results, we derive LP sharp bounds in a general identification setting. Our approach allows applied work to employ previously intractable conditions. We apply our findings to estimating returns to education. To that end, we propose the conditionally monotone IV assumption (cMIV) that tightens the classical monotone IV (MIV) bounds. We argue that cMIV remains unrestrictive relative to MIV and provide a formal test for it. Under cMIV, university education in Colombia is shown to increase the average wage by at least 5.91%, whereas classical conditions fail to produce an informative bound.

---

\*Department of Economics, UCLA. Email: [avoronin@g.ucla.edu](mailto:avoronin@g.ucla.edu). I am grateful to Denis Chetverikov, Andres Santos, Rosa Matzkin, Jinyong Hahn, Bulat Gafarov, Tim Armstrong, Kirill Ponomarev, Shuyang Sheng and Manu Navjeevan as well as to all the participants of the 2024 California Econometrics Conference and the 2024 European Winter Meeting of the Econometric Society for the valuable discussions and criticisms.

# 1. Introduction

In many partial identification frameworks<sup>1</sup>, the sharp bounds on parameters correspond to the values of linear programs (LPs) that depend on identified functionals of the underlying probability measure. Specifically, the bounds take the form  $B(\mathbb{P}) = B(\theta_0(\mathbb{P}))$ , where

$$B(\theta) \equiv \min_{Mx \geq c} p'x, \tag{1}$$

and  $\theta_0(\mathbb{P})$  is the true value of parameter  $\theta = (p', \text{vec}(M)', c)'$ , estimated via a  $\sqrt{n}$ -consistent estimator  $\hat{\theta}_n$ . However, optimization problem (1) exhibits non-regular behavior, particularly when the underlying model is rich enough that some linear functionals of  $x$  are nearly or exactly point-identified over  $\Theta_I = \{x \in \mathbb{R}^d : Mx \geq c\}$ . In such cases, existing estimators of  $B(\mathbb{P})$  are either inconsistent or rate-conservative, creating an undesirable tradeoff in empirical work: richer models provide tighter bounds but complicate their estimation.

To address this issue, we develop a novel estimator of  $B(\mathbb{P})$ . Our estimator takes the form  $\hat{B}(\hat{\theta}_n; w_n)$ , where

$$\hat{B}(\theta; w) \equiv \sup_{x \in \tilde{\mathcal{A}}(\theta; w)} p'x, \quad \tilde{\mathcal{A}}(\theta; w) \equiv \arg \min_{x \in \mathcal{X}} p'x + wl'(c - Mx)^+, \tag{2}$$

and  $w_n \rightarrow \infty$  is the penalty parameter, with  $\frac{w_n}{\sqrt{n}} \rightarrow 0$ . We refer to  $\hat{B}(\hat{\theta}_n; w_n)$  as *the debiased penalty function estimator*. Only assuming that  $\Theta_I$  is non-empty and contained in a known compact  $\mathcal{X}$ , we show that  $\hat{B}(\hat{\theta}_n; w_n)$  is  $\sqrt{n}$ -consistent for any  $w_n$  satisfying the above conditions<sup>2</sup>. In contrast, the plug-in estimator is not generally consistent and may fail to exist with non-vanishing probability, while the alternative set-expansion estimator based on Chernozhukov et al. (2007) is rate-conservative and may fail to exist in finite samples. Figure 1 gives a preview of the comparative performance of these estimators.

We obtain an asymptotically normal version of our estimator via sample-splitting and construct confidence regions with exact asymptotic coverage. By comparison, existing procedures either rely on further conditions (Gafarov, 2024), or lead to asymptotically conservative inference (Cho and Russell, 2023). Notably, the approach most commonly used in applied research—combining plug-in estimation with bootstrap (De Haan (2017), Cygan-Rehm et al. (2017), Siddique (2013), Kreider et al. (2012), Gundersen et al. (2012), Blundell et al. (2007))—may not provide valid confidence intervals even when the underlying model is far from point-identification.

---

<sup>1</sup>Including conditional moment inequalities (Andrews et al., 2023), generalized IV models (Mogstad et al., 2018), revealed preference restrictions (Kline and Tartari, 2016), intersection bounds (Honoré and Lleras-Muney, 2006), dynamic discrete choice panels (Honoré and Tamer, 2006) and shape restrictions Manski and Pepper (2000).

<sup>2</sup>The selection of this parameter and its relation to uniform properties of the estimator is discussed in the Appendix in great detail.

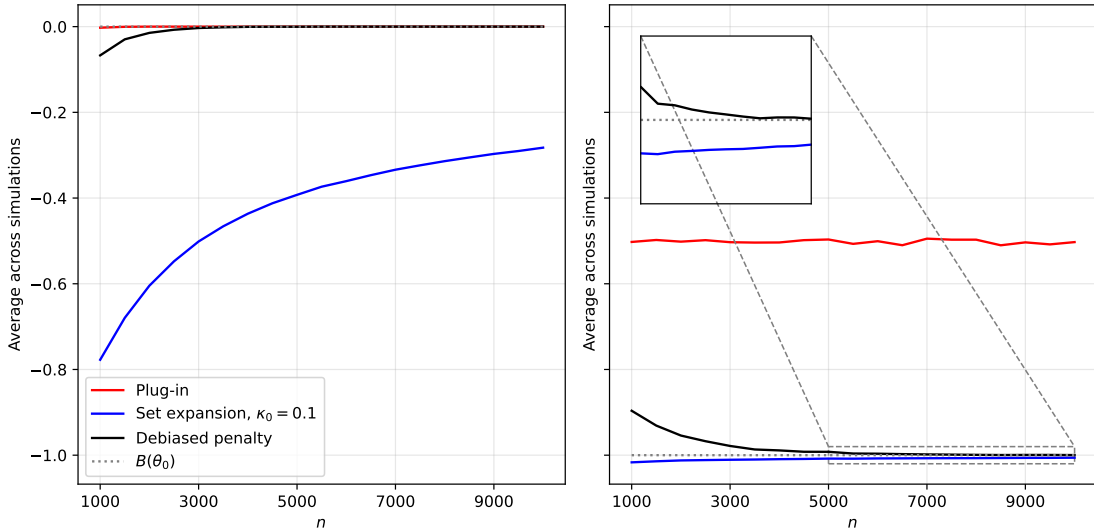


Figure 1: Comparison of estimators for two measures with the true values of 0 and  $-1$ , left to right. Average estimate across 10000 simulations. See Section 2.6 for details.

Using an equivalence between unconstrained piecewise-linear problems and auxiliary LPs, one can compute our estimator in polynomial time with a LP solver. This, combined with the closed-form expression for the asymptotic variance, makes our inference procedure the most computationally efficient in the existing literature<sup>3</sup>.

Turning to uniform asymptotic theory, we first establish a general impossibility result: using Le Cam’s binary testing method, we show that no uniformly consistent estimator exists when the estimated functional is discontinuous in the total variation norm. Applied to LP, this implies that  $B(\mathbb{P})$  cannot be uniformly consistently estimated over the unrestricted set of measures  $\mathcal{P}$ . To make progress, we introduce the ‘ $\delta$ -condition’ that parametrizes  $\mathcal{P}$  by restricting it to the measures at which the smallest singular value of some full-rank submatrix of constraints binding at an optimal vertex is lower-bounded by a  $\delta > 0$ . This condition is minimal in the sense that any measure from  $\mathcal{P}$  satisfies it for some  $\delta$ , ensuring the family of restricted measures’ sets covers  $\mathcal{P}$  as  $\delta$  grows small. Unlike the conditions in Gafarov (2024), it does not exclude economically relevant *problematic* cases, such as point-identification and over-identification, nor does it preclude solution multiplicity. Under the  $\delta$ -condition, our estimator is shown to be uniformly consistent.

To complement our estimation procedure, we develop a general identification framework that results in LP bounds. In particular, we derive novel sharp bounds for a broad class of

<sup>3</sup>Both Gafarov (2024) (BG) and Cho and Russell (2023) (CR) rely on resampling methods, which require to compute one or multiple LPs at each iteration. Computing a confidence interval for a LP with 32 variables takes 16.81 seconds with the approach of BG and 40.65 seconds with the approach of CR, according to the latter work. Our approach requires computing a LP once, which takes around 0.0022 seconds on average.

treatment parameters<sup>4</sup> under arbitrary affine inequalities over conditional moments (AICM), potentially augmented with affine almost sure restrictions and missing data conditions. In the simplest case, AICM identifying restrictions have the form

$$M^*(\mathbb{E}[Y(d)|T = t, Z = z])_{d,t,z} + b^* \geq 0 \quad \text{and} \quad \tilde{M}(Y(d))_d + \tilde{b} \geq 0 \text{ a.s.}, \quad (3)$$

where  $(Y(d))_d$  are continuous potential outcomes corresponding to the legs of treatment  $T$  and  $Z$  are other covariates. Identified matrices  $M^*$ ,  $\tilde{M}$  and vectors  $b^*$ ,  $\tilde{b}$  are chosen by the researcher. In AICM models,  $\theta$  from (1) is usually a function of identified conditional moments  $(\mathbb{E}[Y|T = t, Z = z])_{t,z}$  and the identified joint distribution of  $T, Z$ , while  $x$  collects relevant unobserved conditional moments. Our approach accommodates arbitrary combinations of existing ‘nonparametric bounds’ restrictions, allows to conduct sensitivity analysis, and extends to more complex conditions where sharp bounds were previously unavailable<sup>5</sup>.

Finally, we develop an application of our approach to estimating returns to education in Colombia. To that end, in Section 4 we first introduce a family of conditionally monotone instrumental variables assumptions (cMIV), nested in (3), that impose

$$\mathbb{E}[Y(t)|T \in A, Z = z] - \text{monotone in } z,$$

where  $(Y(t))_{t \in \mathcal{T}}$  are potential log-wage schedules corresponding to education levels  $T$ , and  $Z$  is a proxy for ability based on Saber test scores. Sets  $A$  parametrize different versions of cMIV and are chosen by the researcher<sup>6</sup>. While an explicit form for the sharp bounds under some versions of cMIV may not be feasible, their LP representation follows from our general identification result for (3). The cMIV conditions we consider yield tighter bounds than the classical monotone instrumental variables (MIV) assumption of Manski and Pepper (2000). We argue, however, that they remain unrestrictive in many applications, including ours. While empirical literature (e.g., De Haan (2017)) has visually examined the monotonicity of observed conditional moments, i.e.  $A = \{t\}$ , to justify applying MIV, we show that such monotonicity is instead equivalent to a particular form of cMIV given that MIV holds and under a mild regularity condition. The formal test of cMIV is obtained as an extension of Chetverikov (2019). Using Saber test scores as a cMIV, we find that obtaining a university degree increases average wages by at least 5.91% in Colombia. In contrast, the classical conditions fail to produce an informative bound.

This paper also contributes two auxiliary results. The first one is concerned with an important special case of (3) - the combination of all classical Manski and Pepper (2000) conditions. Since such combination possesses the greatest identifying power out of all classical

<sup>4</sup>Including ATE and CATE, among other typically studied parameters, see Section (3).

<sup>5</sup>For example, cMIV and the mixture of all classical Manski and Pepper (2000) conditions, see below.

<sup>6</sup>We explore various choices, such as letting  $A$  run through all singletons  $\{d\}$  and the full support  $\mathcal{T}$ . See Section 4 for details.

restrictions, empirical work has used it even in the absence of a theoretical justification, obtaining bounds that were either not sharp, or invalid (see [Lafférs \(2013\)](#)). Our method yields sharp bounds and a valid estimation procedure for that setting even under continuous outcomes. Another auxiliary contribution consists in a novel lower bound on the  $\ell_1$ -deviation from a non-empty bounded polytope in terms of Euclidean distance from the polytope. It may provide insights into the behavior of  $\ell_1$ -penalized solutions of systems of linear inequalities studied in the control theory literature (e.g. [Pinar and Chen \(1999\)](#)).

We briefly note the limitations of our approach. On the identification side, the absence of restrictions on treatment selection prevents us from studying more granular parameters, such as marginal treatment responses. Furthermore, our identification results are given for discrete treatment and instrument. An extension to the continuous case is feasible, but is outside the scope of this paper<sup>7</sup>. On the estimation side, while our estimator is pointwise  $\sqrt{n}$ -consistent in general, we only establish  $\sqrt{n}/w_n$ -uniform consistency for a slowly diverging sequence  $w_n$ <sup>8</sup>. We provide further evidence on the uniform rate of consistency in the Appendix. A theoretically  $\sqrt{n}$ -uniformly consistent estimator follows from our analysis, but it depends on an unobserved parameter  $\delta$  that is difficult to estimate, so we do not recommend using it in practice. Finally, while our inference procedure naturally extends to uniform setup under sufficient regularity conditions, exploring this is left for future work.

## Relationship to literature

The strand of literature relevant to the estimation of (1) is concerned with statistical inference in the LP estimation framework. [Semenova \(2023\)](#) considers a LP with an estimated constraint vector  $\hat{c}_n$ , but a known matrix  $M$  and a coefficient vector  $p$ . [Bhattacharya \(2009\)](#) considers a LP with estimated  $\hat{p}_n$  and known  $M, c$ . Methods developed under a known  $M$  assumption cannot be easily extended to the setting when  $M$  is estimated, as will become evident in Section 2. [Mogstad et al. \(2018\)](#) construct a set-expansion estimator and prove its consistency. [Syrngkanis et al. \(2021\)](#) develop a testing procedure for the failure of LP feasibility. [Gafarov \(2024\)](#) develops uniform inference for a LP described by affine inequalities over unconditional moments, provided uniform Linear Independence Constraint Qualification (LICQ) and Slater’s condition (SC) hold. Gafarov’s conditions may be restrictive in some applications - for example, under AICM, see Section 2. [Andrews et al. \(2023\)](#) develop inference in a special case of LP estimation framework, which arises from their model. In their problem, SC holds and  $\theta$  has a particular structure, making their findings hard to generalize. [Cho and Russell \(2023\)](#) add random distortions to  $p$  and introduce random non-vanishing expansions to  $\Theta_I$  to enforce uniform Hadamard differentiability of the perturbed LP. Their approach yields uniformly valid, yet conservative confidence regions for  $B(\mathbb{P})$  when  $\theta$  is affine

---

<sup>7</sup>Even when continuous identification results are available, in practice estimation is still carried out with discretized covariates. This is true for all empirical work referenced below.

<sup>8</sup>Theoretically,  $w_n$  can diverge arbitrarily slowly. Practical guidance on its selection is provided below.

in unconditional population moments, and their practical procedure implicitly assumes that the SC holds<sup>9</sup>. We conduct Monte Carlo simulations to compare our approach with existing methods in Section 2.6.

A number of frameworks that result in bounds of form (1) are described in the review article by Kline and Tamer (2023). Other examples include conditional moment inequalities (Andrews et al., 2023), generalized IV models (Mogstad et al., 2018) and a strand of models nested in (3), pioneered by Manski (1997) and Manski and Pepper (2000, 2009) (MP). The novel LP sharp bounds in Theorem 3.1 coincide with or tighten the bounds in Blundell et al. (2007), Boes (2009), Siddique (2013), Kreider et al. (2012), De Haan (2017) and Cygan-Rehm et al. (2017). To compute the bounds, the above papers use the plug-in estimator  $B(\hat{\theta}_n)$  combined with bootstrap for inference. Our results show that this procedure relies on rather strong assumptions. The exact inference procedure in Theorem 2.3 could be used instead.

AICM approach complements the method of Mogstad et al. (2018), who develop identification theory for generalized IV estimators and obtain bounds in form (1). By virtue of imposing a Heckman and Vytlacil (1999, 2005) treatment selection mechanism and working in the binary treatment case, they accommodate arbitrary a.s. restrictions on the shape of the marginal treatment response functions and produce bounds for a wider family of treatment parameters. Additive separability in treatment selection is equivalent to the Imbens and Angrist (1994) IV conditions under instrument exogeneity (Vytlacil, 2002). Even though (3) nests mean-independence conditions, it appears most useful when an IV is not available. For that reason, a separable selection mechanism is not justified for AICM<sup>10</sup>. AICM is not related to the model in Andrews et al. (2023) other than by virtue of resulting in bounds of form (1). While our inequalities are imposed *over* affine combinations of counterfactual conditional moments, the latter work effectively generalizes the regression framework to moment inequality restrictions on the error term. We are not aware of conditions, similar to those in Imbens and Angrist (1994), that would allow to state linear conditional moment inequalities models in the potential outcomes form.

## Notation

All vectors are column vectors, and  $M'$  denotes the transpose of  $M \in \mathbb{R}^{n \times m}$ . If  $A$  is a set,  $A'$  stands for its complement. A collection  $(x_j)_{j \in J}$  is a column vector.  $2^A$  denotes the powerset of set  $A$ , and  $\overline{m, n}$  is the collection of integers from  $m$  to  $n$ .  $\times$  is a Cartesian product of sets, while  $\otimes$  is the Kronecker product. The sign  $\sqcup$  denotes a disjoint union. Signs  $\wedge$  and  $\vee$  stand for logical ‘and’ and ‘or’ operators respectively. For  $M \in \mathbb{R}^{m \times n}$  and  $A \subseteq \overline{1, m}$ ,  $M_A \in \mathbb{R}^{|A| \times n}$  is the submatrix of the rows

<sup>9</sup>We discuss this in more detail in Section 2.6.

<sup>10</sup>Thus, if one is faced with i) a binary treatment setup, ii) has a valid IV and iii) no outcomes’ data is missing, the method of Mogstad et al. (2018) may be used. If any of these conditions fail, our approach is an alternative.

of  $M$  with indices in  $A$ . If  $j \in \overline{1, n}$ , write  $M_j \equiv M'_{\{j\}}$ .  $\mathcal{R}(M)$  stands for the range of  $M$ , and  $\sigma_d(M)$  is the  $d$ -th largest singular value of  $M$ . In a normed space  $S$ , the distance between  $x \in S$  and  $A \subseteq S$  is written as  $d(x, A) \equiv \inf_{a \in A} \|x - a\|$ , and  $d_H(A, B) \equiv \max\{\sup_{b \in B} d(b, A), \sup_{a \in A} d(a, B)\}$  is the Hausdorff distance between  $A, B \subseteq S$ . For  $A \subseteq S$  the open expansion is  $A^\varepsilon \equiv \{s \in S : d(s, A) < \varepsilon\}$ .  $\text{Int}(A)$  and  $\text{Cl}(A)$  are the interior and closure of  $A \subseteq \mathbb{R}^d$ , while  $\text{Cone}(A)$  is its conical hull. If  $A$  is a matrix,  $\text{Cone}(A)$  is the conical hull of its columns.  $s(x, A) \equiv \max_{a \in A} x'a$  for a compact  $A \subseteq \mathbb{R}^d$  and  $x \in \mathbb{R}^d$  is a support function. For  $v = (v_j)_{j \in \overline{1, d}}$ , define  $v^+ \equiv (\max\{v_j, 0\})_{j \in \overline{1, d}}$ . For  $v, u \in \mathbb{R}^d$  vector inequalities  $v > u$  and  $v \geq u$  mean  $v_i > u_i \forall i \in \overline{1, d}$  and  $v_i \geq u_i \forall i \in \overline{1, d}$  respectively.  $\iota_d \in \mathbb{R}^d$  is a vector of ones,  $I_d \in \mathbb{R}^{d \times d}$  is the identity matrix, and the subscript is dropped occasionally. Operator  $\mathbb{E}_{\mathbb{P}}$  is the expectation under a measure  $\mathbb{P}$ , and the subscript is dropped whenever it does not cause confusion. The statement  $w_n \rightarrow \infty$  w.p.a.1 means that  $\forall M > 0, \lim_{n \rightarrow \infty} \mathbb{P}[w_n > M] = 1$ . We adopt the convention  $\inf \emptyset = +\infty$ , and  $\sup \emptyset = -\infty$ .

## 2. LP estimation framework

In many partial identification settings, bounds on the parameters of interest can be expressed as LP values (see [Kline and Tamer \(2023\)](#) for a review). Readers who prefer to first see an identification framework resulting in such bounds may refer to Section 3, where LP sharp bounds are derived for a general class of AICM models. This section focuses on the estimation theory for such problems. The LP value function is given by

$$B(\theta) \equiv \inf_{Mx \geq c} p'x, \quad (4)$$

where  $M \in \mathbb{R}^{q \times d}$ ,  $c \in \mathbb{R}^q$  and  $p \in \mathbb{R}^d$ . The vector  $\theta \equiv (p', c', \text{vec}(M)')'$  collects parameters of the LP. The estimable value of these parameters at a fixed true measure is denoted by  $\theta_0 \in \mathbb{R}^S$ , with  $S = qd + q + d$ . The value of interest is therefore  $B(\theta_0)$ .

**Remark 2.1.** (4) does not rule out equality constraints, as  $Ax = b \iff Ax \geq b \wedge -Ax \geq -b$ .

We denote the constraint set by  $\Theta_I(\theta) \equiv \{x \in \mathbb{R}^d | Mx \geq c\}$  and omit the argument when  $\theta_0$  is concerned. In the context of Section 3 and other existing applications (e.g. [Mogstad et al. \(2018\)](#)), the set  $\Theta_I$  is the identified set for an unobserved feature  $x$  of the underlying distribution.

**Assumption A0 (Pointwise setup).** *Suppose that at the fixed true parameter  $\theta_0$ : i) The identified set is non-empty,  $\Theta_I(\theta_0) \neq \emptyset$ ; ii)  $\Theta_I(\theta_0) \subseteq \mathcal{X}$  for a known compact  $\mathcal{X} \subseteq \mathbb{R}^d$  and iii) There is an estimator  $\hat{\theta}_n \equiv (\hat{p}'_n, \hat{c}'_n, \text{vec}(\hat{M}_n)')'$ :  $\|\hat{\theta}_n - \theta_0\| = O_p(1/\sqrt{n})$*

Assumption A0 is maintained throughout this section, while the rest of the conditions are stated explicitly. A0.i means that the underlying model *cannot be rejected*. It does not imply that the identifying restrictions are correctly specified. A0.ii is a mild restriction that usually



holds in applications, for example under bounded outcomes in AICM models (see Section 3). The  $\sqrt{n}$ -consistent estimator<sup>11</sup> in A0.iii typically follows from CLT and the Delta-Method.

The following primal and dual solution sets will prove useful in our discussion:

$$\mathcal{A}(\theta) \equiv \arg \min_{Mx \geq c} p'x, \quad \Lambda(\theta) = \arg \max_{M'\lambda = p} c'\lambda.$$

Assumption A0 implies that a finite  $B(\theta_0)$  is attained as a minimum in (4),  $\Theta_I(\theta_0)$  and  $\mathcal{A}(\theta_0)$  are non-empty compacts, and  $\Lambda(\theta_0)$  is non-empty. We now briefly discuss the typically imposed regularity conditions.

**Definition.** Slater’s condition (SC) is the assertion that  $\text{Int}(\Theta_I) \neq \emptyset$ .<sup>12</sup>

SC rules out point-identification of any linear functional of  $x$ . In particular, it precludes exact point-identification of the target  $B(\theta_0)$  and point-identification of  $x$ , i.e. the case when  $|\Theta_I| = 1$ . Most existing methods rely on SC explicitly (Gafarov (2024)) or implicitly (Cho and Russell (2023), Andrews et al. (2023)). Even an ‘approximate’ failure of SC, when the true identified set  $\Theta_I$  becomes ‘thin’, may be problematic for the existing methods in finite samples. Our simulation evidence illustrates this, see Section 2.6. This creates an undesirable tradeoff: high identification power implies poor estimation quality.

**Definition.** Linear independence constraint qualification (LICQ) is the assertion that the submatrix of binding inequality constraints at any  $x \in \mathcal{A}(\theta_0)$  is full-rank.

LICQ precludes the existence of overidentifying constraints at the optimum. It may be hard to justify in ‘bigger’ models, like the one we develop and apply in Sections 4 and 5. These feature a larger number of inequality constraints that may have similar identifying power, so it is not ex-ante clear why there must not be overidentification at the optimum. One may also interpret LICQ as ruling out parameters-on-the-boundary, as the following example clarifies.

**Example 2.1.** Example 2.1 in Fang and Santos (2018) can be restated as a LP:  $B(\theta_0) = \max\{0, \mathbb{E}[X]\} = \min_{t \in \mathbb{R}} t$  s.t.  $t \geq 0, t \geq \mathbb{E}[X]$ . LICQ fails in that program if  $\mathbb{E}[X] = 0$ , which corresponds to the parameter-on-the-boundary case from Andrews (1999, 2000).

**Definition.** No flat faces condition (NFF) is the assertion that  $|\mathcal{A}(\theta_0)| = 1$ .

The notion of flat faces thus corresponds to  $|\mathcal{A}(\theta_0)| \neq 1$ , i.e. the situation in which the bound on the target parameter  $p'x$  is achieved at multiple partially identified features  $x$ .

<sup>11</sup> $\sqrt{n}$  rate straightforwardly generalizes to any  $r_n \rightarrow \infty$ . We focus on  $\sqrt{n}$  throughout for expositional simplicity.

<sup>12</sup>We give simplified versions of assumptions here for simplicity of exposition. In the presence of ‘true’ equalities  $Ax = b$ , SC should be stated in terms of *Relint*, allowing for point-identification along ‘true’ equalities. LICQ should similarly be restated to account for equalities, as in Gafarov (2024).



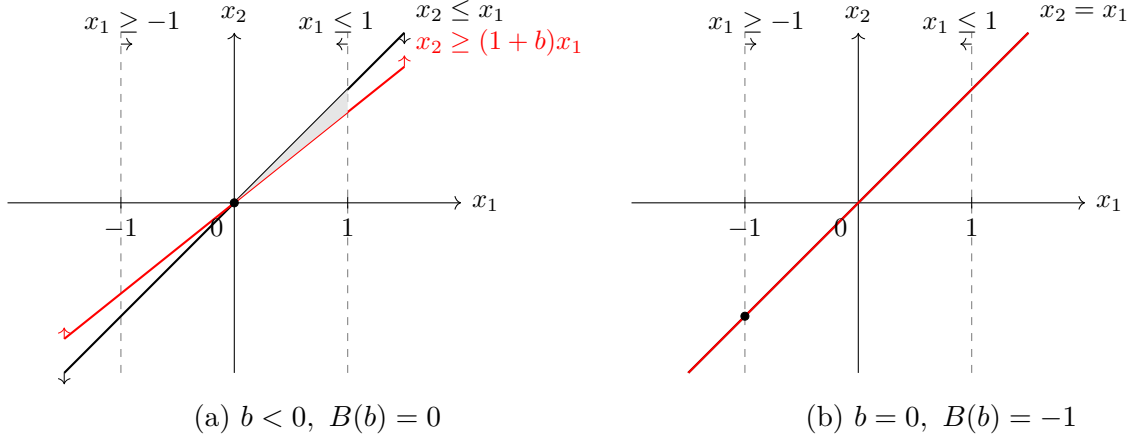


Figure 2: The feasible region in (5) for two values of  $b$ .

Assumption A0 does not impose LICQ or SC, nor does it rule out flat faces. Estimating  $B(\theta_0)$  without these conditions is challenging due to the irregular behavior of  $B(\cdot)$ . If SC fails,  $B(\cdot)$  may be discontinuous at  $\theta_0$ , and the plug-in estimator  $B(\hat{\theta}_n)$  may not be pointwise consistent.

**Proposition 2.1.** *If SC fails for  $\Theta_I(\theta_0)$ ,  $B(\hat{\theta}_n)$  is not, in general, consistent for  $B(\theta_0)$ . Moreover,  $B(\hat{\theta}_n)$  may fail to exist with non-vanishing probability.*

*Proof.* We provide a simple example for the first part. Consider

$$B(b) = \min_x x_1 \quad \text{s.t. : } x_2 \geq (1+b)x_1, \quad x_2 \leq x_1, \quad x_1 \in [-1; 1], \quad (5)$$

where  $b$  is estimated via  $\hat{b}_n = b + \frac{1}{n} \sum_{i=1}^n U_i$  with  $U_i \sim U[-1; 1]$  i.i.d. Suppose in population  $b = 0$ , as in Figure 2b. The true value is then  $B(0) = -1$ , attained at  $x^* = -1$ . The plug-in estimator collapses to

$$B(\hat{b}_n) = -\mathbf{1}\{\hat{b}_n \geq 0\} \rightarrow -1 \text{ in probability.}$$

For the second part, consider the family

$$B(a) = \min_x x_1 \quad \text{s.t. : } x_2 \geq x_1 + a, \quad x_2 \leq x_1, \quad x_1 \in [-1; 1],$$

where  $a$  is estimated via  $\hat{a}_n = a + \frac{1}{n} \sum_{i=1}^n U_i$  with  $U_i$  as before. Suppose  $a = 0$ . If  $\hat{a}_n > 0$ , the plug-in estimator  $B(\hat{a}_n)$  does not exist. This occurs with probability  $1/2$  for any  $n \in \mathbb{N}$ . ■

In some special cases (e.g., [Honoré and Tamer \(2006\)](#)), SC may be argued to hold, ensuring the continuity of  $B(\cdot)$ . However, an additional challenge arises:  $B(\cdot)$  is not necessarily

Hadamard differentiable unless LICQ and NFF also hold. This complicates inference, as Proposition 2.4 in Section 2.2 illustrates.

This section addresses Propositions 2.1 and 2.4. Section 2.1 introduces the penalty function estimator and its debiased version, which we show to be  $\sqrt{n}$ -pointwise consistent under A0. Section 2.2 develops a computationally efficient inference procedure with exact coverage under A0. Turning to uniform properties, Section 2.3 presents a general impossibility result for discontinuous functionals. It implies that the LP value cannot be uniformly consistently estimated under a uniform version of A0 alone. We then characterize a broad class of measures over which a uniformly consistent estimator exists. Sections 2.4 and 2.5 establish the uniform rates of the penalty function estimators over this class. Section 2.6 provides simulation evidence.

## 2.1. Consistency

**2.1.a. Penalty function estimator.** We now develop a consistent estimator that is inspired by the theory of exact penalty functions. The idea is to restate (4) as an unconstrained penalized problem. Define the  $L_1$ -penalized version of the LP objective as

$$L(x; \theta, w) \equiv p'x + w'(c - Mx)^+,$$

and consider the unconstrained problem

$$\tilde{B}(\theta; w) \equiv \min_{x \in \mathcal{X}} L(x; \theta, w), \quad \tilde{A}(\theta; w) \equiv \arg \min_{x \in \mathcal{X}} L(x; \theta, w). \quad (6)$$

We use  $\tilde{B}(\cdot)$  to obtain a preliminary estimator, which we term *the penalty function estimator*. Note that  $L(x; \theta, w) = p'x$  at any  $x \in \Theta_I$ , i.e. the penalized function is equal to the objective function whenever the constraint in (4) holds.

**Assumption A1 (Penalty parameter).** *The penalty vector  $w \in \mathbb{R}^q$  is such that in the initial LP there exists a KKT vector  $\lambda^* \in \Lambda(\theta_0)$  such that  $w > \lambda^*$ .*

**Remark 2.2.** Note that Assumption A1 does not require  $w$  to be component-wise larger than all KKT vectors. In any solvable LP there exists at least one  $\lambda^* < \infty$ , so at any fixed  $\theta_0$  there always exists a large enough  $w$  that satisfies A1.

**Remark 2.3.** If it is known that i)  $B(\theta_0) < K$  for some  $K > 0$  and ii)  $c > \underline{c} > 0$  for some known  $\underline{c} > 0$ , Assumption A1 is satisfied for  $w = \iota K / \underline{c}$ , which is known.

The following Lemma is key to understanding the penalty function approach. It asserts that under Assumption A1 the  $L_1$ -penalty function is *exact* for the LP in (4).

**Lemma 2.1.** For any  $(\theta, w) \in \mathbb{R}^S \times \mathbb{R}_+^q$ , if  $\Theta_I(\theta) \subseteq \mathcal{X}$ , then:

$$\tilde{B}(\theta; w) \leq B(\theta) \tag{7}$$

If the pair  $(\theta_0, w)$  satisfies Assumption A1, then: i) (7) holds with an equality, and ii) optimal solutions coincide:  $\tilde{\mathcal{A}}(\theta_0; w) = \mathcal{A}(\theta_0)$ .

The deterministic result in Lemma 2.1, combined with the observation that the objective function in probability converges uniformly in  $x$  under A0, establish that the penalty function estimator with a fixed  $w$  is consistent under A1:

**Proposition 2.2.** Under Assumption A1,

$$\tilde{B}(\hat{\theta}_n; w) \xrightarrow{p} B(\theta_0).$$

In what follows, the variation across coordinates of  $w$  will not be of interest. We shall thus treat  $w \in \mathbb{R}_+$  as a scalar that induces the penalty vector  $w\iota$ . Based on Lemma 2.1 and Proposition 2.2, it might seem that  $w$  should be selected to be as large as possible. This, however, yields a generally inconsistent estimator if SC fails.

**Remark 2.4.** Consider (5) with  $b = 0$  and suppose  $w > 2$ . If  $\hat{b}_n < 0$ , there exists a sample KKT vector  $\hat{\lambda}_n \in \Lambda(\hat{b}_n)$ , whose largest coordinate is  $\|\hat{\lambda}_n\|_\infty = |\hat{b}_n^{-1}|$ . So, if also  $w > |\hat{b}_n^{-1}|$ , the penalty estimator selects an incorrect optimum  $(0, 0)$  in light of Lemma 2.1. Since  $\hat{b}_n \xrightarrow{p} 0$ , at a large enough sample size  $|\hat{b}_n^{-1}|$  will exceed any fixed  $w$  with high probability and the correct minimum of  $-1$  will be estimated. However, that logic fails in finite samples if  $w$  is ‘large’.

This observation justifies the need to study  $w \rightarrow \infty$  asymptotic theory. We show that the penalty parameter can be allowed to diverge at the rate dominated by  $\sqrt{n}$ .

**Theorem 2.1.** For any  $w_n \rightarrow \infty$  w.p.a.1 with  $\frac{w_n}{\sqrt{n}} \xrightarrow{p} 0$ , we have

$$|\tilde{B}_n(\hat{\theta}_n, w_n) - B(\theta_0)| = O_p\left(\frac{w_n}{\sqrt{n}}\right)$$

The estimator in Theorem 2.1 does not rely on Assumption A1, as the latter is always satisfied at a fixed measure for a large enough  $n$  when  $w_n \rightarrow \infty$ .

**2.1.b. Debiased penalty estimator.** The  $\frac{w_n}{\sqrt{n}}$  rate of convergence in Theorem 2.1 is determined by the slowly vanishing penalty term. This term is a product of the deviation from the true polytope, that vanishes at  $\frac{1}{\sqrt{n}}$ , and an exploding sequence  $w_n$ . It is thus reasonable to

ask whether the  $\sqrt{n}$ -rate could be restored by dropping the penalty term, i.e. *debiasing* the penalty function. We show that this can be done. Before we proceed, let us make the following simplification. Without loss of generality, suppose that

$$\hat{p}_n = p \text{ - non-random.} \quad (8)$$

To see why this is w.l.g., note that one can set  $p = e_1 = (1 \ 0 \dots 0)'$  and add an auxiliary variable for the value of the problem in the first position of  $x$  (see [Gafarov \(2024\)](#)).

The following theorem is one of the main contributions of this paper.

**Theorem 2.2.** Suppose  $\mathcal{A}(\theta_0) \subseteq \text{Int}(\mathcal{X})$ . For any  $w_n \rightarrow \infty$  w.p.a.1 with  $\frac{w_n}{\sqrt{n}} \xrightarrow{p} 0$ ,

$$\sup_{x \in \mathcal{A}(\hat{\theta}_n, w_n)} |p'x - B(\theta_0)| = O_p\left(\frac{1}{\sqrt{n}}\right).$$

Before we discuss the Theorem, let us introduce some notation and jargon. For  $x \in \mathbb{R}^d$ , define the set of constraints that bind at it, when evaluated at  $\tilde{\theta}$ , by  $J(x; \tilde{\theta}) \equiv \{j \in [q] : \tilde{M}_j x = \tilde{c}_j\}$ .

**Definition (Vertex).** We call  $x \in \tilde{\mathcal{A}}(\hat{\theta}_n; w_n)$  a *vertex-solution* if the corresponding matrix of binding constraints,  $\hat{M}_{nJ(x; \hat{\theta}_n)}$ , has full column rank.

**Definition (Nice face).** We say that a set  $A \subseteq [q]$  corresponds to a *nice face*  $F \equiv \{x \in \mathbb{R}^d : M_A x = c_A\}$  if  $p'x = B(\theta_0)$  for any  $x \in F$ .

It should be noted that a nice face  $F$  is not necessarily a valid  $k$ -face of the polytope  $\Theta_I$ .

Intuitively, the proof of [Theorem 2.2](#) proceeds in two steps. First, by anti-concentration arguments we establish that with high probability asymptotically the penalty function estimator manages to select a vertex-solution  $\hat{x}_n \in \tilde{\mathcal{A}}(\hat{\theta}_n; w_n)$  with  $\hat{A}_n = J(\hat{x}_n; \theta_n)$ , such that  $\hat{A}_n$  corresponds to a nice face  $F = \{x \in \mathbb{R}^d : M_{\hat{A}_n} x = c_{\hat{A}_n}\}$ .

Once a nice face has been selected, the  $\sqrt{n}$ -convergence of  $p'\hat{x}_n$  to  $B(\theta_0)$  obtains as a consequence of  $(\hat{M}_{nA}, \hat{c}_{nA})$  converging to  $(M_A, c_A)$  at this rate for a fixed  $A \subseteq [q]$ . To illustrate this idea, suppose for a moment that with high probability asymptotically the penalty function estimator only selects nice faces corresponding to  $\hat{A}_n$ , such that  $M_{\hat{A}_n}$  is also an invertible square matrix. Then,  $\hat{x}_n = \hat{M}_{\hat{A}_n}^{-1} \hat{c}_{\hat{A}_n} + o_p(n^{-0.5})$ , which converges to  $x = M_{\hat{A}_n}^{-1} c_{\hat{A}_n} + o_p(n^{-0.5})$  at rate  $\sqrt{n}$ . Because  $x \in F$ , with  $F$  being a nice face w.p.a.1, one has  $p'x = B(\theta_0) + o_p(n^{-0.5})$ , so that also  $p'\hat{x}_n \rightarrow B(\theta_0)$  at  $\sqrt{n}$  by CMT.

**Remark 2.5.** The result in [Theorem 2.2](#) is uniform over the argmin set, so in the context of lower/upper bound estimation one may use  $\max / \min_{\tilde{\mathcal{A}}(\hat{\theta}_n; w_n)} p'x$  to obtain the tightest bound.

We define *the debiased penalty function estimator* as

$$\hat{B}(\hat{\theta}_n; w_n) \equiv \max_{x \in \tilde{\mathcal{A}}(\hat{\theta}_n; w_n)} p'x.$$

The discussion of uniform asymptotic theory in Section 2.3 sheds light on the role of  $w_n$  and the trade-off involved in its selection. The practical guidance on selecting  $w_n$  is then developed on the basis of our results and random matrix theory, see the Appendix.

**Remark 2.6.** An alternative estimator can be constructed using a set-expansion argument:  $\check{B}_n = \min \hat{p}'_n x$  s.t.  $\hat{M}_n x \geq \hat{c}_n - \sqrt{\frac{\kappa_n}{n}} \iota$ . In the Appendix, we show that the results from Chernozhukov et al. (2007) and the geometry of polytopes imply that  $\check{B}_n$  with an appropriately chosen, diverging  $\kappa_n$ , is consistent for  $B(\theta_0)$ . However, it can be rate-conservative, converging at  $\sqrt{n\kappa_n}^{-1/2}$ . It appears to perform worse than the debiased penalty function estimator  $\hat{B}_n$  in our simulations, see Section 2.6.

## 2.2. Inference

This section develops an inference procedure for a general LP estimator, in which all parameters are inferred from the data. This procedure nests special cases in which some parameters remain fixed, as in Semenova (2023) or Bhattacharya (2009).

**Assumption B0 (Random sample).** Suppose  $\hat{\theta}_n = \hat{\theta}_n(\mathcal{D}_n)$  is a measurable function of the sample  $\mathcal{D}_n \equiv \{W_1, W_2, \dots, W_n\}$ , where  $W_i \in \mathbb{R}^{d_w}$ ,  $i \in [n]$  are i.i.d. random vectors.

We suppose that Assumption B0 holds throughout Section 2.2, whereas the rest of the conditions are imposed explicitly.

**Assumption B1 (Asymptotic normality).** The estimator  $\hat{\theta}_n$  is such that, for  $\Sigma < \infty$ ,

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{L} \mathbb{G}_0 \sim \mathcal{N}(0, \Sigma).$$

Assumption B1 is typically warranted by reference to CLT and the Delta Method when  $\hat{\theta}_n = g(n^{-1} \sum_{i=1}^n W_i)$  for some smooth  $g(\cdot)$ , as in the AICM models (3), see Section 3.

**2.2.a. Bootstrapping the plug-in fails even under SC.** To further justify the need for our inferential procedure, we first examine the properties of the approach that combines bootstrap on  $\hat{\theta}_n$  with the plug-in estimator  $B(\hat{\theta}_n)$ . This method is widely used in empirical literature applying AICM conditions (Blundell et al. (2007), Kreider et al. (2012), Gundersen et al. (2012), Siddique (2013), De Haan (2017), and Cygan-Rehm et al. (2017)).

In light of Proposition 2.1, this approach is inapplicable when SC fails. In practice, researchers attempting to apply it to a LP with a small or empty interior of  $\Theta_I$  will likely encounter frequent bootstrap failures.

In some cases, SC may be established. This is true, for example, if the bound of interest can be expressed as an intersection bound  $B(\theta_0) = \max\{c_1, c_2, \dots, c_q\} = \min_{t \in \mathbb{R}} t$  s.t.  $t \geq c_i$ ,  $i \in [q]$ , where  $\theta_0 = (1 \ 1 \ c)'$  (as in [Chernozhukov et al. \(2013\)](#)).

Yet, even if SC holds, we demonstrate that bootstrap inference based on the plug-in estimator is not valid, unless NFF and LICQ also hold. This observation can be viewed as a generalization of the parameter-on-the-boundary problem of [Andrews \(1999, 2000\)](#).

**Definition.** Let  $\mathbb{D}$  and  $\mathbb{E}$  be Banach spaces, and  $f : \mathbb{D}_f \subseteq \mathbb{D} \rightarrow \mathbb{E}$ . The map  $f$  is said to be Hadamard directionally differentiable at  $v \in \mathbb{D}_f$  tangentially to  $\mathbb{D}_0 \subseteq \mathbb{D}$ , if there is a continuous map  $f'_v : \mathbb{D}_0 \rightarrow \mathbb{E}$ , such that

$$\lim_{n \rightarrow \infty} \left\| \frac{f(v + t_n h_n) - f(v)}{t_n} - f'_v(h) \right\|_{\mathbb{E}} = 0,$$

for all sequences  $\{h_n \in \mathbb{D}\}$  and  $\{t_n\} \subset \mathbb{R}_+$  such that  $t_n \rightarrow 0^+$ ,  $h_n \rightarrow h \in \mathbb{D}_0$  as  $n \rightarrow \infty$  and  $v + t_n h_n \in \mathbb{D}_f$  for all  $n$ . If, moreover,  $f'_v(h)$  is linear in  $h$ , the map  $f$  is said to be fully Hadamard differentiable.

For simplicity of exposition, we abstract from the case of ‘true equalities’ in  $\Theta_I$ . The results extend trivially to this case if SC is defined in terms of relative interior.

**Lemma 2.2.** Under SC,  $B(\cdot)$  is Hadamard directionally differentiable at  $\theta_0$ . The directional derivative is given by

$$B'_{\theta_0}(h) = \inf_{x \in \mathcal{A}(\theta_0)} \sup_{\lambda \in \Lambda(\theta_0)} h'_p x + \sum_{i=1}^Q \lambda_i (h_{c_i} - h'_{M_i} x), \quad (9)$$

where  $h = (h'_p, h'_{M_1}, \dots, h'_{M_q}, h_{c_1}, \dots, h_{c_q})'$  is the direction of the increment in  $\theta$ .

*Proof.* [Duan et al. \(2020\)](#), Theorem 4.1 with second-order terms’ coefficients set to 0. ■

Hadamard directional differentiability of  $B(\theta)$  is sufficient for convergence in law.

**Proposition 2.3.** Under SC and Assumption B1, it follows that

$$\sqrt{n}(B(\hat{\theta}_n) - B(\theta_0)) \xrightarrow{L} B'_{\theta_0}(\mathbb{G}_0)$$

*Proof.* [Fang and Santos \(2018\)](#) Theorem 2.1. combined with Lemma 5. ■

Gaussianity of  $B'_{\theta_0}(\mathbb{G}_0)$  is a necessary condition for bootstrap consistency ([Fang and Santos, 2018](#)). Consequently, the empirical literature using bootstrap with the plug-in estimator has implicitly relied on this assumption. However,  $B'_{\theta_0}(\mathbb{G}_0)$  is not normal unless full Hadamard

differentiability holds, i.e.  $B'_{\theta_0}(h)$  is linear in  $h$ . As (9) suggests, this is not generally the case. Theorem 3.1 in Fang and Santos (2018) establishes that bootstrap is inconsistent for the distribution when  $B'_{\theta_0}(h)$  fails to be linear. The typically applied plug-in and bootstrap combination is then only valid under further restrictive assumptions<sup>13</sup>:

**Proposition 2.4.** *If Assumption B1 holds,  $\text{Supp}(\mathbb{G}_0) = \mathbb{R}^S$  and  $\theta_n^*$  satisfies Assumption 3 in Fang and Santos (2018), bootstrap is consistent for the distribution of  $B(\hat{\theta}_n)$  in the sense that*

$$\sup_{f \in BL_1(\mathbb{R})} \left| \mathbb{E}[f(\sqrt{n}(B(\theta_n^*) - B(\hat{\theta}_n))) | \mathcal{D}_n] - \mathbb{E}[f(B'_{\theta_0}(\mathbb{G}_0))] \right| = o_p(1),$$

*if and only if i) SC, ii) NFF, iii) LICQ all hold at  $\theta_0$ .*

**Remark 2.7.** A consistent estimator for the distribution of the plug-in under SC can be obtained by combining the Functional Delta Method (FDM) of Fang and Santos (2018) with the Numerical Delta Method (NDM) given in Hong and Li (2015), see the Appendix.

**Remark 2.8.** The penalty function estimator is H.d.d. in  $\theta$  (see the Appendix), so FDM + NDM combination yields exact inference for it. This approach still relies on an arbitrarily selected sequence  $\epsilon_n$  and features a fixed  $w$ , and so does not appear satisfactory.

**Remark 2.9.** In the Appendix, we show that the set-expansion estimator has a Lipschitz-bounded bias. A conservative inference procedure based on it can then be obtained using FDM + NDM.

**2.2.b. Exact inference on a debiased estimator.** In this section, we develop our approach to statistical inference on  $B(\theta_0)$ , which achieves exact asymptotic coverage. Our method relies on an asymptotically normal version of the debiased penalty estimator with  $w_n \rightarrow \infty$  and is outlined in Algorithm 1. Before presenting the main result, we introduce auxiliary constructions and discuss our assumptions.

For the true  $\theta_0 = (c', \text{vec}(M)')$  and some subset of indices  $A \subseteq [q]$ , consider

$$\exists x \in \mathcal{A}(\theta_0) : M_A x = c_A, \tag{10}$$

$$p \in \mathcal{R}(M'_A). \tag{11}$$

Equation (10) is satisfied if constraints  $A$  may bind simultaneously at some solutions of the original LP, while (11) holds if the objective function's gradient  $p$  is a linear combination of the gradients of inequalities from  $A$ . For example, if  $A$  is a set of *all binding constraints* at some  $x \in \mathcal{A}(\theta_0)$ , equation (11) follows from KKT conditions.

---

<sup>13</sup>The full-support condition in Proposition 2.4 is imposed for expositional purposes. Sufficiency of conditions i, ii, iii holds generally, whereas necessity obtains whenever the derivative  $B'_{\theta_0}(h)$  is not linear over  $\text{Supp}(\mathbb{G}_0)$  when  $\mathcal{A}(\theta_0), \Lambda(\theta_0)$  are not singletons.



---

**Algorithm 1 (Inference procedure)**


---

Given data  $\mathcal{D}_n$ , estimators  $\hat{\theta}(\mathcal{D}_n)$  and  $\hat{\Sigma}_n = \hat{\Sigma}(\mathcal{D}_n)$ , penalty vector  $w(n, \mathcal{D}_n) \in \mathbb{R}^q$  and constants  $\gamma \in (0; 1)$ ,  $\bar{v} > 0$ , follow the steps below to obtain confidence intervals for  $B(\theta_0)$ .

**Step 1 (Split the sample):**

- 1: Randomly split  $\mathcal{D}_n$  into two folds  $\{\mathcal{D}^{(f)}\}_{f=1,2}$  of sizes  $n_1 = \lfloor \gamma n \rfloor$ ,  $n_2 = n - n_1$
- 2: Compute  $\hat{\theta}^{(f)} \equiv \hat{\theta}(\mathcal{D}^{(f)}) = (\hat{c}^{(f)'}, \text{vec}(\hat{M}^{(f)'})' )'$  for  $f = 1, 2$

**Step 2 (Find the vertex):**

- 1: On the first fold, compute the penalty estimator's arg min as

$$\hat{A} \equiv \arg \min_{x \in \mathbb{R}^d, a \in \mathbb{R}^q} p'x + w(n_1, \mathcal{D}^{(1)})'a, \quad \text{s.t.: } a \geq \hat{c}^{(1)} - \hat{M}^{(1)}x, \quad a \geq 0.$$

- 2: Find the (finite) set of vertex-solutions  $\hat{\mathcal{V}}_x \equiv \{x \in \mathbb{R}^d : (x, a) \in \hat{A}, \text{rk}(\hat{M}_{J(x; \hat{\theta}^{(1)})}^{(1)}) = d\}$
- 3: Find the optimal vertex-solution  $\hat{x} \in \arg \max_{x \in \hat{\mathcal{V}}_x} p'x$
- 4: Compute the set of binding inequalities  $\hat{A} \equiv J(\hat{x}; \hat{\theta}^{(1)})$
- 5: Compute  $\check{v} = \arg \min_{v \in \mathbb{R}^{|\hat{A}|}} \|p - \hat{M}_{\hat{A}}^{(1)' }v\|^2$ , s.t.  $\|v\| \leq \bar{v}$

**Step 3 (Construct the C.I.)**

- 1: Compute  $\hat{\sigma}_n^2 = \sigma^2(\hat{A}, \hat{x}, \hat{v}, \hat{\Sigma}_n)$  using the formula in Lemma 6.6.
- 2: Compute an updated estimate  $\check{B} \equiv \check{v}'(\hat{c}_{\hat{A}}^{(2)} - \hat{M}_{\hat{A}}^{(2)}\hat{x}) + p'\hat{x}$
- 3: The right, two-side and left  $\alpha$ -confidence intervals for  $B(\theta_0)$  are given by

$$\left(\check{B} - \frac{\hat{\sigma}_n}{\sqrt{n_2}}z_{1-\alpha}; +\infty\right), \quad \left(\check{B} - \frac{\hat{\sigma}_n}{\sqrt{n_2}}z_{1-\alpha/2}; \check{B} + \frac{\hat{\sigma}_n}{\sqrt{n_2}}z_{1-\alpha/2}\right), \quad \left(-\infty; \check{B} + \frac{\hat{\sigma}_n}{\sqrt{n_2}}z_{1-\alpha/2}\right)$$


---

Continuing the discussion in Section 2.1.b, we note that subsets  $A$  that satisfy (10) and (11) correspond to the nice faces (see p.12).

**Lemma 2.3.** If  $A \subseteq [q]$  satisfies (10) and (11), then  $F = \{x \in \mathbb{R}^d : M_A x = c_A\}$  is a nice face,

$$B(\theta_0) = p'x, \quad \forall x \in F.$$

Define the set  $\mathbb{A} \equiv \{A \in 2^{[q]} : |A| \geq d, A \text{ satisfies (10) and (11)}\}$ . With probability approaching 1, the penalty function estimator manages to select a vertex-solution  $\hat{x}_n \in \tilde{\mathcal{A}}(\hat{\theta}_n; w_n)$ , determined by the binding constraints  $\hat{A}_n = J(\hat{x}_n; \hat{\theta}_n) \in \mathbb{A}$  that satisfy (10) and (11) and thus correspond to a nice face by Lemma 2.3.

The debiased estimator may hence be understood as a two-stage procedure: one first finds the set of binding inequalities  $\hat{A}_n \in \mathbb{A}$ , and then estimates  $B(\theta_0)$  as  $p'(\hat{M}_{n\hat{A}_n})^\dagger \hat{c}_{n\hat{A}_n}$ . Performing inference on that object directly requires working with the joint distribution of  $\hat{A}_n$ , and  $\hat{\theta}_n$ , leading to a complex and likely non-normal asymptotic distribution.

We address this by ‘disentangling’ the variation in  $\hat{A}_n$  and  $\hat{\theta}_n$  via sample splitting. Intuitively, a vertex is estimated on one part of the sample, while the noise in the parameter estimation comes from the other. We now state our assumptions and present the main result.

**Assumption B2.** *B1 holds, and there exists an estimator  $\hat{\Sigma}_n \xrightarrow{p} \Sigma$ .*

Assumption B2 requires the researcher to possess a consistent estimator of the asymptotic variance of  $\hat{\theta}_n$ . If  $\hat{\theta}_n = g(n^{-1} \sum_{i=1}^n W_i)$  for some smooth and known  $g(\cdot)$ , such estimator can typically be obtained from the estimated covariance matrix of  $W_i$  using via Delta-method. In more complicated scenarios, bootstrap on  $\hat{\theta}_n$  may be employed.

Define the set  $\mathcal{S}_A \equiv \{v \in \mathbb{R}^{|A|} : p = M'_A v\}$  and note that (11) is equivalent to  $\mathcal{S}_A \neq \emptyset$ .

**Assumption B3.** *For a constant  $\bar{v} > 0$ ,  $\max_{A \in \mathbb{A}} \min_{v \in \mathcal{S}_A} \|v\| \leq \bar{v}$ .*

Assumption B3 is a technical condition ensuring that we can find a sequence approaching  $\mathcal{S}_A$  asymptotically, i.e.  $d(\check{v}, \mathcal{S}_A) = o_p(1)$  for  $\check{v}$  defined in (12). Practical guidance on choosing  $\bar{v}$  is provided in the Appendix. Unlike the penalty parameter  $w_n$ , our simulations suggest that the specific choice of  $\bar{v}$  has little impact on inference, as long as it is sufficiently large.

**Definition** (Optimal triplet). We call  $(A, x, v) \in 2^{[q]} \times \mathbb{R}^d \times \mathbb{R}^q$  an optimal triplet if i)  $|A| \geq d$ , ii)  $x \in \mathcal{A}(\theta_0)$ , iii)  $M_A x = c_A$ , iv)  $p = M'_A v_A$ , and v)  $A = \text{Supp}(v)$ .

We randomly split  $\mathcal{D}_n$  into two disjoint, collectively exhaustive folds  $\mathcal{D}_n^{(f)}$  of size  $n_f$  for  $f = 1, 2$ , with  $n_1 = \lfloor \gamma n \rfloor$  and  $n_2 = n - \lfloor \gamma n \rfloor$  for some fixed  $\gamma \in (0, 1)$ . Our inference procedure uses the data from  $\mathcal{D}^{(1)}$  to estimate an optimal triplet  $(\hat{A}, \hat{x}, \hat{v})$ . The vertex<sup>14</sup> is estimated as

$$\hat{x} \in \arg \max_{x \in \tilde{\mathcal{A}}(\hat{\theta}^{(1)}; w_{n_1})} p'x \quad \text{s.t.:} \quad \text{rk}(\hat{M}_{J(x; \hat{\theta}^{(1)})}) = d,$$

the set of binding constraints that define it is denoted by  $\hat{A} \equiv J(\hat{x}; \hat{\theta}^{(1)})$ , and

$$\check{v} \in \arg \min_{\|v\| \leq \bar{v}} \|\hat{M}_{\hat{A}}^{(1)'} v - p\|^2. \quad (12)$$

Our procedure is then based on showing that, for large  $n$ ,

$$\sqrt{n_2} \left( \check{v}' (\hat{c}_{\hat{A}}^{(2)} - \hat{M}_{\hat{A}}^{(2)} \hat{x}) + p' \hat{x} - B(\theta_0) \right) \approx \mathcal{N}(0, \sigma^2(\hat{A}, \hat{x}, \hat{v}, \Sigma)),$$

where  $\sigma^2(\cdot)$  is derived in the Appendix, and  $\hat{v} \in \mathbb{R}^q$  satisfies  $\hat{v}_{\hat{A}} = \check{v}$  and  $\hat{v}_j = 0$  for  $j \notin \hat{A}$ .

**Assumption B4 (Non-degeneracy).** *Suppose  $\sigma(A, x, v, \Sigma) > 0$  for any optimal triplet  $A, x, v$ .*

---

<sup>14</sup>While we assume that  $\tilde{\mathcal{A}}(\hat{\theta}^{(1)}; w_{n_1})$  is estimated precisely, the results do not change if one is only able to estimate a single optimum. This may occur if numerical errors do not allow the LP-solver to find all of the LP solutions. Such optimum will satisfy  $\text{rk}(\hat{M}_{J(x; \hat{\theta}^{(1)})}) = d$  by definition, and so will be a valid vertex-solution.

An inspection of the proof of Theorem 2.3 below reveals that Assumption B4 rules out the scenarios when finding a  $A \in \mathbb{A}$  exactly determines the value  $B(\theta_0)$ . This may occur, for example, if the corresponding  $\hat{c}_A^{(2)}, \hat{M}_A^{(2)}$  are deterministic. In this case, if  $A$  is also unique, meaning  $|\mathbb{A}| = 1$ , the debiased estimator has 0 asymptotic variance, because  $A$  and therefore  $B(\theta_0)$  are correctly estimated with probability approaching 1.

**Theorem 2.3.** Suppose  $\mathcal{A}(\theta_0) \subseteq \text{Int}(\mathcal{X})$  and Assumptions B1, B3, B4 hold. Moreover,

$$\hat{\sigma}_n(A, x, v) \xrightarrow{p} \sigma(A, x, v, \Sigma)$$

for any optimal triplet  $(A, x, v)$  with  $\|v\| \leq \bar{v}$ , which holds for  $\hat{\sigma}_n(A, x, v) = \sigma(A, x, v, \hat{\Sigma}_n)$  under Assumption B2. Then, for any  $\alpha \in (0; 1)$ , and any  $w_n \rightarrow \infty$  w.p.a.1 such that  $w_n = o_p(\sqrt{n})$ ,

$$\mathbb{P} \left[ \frac{\sqrt{n_2}}{\hat{\sigma}_n(\hat{A}, \hat{v}, \hat{x})} \left( \check{v}'(\hat{c}_{\hat{A}}^{(2)} - \hat{M}_{\hat{A}}^{(2)} \hat{x}) + p' \hat{x} - B(\theta_0) \right) \leq z_{1-\alpha} \right] = 1 - \alpha + o(1).$$

**Remark 2.10.** Following Gafarov (2024), one can drop Assumption B4 by using  $\max\{\hat{\sigma}_n(\cdot), \underline{\sigma}\}$  for some small  $\underline{\sigma} > 0$  instead of  $\hat{\sigma}_n(\cdot)$  in Theorem 2.3. In that case, the test would have a correct level, but potentially conservative size.

**Remark 2.11.** If an estimator  $\hat{\Sigma}_n$  is not available, one may alternatively compute the quantiles by performing bootstrap on  $\sqrt{n_2} \hat{v}_{\hat{A}} (\hat{c}_{\hat{A}}^{(2)} - c_{\hat{A}} - (\hat{M}_{\hat{A}}^{(2)} - M_{\hat{A}}) \hat{x})$ , where  $\hat{A}, \hat{x}, \hat{v}$  are fixed, while  $\hat{c}^{(2)}$  and  $\hat{M}^{(2)}$  are bootstrapped by resampling from the second fold  $\mathcal{D}^{(2)}$ .

### 2.3. Uniform asymptotic theory

The optimization problem (1) is challenging to study under no further assumptions, since it may feature instability with respect to arbitrary parameters' perturbations. We now show that this not only leads the plug-in  $B(\hat{\theta}_n)$  to fail, but also precludes the existence of uniformly consistent estimators in general. The following auxiliary result establishes that there exists no uniformly consistent estimator for any real-valued functional from a space of probability measures if it is discontinuous in the total variation norm.

**Lemma 2.4.** Suppose a functional  $V : (\mathcal{P}, \|\cdot\|_{TV}) \rightarrow (\mathbb{R}, |\cdot|)$  is discontinuous at  $\mathbb{P}_0 \in \mathcal{P}$ . Then, there exists no uniformly consistent estimator  $\hat{V}_n = \hat{V}_n(X)$ , which is a sequence of measurable functions of the data  $X \sim \mathbb{P}^n$ . Moreover, if  $\delta > 0$  is the jump at  $\mathbb{P}_0$ , then

$$\inf_{\hat{V}_n} \sup_{\mathbb{P} \in \mathcal{P}} \mathbb{E}_{\mathbb{P}} [ \|V(\mathbb{P}) - \hat{V}_n(X(\mathbb{P}^n))\| ] \geq \frac{\delta}{4}, \quad \forall n \in \mathbb{N},$$

where infimum is taken over all measurable functions of the data.

In this section, we treat the parameter  $\theta_0$  as a functional of the underlying probability measure  $\mathbb{P} \in \mathcal{P}$ . We then make the following assumption on the pair  $\theta_0(\cdot), \mathcal{P}$ :

**Assumption U0 (Uniform setup).** *The functional  $\theta_0(\cdot)$  and the set of probability measures  $\mathcal{P}$  are such that: i)  $\theta_0 : (\mathcal{P}, \|\cdot\|_{TV}) \rightarrow (\mathbb{R}^S, \|\cdot\|_2)$  is continuous; ii)  $\theta_0(\mathcal{P}) = \{y \in \mathbb{R}^S \text{ s.t. } \Theta_I(y) \neq \emptyset, \Theta_I(y) \subseteq \mathcal{X}\}$  for a known and fixed compact  $\mathcal{X}$*

Assumption U0 defines what is meant by ‘the unrestricted set of measures’. U0.i demands that the true parameter be continuous in  $P$ , which holds, for example, in AICM models (see Section 3, Assumption I0). U0.ii assumes that  $\theta_0$  has full support over  $\mathcal{P}$ , meaning that any  $\theta$  generating a non-rejectable model with  $\Theta_I(\theta) \subseteq \mathcal{X}$  is attained at some  $P \in \mathcal{P}$ .

**Theorem 2.4.** Under U0 there exists no uniformly consistent estimator  $\hat{B}_n$  of  $B(\theta_0)$ .

*Proof.* Combining U0, Lemma 2.4 and trivially extending the example in (5) to  $\mathbb{R}^d$ . ■

Given this negative result, it seems natural to seek a minimal restriction on  $\mathcal{P}$  for which a uniformly consistent estimator may exist. We now show that the condition ensuring uniform consistency of the penalty function approach can be considered minimal in the sense to be made precise in Proposition 2.6.

Remark 2.4 illustrates that  $w_n$  cannot be allowed to diverge faster than  $\sqrt{n}$ , as otherwise the penalty approach may fail at measures where SC fails. At the same time, if all KKT vectors  $\lambda^* \in \Lambda(\theta_0)$  grow large, an arbitrarily large  $w$  is needed for Assumption A1 to hold. This occurs when optimal vertices become ‘sharp’, i.e. all relevant full-rank submatrices of binding inequality constraints grow closer to being degenerate. The condition that ensures uniform consistency of the penalty function approach should therefore bound such ‘sharpness’.

We begin with an existence result based on the Caratheodory’s Conical Hull Theorem.

**Proposition 2.5.** *The problem (4) admits a solution  $x^*$  and the associated KKT vector  $\lambda^*$  such that for some index subset  $J^* \subseteq \{1, \dots, q\}$  with  $|J^*| = d$ ,  $M_{J^*}$  is invertible and:*

$$\begin{aligned} x^* &= M_{J^*}^{-1} c_{J^*}, \\ \lambda_{J^*}^* &= M_{J^*}^{-1'} p, \\ \lambda_i^* &= 0 \text{ if } i \notin J^*. \end{aligned}$$

Proposition 2.5 asserts that any finite and feasible LP has an optimal vertex  $x^*$  at which there is a subset  $J^*$  of binding constraints, such that i) the corresponding gradients form a full-rank square matrix, and ii) the objective function gradient belongs to the conical hull formed by the gradients of the constraints from  $J^*$ .

**Assumption U1 ( $\delta$ -condition).** *The class of measures  $\overline{\mathcal{P}}$  satisfies the  $\delta$ -condition for a given  $\delta > 0$ , if*

$$\inf_{\mathbb{P} \in \overline{\mathcal{P}}} \max_{J^* \in \mathcal{J}^*(\theta(\mathbb{P}))} \sigma_d(M_{J^*}(\mathbb{P})) > \delta, \quad (13)$$

where  $\mathcal{J}^*(\theta(\mathbb{P}))$  collects all  $J^*$  defined in Proposition 2.5 at a given  $\theta(\mathbb{P})$ .

The  $\delta$ -condition does not rule out the failure of LICQ, SC or NFF, and is weaker than the conditions usually imposed to establish uniform consistency of LP estimators. To formalize this, let us introduce three families of measures. Firstly, denote the family of measures satisfying U1 for a given  $\delta > 0$  by  $\mathcal{P}^\delta$ . A measure satisfies the Slater's condition if  $\mathbb{P} \in \mathcal{P}^{SC} \equiv \{\mathbb{P} \in \mathcal{P} \mid \text{Int}(\Theta_I(\theta(\mathbb{P}))) \neq \emptyset\}$ . Similarly, a measure satisfies a uniform  $\varepsilon$ -LICQ condition (as in Gafarov (2024)) if  $\mathbb{P} \in \mathcal{P}^{LICQ;\varepsilon}$ , where

$$\mathcal{P}^{LICQ;\varepsilon} \equiv \{\mathbb{P} \in \mathcal{P} \mid M(\mathbb{P})_A \in \mathbb{R}^{d \times d}, \sigma_d(M(\mathbb{P})_A) > \varepsilon \forall A \in \mathcal{V}(\mathbb{P})\},$$

where the set  $\mathcal{V}(\mathbb{P}) \subseteq 2^{[q]}$  consists of sets indices of binding inequalities that define vertices of the polytope  $\Theta_I(\theta(\mathbb{P}))$ .

**Proposition 2.6.** *The following hold:*

1.  $\mathcal{P}^{Slater} \cup \mathcal{P}^{LICQ;0} \subset \mathcal{P} = \bigcup_{\delta > 0} \mathcal{P}^\delta$ , where the inclusion is strict
2.  $\mathcal{P}^{LICQ;\varepsilon} \subset \mathcal{P}^\delta$  for any  $\delta \leq \varepsilon$ , where the inclusion is strict

*Proof.* Part 1 follows from a trivial extension of the example in (5) to  $\mathbb{R}^d$  and Proposition 2.5. Part 2 is true by definition of  $\mathcal{P}^\delta, \mathcal{P}^{LICQ;\varepsilon}$ . ■

Intuitively, the  $\delta > 0$  in Assumption U1 merely parametrizes the degree of irregularity that the researcher is willing to allow for the identified polytope over the considered set of measures. The resulting family of sets ‘covers’ the unconstrained set of measures asymptotically as  $\delta$  decreases to 0. Uniform LICQ and SC, on the contrary, both restrict the set of measures. That is because measures like the one in Figure 2b do not belong to either  $\mathcal{P}^{LICQ;\varepsilon}$  for any  $\varepsilon \geq 0$ , or  $\mathcal{P}^{SC}$ .

**Example 2.1.** Figure 4 plots the range of  $b$ -values that correspond to the measures satisfying the  $\delta$ -condition for a given  $\delta > 0$  in the example (5). The  $\delta$ -condition cuts off an interval of  $b$  along which the optimal vertex becomes ‘too sharp’. Recall that the case  $b = 0$  leads to the failure of SC. However, it satisfies the  $\delta$ -condition for a relatively large  $\delta$ , because at the optimum  $x^* = (-1 \ -1)'$  there is a set  $J^* = \{1, 3\}$  from Proposition 2.5 such that the relevant matrix of binding constraints  $M_{J^*}$  has the smallest singular value  $\sigma_2(M_{J^*}) \approx 0.62 \gg 0$ .

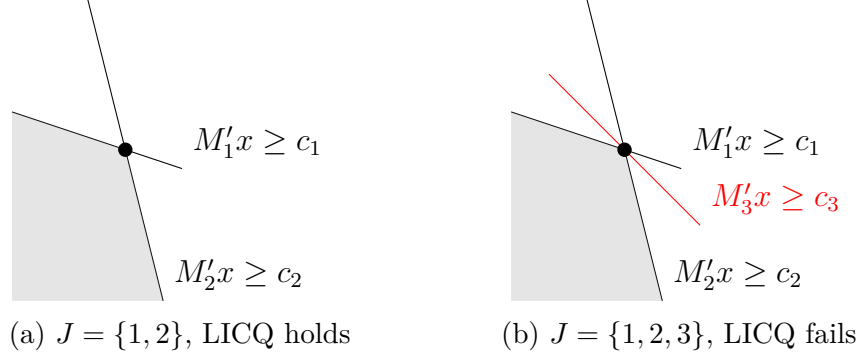


Figure 3: Illustrating the difference between LICQ and  $\mathcal{P}^\delta$ . In (a) and (b) the  $\delta$ -condition holds with the same  $\delta = \sigma_2(M_{\{1,2\}}) \gg 0$ .

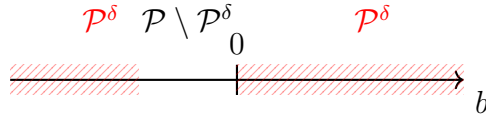


Figure 4: The set of  $b$  in problem (5) satisfying a  $\delta$ -condition.

**Example 2.2.** There are LPs in which SC, LICQ and NFF all fail, but the  $\delta$ -condition is satisfied for a relatively large  $\delta$ . One example is the problem:

$$\min -x_1 + x_2, \quad \text{s.t.} \quad x_2 \leq x_1, x_2 \geq x_1, x_1 \in [-1; 1].$$

The  $\delta$ -condition is satisfied for  $\delta > \sigma_2(M_{\{1,3\}}) \approx 0.62$ , which is large relative to the penalty we suggest (see Section 2.6). There are flat faces, as any pair with  $x_2 = x_1$  and  $x_1 \in [-1; 1]$  is a solution, SC fails, as  $\text{Int}(\Theta_I) = \emptyset$ , and LICQ fails, as at an optimal  $x_2 = x_1 = -1$  there are three binding constraints.

## 2.4. Uniform consistency of penalty function estimator

**Theorem 2.5.** Suppose i)  $\hat{\theta}_n = \hat{\theta}_n(\mathbb{P})$  converges to  $\theta_0(\mathbb{P})$  a.s. uniformly over  $\mathcal{P}$  at rate  $\bar{r}_n \uparrow \infty$ , i.e. for all  $r_n \uparrow \infty$  with  $r_n = o(\bar{r}_n)$  and any  $\varepsilon > 0$ ,

$$\limsup_{n \rightarrow \infty} \sup_{\mathbb{P} \in \mathcal{P}} \mathbb{P}[\sup_{m \geq n} r_m \|\hat{\theta}_m - \theta(\mathbb{P})\| \geq \varepsilon] = 0, \quad (14)$$

and ii)  $w_n(\mathbb{P}) = w_n \rightarrow \infty$  w.p.a.1 a.s. uniformly over  $\mathcal{P}$ , i.e. for any  $M > 0$ ,

$$\lim_{n \rightarrow \infty} \inf_{\mathbb{P} \in \mathcal{P}} \mathbb{P}[\inf_{m \geq n} w_m > M] = 1.$$

Then, for all  $r_n \uparrow \infty$  with  $r_n = o(\frac{\bar{r}_n}{w_n})$  and any  $\varepsilon > 0$ ,

$$\sup_{\delta > 0} \limsup_{n \rightarrow \infty} \sup_{\mathbb{P} \in \mathcal{P}^\delta} \mathbb{P}[\sup_{m \geq n} r_m \|\tilde{B}(\hat{\theta}_m; w_m) - B(\theta_0(\mathbb{P}))\| \geq \varepsilon] = 0. \quad (15)$$

**Example 2.3.** In the AICM models (3),  $\theta_0$  is linear in moments of interactions of  $Y(t)$  with treatment indicators and linear or hyperbolic in probabilities of  $\{T = t, Z = z\}$  (see Section 3). Thus, condition ii) in Theorem 2.5 is established by, firstly, imposing that:

$$\limsup_{C \rightarrow \infty} \sup_{\mathbb{P} \in \mathcal{P}} \mathbb{E}_{\mathbb{P}} \|\mathbb{Y}\| \mathbb{I}\{\|\mathbb{Y}\| > C\} = 0, \quad \mathbb{Y} \equiv (Y(t))_{t \in \mathcal{T}},$$

which holds whenever bounded outcomes are assumed. If also a full-support condition holds:

$$\inf_{\mathbb{P} \in \mathcal{P}, (t,z) \in \mathcal{T} \times \mathcal{Z}} \mathbb{P}[T = t, Z = z] > C$$

for some  $C > 0$ , then  $\theta(\cdot)$  is uniformly continuous in population moments. Combining this with the LLN uniform in probability measure (see Proposition A.5.1 on p. 456 of Van Der Vaart et al. (1996)) yields condition (ii). If one additionally assumes that:

$$\limsup_{C \rightarrow \infty} \sup_{\mathbb{P} \in \mathcal{P}} \mathbb{E}_{\mathbb{P}} \|\mathbb{Y}\|^2 \mathbb{I}\{\|\mathbb{Y}\| > C\} = 0,$$

the rate in (14) is  $r_n = \sqrt{n}$  (see Proposition A.5.2 on p. 457 of Van Der Vaart et al. (1996)).

**Remark 2.12.** The existence of a uniformly consistent estimator over  $\mathcal{P}$  implies that  $B(\cdot)$  is continuous over  $\theta_0(\mathcal{P})$ . However, as Proposition (2.1) illustrates, it is not necessarily continuous over the support of  $\hat{\theta}_n$ . It is straightforward to see that the case  $b = 0$  in (5) satisfies the  $\delta$ -condition for a relatively large  $\delta$ , but the plug-in estimator is still pointwise inconsistent at such measure.

In light of the findings in this section, we develop an approach to selecting the penalty parameter  $w_n$ . It leverages random matrix theory and is given in Appendix.



## 2.5. Uniform consistency of the debiased penalty estimator

We need to introduce the following two objects, which we term the face and polytope condition numbers.

**Definition** (Face condition number). For a  $k$ -face of a polytope,  $f$ , which is described by binding constraints  $A \subseteq \overline{1, q}$  with  $|A| \geq d - k$  such that  $\text{rk}(M_A) = d - k$ , define the face condition number to be:

$$\tilde{\kappa}(f) \equiv \min_{B \subseteq A: \text{rk}(M_B) = d - k} \sigma_{d-k}(M_B)$$

**Definition** (Polytope condition number). For a polytope  $\Theta$ , define the polytope condition number as:

$$\kappa(\Theta) = \min_{f \text{-face of } \Theta} \tilde{\kappa}(f) = \min_{f \text{-vertex of } \Theta} \tilde{\kappa}(f)$$

To see why the second equality above holds, note that any full-rank matrix at a  $k$ -face  $f$  for  $k > 0$  can be obtained by removing  $k$  vectors from a full-rank matrix at some vertex (0-face)  $f^* \subseteq f$ , so the condition number of  $f$  is greater or equal than that of  $f^*$  by the interlacing inequality for singular values.

**Assumption U2 (Polytope  $\delta$ -condition).** *The class of measures  $\overline{\mathcal{P}}$  satisfies the polytope  $\delta$ -condition for a given  $\delta > 0$ , if:*

$$\inf_{\mathbb{P} \in \overline{\mathcal{P}}} \kappa(\Theta_I(\mathbb{P})) \geq \delta$$

The polytope  $\delta$ -condition thus lower bounds the smallest singular values of all full-rank matrices that can be constructed from the vertices of the polytope. Similarly to Assumption U1, it parametrizes the unconstrained set of probability measures, since at any fixed  $\mathbb{P} \in \mathcal{P}$  we have  $\kappa(\Theta_I(\mathbb{P})) > 0$  by definition. Proposition 2.6 continues to hold for U2, if one substitutes  $\mathcal{P}^\delta$  with  $\mathcal{P}_p^\delta$  - the family of measures satisfying U2 for a given  $\delta > 0$ .

The following ‘anticoncentration’ Lemma appears to be mathematically novel.

**Lemma 2.5.** For any non-empty and bounded polytope  $\Theta = \{x \in \mathbb{R}^d | Mx \geq c\}$ :

$$l'(c - Mx)^+ \geq \frac{d(x, \Theta)\kappa(\Theta)}{d}$$

By Theorem 2.5, the biased penalty estimator is uniformly consistent at rate  $\frac{\sqrt{n}}{w_n}$  under the  $\delta$ -condition. The following Theorem asserts that the debiased estimator converges at least at the same rate.

**Theorem 2.6.** Suppose  $\hat{\theta}_n$  converges to  $\theta(\mathbb{P})$  a.s. uniformly over  $\mathcal{P}$ , i.e. for any  $\varepsilon > 0$ :

$$\limsup_{n \rightarrow \infty} \sup_{\mathbb{P} \in \mathcal{P}} \mathbb{P}[\sup_{m \geq n} r_m \|\hat{\theta}_m - \theta(\mathbb{P})\| \geq \varepsilon] = 0$$

Then:

$$\sup_{\delta > 0} \limsup_{n \rightarrow \infty} \sup_{\mathbb{P} \in \mathcal{P}_p^\delta} \mathbb{P}[\sup_{m \geq n} \frac{r_m}{w_m} \|\hat{B}(\hat{\theta}_m; w_m) - B(\theta(\mathbb{P}))\| \geq \varepsilon] = 0$$

It is unclear if the  $\frac{\sqrt{n}}{w_n}$  rate is uniformly sharp. Our simulation evidence (see the Appendix) suggests that this may be the case, but only along the sequences of measures for which SC, LICQ and NFF all fail in the limit. Once such sequences are ruled out, the rate of  $\sqrt{n}$  appears to be restored uniformly.

## 2.6. Monte Carlo

Consistency and inference simulations use  $10^4$  and  $10^3$  repetitions respectively.

**2.6.a. Consistency.** We first describe the example from Proposition 2.1 in more detail. Consider a setup with  $d = 2$  variables and  $q = 4$  constraints. Recall that  $\theta = (p', \text{vec}(M)', c)'$ . In this case, the parameters are defined as

$$p = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \quad M = \begin{pmatrix} -(1+b) & 1 \\ 1 & -1 \\ 1 & 0 \\ -1 & 0 \end{pmatrix}, \quad c = \begin{pmatrix} 0 \\ 0 \\ -1 \\ -1 \end{pmatrix}. \quad (16)$$

The sample analogues  $\hat{M}_n, \hat{c}_n$  of parameters in (16) are obtained by substituting  $b$  with its estimator  $\hat{b}_n$ . In our simulation study, that estimator is given by  $\hat{b}_n = b + n^{-1} \sum_{i=1}^n U_i^b$ , where  $U_i^b \sim U[-1; 1]$  for  $i \in [n]$  are i.i.d.

The true parameter  $b$  is in one-to-one correspondence with the underlying measure, and  $\theta$  is completely described by it:  $\theta(\mathbb{P}) = \theta(b(\mathbb{P}))$ . Consequently, it is sufficient to index the parameter  $\theta$ , the identified polytope  $\Theta_I$  and the program's value  $B$  by  $b$  only.

Observe that the example (16) is so engineered that  $\Theta_I(\hat{b}_n)$  would never be empty or unbounded, ensuring the existence of the plug-in estimator  $B(\hat{b}_n)$ . A slight modification to that setup, that we consider in the next subsection, would lead the identified polytope to be empty with a potentially non-vanishing probability.

**Measures** We consider two values of  $b$ :  $b = 0$  and  $b = -0.05$ . At  $b = 0$ , SC fails, because  $\Theta_I(0)$  has an empty interior. LICQ also fails at  $b = 0$ , as at the optimum  $x = (-1, -1)$ ,

inequalities 1, 2, 4 are binding. There are no flat faces at  $b = 0$ . At  $b = -0.05$ , SC, LICQ and NFF all hold. The values  $b = 0$  and  $b = -0.05$  result in the smallest singular values of  $M_{J^*(b)}(b)$  matrices that correspond to the 75–th and the 19–th percentiles of the Tao-Vu distribution given in Theorem 6.1.

If  $b < 0$ , the norm of the ‘smallest’ KKT vector  $\lambda$  in the true LP corresponding to (16) is proportional to  $|b^{-1}|$ . So, for a small negative  $b$ , the  $\delta$ -condition is only satisfied for a small value of  $\delta > 0$ . In this case the ‘optimal’  $w$  is large. By that logic, the plug-in estimator, which in this case obtains as the limit  $w \rightarrow \infty$ , should perform well at such  $b$ , and potentially outperform the debiased penalty estimator. This is partly an artifact of the setup in (16): additional noise in one of the slanted lines’ intercepts would render the problem infeasible with positive probability, which would worsen the performance of the plug-in estimator in finite samples (see Figure 8).

At  $b \geq 0$ , the  $\delta$ -condition is satisfied with a relatively large  $\delta = \frac{\sqrt{5}-1}{2}$ , and so the ‘optimal’  $w$  is relatively small. If  $b = 0$ , in 50% of the cases,  $\hat{b}_n < 0$  is estimated. If, moreover,  $w_n > \text{const} \times |\hat{b}_n^{-1}|$ , the debiased penalty estimator would select the incorrect maximum of 0 in such cases. A larger  $w$  at  $b = 0$  thus hampers the performance of the debiased penalty estimator.

**Parameters** We set  $w_n = \delta_{0.2}^{-1}(d) \|p\| \frac{\ln \ln n}{\ln \ln 100}$ , where  $\delta_\alpha(d)$  is the  $\alpha$ -quantile of the  $d$ -dimensional Tao-Vu distribution given in Theorem 6.1 in the Appendix. To ensure the same expansion rate for the set-expansion estimator, we set  $\sqrt{\kappa_n} = \kappa_0 \times \ln \ln n$ . There is no guidance as to the selection of  $\kappa_0$ . Our baseline is  $\kappa_0 = 0.1$ , and we explore other values in Figure 7.

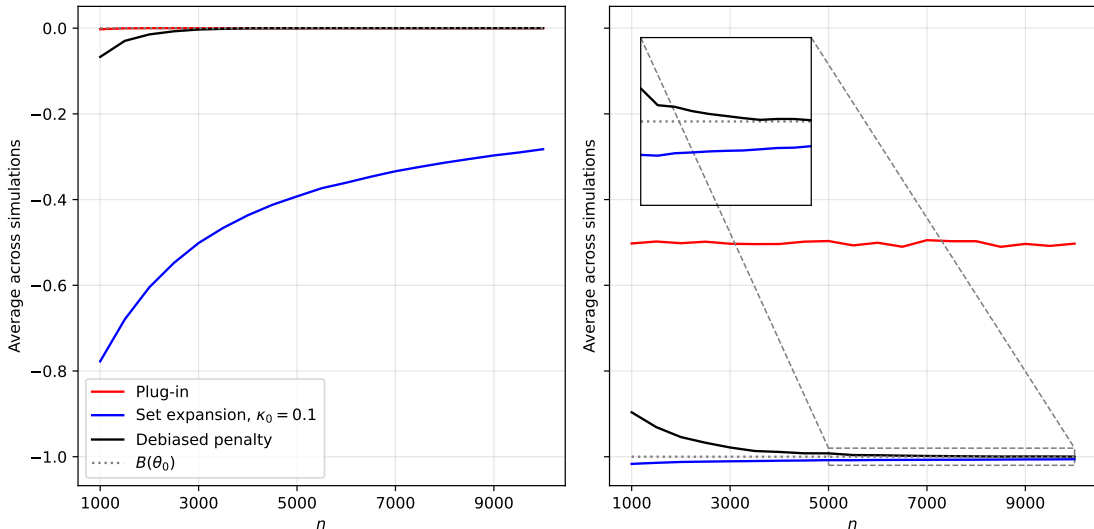


Figure 5: Simulation of example in (16) for  $b = -0.05$  and  $b = 0$  respectively.

**Discussion** Consider the right panel of Figure 5, corresponding to  $b = 0$ . The plug-in estimator is inconsistent, while the set-expansion estimator performs well, approaching  $-1$  from below for larger  $n$ . The debiased penalty-function estimator is slightly upward-biased for smaller  $n$ , but yields the value of almost exactly  $-1$  in larger samples. In contrast, the set-expansion estimator has a conservative rate, and remains slightly downward-biased even in larger samples.

The case of  $b = -0.05$  is depicted in the left panel of Figure 5. The plug-in estimator is consistent and appears to be the best estimator out of the three. The set-expansion estimator is severely downward-biased. This is because when the optimal vertex has a ‘sharp’ angle, a small expansion of the inequalities’ RHS may lead to a large shift of the vertex. To see that, consider shifting both inequalities outwards in Figure 12 for a small and a large absolute value of  $b$ , when  $b$  is negative. Once the expansion grows smaller, the set-expansion estimator slowly converges to the true value of 0 from below. While selecting a smaller  $\kappa_0$  parameter would improve the performance of the set-expansion estimator at  $b = -0.05$ , in the next part of our analysis we demonstrate that in this example  $\kappa_0 = 0.1$  is close to being optimal in the uniform sense, because smaller  $\kappa_0$  worsens the estimator’s performance at  $b = 0$ . The debiased penalty estimator, in contrast, converges rather quickly. It is slightly conservative at smaller  $n$ , as it selects the incorrect vertex of  $-1$  whenever  $\hat{b}_n < 0$  and  $|\hat{b}_n^{-1}| > \text{const} \times w_n$ , i.e. when the penalty parameter is not large enough.

**Robustness of the debiased penalty estimator** For  $b < 0$ , the debiased estimator may be expected to perform better at measures with a larger  $|b|$ . Figure 7 illustrates that point:

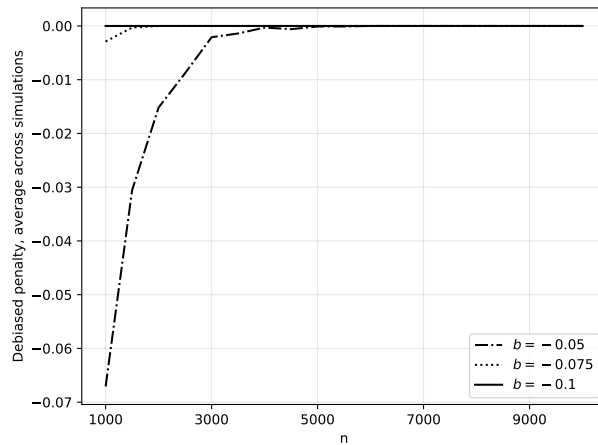


Figure 6: Performance of the debiased penalty estimator for different  $b$  in example (16).

Researchers applying any LP estimator should exercise caution when operating in smaller samples due to irregularity inherent in 4. In example (16), the debiased penalty estimator exhibits desirable behavior even along highly irregular measures for sample sizes of order

$n = 5000$ . Such sample sizes are not uncommon in partially identified settings. In our application, estimation is performed on 664633 observations.

**Alternative  $\kappa_0$**  We now investigate to which extent the behavior of the set-expansion estimator can be improved by selecting an alternative  $\kappa_0$  parameter. Clearly, a larger  $\kappa_n$  makes the set-expansion estimator more conservative. However,  $\kappa_n$  cannot be selected as small is possible. If the expansion is too small, it may be insufficient to counteract the noise involved in estimating  $\Theta_I$ , leading to poor performance of the estimator at measures where SC fails to hold.

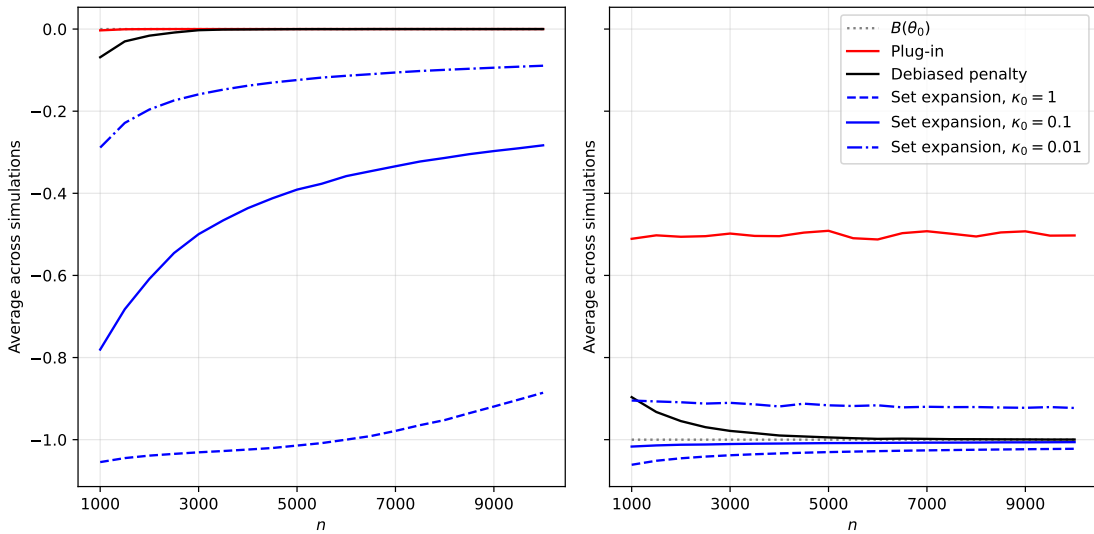


Figure 7: Simulation of (16) for  $b = -0.05$  and  $b = 0$  resp., different  $\kappa_0$ .

That tradeoff is very clear in Figure 7. Compare the performance of the estimator with  $\kappa_0 = 0.01$  to the baseline of  $\kappa_0 = 0.1$ . Decreasing  $\kappa_0$  down to 0.01 makes the estimator less conservative at  $b = -0.05$ , but results in an upward bias at  $b = 0$ . This occurs, because at  $\kappa_0 = 0.01$  the resulting set-expansion sequence  $\kappa_n$  is too small to counteract the estimation noise for the considered sample sizes. This logic is ‘monotone’, meaning that selecting an even smaller  $\kappa_0$  would worsen the performance at  $b = 0$  further. Even at  $\kappa_0 = 0.01$ , the set-expansion estimator is quite conservative for the measure  $b = -0.05$ , and the parameter can clearly not be reduced any further without affecting the validity of the estimator at  $b = 0$ . It appears that the baseline choice of  $\kappa_0 = 0.1$  is close to being optimal in our example. Therefore, the conservative behavior of the set-expansion estimator at  $b = -0.05$  is not explained by a poor choice of the tuning parameter, but rather is a feature of the estimator itself.

**Noise in  $\hat{c}_n$**  We also report the result of consistency simulations for the DGP described in (17). This DGP features noise on the right-hand-side of the inequalities that describe the polytope, which allows the estimated polytope to be empty.

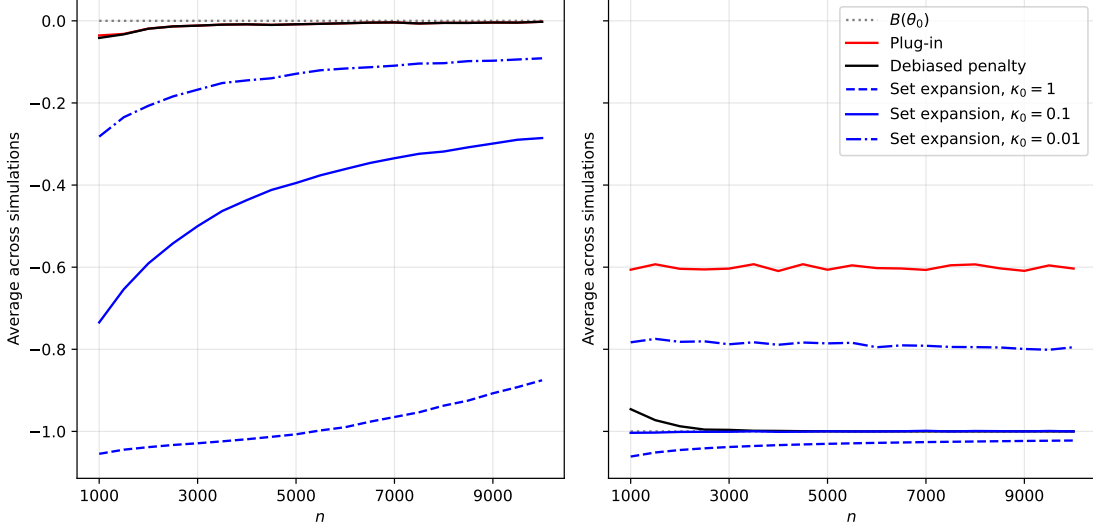


Figure 8: Simulation of (17) for  $b = -0.05$  and  $b = 0$  resp., different  $\kappa_0$ . Averages for the plug-in and set-expansion estimators ignore failed iterations.

In Figure 8, the plug-in and the debiased penalty estimators perform equally well at  $b = -0.05$ . The rest of the conclusions are qualitatively unchanged.

**2.6.b. Inference.** In this section, we assess the performance of our inferential procedure. We compare it to the performance of two recently developed methods, Cho and Russell (2023) (CR) and Gafarov (2024) (BG). We consider a slightly modified version of example (16):

$$M = \begin{pmatrix} -(1+b) & 1 \\ 1+\zeta & -1 \\ 1 & 0 \\ -1 & 0 \end{pmatrix}, \quad c = \begin{pmatrix} \nu \\ \zeta \\ -1-\nu \\ -1 \end{pmatrix} \quad (17)$$

The sample analogues  $\hat{M}_n, \hat{c}_n$  of parameters in (17) are obtained by substituting  $b, \zeta, \nu$  with their estimators:

$$\hat{b}_n = b + n^{-1} \sum_{i=1}^n U_i^b, \quad \hat{\zeta}_n = n^{-1} \sum_{i=1}^n U_i^\zeta, \quad \hat{\nu}_n = n^{-1} \sum_{i=1}^n U_i^\nu,$$

where  $U_i^k \sim U[-0.5; 0.5]$ , i.i.d. across  $i \in [n]$  and  $k \in \{b, \zeta, \nu\}$ . We consider measures, for which the true values of  $\zeta = \nu = 0$ , whereas we still vary  $b$  as in the example before. As before, we mainly consider  $b = -0.05$  and  $b = 0$ .

Note that in (17) the plug-in estimator of the polytope  $\hat{\Theta}_I$  may be empty for some realizations of  $\{U_i^k\}_{i,k}$ . That is the case, for example, if  $\hat{b}_n = \hat{\zeta}_n = 0$  and  $\hat{\nu}_n > 0$ .

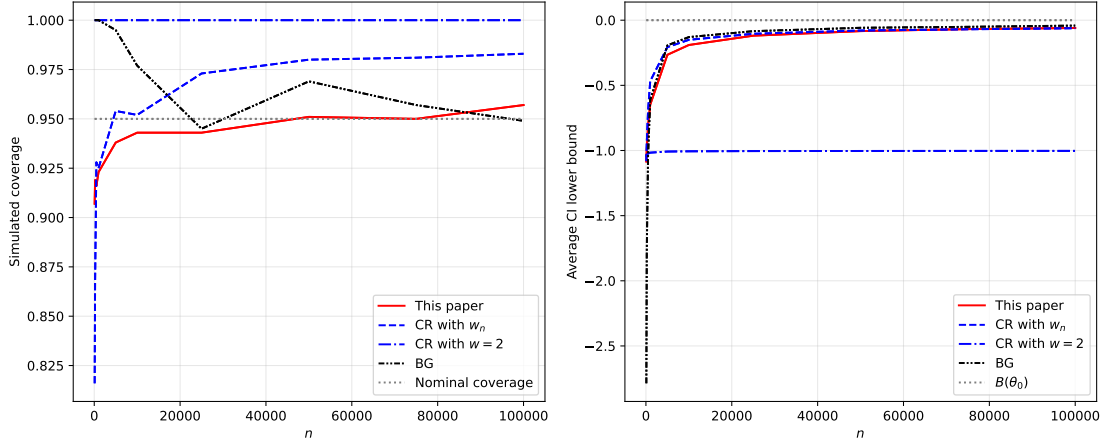
**Other methods** In case SC fails, the BG estimator is not applicable. It fails to exist in around 25% of cases at  $b = 0$ . The CR estimator in the main text of the paper also relies on  $B(\hat{\theta}_n)$  and is not applicable if SC fails. The procedure described on p.47 of the Supplementary Appendix in Cho and Russell (2023) would not be practically applicable without the results contained in the present paper. It can be shown that the CR augmented procedure combines a random, non-vanishing set expansion and objective function perturbation with the penalty function approach<sup>15</sup>. The authors, however, treat the analogue of the penalty parameter as ‘some [fixed] large value’. They proceed to argue that the inequalities, which establish the size of their inference procedure, hold for any value of  $w$ , which appears to suggest that its selection is unimportant. There is no practical guidance on  $w$  selection, and CR do not implement the augmented estimator in their simulations. In this paper, we studied penalized estimation in great detail, and our results suggest that the appropriate choice of  $w$  is critical. Both our previous findings and simulations in Figure 9 demonstrate that for different values of  $w$  the performance of the CR procedure ranges from highly conservative to invalid. Selecting ‘a large value’ of  $w$  does not yield a valid procedure in finite samples. Our simulations also suggest that CR augmented approach can perform relatively well if combined with our results on the rate and level of  $w_n$ . Unlike our approach, however, it remains asymptotically conservative due to the use of non-vanishing random expansions.

**Implementation details** As mentioned in Section 2.2, one can estimate the as. covariance matrix  $\hat{\Sigma}_n$  of  $\hat{\theta}_n$  with resampling. One then plugs it into the expression for  $\sigma(\cdot)$  to obtain the required s.e. An even simpler approach is to compute  $\sigma(\cdot)$  directly by bootstrapping the quantities estimated on the second fold, while keeping the first fold quantities fixed. We have verified that the performance of the two approaches is similar to using the closed-form estimator for  $\Sigma$ , namely  $\hat{\Sigma}_n = G\hat{\Omega}_nG'$ , where  $G \equiv \frac{\partial\theta}{\partial(b \ \zeta \ \nu)}$  and  $\hat{\Omega}_n$  is the sample covariance matrix of  $(b_i \ \zeta_i \ \nu_i)'$ . We employ the latter estimator in our simulations. When implementing the procedure in Cho and Russell (2023) (CR), we use uniform noise with the support size of 0.001, as recommended in the paper. Note that we refer to their parameter  $M$  as  $w$ . The estimator in Gafarov (2024) (BG) was implemented using the code kindly provided by the author.

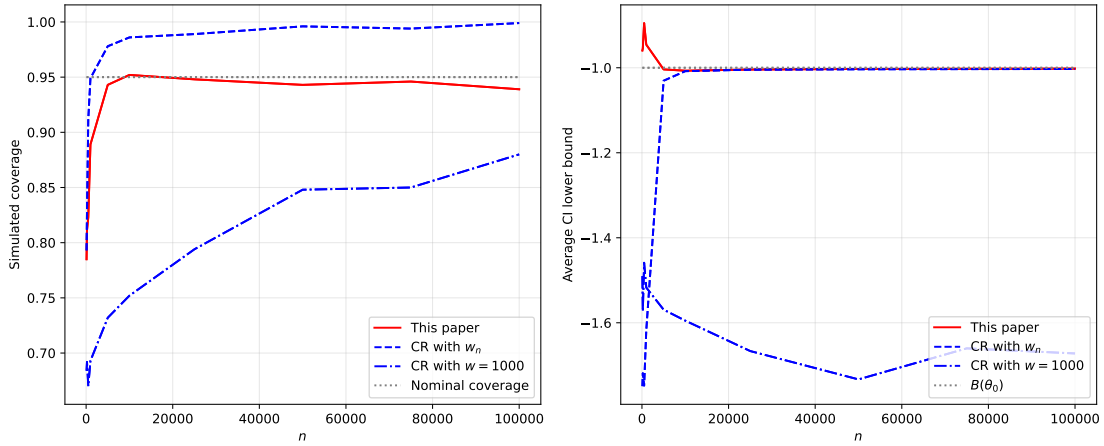
---

<sup>15</sup>Because the penalty function estimator can be rewritten in the form of an equivalent problem w.p.a.1.





(a)  $b = -0.05$ . For BG, the number of failed simulations was 206, 158, 46, 6 at  $n = 100, 200, 500, 1000$  respectively, none at larger  $n$ .



(b)  $b = 0$ . BG omitted, as around 250 simulations fail.

Figure 9: Performance of different inferential procedures over 1000 simulations of (17). Left panel - estimated coverage of a 95% one-sided C.I.; right panel - average lower confidence bound.

**Discussion** Results of our simulations are given in Figure 9. Overall, it appears that our estimator has the correct nominal level even at smaller sample sizes, whereas the CR with penalty  $w_n$  overcovers asymptotically. BG estimator achieves nominal coverage asymptotically, although it over-covers in smaller samples and may yield a very conservative left confidence bound, as is evident from the right panel. It fails at  $b = 0$ , because the SC fails<sup>16</sup>. For different fixed values of  $w$  the Cho-Russell procedure's performance may range from highly conservative to invalid. We illustrate this by adding lines with  $w = 2$  and  $w = 1000$  to the figures corresponding to  $b = -0.05$  and  $b = 0$  respectively.

<sup>16</sup>We have still run the corresponding simulations, but are not displaying them. BG fails to exist in around 25% of cases, while the remaining simulations result in highly conservative bounds with incorrect coverage

### 3. Special case of LP bounds

Let  $Y \in \mathbb{R}$  denote the outcome of interest<sup>17</sup>,  $T \in \mathbb{R}$  stand for the treatment, and  $Z \in \mathbb{R}^{dz}$  be the candidate instrument. Here *treatment* is any variable which effect on  $Y$  we attempt to infer, whereas the term *instrument* refers to an auxiliary variable that allows us to partially identify the treatment effect of interest. Our approach nests the case when  $Z$  is the usual IV. The reader seeking an economic intuition may refer to the classical example from Manski and Pepper (2000), where  $Y$  is the wage,  $T$  is an indicator of educational degree and  $Z$  is the level of ability.

Denote the supports  $Y, T, Z$  as  $\mathcal{Y}, \mathcal{T}$  and  $\mathcal{Z}$  respectively. Throughout this section we consider the case of continuous outcomes and discrete treatment and instrument, namely  $\mathcal{Y}$  is uncountable, while  $N_T \equiv |\mathcal{T}| < \infty$  and  $N_Z \equiv |\mathcal{Z}| < \infty$ . In non-parametric bounds literature it is rather conventional to employ a discrete instrument at the estimation stage (see Manski and Pepper (2009)). While the main identification result could be extended to continuous  $Z$ , we make the discreteness assumption early on to avoid unnecessary technical complications.

Our setup accommodates missing observations of the dependent variable. Namely, we split the set of treatments into two disjoint subsets  $\mathcal{T} = \mathcal{O} \sqcup \mathcal{U}$ . Whenever  $T \in \mathcal{O}$ , the researcher observes  $Y, T, Z$ , whereas if  $T \in \mathcal{U}$ , only the covariates  $T, Z$  are observed. For example, in Blundell et al. (2007) the wage is observed only if an individual is employed. Corresponding to the legs of the treatment are the potential outcomes  $Y(t)$ ,  $t \in \mathcal{T}$ :

$$Y = \sum_{t \in \mathcal{O}} \mathbb{I}\{T = t\}Y(t) + \sum_{t \in \mathcal{U}} \mathbb{I}\{T = t\}Y(t)$$

Continuing the wages and education example, the value of  $Y(t)$  for a fixed  $t \in \mathcal{T}$  may then correspond to the potential wage that an individual with the associated random characteristics would get, had she obtained education  $t$ .

Let us collect the potential outcomes in the vector  $\mathbb{Y} \equiv (Y(t))_{t \in \mathcal{T}} \in \mathbb{R}^{N_T}$ . Variables  $(\mathbb{Y}, T, Z, Y)$  are jointly defined on the true probability space  $(\mathbb{P}, \Omega, \mathcal{S})$  and we let  $\mathcal{P}$  denote the considered collection of probability measures on  $(\Omega, \mathcal{S})$ , such that  $\mathbb{P} \in \mathcal{P}$ . We impose the following conditions on the set of considered measures throughout this section:

**Assumption I0 (Conditions on  $\mathcal{P}$ ).**  $\mathcal{P}$  is such that  $P \in \mathcal{P}$  if: i)  $P$  generates  $F_{T,Z}(\cdot)$  and  $\{F_{Y|T=t,Z}(\cdot)\}_{t \in \mathcal{O}}$ ; ii)  $P[T = d, Z = z] > 0 \forall d, z \in \mathcal{T} \times \mathcal{Z}$  and iii)  $|\mathbb{E}_P[Y(t)|T = d, Z = z]| < \infty$  for all  $z \in \mathcal{Z}$  and  $t, d \in \mathcal{T}$

Part i) of the Assumption I0 formalizes the assumed identification pattern. It says that the joint distribution of  $T, Z$  is always identified and the researcher also observes the joint

---

<sup>17</sup>Univariate case is considered for simplicity of exposition, but the extension to multivariate outcomes is immediate.

distribution of  $Y, T, Z$  whenever  $T \in \mathcal{O}$ . Parts ii) and iii) of I0 ensure that all conditional expectations and probabilities are well-defined and finite<sup>18</sup>.

**Remark 3.1.** Under no missing data, i.e.  $\mathcal{T} = \mathcal{O}$ , condition i) is equivalent to  $P$  generating the identified joint distribution  $F_{Y,T,Z}(\cdot)$ .

We define the vector  $m$  collecting all elementary conditional moments as

$$m(P) \equiv (\mathbb{E}_P[\mathbb{Y}|T = d, Z = z])_{d \in \mathcal{T}, z \in \mathcal{Z}} \in \mathbb{R}^{N_T^2 N_Z},$$

and suppose that the researcher is interested in the target parameter  $\beta^*$ , given by

$$\beta^* = \mu^*(\mathbb{P})' m(\mathbb{P}), \quad (18)$$

where  $\mu^* : \mathcal{P} \rightarrow \mathbb{R}^{N_T^2 N_Z}$  is *identified* and *chosen* by the researcher. It parametrizes the choice of the outcome of interest, as the following remark clarifies.

**Remark 3.2.** The form (18) nests i)  $\mathbb{E}[Y(t)]$ , ii)  $ATE_{td} = \mathbb{E}[Y(t) - Y(d)]$  and iii)  $CATE_{td,A,B} = \mathbb{E}[Y(t) - Y(d)|T \in A, Z \in B]$ .

### 3.1. Affine inequalities over conditional moments

We now introduce the general class of identifying conditions described by affine inequalities over conditional moments, potentially augmented with affine a.s. restrictions. These restrict the set of admissible measures to  $\mathcal{P}^*$ :

$$\mathcal{P}^* \equiv \{P \in \mathcal{P} | (M^* m + b^*)(P) \geq 0 \wedge (\tilde{M} \mathbb{Y} + \tilde{b})(P) \geq 0 \text{ } P\text{-a.s.}\}, \quad (19)$$

where  $b^* : \mathcal{P} \rightarrow \mathbb{R}^R$ ,  $M^* : \mathcal{P} \rightarrow \mathbb{R}^{R \times N_T^2 N_Z}$ ,  $\tilde{b} : \mathcal{P} \rightarrow \mathbb{R}^{\tilde{R}}$  and  $\tilde{M} : \mathcal{P} \rightarrow \mathbb{R}^{\tilde{R} \times N_T}$  are *identified* parameters, *chosen* by the researcher. These parametrize the choice of  $R \in \mathbb{N}$  identifying inequalities on conditional moments of potential outcomes as well as  $\tilde{R} \in \mathbb{N}$  almost sure inequalities on the potential outcomes. In general,  $\psi(P) \equiv (\mu^*(P), b^*(P), M^*(P), \tilde{b}(P), \tilde{M}(P))$  and  $m(P)$  are functionals of  $P$ . We omit this dependence whenever it does not cause confusion.

The family of models that can be written in the form (19) is very rich, as illustrated by the following examples.

**Example 3.1.** MIV with  $Z \in \mathbb{R}$  (Manski and Pepper, 2000) imposes that for each  $t \in \mathcal{T}$  and  $z, z' \in \mathcal{Z}$ ,  $z' \geq z \implies \mathbb{E}[Y(t)|Z = z'] \geq \mathbb{E}[Y(t)|Z = z]$ . It is nested for an appropriate choice of matrix  $M^* = M_{MIV}$  and  $b^* = 0$ . MTS from (Manski and Pepper, 2000) obtains when  $Z = T$ .

---

<sup>18</sup>Similar identification results can still be obtained if one relaxes the full-support condition for some known pairs from  $\mathcal{Z} \times \mathcal{T}$ . Note that it can also be verified in the data.

**Example 3.2.** IV with  $Z \in \mathbb{R}$  imposes that for each  $t \in \mathcal{T}$  and  $z, z' \in \mathcal{Z}$ ,  $\mathbb{E}[Y(t)|Z = z'] = \mathbb{E}[Y(t)|Z = z]$ . It is nested for an appropriate choice of matrix  $M_{IV}$  that can, for example, be constructed as  $M^* = M_{IV} = \begin{pmatrix} M_{MIV} \\ M_{MIV} \end{pmatrix}$  and  $b^* = 0$ .

**Example 3.3.** MTR (Manski and Pepper, 2000) imposes that for each  $t, t' \in \mathcal{T}$ :  $t' > t$ ,  $Y(t') \geq Y(t)$  a.s. It is nested for an appropriate choice of matrix  $\tilde{M} = \tilde{M}_{MTR}$  with  $\tilde{b} = 0$ .

**Example 3.4.** Roy model (Laffers, 2019) imposes that for each  $t \in \mathcal{T}$ , the individual's choice is, on average, optimal  $\mathbb{E}[Y(t)|T = t, Z = z] = \max_{d \in \mathcal{T}} \mathbb{E}[Y(d)|T = t, Z = z]$ . It is nested for an appropriate choice of matrix  $M^* = M_{ROY}$  and  $b^* = 0$ .

**Example 3.5.** Missing data. Blundell et al. (2007) derives bounds on  $F(w|x)$  - the cdf of wages evaluated at some  $w$ , with the wages observed if the individual is employed,  $E = 1$ , and unobserved otherwise, if  $E = 0$ . Introduce  $\mathcal{O} = \{1\}$  and  $\mathcal{U} = \{0\}$ . Let  $Y(t) \equiv \mathbb{I}\{W \leq w\}$ , so that  $\mathbb{E}[Y(t)|X = x] = F(w|x)$ . Our approach allows to accommodate all identifying conditions in the original paper by appropriately choosing  $M^*, b^*$  and  $\tilde{M}, \tilde{b}$ .

**Remark 3.3.** Combinations of assumptions are obtained by stacking the respective matrices, as in Example 3.2. Sensitivity analysis can be performed via relaxations  $b_\ell^* = b^* + \ell$ , or  $\tilde{b}_\ell = \tilde{b} + \ell$  for some  $\ell \geq 0$ . For example, given some  $\ell = \{\ell(t, z, z')\}_{t, z, z'} \geq 0$ , inequalities  $\mathbb{E}[Y(t)|Z = z'] - \mathbb{E}[Y(t)|Z = z] \geq -\ell(t, z, z')$  for  $z, z' \in \mathcal{Z}$  with  $z' > z$  and  $t \in \mathcal{T}$  yield a relaxation of MIV. In De Haan (2017) the shape of observed moments may suggest a failure of monotonicity near the boundaries of  $\text{Supp}(Z)$ . Selecting positive  $\ell(z, z')$  for values of  $z, z'$  close to the boundaries could constitute a meaningful robustness check.

### 3.2. Linear programming bounds

We now provide the general identification result for the models described by  $\mathcal{P}^*$ . Let us construct  $x$  that collects unobserved pointwise-conditional moments and the vector  $\bar{x}$  of those pointwise-conditional moments that are identified:

$$\begin{aligned} x &\equiv (\mathbb{E}[Y(t)|T = d, Z = z])_{z \in \mathcal{Z} \wedge (t, d \in \mathcal{T}: t \neq d \vee t, d \in \mathcal{U}: t = d)}, \\ \bar{x} &\equiv (\mathbb{E}[Y(t)|T = t, Z = z])_{z \in \mathcal{Z}, t \in \mathcal{O}}. \end{aligned}$$

For *known* selector matrices  $P_m, \bar{P}_m$ , one can then decompose  $m$  as

$$m = \underbrace{P_m x}_{\text{partially identified}} + \underbrace{\bar{P}_m \bar{x}}_{\text{identified}}.$$

It is also straightforward to observe that  $\tilde{M}Y + \tilde{b} \geq 0$  a.s. implies

$$(I_{N_T N_Z} \otimes \tilde{M})m + \iota_{N_T N_Z} \otimes \tilde{b} \geq 0. \tag{20}$$

Before we state the main identification result, let us construct the matrix  $M^{**}$  and the vector  $b^{**}$  that combine the conditional restrictions with the implications of the almost sure restrictions in (20), as

$$M^{**} \equiv \begin{pmatrix} I_{N_T N_Z} \otimes \tilde{M} \\ M^* \end{pmatrix}, \quad b^{**} \equiv \begin{pmatrix} \iota_{N_T N_Z} \otimes \tilde{b} \\ b^* \end{pmatrix}. \quad (21)$$

**Theorem 3.1.** Suppose Assumption I0 holds. For any  $\psi$ , the sharp identified set  $\mathcal{B}^*$  for  $\beta^*$  satisfies

$$\mathcal{B}^* \equiv (\mu^{*'} m)(\mathcal{P}^*) \subseteq \{\beta \in \mathbb{R} \mid \inf_{x: Mx \geq b} p'x \leq \beta - \bar{p}'\bar{x} \leq \sup_{x: Mx \geq b} p'x\}, \quad (22)$$

where

$$\bar{p} \equiv \bar{P}'_m \mu^*, \quad p \equiv P'_m \mu^*, \quad M \equiv M^{**} P_m, \quad b \equiv -b^{**} - M^{**} \bar{P}_m \bar{x}.$$

The converse inclusion in (22) holds if  $\tilde{M} = 0'_{N_T}$ ,  $\tilde{b} = 0$ , or:

1. MTR holds, i.e.  $Y(t_1) \geq Y(t_0)$  a.s.  $\forall t_1, t_0 \in \mathcal{T}$  s.t.  $t_1 > t_0$ :

$$\tilde{M} = \tilde{M}_{MTR}, \quad \tilde{b} = 0_{N_T-1}$$

2. Outcomes are bounded,  $Y(t) \in [K_0; K_1]$ ,  $\forall t \in \mathcal{T}$  a.s. for known  $K_1 > K_0$ :

$$\tilde{M} = \tilde{M}_b \equiv \begin{pmatrix} I_{N_T} \\ -I_{N_T} \end{pmatrix}, \quad \tilde{b} = \tilde{b}_b \equiv \begin{pmatrix} -K_0 \cdot \iota_{N_T} \\ K_1 \cdot \iota_{N_T} \end{pmatrix}. \quad (23)$$

3. MTR holds, outcomes are bounded and  $(\Omega, \mathcal{S})$  can support a  $U[0; 1]$  r.v.:

$$\tilde{M} = \begin{pmatrix} \tilde{M}_{MTR} \\ \tilde{M}_b \end{pmatrix}, \quad \tilde{b} = \begin{pmatrix} \tilde{b}_{MTR} \\ 0_{N_T-1} \end{pmatrix}. \quad (24)$$

The matrix  $\tilde{M}_{MTR}$  is defined in the Appendix.

Theorem 3.1 postulates that bounds on the target parameter  $\beta^*$  under  $\mathcal{P}^*$  can be obtained by solving two linear programs. The LP bounds are sharp if there are no a.s. inequalities in the model, or if the a.s. inequalities parametrize three special cases that are typically used in the literature. Otherwise, the LP bounds may not be sharp, as we show in the Appendix. This is because, in general, the entire distribution of  $\mathbb{Y}, T, Z$  is relevant for  $\beta^*$  under a.s. restrictions. The naive approach of searching over such joint distributions, however, would

involve infinite-dimensional optimization, because  $|\mathcal{Y}| = \infty$ .

**Remark 3.4.** We are not aware of affine a.s. restrictions  $\tilde{M}, \tilde{b}$  used in applied work that are not special cases 1-3 in Theorem 3.1. The special cases 1-3 appear in Blundell et al. (2007), Kreider et al. (2012), Gundersen et al. (2012), Siddique (2013).

**Remark 3.5.** The empirical literature has extensively relied on the MIV + MTR + MTS combination of Manski and Pepper (2000) assumptions, as it yields the tightest bounds out of all classical conditions. In the absence of a theoretical justification, this has led to errors (Laffers, 2013). Theorem 3.1 provides the first available sharp bounds under this combination when  $|\mathcal{Y}| = \infty$ .

## 4. Conditional monotonicity assumptions

A particular family of identifying conditions that can be written in the form (19) is the *conditional monotonicity* class of assumptions. These impose that potential outcomes are mean-monotone in the instrument even within some treatment subgroups. While more restrictive than the conventional MIV, conditionally monotone instrumental variables (cMIV) allow to sharpen the bounds on the outcomes of interest. Throughout this section, we assume that outcomes are bounded  $Y(t) \in [K_0; K_1]$  a.s. for known  $K_0, K_1 \in \mathbb{R}$ ,  $K_0 < K_1$ . We also suppose that there are no missing data<sup>19</sup>, i.e.  $\mathcal{T} = \mathcal{O}$ .

We argue that cMIV assumptions are reasonable in classical applications, discuss the difference between MIV and cMIV and develop a formal testing strategy for a particular version of cMIV. This testing procedure relies on the observed outcomes' monotonicity, which has been typically used in applied work to justify applying MIV. Our results imply that if such monotonicity is observed and the researcher is comfortable with MIV, the cMIV assumption is *inexpensive*, and can be applied to sharpen the bounds on the outcomes of interest. In some applications, as is the case in Section 5, cMIV yields informative bounds even when the classical conditions fail to do so.

While we only discuss three variations of cMIV, the class of such assumptions is potentially richer<sup>20</sup>, and Theorem 3.1 applies in any such framework.

**Assumption cMIV-s.** Suppose that for any  $t \in \mathcal{T}$ ,  $A \subseteq \mathcal{T} : A \neq \{t\}$  and  $z, z' \in \mathcal{Z}$  s.t.  $z' > z$  we have:

$$\mathbb{E}[Y(t)|T \in A, Z = z'] \geq \mathbb{E}[Y(t)|T \in A, Z = z],$$

<sup>19</sup>Although it is hopefully clear from our general approach how cMIV conditions extend to the missing data case.

<sup>20</sup>One can consider the class of conditional restrictions  $\mathbb{E}[Y(t)|T \in A, Z = z'] \geq \mathbb{E}[Y(t)|T \in A, Z = z]$ ,  $\forall A \in \mathcal{F}_t$  for all  $t \in T$  where subcollections  $\mathcal{F}_t \subseteq \mathcal{T}$  are chosen by the researcher.

*i.e. the potential outcomes are, on average, non-decreasing in  $Z$  for any treatment subgroup.*

The strong conditional monotonicity assumption possesses the greatest identifying power across all considered cMIV conditions. To see that cMIV-s implies MIV, set  $A = \mathcal{T}$  in the definition above.

**Assumption cMIV-w.** *Suppose MIV holds and for any  $t \in \mathcal{T}$  and  $z, z' \in \mathcal{Z}$  s.t.  $z' > z$  we have:*

$$\mathbb{E}[Y(t)|T \neq t, Z = z'] \geq \mathbb{E}[Y(t)|T \neq t, Z = z],$$

*i.e. the potential outcomes are, on average, non-decreasing in  $Z$  for the non-treated and the whole population.*

The weak conditional monotonicity assumption allows for closed-form expressions for sharp bounds that are easy to compute and perform inference on, see the Appendix.

**Assumption cMIV-p.** *Suppose MIV holds and for any  $t \in \mathcal{T}, d \in \mathcal{T} \setminus \{t\}$  and  $z, z' \in \mathcal{Z}$  s.t.  $z' > z$  we have:*

$$\mathbb{E}[Y(t)|T = d, Z = z'] \geq \mathbb{E}[Y(t)|T = d, Z = z],$$

*i.e. the potential outcomes are, on average, non-decreasing in  $Z$  conditional on any counterfactual level of treatment.*

The pointwise conditional monotonicity assumption is directly testable under a mild homogeneity condition, see Section 4.2.

Conditional monotonicity restrictions differ in the collection of treatment subsets over which monotonicity in the instrument is assumed. The strong conditionally monotone instruments are such that, among individuals from any given counterfactual treatment subgroup, higher values of  $Z$  are, on average, associated with higher potential outcomes. The weak conditional monotonicity restriction only imposes the same mean-monotonicity on the whole population and on the untreated, whereas the pointwise form assumes it over the entire population as well as conditional on each counterfactual level of treatment.

**Remark 4.1.** All cMIV assumptions imply MIV. Moreover, cMIV-w, cMIV-p are implied by cMIV-s. If treatment is binary, cMIV-s, cMIV-w and cMIV-p are equivalent.

While it is possible for the general approach of form (19), cMIV conditions avoid assuming monotonicity over the observed treatment subset  $\{T = t\}$ . This is because such monotonicity is identified. If it holds, it should not add any identifying power to our conditions in theory. On the other hand, large violations of the observed outcomes' monotonicity will lead the test developed in Section 4.2 to reject cMIV-p and cMIV-s.

The following observation motivates the use of cMIV assumptions.

**Proposition 4.1.** *Manski and Pepper (2000) MIV bounds are not sharp under either cMIV-s, cMIV-w or cMIV-p.*

*Proof.* Consider a binary treatment  $T$ , three levels of the instrument  $Z \in \{z_0, z_1, z_2\}$  with  $z_0 < z_1 < z_2$  and  $-K_0 = K_1 = 1$ . Suppose for a fixed  $t \in \{0, 1\}$ , we have  $\mathbb{E}[Y(t)|T = t, Z = z_i] = 0$ , with  $P[T \neq t|Z = z_0] = 0.125$ ,  $P[T \neq t|Z = z_1] = 0.5$ ,  $P[T \neq t|Z = z_2] = 0.25$ . The no-assumptions lower bounds on  $\mathbb{E}[Y(t)|Z = z_i]$  are  $(-0.125, -0.5, -0.25)$ . MIV ‘irons’ the no-assumptions bounds to  $(-0.125, -0.125, -0.125)$ , which also implies the lower bounds on  $\mathbb{E}[Y(t)|T \neq t, Z = z_i]$ :  $(-1, -0.25, -0.5)$ . Under cMIV, one can further ‘iron’ these to improve the lower bound for  $z_2$  up to  $-0.25$ , so that the lower bound on  $\mathbb{E}[Y(t)|Z = z_2]$  becomes  $-1/16 > -1/8$ . ■

Sharp bounds for all versions for cMIV follow from Theorem 3.1. We also show that under cMIV-w the bounds can be characterized explicitly, which is especially convenient if the treatment is binary, so that all cMIV assumptions coincide. For didactic purposes, we provide the detailed construction of the triplet  $M, c, p$  from Theorem 3.1 under cMIV-s and cMIV-p. All details on the identification under cMIV are provided in the Appendix.

#### 4.1. Discussion of cMIV

This section illustrates the difference between MIV and cMIV by considering two parametric examples with classical applications.

**4.1.a. Education selection.** Consider the following empirical setup. Suppose  $T$  is an indicator of whether or not an individual has a university degree,  $Y(t)$  are potential log wages and  $Z$  is an observed indicator of ability.

MIV assumption on  $Z$  implies that more able individuals can do better both with and without a college degree on average:  $\mathbb{E}[Y(t)|Z = z]$  - monotone in  $z$ . cMIV additionally imposes that: i) among those who have a college degree, a *smarter* individual could have done relatively better on average than their counterpart if both did not have it:  $\mathbb{E}[Y(0)|Z = z, T = 1]$  - monotone in  $z$ ; and ii) among those who do not have a college degree, a *smarter* individual could have done relatively better on average than their counterpart if both had it:  $\mathbb{E}[Y(1)|Z = z, T = 0]$  - monotone in  $z$ .

We now consider a parametric example. Suppose that  $\eta$  measures how *diligent* one is from birth and is ex-ante mean-independent of  $Z$ . While  $Z$  is observed by both the employers and the econometrician (e.g. an IQ score), the employer additionally observes the employee effort level  $\eta + \varepsilon$  with  $\varepsilon \perp\!\!\!\perp (Z, T, \eta)$ . Suppose  $Var(Z) = Var(\eta) = 1$  and  $\mathbb{E}[Z] = \mathbb{E}[\eta] = \mathbb{E}[\varepsilon] = 0$ . Suppose that, on average, employees choose  $T$  to maximize their expected earnings. This



motivates a stylized Roy selection model with:

$$Y(t) = \beta_0(t) + \beta_1(t)Z + \beta_2(t)\eta + \varepsilon(t), \quad T = \mathbb{1}\{\mathbb{E}[Y(1) - Y(0)|Z, \eta] + \nu \geq 0\},$$

where  $\nu \perp\!\!\!\perp (Z, \eta, \varepsilon)$  is remaining heterogeneity, and  $\varepsilon(t) \equiv \beta_2(t)\varepsilon$ . MIV demands that:

$$\text{(MIV)} : \beta_1(t) \geq 0$$

MIV postulates that the direct effect of ability on potential earnings is positive. It seems reasonable to suppose that  $\beta_i(t) \geq 0$ ,  $i = 1, 2$ ,  $t = 0, 1$ , i.e. both effort and ability increase potential wages. Letting  $\delta_z \equiv \beta_1(1) - \beta_1(0)$  and  $\delta_\eta \equiv \beta_2(1) - \beta_2(0)$  denote the differentials in the effects of ability and effort respectively, the additional requirement of cMIV is that:

$$\underbrace{\beta_1(0)z}_{\text{direct effect}} + \underbrace{\beta_2(0)\mathbb{E}[\eta|\delta_z z + \delta_\eta \eta + \tilde{\nu} \geq 0]}_{\text{selection given } T = 1} \text{--increasing} \quad (25)$$

$$\underbrace{\beta_1(1)z}_{\text{direct effect}} + \underbrace{\beta_2(1)\mathbb{E}[\eta|\delta_z z + \delta_\eta \eta + \tilde{\nu} \leq 0]}_{\text{selection given } T = 0} \text{--increasing,} \quad (26)$$

where  $\tilde{\nu} \equiv \beta_0(1) - \beta_0(0) + \nu$ .

Notice that if  $\delta_z$  and  $\delta_\eta$  are of different signs, for example because the jobs that one may apply for with a college degree are more ability-intensive ( $\delta_z > 0$ ), whereas those which are available otherwise are more skill-intensive ( $\delta_\eta < 0$ ), the additional conditional monotonicity requirements (25)-(26) are less strict than MIV. This is because, *conditional* on both having a degree and not having it, ability and effort are *positively* associated.

Intuitively, among those who do not have a degree ( $T = 0$ ), people of higher ability must have had stronger incentives to forgo college. This should have been because a higher level of diligence gives them a comparative advantage in effort-intensive jobs. Among those with a degree, higher ability implies a comparative advantage in ability-intensive occupations, which explains their willingness to select into this option ( $T = 1$ ). It does not, therefore, signal as low an effort level as it would for a less capable individual.

Now consider the same setup with<sup>21</sup>:

$$T = \mathbb{I}\{\eta + Z \geq 0\}$$

This selection mechanism can be explained by the fact that to get a degree one needs to be either hard-working or of high ability. The requirement of MIV is unchanged, and cMIV

---

<sup>21</sup>Setting  $\delta_z = \delta_\eta > 0$  and  $\tilde{\nu} = 0$ .

necessitates that:

$$\beta_1(0)z + \beta_2(0)\mathbb{E}[\eta|\eta \geq -z] - \text{increasing} \quad (27)$$

$$\beta_1(1)z + \beta_2(1)\mathbb{E}[\eta|\eta \leq -z] - \text{increasing} \quad (28)$$

In this case, conditional on each level of education, effort level  $\eta$  and ability  $Z$  are negatively associated, so the conditional selection terms in (27)-(28) make cMIV a stricter assumption than MIV. Intuitively, a more able individual with a college degree did not need to work as hard to get it as her counterpart with a lower ability. Similarly, if an individual is capable, but does not have a degree, she has to be of lower effort as otherwise she would have selected into education.

Even if MIV holds, cMIV can thus fail if employer prefers effort over ability to the extent that the conditional negative association between the two outweighs the direct impact of ability on wages as well as any ex-ante positive correlation between the employer-observed signal of diligence and the ability.

An examination of equations (25) and (26) suggests that cMIV is more likely to hold whenever  $\delta_z$  is small relative to  $\delta_\eta$ , while  $\beta_1(\cdot)$  is large relative to  $\beta_2(\cdot)$ . This means that  $Z$  should be *relatively weak* in the parlance of the classical IV models, and *strongly monotone*.

Overall, it seems reasonable to use a proxy for the level of ability as a conditionally monotone instrument in the estimation of returns to schooling. One would be inclined to think that while  $Z$  does enter selection, it affects the potential outcomes directly and strongly enough, so that there are no subgroups by schooling for which a higher value of ability would correspond to lower potential wages on average.

**4.1.b. Simultaneous equations.** As some aspects of mathematical intuition may be muted in discrete models, we also consider a simple continuous setup to confirm the insights derived from the previous analysis. For illustrative purposes, drop the boundedness and discreteness assumptions and consider the demand and supply simultaneous equations:

$$q^k(p) = \alpha^k(p) + \beta^k(p)Z + \gamma^k(p)\eta + \kappa^k(p)\varepsilon^k, \quad k \in \{s, d\}$$

The observed log-price  $P$  clears the market:

$$P \in \{p \in \mathbb{R} | \mathbb{E}[q^s(p)|Z, \eta] = \mathbb{E}[q^d(p)|Z, \eta]\}, \quad (29)$$

where  $\eta, Z$  are continuous unobserved and observed random variables respectively, with  $\mathbb{E}[\eta|Z = z] = 0$ <sup>22</sup> and  $\mathbb{E}[\varepsilon^k] = 0$  with  $\varepsilon^k \perp\!\!\!\perp (\eta, Z, \varepsilon^{-k})$  for  $k \in \{s, d\}$ . Further assume that all functions of  $p$  are continuous.

---

<sup>22</sup>Note that, once again, mean independence is not restrictive, as otherwise we could always redefine the data generating process in an observationally equivalent way.

Potential price  $p$  indexes the potential outcomes, giving rise to the demand and supply *schedules*. Suppose we aim to identify the elasticity of supply,  $\mathbb{E}[\frac{q^s(p_1) - q^s(p_0)}{p_1 - p_0}]$  for some  $p_1 > p_0$ , and  $Z$  is a monotone instrument for  $q^s(p)$ , while  $P$  can be interpreted as treatment.  $\eta$  is unobserved heterogeneity and  $\varepsilon^k$  are random violations from the market clearing condition or measurement errors independent of the rest of the model. For an individual realization of market clearing an econometrician observes  $\{P, \{q^k(P)\}_k, Z\}$ , but does not observe the schedules at other prices  $\{q^k(p)\}_k$  for  $p \neq P$ , nor disturbances  $\{\eta, \{\varepsilon^k\}_k\}$ .

Define  $\delta_z(p) \equiv \beta^s(p) - \beta^d(p)$  and similarly for  $\eta$ , with  $\delta_p(p) \equiv \alpha^s(p) - \alpha^d(p)$ . As stated, the model is potentially *incomplete* or *incoherent*, as for a given vector  $(Z, \eta)$  equation (29) may have multiple or no solutions. To avoid that, so long as that the support of  $Z, \eta, \varepsilon^k$  is full, it is necessary that  $\delta_z(p), \delta_\eta(p)$  be constant. We shall assume that for simplicity. Provided that  $\delta_p(p)$ , which determines the *excess supply* at fixed  $(Z, \eta)$ , is strictly increasing and has full image, the model is *complete* and *coherent* and:

$$P = \delta_p^{-1}(-\delta_z Z - \delta_\eta \eta) \quad (30)$$

Equation (30) introduces a deterministic linear relationship between  $Z$  and  $\eta$  conditional on each given value of  $P$ . As we saw in the previous example, this constitutes the worst-case scenario for cMIV, if  $\delta_z$  and  $\delta_\eta$  have the same sign. A noisier selection mechanism would relax the conditional link between  $Z$  and  $\eta$ , and would thus weaken the conditional selection channel.

Note that the reduced-form error is  $u \equiv \gamma^s(P)\eta + \kappa^s(P)\varepsilon^s$  and there is a simultaneity bias:

$$\mathbb{E}[Pu] = \mathbb{E}[P\gamma^s(P) \underbrace{\mathbb{E}[\eta|\delta_z Z + \delta_\eta \eta = P]}_{\text{simultaneity/omitted variable}}] \neq 0$$

In this setup, MIV requires:

$$(MIV) : \beta^s(p) \geq 0, \forall p \in \mathbb{R}$$

Whereas cMIV-p additionally imposes that:

$$\beta^s(p)z + \gamma^s(p)\mathbb{E}[\eta|\delta_z z + \delta_\eta \eta = -\delta_p(d)] \geq 0 - \text{increasing in } z, \forall p, d \in \mathbb{R} : d \neq p \quad (31)$$

Suppose that  $\delta_z, \delta_\eta \neq 0$  to rule out uninteresting cases. (31) rewrites as:

$$\beta^s(p) \geq \gamma^s(p) \frac{\beta^s(p) - \beta^d(p)}{\gamma^s(p) - \gamma^d(p)} \quad (32)$$

For concreteness, consider two positive supply shocks, i.e.  $\beta^s(p), \gamma^s(p) > 0$ . Equation (32)

then says that either  $\eta$  and  $Z$  affect the reduced-form equilibrium price in different directions (recall the comparative advantage example), or the effect of  $Z$  on the equilibrium price relative to its effect on the supply schedule is smaller than that of  $\eta$ :

$$\text{sgn}(\delta_\eta) \neq \text{sgn}(\delta_z) \quad \text{or} \quad \left| \frac{\beta^s(p) - \beta^d(p)}{\beta^s(p)} \right| \leq \left| \frac{\gamma^s(p) - \gamma^d(p)}{\gamma^s(p)} \right| \quad (33)$$

Under  $\text{sgn}(\delta_\eta) = \text{sgn}(\delta_z)$ , equation (33) once again requires that  $Z$  be *strongly monotone* and *relatively weak*. The logic we described may help the researcher navigate the potential economic forces in a given application to decide whether cMIV-p is a suitable assumption.

For example, consider estimating the supply elasticity in the market for plane tickets in the early days of Covid-19 pandemic. Suppose  $Z$  is an inverse Covid-stringency index for the economy, while  $\eta$  may be interpreted as residual cost shocks, defined to be mean-independent of  $Z$ . It is likely that  $\delta_\eta \approx \gamma^s$ , i.e. residual cost shocks affect mainly the supply in that sector, and not the demand. It is also likely that either supply is less responsive to  $Z$  than demand (so that cMIV is implied by MIV), or the effects are of the same order of magnitude.  $Z$  is therefore likely to be a conditionally monotone instrument.

## 4.2. Testing cMIV

One could argue against cMIV conditions whenever  $\mathbb{E}[Y(t)|T = t, Z = z]$  fail to be monotone in the data. In general, the power and size of that test are unclear. There is, however, a special case when cMIV can be tested directly, given that the researcher believes in MIV. In some applications one may conjecture that the potential outcomes' functions  $Y(t)$ , either in the reduced or in the structural form, are such that the relative effects of  $Z$  and the unobserved variable(s)  $\eta$ , potentially correlated with  $Z$ , are unchanged across outcome indices  $t$ .

Researchers often impose even stricter versions of this homogeneity assumption. For example, Manski and Pepper (2009) discuss MIV identification under HLR condition:  $Y(t) = \beta t + \eta$ . Conditions in Proposition 4.2 relax HLR to an arbitrary shape of response of a potential outcome to treatment and allow for a generally heteroscedastic/treatment-specific response to unobserved variables and instrument, so long as the relative effects are unchanged across potential outcomes.

**Proposition 4.2.** *Suppose that a): i)  $Y(t) = g(t, \xi) + h(t)\psi(Z, \eta)$ ,  $h(t) \neq 0$  with  $\xi \perp\!\!\!\perp (T, Z, \eta)$  and ii) MIV holds, strictly for some  $z, z' \in \mathcal{Z}$  with  $z' > z$ ; or b): i)  $Y(t) = g(t, \xi, T) + h(t)\psi(Z, \eta)$  with  $\xi \perp\!\!\!\perp (T, Z, \eta)$ , ii)  $\frac{h(t)}{h(d)} > 0 \forall t, d \in \mathcal{T}$ ; and iii) MIV holds. Then Assumption cMIV-p holds iff  $\mathbb{E}[Y(t)|T = t, Z = z]$  are all monotone.*

Note that whether or not  $h(t) \neq 0$  is observable in the data for case (a) and whether or not  $h(t)/h(d) > 0$  is also identified for (b).

**Remark 4.2.** The monotonicity of observed outcomes has been routinely used in applied work to motivate the use of MIV condition (e.g. De Haan (2017)). We show that, given that MIV holds and under a homogeneity condition, the observed monotonicity is instead equivalent to cMIV-p.

**Remark 4.3.** cMIV is testable in the Example 3.2.2, because the reduced form expression has the form  $b) : i)$ . It also becomes testable in the Example 3.2.1 if instead of separately observing  $\eta, Z$ , employers on average observe a mixed signal of ability and effort,  $s \equiv aZ + b\eta$  for some  $a, b \in \mathbb{R}$ .

A test of cMIV-p is thus the test of all  $f_t(z) \equiv \mathbb{E}[Y(t)|T = t, Z = z]$  being monotone:

$$\mathcal{H}_0 : f_t(z) - \text{increasing in } z, \forall t \in \mathcal{T}$$

To obtain such a test, we may extend the procedure in Chetverikov (2019)<sup>23</sup>. Denote the set of all observations with treatment level  $t$  as  $\mathcal{I}_t \equiv \{i \in \overline{1, n} : T_i = t\}$  with  $n_t \equiv |\mathcal{I}_t|$ . Suppose  $\phi_{n_t}^t$  is the corresponding Chetverikov's regression monotonicity test (or a corresponding parametric test for discrete  $Z$ ) with the confidence level  $\alpha_t \in (0; 0.5)$ . We define the joint test as:

$$\phi_n \equiv \max_{t \in \mathcal{T}} \phi_{n_t}^t$$

Denote  $\mathcal{P}^C$  to be the set of probability measures, such that for all  $P \in \mathcal{P}^C$  and all  $t \in \mathcal{T}$  the conditional probability measure given  $T = t$  that  $P$  generates satisfies the regularity conditions in Theorem 3.1 in Chetverikov (2019). Similarly, let  $\mathcal{P}_t^C$  be the set of all the conditional probability measures given  $T = t$  that measures from  $\mathcal{P}^C$  generate.

**Proposition 4.3.** *If  $\prod_{t \in \mathcal{T}} (1 - \alpha_t) \geq 1 - \alpha$ , then:*

$$\inf_{P \in \mathcal{P}^C \cap \mathcal{H}_0} P[\phi_n = 0] \geq 1 - \alpha + o(1) \tag{34}$$

as  $n \rightarrow \infty$ .

*Proof.* Notice that each  $\phi_{n_t}^t$  is a function of the observations from  $\mathcal{I}_t$  only. Since  $\mathcal{I}_t$  are mutually exclusive by construction and because the data are i.i.d., we have  $P[\phi_n = 0] = \prod_{t \in \mathcal{T}} P[\phi_{n_t}^t = 0]$ .

---

<sup>23</sup>This test is developed for continuous  $Z$ , which is used in our application. Although the instrument is discretized at the estimation stage, the monotonicity of  $\mathbb{E}[Y(t)|T = t, Z = z]$  for continuous  $Z$  clearly implies the monotonicity of the discretized moments. The procedure we describe straightforwardly accommodates testing discrete instruments. As noted in Chetverikov (2019), for discrete conditioning variable the test is a simple parametric problem, since the conditional moment function can be  $\sqrt{n}$ -consistently estimated at each point from the support.

By the standard optimization argument:

$$\Pi_{t \in \mathcal{T}} \inf_{P \in \mathcal{P}_t^C} P[\phi_{n_t}^t = 0] \leq \inf_{P \in \mathcal{P}^C} \Pi_{t \in \mathcal{T}} P[\phi_{n_t}^t = 0] \quad (35)$$

Theorem 3.1 from Chetverikov (2019) and  $\Pi_{t \in \mathcal{T}}(1 - \alpha_t) \geq 1 - \alpha$  then yield the result. ■

**Remark 4.4.** One may set  $\alpha_t = 1 - (1 - \alpha)^{1/N_T}$  as a baseline. If the domain knowledge suggests that for some treatments monotonicity is more likely to hold, one can set a higher  $\alpha_t$  for them, so long as  $\Pi_{t \in \mathcal{T}}(1 - \alpha_t) \geq 1 - \alpha$ . This may improve the power of the test.

## 5. Returns to education in Colombia

Our data is comprised of 861492 observations from Colombian labor force. The sample represents a snapshot of those individuals who could be matched across the educational, formal employment and census datasets in 2021<sup>24</sup>. For 664633 individuals from this dataset we observe their average lifetime wages, education level and Saber 5 or Saber 11 scores for Mathematics and Spanish language tests<sup>25</sup>.

The outcome variable we consider ( $Y_i$ ) is a log-wage, and  $T_i$  is the education level. We distinguish four education levels: primary, secondary and high school as well as 'university'<sup>26</sup>. Our measure of ability is constructed as a CES aggregator, which is then split into deciles:

$$Z_i \equiv (MATH_i^{1/2} + SPANISH_i^{1/2})^2$$

---

<sup>24</sup>Educational dataset was assembled by the testing authority Instituto Colombiano para la Evaluación de la Educación (ICFES), formal employment dataset comes from social security data based on Planilla Integrada de Liquidación de Aportes (PILA), whereas census data is handled by Departamento Administrativo Nacional de Estadística (DANE). The data was merged and anonymized by ICFES.

<sup>25</sup>Saber 5 and 11 tests are taken at different ages, but designed to be comparable between each other, which justifies merging them.

<sup>26</sup> $T_i$  is based on the number of years of schooling,  $S_i$ . If  $S_i < 9$ , set  $T_i \equiv 0$  meaning the individual only graduated from primary school.  $S_i \in [9; 11)$  and  $T_i \equiv 1$  correspond to completing compulsory education (secondary school),  $S_i = 11$  and  $T_i \equiv 2$  means that the individual is a high-school graduate, whereas  $S_i > 11$  with  $T_i \equiv 3$  means university education. Unfortunately,  $S_i$  is capped at 17 years in our sample, making it impossible to distinguish between those who continued to graduate education and those who just finished the 6–years degree.

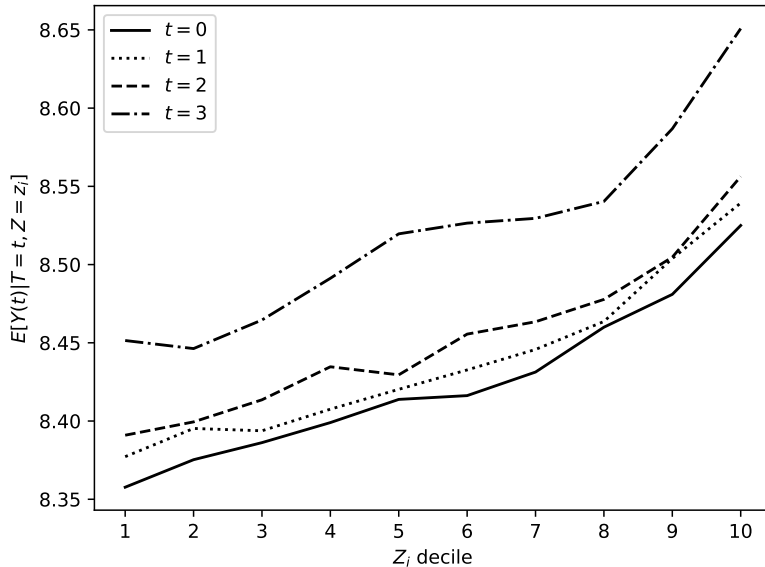


Figure 10: Estimated conditional moments of log-wages given ability and education level.

We first test whether cMIV is a reasonable assumption in our setup by implementing the test discussed in Section 3.3. To that end, we use the parameters and kernel functions recommended by Chetverikov (2019) and focus on the theoretically most powerful procedure, the step-down approach. The estimated  $p$ -value of the test is 0.29, see Table 1. We thus conclude that cMIV-p is a credible assumption provided that MIV holds.

$t$	$R_t^{st}$	$R_{t,0.1}^{crit}$	$p$ -value	$n_t$
0	0.98	2.33	0.34	274295
1	-1.17	2.17	0.95	143299
2	-1.51	2.30	1.00	216336
3	1.86	2.38	0.08	30703

Table 1: Results of the monotonicity test, see Section 4.2. Second column gives the estimated Chetverikov (2019) test-statistic, third column contains the  $\alpha = 0.1$  critical values, corresponding to  $\alpha_t = 1 - (1 - 0.1)^{1/4} \approx 0.026$  individual critical value. The last column gives a  $p$ -value against the individual null for each  $t$ . The overall  $p$ -value is 0.29.

The data we study is rather noisy. One would expect a considerable measurement error in the construction of both treatment levels and the outcome variable<sup>27</sup>. In line with that, the strongest form of cMIV is not sufficient to provide identification in the absence of further

<sup>27</sup>In particular, age is self-reported when filling an online questionnaire and appears to be of low quality, so we are forced to merge multiple cohorts.

assumptions. While the resulting bounds are tighter than that under MIV, they remain uninformative.

To achieve identification, we augment our assumptions with the MTR condition. While MIV and cMIV-w remain uninformative, both cMIV-p and cMIV-s result in positive lower bounds on the ATEs. Under cMIV-p the effect of obtaining a 'university education' is estimated to be at least as large as 3.62%, and 5.91% under cMIV-s. This is consistent with previous evidence. Causal estimates for the US (Card (1993), Brand and Xie (2010) and Angrist and Chen (2011)) report the return of at least 10% for a 4-year college degree. Recent evidence suggests that this number may be substantially lower for Colombia (Gomez, 2022).

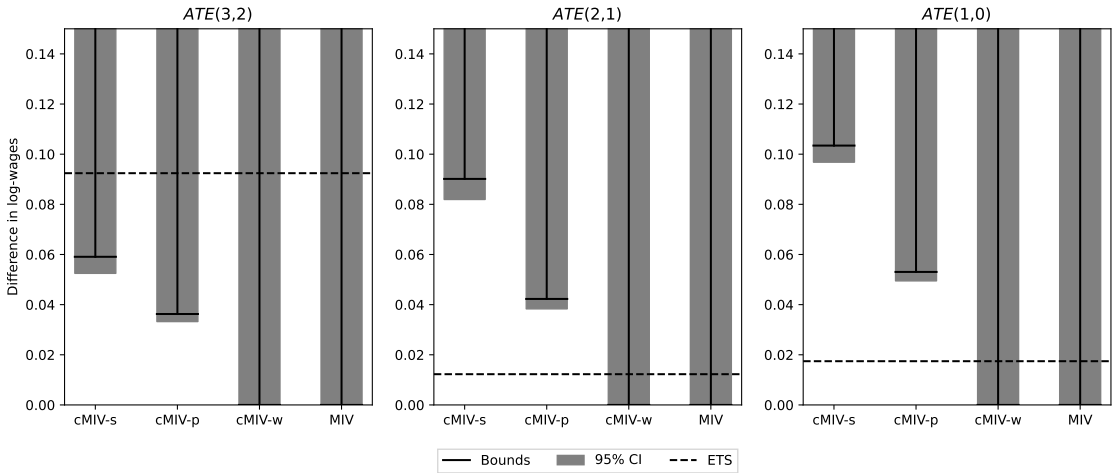


Figure 11: Estimation results for the monotonicity assumptions augmented with MTR. CI constructed according to Proposition 11. The exogenous treatment selection estimates (ETS) are  $ATE_{t,d}^{ETS} \equiv \mathbb{E}_n[Y(t)|T = t] - \mathbb{E}_n[Y(d)|T = d]$

We also find significantly positive effects at other education stages, see Figure 11. Further details on data construction and estimation as well as robustness checks are available in the Appendix.



## References

- ALIPRANTIS, C. AND K. BORDER (2007): *Infinite Dimensional Analysis: A Hitchhiker's Guide*, Springer.
- ANDREWS, D. W. (1999): "Estimation when a parameter is on a boundary," *Econometrica*, 67, 1341–1383.
- ANDREWS, D. W. K. (2000): "Inconsistency of the Bootstrap When a Parameter is on the Boundary of the Parameter Space," *Econometrica*, 68, 399–405.
- ANDREWS, I., J. ROTH, AND A. PAKES (2023): "Inference for Linear Conditional Moment Inequalities," *The Review of Economic Studies*, 90, 2763–2791.
- ANGRIST, J. D. AND S. H. CHEN (2011): "Schooling and the Vietnam-era GI Bill: Evidence from the draft lottery," *American Economic Journal: Applied Economics*, 3, 96–118.
- ANGRIST, J. D., G. W. IMBENS, AND D. B. RUBIN (1996): "Identification of causal effects using instrumental variables," *Journal of the American statistical Association*, 91, 444–455.
- BERESTEANU, A. AND F. MOLINARI (2008): "Asymptotic properties for a class of partially identified models," *Econometrica*, 76, 763–814.
- BERTSEKAS, D. P. (1975): "Necessary and sufficient conditions for a penalty method to be exact," *Mathematical programming*, 9, 87–99.
- BHATTACHARYA, D. (2009): "Inferring optimal peer assignment from experimental data," *Journal of the American Statistical Association*, 104, 486–500.
- BLUNDELL, R., A. GOSLING, H. ICHIMURA, AND C. MEGHIR (2007): "Changes in the distribution of male and female wages accounting for employment composition using bounds," *Econometrica*, 75, 323–363.
- BOES, S. (2009): "Bounds on counterfactual distributions under semi-monotonicity constraints," .
- BRAND, J. E. AND Y. XIE (2010): "Who Benefits Most from College?: Evidence for Negative Selection in Heterogeneous Economic Returns to Higher Education," *American Sociological Review*, 75, 273–302, PMID: 20454549.
- CARD, D. (1993): "Using geographic variation in college proximity to estimate the return to schooling," .
- CHERNOZHUKOV, V., H. HONG, AND E. TAMER (2007): "Estimation and Confidence Regions for Parameter Sets in Econometric Models," *Econometrica*, 75, 1243–1284.
- CHERNOZHUKOV, V., S. LEE, AND A. M. ROSEN (2013): "Intersection Bounds: Estimation and Inference," *Econometrica*, 81, 667–737.
- CHETVERIKOV, D. (2019): "TESTING REGRESSION MONOTONICITY IN ECONOMETRIC MODELS," *Econometric Theory*, 35, 729–776.
- CHO, J. AND T. M. RUSSELL (2023): "Simple Inference on Functionals of Set-Identified Parameters Defined by Linear Moments," *Journal of Business & Economic Statistics*, 0, 1–16.
- CYGAN-REHM, K., D. KUEHNLE, AND M. OBERFICHTNER (2017): "Bounding the causal effect of unemployment on mental health: Nonparametric evidence from four countries," *Health Economics*, 26, 1844–1861.
- DE HAAN, M. (2017): "The Effect of Additional Funds for Low-ability Pupils: A Non-parametric

- Bounds Analysis,” *The Economic Journal*, 127, 177–198.
- DEMYANOV, V. F. (2009): *Minimax: directional differentiability* *Minimax: Directional Differentiability*, Boston, MA: Springer US, 2075–2079.
- DUAN, Q., M. XU, L. ZHANG, AND S. ZHANG (2020): “Hadamard directional differentiability of the optimal value of a linear second-order conic programming problem,” *Journal of Industrial and Management Optimization*, 17, 3085–3098.
- FANG, Z. AND A. SANTOS (2018): “Inference on Directionally Differentiable Functions,” *The Review of Economic Studies*, 86, 377–412.
- GAFAROV, B. (2024): “Simple subvector inference on sharp identified set in affine models,” *arXiv e-prints*, arXiv-1904, conditionally Accepted at Journal of Econometrics, 2024.
- GOMEZ, N. (2022): “Returns to college education in Colombia,” *Higher Education Policy*, 35, 692–708.
- GUNDERSEN, C., B. KREIDER, AND J. PEPPER (2012): “The impact of the National School Lunch Program on child health: A nonparametric bounds analysis,” *Journal of Econometrics*, 166, 79–91, annals Issue on “Identification and Decisions”, in Honor of Chuck Manski’s 60th Birthday.
- HANSEN, B. E. (2017): “Regression kink with an unknown threshold,” *Journal of Business & Economic Statistics*, 35, 228–240.
- HECKMAN, J. J. AND E. VYTLACIL (2005): “Structural equations, treatment effects, and econometric policy evaluation 1,” *Econometrica*, 73, 669–738.
- HECKMAN, J. J. AND E. J. VYTLACIL (1999): “Local instrumental variables and latent variable models for identifying and bounding treatment effects,” *Proceedings of the national Academy of Sciences*, 96, 4730–4734.
- HONG, H. AND J. LI (2015): “The numerical delta method and bootstrap,” Tech. rep., Working paper.
- HONORÉ, B. E. AND A. LLERAS-MUNEY (2006): “Bounds in competing risks models and the war on cancer,” *Econometrica*, 74, 1675–1698.
- HONORÉ, B. E. AND E. TAMER (2006): “Bounds on parameters in panel dynamic discrete choice models,” *Econometrica*, 74, 611–629.
- IMBENS, G. W. AND J. D. ANGRIST (1994): “Identification and Estimation of Local Average Treatment Effects,” *Econometrica*, 62, 467–475.
- IMBENS, G. W. AND C. F. MANSKI (2004): “Confidence Intervals for Partially Identified Parameters,” *Econometrica*, 72, 1845–1857.
- KLINE, B. AND E. TAMER (2023): “Recent Developments in Partial Identification,” *Annual Review of Economics*, 15, 125–150.
- KLINE, P. AND M. TARTARI (2016): “Bounding the Labor Supply Responses to a Randomized Welfare Experiment: A Revealed Preference Approach,” *American Economic Review*, 106, 972–1014.
- KREIDER, B., J. V. PEPPER, C. GUNDERSEN, AND D. JOLLIFFE (2012): “Identifying the effects of SNAP (food stamps) on child health outcomes when participation is endogenous and misreported,” *Journal of the American Statistical Association*, 107, 958–975.
- LAFFÉRS, L. (2019): “Bounding average treatment effects using linear programming,” *Empirical*

- economics*, 57, 727–767.
- LAFFÉRS, L. (2013): “A note on bounding average treatment effects,” *Economics Letters*, 120, 424–428.
- LI, W. (1993): “The sharp Lipschitz constants for feasible and optimal solutions of a perturbed linear program,” *Linear algebra and its applications*, 187, 15–40.
- MANSKI, C. F. (1997): “Monotone Treatment Response,” *Econometrica*, 65, 1311–1334.
- MANSKI, C. F. AND J. V. PEPPER (2000): “Monotone Instrumental Variables: With an Application to the Returns to Schooling,” *Econometrica*, 68, 997–1010.
- (2009): “More on monotone instrumental variables,” *The Econometrics Journal*, 12, S200–S216.
- MASTEN, M. A. AND A. POIRIER (2018): “IDENTIFICATION OF TREATMENT EFFECTS UNDER CONDITIONAL PARTIAL INDEPENDENCE,” *Econometrica*, 86, 317–351.
- MEYER, R. (1979): “Continuity properties of linear programs,” Tech. rep., University of Wisconsin-Madison Department of Computer Sciences.
- MOGSTAD, M., A. SANTOS, AND A. TORGOVITSKY (2018): “Using instrumental variables for inference about policy relevant treatment parameters,” *Econometrica*, 86, 1589–1619.
- PINAR, M. Ç. AND B. CHEN (1999): “1 1 solution of linear inequalities,” *IMA journal of numerical analysis*, 19, 19–37.
- RICHEY, J. (2016): “An odd couple: Monotone instrumental variables and binary treatments,” *Econometric Reviews*, 35, 1099–1110.
- ROCKAFELLAR, R. T. (1970): *Convex Analysis*, Princeton: Princeton University Press.
- SEMENOVA, V. (2023): “Adaptive Estimation of Intersection Bounds: a Classification Approach,” .
- SHAPIRO, A. (1990): “On concepts of directional differentiability,” *Journal of optimization theory and applications*, 66, 477–487.
- SIDDIQUE, Z. (2013): “Partially Identified Treatment Effects Under Imperfect Compliance: The Case of Domestic Violence,” *Journal of the American Statistical Association*, 108, 504–513.
- SYRGKANIS, V., E. TAMER, AND J. ZIANI (2021): “Inference on auctions with weak assumptions on information,” *arXiv preprint arXiv:1710.03830*.
- TAO, T. AND V. VU (2010): “Random matrices: The distribution of the smallest singular values,” *Geometric And Functional Analysis*, 20, 260–297.
- VAN DER VAART, A. W., J. A. WELLNER, A. W. VAN DER VAART, AND J. A. WELLNER (1996): *Weak convergence*, Springer.
- VYTLACIL, E. (2002): “Independence, monotonicity, and latent index models: An equivalence result,” *Econometrica*, 70, 331–341.
- WAINWRIGHT, M. J. (2019): *High-dimensional statistics: A non-asymptotic viewpoint*, vol. 48, Cambridge university press.
- WRIGHT, S. J. (1997): *Primal-dual interior-point methods*, SIAM.
- YAKUSHEVA, O. (2010): “Return to college education revisited: Is relevance relevant?” *Economics of Education Review*, 29, 1125–1142.

## 6. Appendix

### 6.1. Proof of Lemma 2.1

*Proof.* Recall that for any sets  $A, B$ , we have  $\inf A \cup B = \min\{\inf A, \inf B\}$ . Fix any  $(\theta, w) \in \mathbb{R}^S \times \mathbb{R}^q$  such that  $\Theta_I(\theta) \subseteq \mathcal{X}$  and take  $A \equiv L(\Theta_I(\theta); \theta, w)$ ,  $B \equiv L(\mathcal{X} \cap \Theta_I(\theta)'; \theta, w)$ . Note that  $A \cup B = L(\mathcal{X}; \theta, w)$ . Substituting the definitions, it follows that

$$\begin{aligned} \tilde{B}(\theta; w) &= \min\left\{ \inf_{x \in \Theta_I(\theta)} p'x, \inf_{x \in \mathcal{X} \cap \Theta_I(\theta)'} p'x + w'(c - Mx)^+ \right\} \\ &= \min\left\{ B(\theta), \inf_{x \in \mathcal{X} \cap \Theta_I(\theta)'} p'x + w'(c - Mx)^+ \right\} \leq B(\theta), \end{aligned}$$

which establishes (7).

For the second part, fix any pair  $(\theta_0, w)$  that satisfies Assumption A0. We write  $\theta_0 = (p', c', \text{vec}(M)')'$ . Let  $\lambda^*$  be the KKT vector from Assumption A1. The definition of KKT vector (see Section 28 in Rockafellar (1970)) requires that

$$B(\theta_0) = \inf_{x \in \mathbb{R}^d} p'x + \lambda^{*'}(c - Mx) \quad (36)$$

Note that, for any  $x \in \mathbb{R}^d$ ,

$$B(\theta_0) \leq p'x + \lambda^{*'}(c - Mx) \leq p'x + \lambda^{*'}(c - Mx)^+ \leq p'x + w'(c - Mx)^+, \quad (37)$$

where the first inequality follows from (36), the second inequality follows from  $\lambda^* \geq 0$  and  $(t)^+ \geq t$  for any  $t \in \mathbb{R}$ , and the third inequality follows by Assumption A1. Taking infimum over  $x \in \mathcal{X}$  on both sides of (37) and combining with (7) yields

$$\tilde{B}(\theta_0; w) = B(\theta_0) \quad (38)$$

We now wish to show that  $\mathcal{A}(\theta_0) = \tilde{\mathcal{A}}(\theta_0; w)$ . From (38), the fact that  $L(x; \theta_0, w) = p'x$  for  $x \in \Theta_I(\theta_0)$  and  $\mathcal{A}(\theta_0) \subseteq \Theta_I(\theta_0)$ , it follows that:

$$\mathcal{A}(\theta_0) \subseteq \tilde{\mathcal{A}}(\theta_0; w) \quad (39)$$

To establish the other direction, we proceed by contradiction. Suppose  $\exists x^* \in \tilde{\mathcal{A}}(\theta_0; w) \cap \mathcal{A}(\theta_0)'$ . Suppose  $x^* \in \Theta_I(\theta_0)$ . Since  $x^* \notin \mathcal{A}(\theta_0)$ , it must then be that  $p'x^* > B(\theta_0)$ , but  $\tilde{B}(\theta_0; w) = L(x^*; \theta_0, w) = p'x^*$ , which yields a contradiction with (38). So,  $x^* \notin \Theta_I(\theta_0)$ . Consider

$$\tilde{B}(\theta_0; w) = p'x^* + w'(c - Mx^*)^+ > p'x^* + \lambda^{*'}(c - Mx^*)^+ \geq p'x^* + \lambda^{*'}(c - Mx^*) \quad (40)$$

where the first inequality follows from Assumption A1 and the fact that  $x^* \notin \Theta_I(\theta_0) \implies$

$\exists j : c_j - M_j x^* > 0$ . Combining (40) with (36), one gets  $\tilde{B}(\theta_0; w) > B(\theta_0)$ , which yields a contradiction with (38). So,  $\mathcal{A}(\theta_0) \supseteq \tilde{\mathcal{A}}(\theta_0; w)$ . Combining with (39) establishes

$$\mathcal{A}(\theta_0) = \tilde{\mathcal{A}}(\theta_0; w).$$

This concludes the proof of the lemma. ■

## 6.2. Proof of Theorem 2.1

*Proof.* Fix the true  $\theta_0 = (p', c', \text{vec}(M))'$ . Recall that for any  $g_1, g_2 \in \mathcal{C}(\mathcal{X})$ , we can bound

$$|\inf_{\mathcal{X}} g_1 - \inf_{\mathcal{X}} g_2| \leq \sup_{\mathcal{X}} |g_1 - g_2|. \quad (41)$$

Clearly,  $L(\cdot; \theta, w) \in \mathcal{C}(\mathcal{X})$  for any  $(\theta, w) \in \mathbb{R}^S \times \mathbb{R}_+$ . Using this, the bound (41) and the definition of  $\tilde{B}(\cdot)$ , one gets

$$|\tilde{B}(\hat{\theta}_n; w_n) - \tilde{B}(\theta_0; w_n)| \leq \sup_{x \in \mathcal{X}} |L(x; \hat{\theta}_n, w_n) - L(x; \theta_0, w_n)| \quad (42)$$

Using the definition of  $L(\cdot)$ , triangle and Cauchy-Schwarz inequalities, for every  $x \in \mathcal{X}$

$$|L(x; \hat{\theta}_n, w_n) - L(x; \theta_0, w_n)| \leq \|\hat{p}_n - p\| \cdot \|x\| + w_n \sqrt{q} \|(\hat{c}_n - \hat{M}_n x)^+ - (c - Mx)^+\|. \quad (43)$$

It is straightforward to observe that for any  $v_1, v_2 \in \mathbb{R}^q$

$$\|v_1^+ - v_2^+\| \leq \|v_1 - v_2\| \quad (44)$$

Further, recall that for any  $A \in \mathbb{R}^{q \times d}$  and  $x \in \mathbb{R}^d$

$$\|Ax\| \leq \|x\| \sup_{\|y\| \leq 1} \|Ay\| = \|x\| \cdot \|A\|_2, \quad (45)$$

where  $\|A\|_2$  is the spectral norm of  $A$ . Also recall that if  $\|\cdot\|_F$  is the Frobenius norm,

$$\|A\|_2 \leq \|A\|_F = \|\text{vec}(A)\|. \quad (46)$$

Combining (44), (45), (46) and using triangle inequality, one gets

$$\begin{aligned} |L(x; \hat{\theta}_n, w_n) - L(x; \theta_0, w_n)| &\leq \|\hat{p}_n - p\| \cdot \|x\| \\ &+ w_n \sqrt{q} \left( \|\hat{c}_n - c\| + \|\text{vec}(\hat{M}_n) - \text{vec}(M)\| \cdot \|x\| \right), \end{aligned}$$

where taking sup on both sides and using (42) yields

$$\begin{aligned} |\tilde{B}(\hat{\theta}_n; w_n) - \tilde{B}(\theta_0; w_n)| &\leq \|\hat{p}_n - p\| \cdot \|x\|_\infty \\ &+ w_n \sqrt{q} \left( \|\hat{c}_n - c\| + \|\text{vec}(\hat{M}_n) - \text{vec}(M)\| \cdot \|x\|_\infty \right) = O_p \left( \frac{w_n}{\sqrt{n}} \right), \end{aligned} \quad (47)$$

where  $\|x\|_\infty = \sup_{x \in \mathcal{X}} \|x\| < \infty$  by Assumption A0.ii, and the last equality follows from Assumption A0.iii.

Finally, since  $\theta_0$  satisfies A0, there exists some  $\lambda^*$  in  $\Lambda(\theta_0)$  with  $\|\lambda^*\|_\infty < \infty$ . Let  $E_n \equiv \{w_n > \|\lambda^*\|_\infty\}$ . By Lemma 2.1,  $E_n \subseteq \{\tilde{B}(\theta_0; w_n) = B(\theta_0)\}$ , so, for any deterministic  $\{r_n\}_{n \in \mathbb{N}}$  and any  $\varepsilon > 0$ ,

$$\begin{aligned} \mathbb{P} \left[ r_n |\tilde{B}(\hat{\theta}_n; w_n) - B(\theta_0)| > \varepsilon \right] &= \\ \mathbb{P} \left[ r_n |\tilde{B}(\hat{\theta}_n; w_n) - \tilde{B}(\theta_0; w_n)| > \varepsilon, E_n \right] &+ \mathbb{P} \left[ r_n |\tilde{B}(\hat{\theta}_n; w_n) - B(\theta_0)| > \varepsilon, E_n' \right] \leq \\ \mathbb{P} \left[ r_n |\tilde{B}(\hat{\theta}_n; w_n) - \tilde{B}(\theta_0; w_n)| > \varepsilon \right] &+ \mathbb{P}[E_n'] \end{aligned} \quad (48)$$

By definition,  $w_n \rightarrow \infty$  w.p.a.1 implies  $\mathbb{P}[w_n > \|\lambda^*\|_\infty] \rightarrow 1$ , so  $\mathbb{P}[E_n'] \rightarrow 0$ . Combining this, (47) and (48), one obtains:

$$\tilde{B}(\hat{\theta}_n; w_n) - B(\theta_0) = O_p \left( \frac{w_n}{\sqrt{n}} \right)$$

This concludes the proof of the Theorem. ■

**6.2.a. Set expansions approach.** Proposition 2.1 highlights that the plug-in estimator fails whenever the constraint set has an empty interior. For completeness of our argument, we develop a natural alternative to the penalty function estimator - the *set-expansion approach*. The idea here is to enlarge  $\Theta_I$  by relaxing each inequality constraint with a sequence  $\kappa_n$ <sup>28</sup>. The resulting estimator has the flavor of the approach in Chernozhukov et al. (2007). Intuitively, it enforces SC at the cost of producing a potentially conservative estimate. We show that, in general, this estimator can indeed have a conservative rate, and thus we do not advocate its use in practice.

The approach in this section is first to prove that the appropriately extended identified set converges to the population identified set in Hausdorff distance, and then use uniform continuity of the criterion function as well as its resemblance to the support function to establish the convergence of the estimator itself.

<sup>28</sup>In the presence of ‘true equality’ constraints  $Ax = b$ , the corresponding inequalities need not be expanded.

Consider the following criterion function and its sample analogue:

$$Q(x) \equiv \|(Mx - c)^-\|^2, \quad \hat{Q}_n(x) \equiv \|(\hat{M}_n x - \hat{c}_n)^-\|^2$$

Denote the identified set as  $\Theta_I \equiv \{x \in \mathcal{X} | Q(x) = 0\} = \{x \in \mathcal{X} | Mx - c \geq 0\}$ .

**Lemma 6.1.**  $\|\hat{Q}_n(x) - Q(x)\|_\infty \xrightarrow{p} 0$ , where  $\|\cdot\|_\infty$  is over  $\Theta_I$ .

*Proof.*

$$|\hat{Q}_n(x) - Q(x)| = \left| \sum_j \left( [\hat{M}_n x - \hat{c}_n]_j^- \right)^2 - \left( [Mx - c]_j^- \right)^2 \right| = \quad (49)$$

$$= \left| \sum_j ([\hat{M}_n x - \hat{c}_n]_j^- - [Mx - c]_j^-) ([\hat{M}_n x - \hat{c}_n]_j^- + [Mx - c]_j^-) \right| \leq \quad (50)$$

$$\leq \sum_j |[\hat{M}_n x - \hat{c}_n]_j^- - [Mx - c]_j^-| \cdot |[\hat{M}_n x - \hat{c}_n]_j^- + [Mx - c]_j^-| \leq \quad (51)$$

$$\leq \left( \max_j [\hat{M}_n x - \hat{c}_n]_j^- + [Mx - c]_j^- \right) \sum_j \left| [(\hat{M}_n - M)x + c - \hat{c}_n]_j \right| \quad (52)$$

Where (52) uses the fact that  $|(y_0)^- - (y_1)^-| = |\max\{0, -y_0\} - \max\{0, -y_1\}| \leq |y_0 - y_1| \forall y_0, y_1 \in \mathbb{R}$ . We now show that the last line converges to 0 is supremum over  $x \in \mathcal{X}$ . Note that, since  $\hat{M}_n \xrightarrow{p} M$ ,  $\hat{c}_n \xrightarrow{p} c$ , the estimator asymptotically lies in any  $\delta$ -vicinity of the true population parameter. In other words,  $\forall \delta > 0$ , we have  $(\hat{c}'_n, \text{vec}(\hat{M}_n)')' \in B_\delta((c', \text{vec}(M)'))'$  w.p. 1 asymptotically.

Since  $\mathcal{X}$  is a compact and because of the former result, both  $\hat{M}_n x - \hat{c}_n$  and  $Mx - c$  are bounded w.p. 1 asymptotically, so there exists  $K > 0$  - large enough:<sup>29</sup>

$$\sup_{x \in \mathcal{X}} \max_j [\hat{M}_n x - \hat{c}_n]_j^- + [Mx - c]_j^- \leq K + o_p(1) \quad (53)$$

Note that by Cauchy-Schwarz,  $\sum_j \left| [(\hat{M}_n - M)x + c - \hat{c}_n]_j \right| \leq |\mathcal{F}_t| \cdot \|(\hat{M}_n - M)x + c - \hat{c}_n\|$ . Further using (52), (53) and noting that for nonnegative  $f, g$  one has  $\sup_A fg \leq \sup_A f \cdot \sup_A g$ , we get:

$$\|\hat{Q}_n(x) - Q(x)\|_\infty \leq (K + o_p) \cdot |\mathcal{F}_t| \cdot \sup_{x \in \mathcal{X}} \|(\hat{M}_n - M)x + c - \hat{c}_n\| \leq \quad (54)$$

$$\leq (\tilde{K} + o_p) \cdot (\|\hat{M}_n - M\| \cdot \|x\|_\infty + \|c - c_n\|) = o_p(1) \quad (55)$$

The proof is complete. ■

<sup>29</sup>In the cMIV setup all terms of  $\hat{M}_n$ ,  $\hat{c}_n$  are known to be bounded, so asymptotic arguments are not necessary. We consider a more general case here.

Analogously to the proof of Lemma 3, one shows that because both  $\hat{c}_n$  and  $\hat{M}_n$  are  $\sqrt{n}$ -consistent from A0, we have:

$$\sup_{\mathcal{X}}(Q - \hat{Q}_n)^+ = O_p(1/\sqrt{n}), \quad \sup_{\Theta_I} \hat{Q}_n = O_p(1/n).$$

The plug-in estimator of the identified set,  $\{x \in \mathcal{X} | \hat{Q}_n = 0\} = \Theta_I(\hat{\theta}_n)$ , may not ‘cover’ the true asymptotically, as discussed in Chernozhukov et al. (2007) (CHT). To address that, consider the following class of set estimators:

$$\{x \in \mathcal{X} | n\hat{Q}_n(x) \leq \kappa_n\}$$

Fix  $\kappa_n$  such that  $P[\kappa_n \geq \sup_{\Theta_I} n\hat{Q}_n] \rightarrow 1$  and  $\frac{\kappa_n}{n} \xrightarrow{p} 0$ . Let  $\hat{\Theta}_n \equiv \{x | \hat{M}_n x - \hat{c}_n \geq -\frac{\sqrt{\kappa_n}}{\sqrt{n}}\iota\}$ . It is the set that we want to prove consistent for the population identified set.

The issue is that the set  $\hat{\Theta}_n$  is not a criterion-based set, so the results in CHT is not directly applicable. However, we can define  $\underline{\Theta}_n \equiv \{x | \hat{Q}_n(x) \leq \frac{\kappa_n}{n}\} \subseteq \hat{\Theta}_n$  and  $\bar{\Theta}_n \equiv \{x | \hat{Q}_n(x) \leq q \frac{\kappa_n}{n}\} \supseteq \hat{\Theta}_n$ .

We then wish to ‘sandwich’  $\hat{\Theta}_n$  between a smaller set that asymptotically covers  $\Theta_I$  and a bigger set that is asymptotically covered by  $\Theta_I$ . The following simple lemma is an analogue of ‘sandwich theorem’ for sets.

**Lemma 6.2.** Consider  $\Theta_I \subseteq \mathcal{X}$  and suppose the random set  $\hat{\Theta}_n \subseteq \Theta$  can be sandwiched between two sets:  $\underline{\Theta}_n \subseteq \hat{\Theta}_n \subseteq \bar{\Theta}_n$ , such that:

$$\begin{aligned} \sup_{x \in \bar{\Theta}_n} d(x, \Theta_I) &= o_p(1) \\ \sup_{x \in \Theta_I} d(x, \underline{\Theta}_n) &= o_p(1) \end{aligned}$$

Then:

$$d_H(\hat{\Theta}_n, \Theta_I) = o_p(1)$$

*Proof.* Writing out the definitions and applying CMT yields the result. ■

The only thing that remains to show consistency of the set-estimator is to prove that the inequalities in Lemma 6.2 hold in our case. The derivation below follows the usual CHT logic. The first equality is established through:

$$\begin{aligned} P[\sup_{x \in \bar{\Theta}_n} d(\theta, \Theta_I) \leq \varepsilon] &= P[\bar{\Theta}_n \subseteq \Theta_I^\varepsilon] = \\ &P[\bar{\Theta}_n \cap \mathcal{X} \setminus \Theta_I^\varepsilon = \emptyset] \geq P[\sup_{x \in \bar{\Theta}_n} Q(\theta) < \inf_{x \in \mathcal{X} \setminus \Theta_I^\varepsilon} Q(\theta)] \end{aligned} \tag{56}$$



Then, by uniform continuity and by the construction of  $\bar{\Theta}_n$ :

$$\sup_{x \in \bar{\Theta}_n} Q(\theta) = \sup_{x \in \bar{\Theta}_n} \hat{Q}_n(\theta) + o_p(1) = q \frac{\kappa_n}{n} + o_p(1) = o_p(1)$$

By construction of  $\Theta_I$  and continuity of  $Q(\theta)$ ,  $\exists \delta > 0$ :  $\inf_{x \in \mathcal{X} \setminus \Theta_I^\varepsilon} Q(\theta) > \delta$ . Thus, the RHS of (56) goes to 1. So,  $\sup_{x \in \bar{\Theta}_n} d(x, \Theta_I) = o_p(1)$ .

The other side follows, as by construction  $\sup_{x \in \Theta_I} \hat{Q}_n(x) \leq \frac{\kappa_n}{n} \implies \Theta_I \subseteq \underline{\Theta}_n$ . So,

$$P[\sup_{x \in \Theta_I} d(\theta, \underline{\Theta}_n) \leq \varepsilon] \geq P[\sup_{x \in \Theta_I} \hat{Q}_n(x) \leq \frac{\kappa_n}{n}] \xrightarrow{p} 1$$

Therefore, using Lemma 3, we conclude that:

$$d_H(\hat{\Theta}_n, \Theta_I) \xrightarrow{p} 0$$

The next step is to recall that if we have two convex, compact sets,  $A, B$ , the following holds:

$$d_H(A, B) = \max_{\|y\| \leq 1} |s(y, A) - s(y, B)|,$$

where  $s(y, S) \equiv \max_{t \in S} y't$  - the support function.

Using uniform convergence of the value function and combining all the results:

$$\begin{aligned} & |\min_{x \in \hat{\Theta}_n} \hat{p}'_n x - \min_{x \in \Theta_I} p'x| = |\min_{x \in \hat{\Theta}_n} p'x - \min_{x \in \Theta_I} p'x| + o_p(1) = \\ & = |s(-p, \Theta_I) - s(-p, \hat{\Theta}_n)| + o_p(1) \leq \|p\| d_H(\Theta_I, \hat{\Theta}_n) + o_p(1) \xrightarrow{p} 0 \end{aligned}$$

This establishes the following proposition:

**Proposition 6.1.** *Let  $\kappa_n : P[\kappa_n \geq \sup_{\Theta_I} n\hat{Q}_n] \rightarrow 1$  and  $\frac{\kappa_n}{n} \xrightarrow{p} 0$ . Then the following estimator is consistent for the sharp lower bound:*

$$\check{B}_n \equiv \min_{\hat{M}_n x - \hat{c}_n \geq -\sqrt{\frac{\kappa_n}{n}}} \hat{p}'_n x \xrightarrow{p} \min_{Mx - c \geq 0} p'x$$

In practice, Chernozhukov et al. (2007) suggest to select some  $\kappa_n$  that diverges sufficiently slowly with the sample size. We use  $\sqrt{\kappa_n} \propto \ln \ln n$  in the simulations in Section 2.6. Under the Slater's condition the naive estimator is consistent, i.e. one could set  $\kappa_n = 0$ .

Although it seems intuitive that  $\check{B}_n$  should converge at the rate  $\sqrt{n\kappa_n^{-1}}$ , deriving that result is outside the scope of this paper, because we do not advocate its use. It is immediate to see, however, that  $\check{B}_n$  can converge as slowly as  $\sqrt{n\kappa_n^{-1}}$ . For that, consider the example in

Proposition 2.1 without the inequality  $x_2 \leq x_1$  and setting  $\hat{b}_n = 0$ . The minimum is attained at  $-1 - \sqrt{\frac{\kappa_n}{n}}$ . Our simulation evidence suggests that the set-expansion estimator can be quite conservative in practice, see Section 2.6.

### 6.3. Proof of Proposition 5

First, notice that the iff result for the following conditions: i) SC holds and ii)  $\mathcal{A}(\theta_0)$  and  $\Lambda(\theta_0)$  are both singletons follows from Theorem 3.1 in Fang and Santos (2018) combined with Lemma 2.2 and Proposition 2.3. Then, observe that  $\mathcal{A}(\theta_0)$  and  $\Lambda(\theta_0)$  are both singletons if and only if both LICQ hold and there are no flat faces.

### 6.4. Inference for LP under SC and the biased penalty function estimator

**6.4.a. LP under Slater's condition.** We now consider inference for the original LP estimator under Slater's condition. Proposition 5 showed that, in the absence of this condition, the plug-in may fail to be consistent, because the value function is not continuous in the parameter  $\theta_0$ <sup>30</sup>.

**Assumption B2 (Slater's condition).**  $\text{Int}(\Theta_I) \neq \emptyset$

The following lemma establishes Hadamard directional differentiability of a linear program under Assumption B2.

There is no apparent reason to suppose that ii) should hold in practice, and therefore we do not endorse applying Proposition 9. Instead, it is intended to illustrate the difficulty with inference on general LP value functions even under the Slater's condition.

**Corollary 1.** Under assumption B2, Proposition 6 holds with  $\kappa_n = 0$ , i.e. the naive estimator is consistent.

One way to obtain a consistent estimator is to employ the procedure developed in Hong and Li (2015). Let:

$$\tilde{B}'_n(\mathbb{Z}_n^*) \equiv \frac{B(\hat{\theta}_n + \epsilon_n \mathbb{Z}_n^*) - B(\hat{\theta}_n)}{\epsilon_n} \quad (57)$$

For  $\epsilon_n \rightarrow 0$  with  $r_n \epsilon_n \rightarrow \infty$ , we have the following proposition:

---

<sup>30</sup>Although Slater's condition is not necessary for duality in the case of LP, its failure allows for unbound- edness of the dual solution set at  $\theta_0$ , see Wright (1997). As shown in Meyer (1979), this will imply that the optimal value function is not continuous with respect to perturbations in  $c$ .

**Proposition 6.2.** *If Assumptions B1, B2 hold, and the bootstrapped  $\mathbb{Z}_n^*$  satisfies the measurability conditions in Hong and Li (2015):*

$$\sup_{f \in BL_1(\mathbb{R})} |\mathbb{E}[f(\tilde{B}'_n(\mathbb{Z}_n^*)) | \{X_i\}_{i=1}^n] - \mathbb{E}[f(B'_{\theta_0}(\mathbb{G}_0))]| = o_p(1) \quad (58)$$

Assumption B2 is rather strong, and one may not be comfortable imposing it directly. This is especially true in cases where many inequality restrictions are involved, such as under cMIV-s, because one would be concerned that the defined system may be close to point-identification. An even more serious problem in practice is that, even if an open ball is contained in  $\Theta_I$  at  $\theta_0$ , the radius of that ball is not inconsequential in finite samples. A thinner identified set leads the bootstrap iterations of the N.D.M. to fail more often, as the constraint set turns empty at perturbed parameter values. Dropping the failed iterations introduces an unknown bias to the estimates, and so is not advised.

One potential solution would be to use the set-expansion estimator as in Section 4.2. Indeed, as long as the true system is feasible, expanding the set from the RHS renders the Slater's condition true, and the procedure described in this section becomes applicable. The bias of such expansion would be controlled as follows:

$$\min_{\Theta_I} p'x - \|p\|d_H(\Theta_I, \tilde{\Theta}_I) \leq \min_{\Theta_I} p'x \leq \min_{\Theta_I} p'x \quad (59)$$

Moreover, by Lipschitz continuity of systems of linear inequalities,  $d_H(\Theta_I, \tilde{\Theta}_I) \leq C|\kappa|$  for some  $C > 0$  depending on  $\theta_0$ , where the vector  $\kappa > 0$  is the RHS-expansion.

This estimator, however, would still be problematic both because it is conservative even in terms of the convergence rate, and because it relies on an arbitrarily selected set expansion. Since a larger expansion leads to a more conservative lower bound, in applied work the researcher would be tempted to select the minimal value that ensures the bootstrap iterations do not fail. The statistical properties of that approach are unclear.

**6.4.b. Inference for the biased penalty.** We now consider the penalty function estimator  $\tilde{B}(\cdot)$  defined in Section 4.1. The main difficulty when conducting inference for it consists of proving its Hadamard directional differentiability.

Observe that we can write  $\tilde{B} \equiv \phi \circ \tilde{L}$ , where  $\tilde{L}(\theta) \equiv L(\cdot; \theta)$  is a map  $\tilde{L} : \mathbb{R}^S \rightarrow \ell^\infty(\mathcal{X})$ , and  $\phi : \ell^\infty(\mathcal{X}) \rightarrow \mathbb{R}$  is given by:

$$\phi(q) \equiv \inf_{x \in \mathcal{X}} q(x),$$

and where we equip  $\ell^\infty(\mathcal{X})$  with the sup norm. By Lemma S.4.9 in the Online Appendix of Fang and Santos (2018),  $\phi$  is Hadamard directionally differentiable. It is therefore tempting to apply the chain rule to find the derivative of  $\tilde{B}$ , which only requires that  $\tilde{L}$  is H.d.d.

However, in the spirit of the example from Hansen (2017), this is not the case. The following remark illustrates that issue.

**Remark 6.1.**  $g(y)(x) \equiv (x + y)^+$  viewed as a map  $g : \mathbb{R} \rightarrow \ell^\infty(A)$  for  $x \in A \equiv [-C; C]$  for some  $C > 0$  is **not** Hadamard directionally differentiable for any fixed  $y \in [-C/2; C/2]$ :

$$\lim_{t_n \rightarrow 0^+, h_n \rightarrow h} \left\| \frac{(y + x + t_n h_n)^+ - (y + x)^+}{t_n} - f(h)(x) \right\|_\infty \neq 0$$

for any continuous  $f(h)(x)$ . To see that, note that the first term converges pointwise to  $\mathbb{I}\{y + x = 0\}h^+ + \mathbb{I}\{y + x > 0\}h$ . Suppose that  $h < 0$  and consider:  $x_n = -y - \frac{t_n}{2}h_n$ , we have:

$$\begin{aligned} \left| \frac{(y + x_n + t_n h_n)^+ - (y + x_n)^+}{t_n} - \mathbb{I}\{y + x_n = 0\}h^+ + \mathbb{I}\{y + x_n > 0\}h \right| &= \\ &= o(1) - \frac{h}{2} \neq o(1) \end{aligned}$$

In light of this finding, it should be almost surprising that  $\tilde{B}(\cdot)$  is still Hadamard directionally differentiable, as we now demonstrate. Instead of using the chain rule, which is of course only a sufficient condition for differentiability, we notice that  $\tilde{B}$  can be rewritten as a new linear program that has a non-empty interior of the constraint set<sup>31</sup>.

**Proposition 6.3.** *The penalty function estimator,  $\tilde{B}(\theta; w)$  is Hadamard directionally differentiable in  $\theta$  at  $\theta_0$  if either i)  $\mathcal{X}$  is a polytope with  $\text{Int}(\mathcal{X}) \neq \emptyset$ , or ii)  $\exists x \in \tilde{\mathcal{A}}(\theta_0; w) \cap \text{Int}(\mathcal{X})$ . The H.d.d. is given by:*

$$\tilde{B}'_{\theta_0}(h; w) = \inf_{x \in \tilde{\mathcal{A}}(\theta_0; w)} \sup_{\lambda \in \tilde{\Lambda}(\theta_0; w)} h'_p x + \sum_{j=1}^{2^q} \lambda_j \sum_{i \in \Pi_j} w_i (h_{c_i} - h'_{M_i} x) \quad (60)$$

where  $h = (h'_p, h_{c_1}, \dots, h_{c_q}, h'_{M_1}, \dots, h'_{M_q})$  is the direction and an upper-hemicontinuous correspondence  $\tilde{\Lambda} : \mathbb{R}^S \rightarrow 2^{1, 2^q}$  is as defined in the proof.

*Proof.* Throughout this proof  $w$  is taken to be fixed, therefore some dependencies on it are omitted in notation for brevity. We proceed in four steps:

1. Notice that  $L(x; \theta, w)$  is a convex piecewise-linear function and it has the following representation:

$$L(x; \theta, w) = \max_{j \in \{1, 2^q\}} \left\{ p'x + \sum_{i \in \Pi_j} w_i (c_i - M'_i x) \right\}, \quad (61)$$

<sup>31</sup>Clearly, this new LP is not equivalent to the original one point-by-point, as that would mean that the plug-in,  $B(\cdot)$ , is always H.d.d., contradicting Proposition 5.

where  $\{\Pi_j\}_{j=1}^{2^q} = \overline{1, 2^q}$ , so that  $\Pi_j$  for different  $j$  contain indices of all possible combinations of positive penalty term. At a given  $x$  these can be interpreted as the sets of violated constraints. Let  $g_j(x, \theta) \equiv p'x + \sum_{i \in \Pi_j} w_i(c_i - M_i'x)$  for  $j \in \overline{1, 2^q}$ .

The initial estimator can then be represented as:

$$\tilde{B}(\theta; w) = \min_{x \in \mathcal{X}} \max_{j \in \overline{1, 2^q}} g_j(x, \theta) \quad (62)$$

2. Assumptions i) or ii) allow us to impose w.l.g. that the known compact set  $\mathcal{X}$  is a fixed, non-empty and bounded polyhedron. To see that for ii), note that the program is convex and therefore the sets of local and global minima coincide. If there exists an interior local minimum, it means that expanding the constraint set does not change the value, and therefore we can set  $\mathcal{X}$  to be some compact and non-empty polyhedron that contains the original set. Then, another representation of the considered problem follows:

$$\tilde{B}(\theta; w) = \min_{t, x} t \quad \text{s.t.} \quad \begin{cases} t \in [\underline{t}, \bar{t}] \\ x \in \mathcal{X} \\ g_j(x, \theta) \leq t, j \in \overline{1, 2^q} \end{cases} \quad (63)$$

For some sufficiently wide  $[\underline{t}, \bar{t}]$ , given  $\theta$  is close to  $\theta_0$  and such that  $\tilde{B}(\theta_0; w) \in (\underline{t}, \bar{t})$ . This is justified because  $\tilde{B}(\theta; w)$  is continuous in  $\theta$ , as shown in the proof of Proposition 5.

3. Note that the constraint set of (63) is compact, non-empty at  $\theta = \theta_0$  and, moreover, it contains an open set. To see that, consider some pair  $x(\theta_0), t(\theta_0)$  from the argmin of the problem, where  $x(\theta_0) \in \tilde{A}(\theta_0; w) \subseteq \mathcal{X}$  and  $t(\theta_0) \equiv \tilde{B}(\theta_0; w)$ . Consider  $\varepsilon \equiv \bar{t} - t(\theta_0)$  and take  $t^* \equiv t(\theta_0) + \frac{\varepsilon}{2}$ . Note that by definition  $t(\theta_0) \geq \max_j g_j(x(\theta_0))$ . By continuity of  $g_j(x, \theta_0)$  in  $x$  for all  $j \in \overline{1, 2^q}$ ,  $\exists \delta > 0$  such that  $t \geq \max_j g_j(x) \forall t \in (t^* - \frac{\varepsilon}{4}; t^* + \frac{\varepsilon}{4}), \forall x \in \mathbb{B}_\delta(x(\theta_0))$ . By either i) or ii)  $\text{Int}(\mathcal{X}) \neq \emptyset$  and as  $x(\theta_0) \in \mathcal{X}$  it follows that  $\text{Int}(\mathcal{X}) \cap \mathbb{B}_\delta(x(\theta_0))$  is non-empty. It is also open as an intersection of two open sets. Therefore, the open set  $\mathcal{O} \equiv (t^* - \frac{\varepsilon}{4}; t^* + \frac{\varepsilon}{4}) \times (\mathbb{B}_\delta(x(\theta_0)) \cap \text{Int}(\mathcal{X}))$  is contained in the constraint set of the induced LP at  $\theta_0$ . That is, the problem at  $\theta_0$  satisfies the Slater's condition and Lemma 6 applies.
4. Suppose  $\check{\Lambda}(\theta_0)$  is the set of Lagrange multipliers of (63) at  $\theta = \theta_0$ , and  $\tilde{\Lambda}(\theta_0)$  is its projection on the coordinates corresponding to the constraints of form  $g_j(x; \theta_0) \leq t$  for all  $j \in \overline{1, 2^q}$ . A typical element of  $\tilde{\Lambda}(\theta_0)$  will be written as  $\lambda = (\lambda_j)_{j=1}^{2^q}$ . Recall that for  $\theta$  in some small open neighbourhood of  $\theta_0$  the value function of (63) is equal to  $\tilde{B}(\theta; w)$  and, moreover, the problems are equivalent, so if  $\check{\Lambda}(\theta)$  is the arg min of (63),

then  $\check{\mathcal{A}}(\theta) = \{\check{B}(\theta; w)\} \times \check{\mathcal{A}}(\theta; w)$ . Using the conclusion of Step 3, direct application of Lemma 6 to (63) yields:

$$\check{B}'_{\theta_0}(h; w) = \inf_{x \in \mathcal{A}(\theta_0; w)} \sup_{\lambda \in \check{\Lambda}(\theta_0)} \sum_{j=1}^{2^q} \lambda_j \left( h'_p x + \sum_{i \in \Pi_j} w_i (h_{c_i} - h'_{M_i} x) \right), \quad (64)$$

where note that there are no terms corresponding to the objective function and the constraints  $t \in [\underline{t}, \bar{t}]$  and  $x \in \mathcal{X}$ , because there are no corresponding increments. Moreover, differentiating the Lagrangean of (63) and recalling that  $t(\theta_0) \in (\underline{t}, \bar{t})$ , so the constraints  $t \in [\underline{t}, \bar{t}]$  do not bind and the corresponding multipliers are 0, one gets that  $\forall \lambda \in \check{\Lambda}(\theta_0)$ , we have  $\sum_{j=1}^{2^q} \lambda_j = 1$ , establishing (60). ■

**Remark 6.2.** By Lemma 2, Assumption A1 ensures ii) in Proposition 11 if  $\Theta_I \subseteq \text{Int}(\mathcal{X})$ .

Assuming A1 holds, exact pointwise inference is then obtained via Proposition 10. It is also straightforward to show that if A1 does not hold, but conditions i) or ii) in Proposition 11 are otherwise satisfied, this inference is asymptotically conservative.

Computational considerations may be important in practice, especially as bootstrap is involved. In Appendix we further show that the penalty function estimator may be computed as a value of a simple LP. If there are  $k$  constraints defining  $\mathcal{X}$  and  $q$  constraints for  $\Theta_I$ , with  $d$  variables, the penalty-induced LP will feature  $d + q$  variables and  $2q + k$  constraints, which makes it almost as simple computationally as the usual plug-in estimator with  $d$  variables and  $q + k$  constraints.

## 6.5. Proof of Theorem 2.2

*Proof.* Fix  $\theta = (p', \text{vec}(M)', c')' = \theta_0$ . We proceed in six steps, first proving the following lemma:

**Lemma 6.3.** Consider  $B \equiv \arg \min_{x \in A} f(x)$  and  $c \equiv f(x^*)$  for any  $x^* \in B$ , where  $f(\cdot)$  is continuous and  $A$  is a non-empty compact. Then, for any measurable random sequence  $\{x_n\} \subseteq A$  such that  $f(x_n) \xrightarrow{p} c$ , there exists a measurable random sequence  $\{x_n^*\} \subseteq B$  such that  $\|x_n^* - x_n\| \xrightarrow{p} 0$ .

*Proof.* Under the assumptions of the Lemma, Berge's maximum theorem implies that  $B$  is a non-empty compact. Because the distance is continuous, the projection  $x_n^*$  of  $x_n$  onto  $B$  is always well-defined for each  $n$ . If it is not unique, we select one of the values that yield the minimum distance. Measurability of at least one such selection is established by reference to Theorem 18.19 in Aliprantis and Border (2007). We then proceed by contradiction. Suppose

that  $\exists \varepsilon > 0$ :

$$\mathbb{P}[\|x_n^* - x_n\| > \varepsilon] \not\rightarrow 0 \quad (65)$$

Then, there exists a  $\delta > 0$  and a subsequence  $\{n_k\}_{k=1}^\infty$  such that, for all  $k \in \mathbb{N}$ :

$$\mathbb{P}[\|x_{n_k}^* - x_{n_k}\| > \varepsilon] > \delta \quad (66)$$

Consider the following problem:

$$\min_{x \in A, d(x, B) \geq \varepsilon} f(x) \quad (67)$$

Notice that the constraint set is compact. It is also non-empty, as for any  $k$  some of the realisations of  $x_{n_k}$  are in it by (66). Therefore the minimum is attained at some  $\tilde{x}$ . Suppose that the minimum is equal to  $f(\tilde{x}) = \tilde{c}$ . If  $\tilde{c} = c$ , it follows that  $\tilde{x} \in B$ , which is not possible as  $d(\tilde{x}, B) \geq \varepsilon$ . Clearly,  $\tilde{c} < c$  is also infeasible as the constraint set of that problem is smaller than that of the original one. Therefore,  $\tilde{c} - c = K > 0$ . Then, note that for any  $k \in \mathbb{N}$ :

$$\|x_{n_k}^* - x_{n_k}\| > \varepsilon \implies f(x_{n_k}) \geq f(\tilde{x}) = c + K > c \quad (68)$$

So,

$$\mathbb{P}[f(x_{n_k}) - f(x^*) \geq K] \geq \mathbb{P}[\|x_{n_k}^* - x_{n_k}\| > \varepsilon] > \delta > 0, \quad (69)$$

where the LHS goes to 0 as  $k \rightarrow \infty$ , since  $f(x_{n_k}) \xrightarrow{p} f(x^*)$  by assumption of the Lemma. This yields a contradiction. Therefore,  $\|x_n^* - x_n\| \xrightarrow{p} 0$ .  $\blacksquare$

1. We first prove that  $\exists \{\delta_n\} \rightarrow 0^+$  such that  $\mathcal{A}(\hat{\theta}_n, w_n) \subseteq \mathcal{A}(\theta_0)^{\delta_n}$  w.p. 1 asymptotically. For this purpose, recall that by Theorem 3 for any sequence  $x_n \in \mathcal{A}(\hat{\theta}_n, w_n)$  for all  $n$  and for any  $x^* \in \mathcal{A}(\theta_0)$ , we have:

$$p'x_n + w_n \iota'(\hat{c}_n - \hat{M}_n x_n)^+ - p'x^* = o_p(1) \quad (70)$$

Furthermore, since  $w_n = o_p(\sqrt{n})$ , we have:

$$w_n \|\hat{c}_n - \hat{M}_n x - c + Mx\|_\infty = o_p(1) \quad (71)$$

Because the argmin is contained in a compact,  $\mathcal{A}(\hat{\theta}_n, w_n) \subseteq \mathcal{X}$ , the first term in (70) is bounded in probability:  $p'x_n = O_p(1)$ , thus, from (70), it also follows that  $w_n \iota'(\hat{c}_n - \hat{M}_n x_n)^+ = O_p(1)$ . By triangle inequality and using with (71), we therefore

conclude:

$$w_n \iota'(c - Mx_n)^+ = O_p(1) \quad (72)$$

As  $w_n \rightarrow \infty$ , it further follows that:

$$(c - Mx_n)^+ = o_p(1) \quad (73)$$

We shall now consider  $\tilde{x}_n$  - a projection of  $x_n$  onto  $\{x \in \mathbb{R}^d | Mx \geq c\}$ . Note that it exists, because distance is a continuous function and the set is a non-empty compact. Note that (73) implies that, for some random  $\kappa_n \geq 0$  for all  $n$ :

$$c - Mx_n \leq \iota \kappa_n \quad (74)$$

where  $\kappa_n = o_p(1)$ . We get:

$$\|x_n - \tilde{x}_n\| = d(x_n, \{x \in \mathbb{R}^d | Mx \geq c\}) \leq \quad (75)$$

$$\leq d_H(\{x \in \mathbb{R}^d | Mx \geq c - \kappa_n\}, \{x \in \mathbb{R}^d | Mx \geq c\}) \leq C\kappa_n, \quad (76)$$

where  $C > 0$  is some fixed constant. The first equality is by definition of projection, the second inequality follows from the definition of the Hausdorff distance and (74) as well as:

$$d(x_n, \{x \in \mathbb{R}^d | Mx \geq c\}) \leq \sup_{x \in \{x \in \mathbb{R}^d | Mx \geq c - \kappa_n\}} d(x, \{x \in \mathbb{R}^d | Mx \geq c\}) \quad (77)$$

The final inequality is implied by Lipschitz-continuity of polytopes in Hausdorff distance with respect to RHS expansions (see Li (1993)). Therefore:

$$\tilde{x}_n - x_n \xrightarrow{p} 0 \quad (78)$$

We now wish to show that  $p'x_n \xrightarrow{p} p'x^*$ , where  $x^*$  is some value from  $\mathcal{A}(\theta_0)$ . For arbitrary  $\varepsilon > 0$  note that:

$$\mathbb{P}[|p'x_n + w_n \iota'(\hat{c}_n - \hat{M}_n x_n) - p'x^*| > \varepsilon] \geq \quad (79)$$

$$\mathbb{P}[p'x_n > p'x^* + \varepsilon - w_n \iota'(\hat{c}_n - \hat{M}_n x_n)] \geq \quad (80)$$

$$\mathbb{P}[p'x_n > p'x^* + \varepsilon] \quad (81)$$

As the LHS goes to 0 by (70), we have:

$$\mathbb{P}[p'x_n > p'x^* + \varepsilon] \rightarrow 0 \quad (82)$$



To prove the other side, note that, as  $\tilde{x}_n \in \Theta_I(\theta_0)$ , by definition of  $x^*$ , it must be that  $p'\tilde{x}_n \geq p'x^*$ . Therefore,

$$\mathbb{P}[p'x_n < p'x^* - \varepsilon] \leq \mathbb{P}[p'x_n < p'\tilde{x}_n - \varepsilon] \rightarrow 0, \quad (83)$$

where the RHS converges to 0 by (78) and CMT. We thus conclude that  $p'x_n \xrightarrow{p} p'x^*$  and, moreover,  $p'\tilde{x}_n \xrightarrow{p} p'x^*$ .

Notice that by Lemma 2, for a fixed, large enough  $w$  satisfying Assumption A1 Lemma 6.3 applies, where one sets  $f(x) = L(x; \theta_0, w)$ ,  $B = \mathcal{A}(\theta_0)$  with  $f(x^*) = p'x^*$  for any  $x^* \in \mathcal{A}(\theta_0)$ . Thus,  $\exists x_n^* \in \mathcal{A}(\theta_0)$  such that  $\|x_n - x_n^*\| \xrightarrow{p} 0$ . Therefore,  $\exists \delta_n \rightarrow 0^+$  such that:

$$\mathbb{P}[\|x_n - x_n^*\| < \delta_n] \rightarrow 1 \quad (84)$$

Recall that the sequence  $x_n$  was arbitrarily selected from  $\mathcal{A}(\hat{\theta}_n, w_n)$ , and we can, for example, select a measurable  $\{x_n\}_{n=1}^\infty$  (by Theorem 18.19 in Aliprantis and Border (2007)):

$$x_n \in \arg \max_{x \in \mathcal{A}(\hat{\theta}_n, w_n)} d(x, \mathcal{A}(\theta_0)) \quad (85)$$

For such  $x_n$ , we get:

$$\|x_n - x_n^*\| < \delta_n \implies d(x, \mathcal{A}(\theta_0)) < \delta_n \quad \forall x \in \mathcal{A}(\hat{\theta}_n, w_n) \quad (86)$$

So:

$$\mathbb{P}[\mathcal{A}(\hat{\theta}_n, w_n) \subseteq \mathcal{A}(\theta_0)^{\delta_n}] \geq \mathbb{P}[\|x_n - x_n^*\| < \delta_n] \rightarrow 1 \quad (87)$$

This establishes the existence of a deterministic  $\delta_n \rightarrow 0^+$  such that  $\mathcal{A}(\hat{\theta}_n, w_n) \subseteq \mathcal{A}(\theta_0)^{\delta_n}$  w.p.a.1.

2. By (87) and using the representation found in Proposition 6.3 we have that:

$$\inf_{x \in \mathcal{X}} L(x; \hat{\theta}_n, w_n) = \inf_{x \in \mathcal{A}(\theta_0)^{\delta_n}} L(x; \hat{\theta}_n, w_n) + o_p(1) \quad (88)$$

$$= \min_{x \in \mathcal{A}(\theta_0)^{\delta_n}} p'x + w_n \max_{j \in \{1, 2^q\}} \left\{ \sum_{i \in \Pi_j} (\hat{c}_{ni} - \hat{M}'_{ni}x) \right\} + o_p(1), \quad (89)$$

where  $o_p(1)$  encompasses realizations at which  $\mathcal{A}(\hat{\theta}_n, w_n) \not\subseteq \mathcal{A}(\theta_0)^{\delta_n}$  or where  $\hat{\theta}_n$  is not in a fixed open vicinity of  $\theta_0$  that was argued to exist in Proposition 6.3. Suppose that at  $\theta_0$  the constraints that do not bind at any  $x \in \mathcal{A}(\theta_0)$  are given by  $I \subseteq \{1, 2, \dots, q\}$ .

By continuity, it follows that  $\exists \delta > 0$  and  $\varepsilon > 0$  such that:

$$c_i - M_i x < -\varepsilon, \forall i \in I \quad (90)$$

for any  $x \in \mathcal{A}(\theta_0)^\delta$ . From (87) it then also follows that:

$$\inf_{x \in \mathcal{X}} L(x; \hat{\theta}_n, w_n) = \min_{x \in \mathcal{A}(\theta_0)^{\delta_n}} p'x + w_n \max_{\Pi \in 2^{\overline{1, q}} \setminus I} \left\{ \sum_{i \in \Pi_j} (\hat{c}_{ni} - \hat{M}'_{ni} x) \right\} + o_p(1) \quad (91)$$

3. Consider the problem in the linear programming representation found in Proposition 6.3, which it admits w. p. 1 as.:

$$\inf_{x \in \mathcal{X}} L(x; \hat{\theta}_n; w_n) = \min_{t, x} t \quad \text{s.t.:} \quad \begin{cases} t \in [\underline{t}; \bar{t}] \\ x \in \mathcal{X} \\ p'x + \sum_{i \in \Pi_j} w_n (\hat{c}_{ni} - \hat{M}'_{ni} x) \leq t, j \in \overline{1, 2^q} \end{cases} \quad (92)$$

The Lagrangian reads as:

$$\mathcal{L} = t + \sum_{\Pi \in 2^{\overline{1, q}}} \lambda_\Pi \left( p'x - t + w_n \sum_{j \in \Pi} \hat{c}_{nj} - \hat{M}'_{nj} x \right), \quad (93)$$

Where the constraints  $x \in \mathcal{X}$  and  $t \in [\underline{t}; \bar{t}]$  are omitted, as they are not binding with probability 1 as. This holds, as  $\mathcal{A}(\theta_0) \subseteq \text{Int}(\mathcal{X})$  and  $B(\theta_0) \in \text{Int}([\underline{t}; \bar{t}])$  by assumption. Because  $\mathcal{A}(\theta_0)$  is compact, there further exists<sup>32</sup> a  $\bar{\delta} > 0$ :  $\mathcal{A}(\theta_0)^{\bar{\delta}} \subseteq \text{Int}(\mathcal{X})$  and as  $\tilde{\mathcal{A}}(\hat{\theta}_n; w_n) \subseteq \mathcal{A}(\theta_0)^{\delta_n}$  w.p. 1 as. for some  $\delta_n \rightarrow 0^+$ , it follows that w.p.a.1  $\mathcal{A}(\hat{\theta}_n; w_n) \subseteq \text{Int}(\mathcal{X})$ . Similar argument establishes that  $t_n^* \in \text{Int}([\underline{t}; \bar{t}])$  w.p.a.1. In what follows, we will simply call such optimal pairs *interior*.

Differentiating with respect to  $t$ , one notes that:

$$\sum_{\Pi} \lambda_\Pi = 1 \quad (94)$$

Next, at any *interior* optimal  $t, x$ :

$$t = p'x + w_n \max_{\Pi} \sum_{j \in \Pi} (\hat{c}_{nj} - \hat{M}'_{nj} x) \quad (95)$$

---

<sup>32</sup>To see that, consider  $A, B \subseteq \mathbb{R}^d$  such that  $A$  is compact,  $B$  is open and  $A \subseteq B$ . Since  $B$  is open, for any  $b \in B \exists \varepsilon > 0 : B_\varepsilon(b) \subseteq B$ . This defines an open cover of  $A$ , as  $A \subseteq \bigcup_{b \in B} B_{\varepsilon_b/2}(b)$ . Since  $A$  is compact, for any cover there exists a finite subcover, i.e.  $\exists (b_k, \varepsilon_{b_k}/2)_{k=1}^K$  such that  $b_k \in B$  and  $A \subseteq \bigcup_{k=1}^K B_{\varepsilon_{b_k}/2}(b_k)$ . Take  $\delta = \min_k \varepsilon_{b_k}/2$ . Then, pick any  $x \in A^\delta$ . It follows that  $\exists y \in A$ :  $\|x - y\| < \delta$ . Because  $y \in A$ , there further  $\exists k$ :  $\|y - b_k\| \leq \varepsilon_{b_k}/2$ . Thus,  $\|x - b_k\| \leq \|y - b_k\| + \|x - y\| < \varepsilon_{b_k}/2 + \delta \leq \varepsilon_{b_k}$ , and so  $x \in B_{\varepsilon_{b_k}}(b_k) \subseteq B$ .

To see that, note that by contradiction, if:

$$t > p'x + w_n \max_{\Pi} \sum_{j \in \Pi} (\hat{c}_{nj} - \hat{M}'_{nj}x) \quad (96)$$

Then, as we assumed that the pair  $(t, x)$  is *interior*, there exists  $\tilde{t} < t$  such that the pair  $(\tilde{t}, x)$  satisfies all the constraints. Therefore,  $(t, x)$  is not optimal. The other direction of the inequality is infeasible, and so the equality must hold. Moreover, since  $\Pi$  may be empty, we also have at any optimal  $x$ :

$$t \geq p'x \quad (97)$$

Furthermore, the problem has a solution w.p.a.1, and therefore it has a vertex-solution, i.e. a solution that is pinned down by a matrix of binding constraints of full column-rank. Because w.p.a.1 any solution is *interior*, any such matrix w.p.a.1 does not feature constraints  $x \in \mathcal{X}, t \in [\underline{t}, \bar{t}]$ . The only constraints that can be satisfied at such vertex-solution with an equality are of the following type:

$$p'x - t = w_n \sum_{j \in \Pi_k} \hat{c}_{nj} - \hat{M}'_{nj}x, \quad k \in \tilde{J} \quad (98)$$

for some  $\tilde{J} \subseteq \overline{1, q} : |\tilde{J}| \geq d + 1$ , where the latter inequality holds by definition of a vertex of a linear program<sup>33</sup>. One can write the complete set of the binding constraints (98) as:

$$\hat{R}_{\tilde{J}n} \begin{pmatrix} t \\ x \end{pmatrix} = \hat{r}_{\tilde{J}n}, \quad (99)$$

where the  $|\tilde{J}| \times (d + 1)$  matrix  $\hat{R}_{\tilde{J}n}$  is of full column rank and the system yields a unique solution  $t_n^*, x_n^*$ .

4. Denote the set of all vertices  $(t^*, x^*)$  that satisfy (98) with  $|\tilde{J}| \geq d + 1$  and a full-column-rank  $\hat{R}_{\tilde{J}n}$  at a given  $\hat{\theta}_n$  by  $\mathcal{V}^*(\hat{\theta}_n)$ . From the previous arguments it follows that  $\mathcal{V}^*(\cdot)$  is non-empty w.p.a.1 and finite, because any finite-dimensional polytope has finitely many vertices and therefore the corresponding LP has finitely many optimal vertices. We will write  $\mathcal{V}_x^*(\hat{\theta}_n)$  for the projection of that set on the  $x$ -coordinates. For any vertex-solution  $(t^*, x^*) \in \mathcal{V}^*(\hat{\theta}_n)$ , suppose constraints  $V^* \subseteq \{1, \dots, q\}$  are violated at it, meaning that:

$$V^*(\hat{\theta}_n, x^*) \equiv \{j \in \overline{1, q} | \hat{c}_{nj} - \hat{M}'_{nj}x^* > 0\} \quad (100)$$

---

<sup>33</sup>Any finite feasible LP has a vertex-solution, at which the matrix of binding constraints has full rank, so that its dimension is at least that of  $(t \ x)'$ .

For brevity, we will write  $V_n^* \equiv V^*(\hat{\theta}_n, x_n^*)$  where  $t_n^*, x_n^* \in \mathcal{V}^*(\hat{\theta}_n)$  is some (measurable) sequence of optimal vertices. Note that:

$$t_n^* = p'x_n^* + w_n \max_{\Pi} \sum_{j \in \Pi} (\hat{c}_{nj} - \hat{M}'_{nj}x_n^*) = p'x_n^* + w_n \sum_{j \in V_n^*} (\hat{c}_{nj} - \hat{M}'_{nj}x_n^*) \quad (101)$$

Consider (98) and suppose  $\tilde{J}_n = \tilde{J}(t_n^*, x_n^*)$  with  $|\tilde{J}_n| \geq d+1$  is the set of the corresponding subsets, i.e.:

$$t_n^* = p'x_n^* + w_n \sum_{j \in \Pi_i} (\hat{c}_{nj} - \hat{M}'_{nj}x_n^*) \quad \forall i \in \overline{1, k} \quad (102)$$

It must be that  $V_n^* \subseteq \Pi_i \quad \forall i \in \tilde{J}_n$ , because  $j \notin V_n^* \implies (\hat{c}_{nj} - \hat{M}'_{nj}x_n^*) \leq 0$ , and so we have:

$$\sum_{j \in V_n^*} (\hat{c}_{nj} - \hat{M}'_{nj}x_n^*) = \sum_{j \in \Pi_i} (\hat{c}_{nj} - \hat{M}'_{nj}x_n^*) \leq \sum_{j \in \Pi_i \cap V_n^*} (\hat{c}_{nj} - \hat{M}'_{nj}x_n^*), \quad (103)$$

where the first equality follows from (102) and (101). We now proceed by contradiction. Suppose that  $\exists j : j \in V_n^* \cap \Pi'_i$  (where the complement is taken with respect to  $\overline{1, q}$ ), then:

$$\begin{aligned} \sum_{j \in \Pi_i \cap V_n^*} (\hat{c}_{nj} - \hat{M}'_{nj}x_n^*) &< \sum_{j \in \Pi_i \cap V_n^*} (\hat{c}_{nj} - \hat{M}'_{nj}x_n^*) + \sum_{j \in \Pi'_i \cap V_n^*} (\hat{c}_{nj} - \hat{M}'_{nj}x_n^*) = \\ &= \sum_{j \in V_n^*} (\hat{c}_{nj} - \hat{M}'_{nj}x_n^*), \end{aligned}$$

which yields a contradiction with (103), so there can be no such  $j$ . In light of (103) it then also follows that  $\forall i \in \tilde{J}_n$  and  $\forall j \in \Pi_i \cap V_n^*$  it must be that:

$$\hat{c}_{nj} - \hat{M}'_{nj}x_n^* = 0 \quad \forall j \in \Pi_k \setminus V_n^* \quad (104)$$

Therefore, the complete system described by equation (102), is equivalent to:

$$\begin{cases} \hat{c}_{nj} - \hat{M}'_{nj}x_n^* = 0 \quad \forall i \in \tilde{J}_n : \Pi_i \neq V_n^*, \quad \forall j \in \Pi_i \setminus V_n^* \\ t_n^* = p'x_n^* + w_n \sum_{j \in V_n^*} \hat{c}_{nj} - \hat{M}'_{nj}x_n^* \end{cases} \quad (105)$$

From the representation (99), we know that the matrix corresponding to system (105) must be of full column rank,  $d+1$ . Dropping the equation defining  $t_n^*$ , it implies that there exists at least  $d$  linearly independent equations of form:

$$\hat{c}_{nj} - \hat{M}'_{nj}x_n^* = 0$$

We denote the set of all binding constraints by  $\Pi^*(\hat{\theta}_n, x_n^*) \equiv \{j \in \overline{1, q} \mid \hat{c}_{nj} - \hat{M}'_{nj} x_n^* = 0\}$ , which we shall occasionally write as  $\Pi_n^*$  for brevity. We thus have:

$$|\Pi_n^*| \geq d, \quad \text{rk}(\hat{M}_{\Pi_n^*}) = d \quad (106)$$

5. Consider two collections of sets:

$$\mathcal{E} \equiv \{A \subseteq 2^{[q]} : M_A x \neq c_A \ \forall x \in \mathcal{A}(\theta_0)\} \quad (107)$$

$$\mathcal{F} \equiv \{A \subseteq 2^{[q]} : p \notin \mathcal{R}(M'_A)\} \quad (108)$$

We shall now consider two events  $E_n$  and  $F_n$ :

$$E_n \equiv \{\Pi_n^* \in \mathcal{E}\}, \quad F_n \equiv \{\Pi_n^* \in \mathcal{F}\} \quad (109)$$

$$(110)$$

We wish to show that  $\mathbb{P}[E_n] \rightarrow 0$  and  $\mathbb{P}[F_n] \rightarrow 0$  and therefore  $\mathbb{P}[E'_n \cap F'_n] \rightarrow 1$ .

a) Let us consider  $E_n$  first. Since  $\mathcal{A}(\theta_0)$  is compact, for a fixed set  $A \in \mathcal{E}$ , the condition  $M_A x \neq c_A \ \forall x \in \mathcal{A}(\theta_0)$  implies that there exists  $\varepsilon(A) > 0$ :

$$\inf_{x \in \mathcal{A}(\theta_0)} \|M_A x - c_A\| > \varepsilon(A) \quad (111)$$

Because  $E$  is a finite collection of sets, we can pick  $\varepsilon = \min_{A \in E} \varepsilon(A)$ , so that:

$$\min_{A \in \mathcal{E}} \inf_{x \in \mathcal{A}(\theta_0)} \|M_A x - c_A\| > \varepsilon \quad (112)$$

By continuity of the objective function in  $x$ , there further  $\exists \kappa > 0$ , such that:

$$\min_{A \in \mathcal{E}} \inf_{x \in \mathcal{A}^\kappa(\theta_0)} \|M_A x - c_A\| > \frac{\varepsilon}{2} \quad (113)$$

We now consider:

$$\mathbb{P}[E_n] \leq \mathbb{P} \left[ \|\hat{M}_{\Pi_n^*} x_n^* - \hat{c}_{\Pi_n^*}\| = 0, \inf_{x \in \mathcal{A}^\kappa(\theta_0)} \|M_{\Pi_n^*} x - c_{\Pi_n^*}\| > \frac{\varepsilon}{2} \right] \quad (114)$$

Observe that for any non-empty  $A \subseteq [q]$ , by Cauchy-Schwartz and triangle inequali-

ties:

$$\begin{aligned} & \|(\hat{M}_{nA}x_n^* - \hat{c}_{nA})\| = \\ & \left\| (M_Ax_n^* - c_A) - \left( (\hat{c}_{nA} - c_A) + (M_A - \hat{M}_{nA})x_n^* \right) \right\| \geq \\ & \|M_Ax_n^* - c_A\| - \left\| \hat{M}_{nA} - M_A \right\| \|x\|_\infty - \|\hat{c}_{nA} - c_A\| \end{aligned}$$

We can thus further rewrite:

$$\begin{aligned} & \mathbb{P} \left[ \left\| \hat{M}_{\Pi_n^*}x_n^* - \hat{c}_{\Pi_n^*} \right\| \leq 0, \inf_{x \in \mathcal{A}^\kappa(\theta_0)} \|M_{\Pi_n^*}x - c_{\Pi_n^*}\| > \frac{\varepsilon}{2} \right] \leq \\ & \mathbb{P} \left[ \left\| M_{\Pi_n^*}x_n^* - c_{\Pi_n^*} \right\| \leq \eta_n, \inf_{x \in \mathcal{A}^\kappa(\theta_0)} \|M_{\Pi_n^*}x - c_{\Pi_n^*}\| > \frac{\varepsilon}{2} \right], \end{aligned}$$

where  $\eta_n \equiv \left\| \hat{M}_{\Pi_n^*} - M_{\Pi_n^*} \right\| \|x\|_\infty + \left\| \hat{c}_{\Pi_n^*} - c_{\Pi_n^*} \right\| = o_p(1)$ . Finally, using  $\mathbb{P}[A \cap B'] + \mathbb{P}[A \cap B] = \mathbb{P}[A]$ :

$$\begin{aligned} & \mathbb{P} \left[ \left\| M_{\Pi_n^*}x_n^* - c_{\Pi_n^*} \right\| \leq \eta_n, \inf_{x \in \mathcal{A}^\kappa(\theta_0)} \|M_{\Pi_n^*}x - c_{\Pi_n^*}\| > \frac{\varepsilon}{2} \right] = \\ & \mathbb{P} \left[ \left\| M_{\Pi_n^*}x_n^* - c_{\Pi_n^*} \right\| \leq \eta_n, \inf_{x \in \mathcal{A}^\kappa(\theta_0)} \|M_{\Pi_n^*}x - c_{\Pi_n^*}\| > \frac{\varepsilon}{2}, x_n^* \in \mathcal{A}^\kappa(\theta_0) \right] + \\ & + \mathbb{P} \left[ \left\| M_{\Pi_n^*}x_n^* - c_{\Pi_n^*} \right\| \leq \eta_n, \inf_{x \in \mathcal{A}^\kappa(\theta_0)} \|M_{\Pi_n^*}x - c_{\Pi_n^*}\| > \frac{\varepsilon}{2}, x_n^* \notin \mathcal{A}^\kappa(\theta_0) \right], \end{aligned}$$

where the second term is  $o(1)$  by Step 1 of the proof. Finally, we note that  $x_n^* \in \mathcal{A}^\kappa(\theta_0) \implies \left\| M_{\Pi_n^*}x_n^* - c_{\Pi_n^*} \right\| > \varepsilon/2$ , which, combined with  $\left\| M_{\Pi_n^*}x_n^* - c_{\Pi_n^*} \right\| \leq \eta_n$ , further implies that  $\eta_n > \varepsilon/2 > 0$ , so that:

$$\begin{aligned} & \mathbb{P} \left[ \left\| M_{\Pi_n^*}x_n^* - c_{\Pi_n^*} \right\| \leq \eta_n, \inf_{x \in \mathcal{A}^\kappa(\theta_0)} \|M_{\Pi_n^*}x - c_{\Pi_n^*}\| > \frac{\varepsilon}{2}, x_n^* \in \mathcal{A}^\kappa(\theta_0) \right] \leq \\ & \mathbb{P} \left[ \frac{\varepsilon}{2} \leq \eta_n \right] = o(1) \end{aligned}$$

This concludes the proof.

- b) We now consider  $F_n$ . To do so, it is convenient to observe that the penalty function estimator and problem (92) are equivalent to yet another LP:

$$B(\hat{\theta}_n) + o_p(1/\sqrt{n}) = \min_{x,a} p'x + w_n t'a \quad \text{s.t. :} \begin{cases} a \geq 0 \\ a \geq \hat{c}_n - \hat{M}_n x \end{cases} \quad (115)$$

Note that we drop the constraints corresponding to  $x \in \mathcal{X}$  in (115), and  $o_p(1/\sqrt{n})$

accommodates the potential non-existence of the interior solution. Write Lagrangian:

$$\mathcal{L} = p'x + w_n t' a + \mu'(\hat{c}_n - \hat{M}_n x - a) - \omega' a$$

The KKT conditions at an interior optimum are:

$$p = \hat{M}'_n \mu \quad (116)$$

$$w_n = \omega + \mu \quad (117)$$

$$\omega' a = 0 \quad (118)$$

$$\mu'(\hat{c}_n - \hat{M}_n x - a) = 0 \quad (119)$$

$$a \geq \hat{c}_n - \hat{M}_n x \quad (120)$$

$$a \geq 0, \omega \geq 0, \mu \geq 0 \quad (121)$$

Analyzing the above system, one observes that if at  $x_n^* \in \mathcal{V}_x^*(\theta_n)$  a constraint is violated,  $j \in V_n^*$ , then  $a_j > 0$ , and so  $\omega_j = 0$ , which implies  $\mu_j = w_n$ . If  $\hat{M}_{nj}x_n^* - \hat{c}_{nj} > 0$ , then  $\hat{c}_{nj} - \hat{M}_{nj}x_n^* - a_j < 0$ , and so  $\mu_j = 0$ . Finally, if  $j \in \Pi_n^*$ , then  $\mu_j \in [0; w_n]$ . Therefore, (116) rewrites as:

$$p = w_n \sum_{j \in V_n^*} \hat{M}'_{nj} + \sum_{j \in \Pi_n^*} \hat{M}'_{nj} \mu_j \quad (122)$$

Since  $\mu_j \leq w_n$  and as  $\hat{M}_n - M = O_p(1/\sqrt{n})$ , we have:

$$p = w_n \sum_{j \in V_n^*} M'_j + \sum_{j \in \Pi_n^*} M'_j \mu_j + O_p\left(\frac{w_n}{\sqrt{n}}\right) \quad (123)$$

Consider a projection  $P_{\Pi_n^*}$  from  $\mathbb{R}^d$  onto  $\mathcal{R}(M'_{\Pi_n^*})$ . For example, one can construct it as  $M'_{\Pi_n^*} (M'_{\Pi_n^*})^\dagger$ , where  $\dagger$  denotes a Moore-Penrose pseudoinverse. We can write:

$$p - O_p\left(\frac{w_n}{\sqrt{n}}\right) = w_n (I - P_{\Pi_n^*}) \sum_{j \in V^*} M'_j + \underbrace{w_n P_{\Pi_n^*} \sum_{j \in V^*} M'_j + \sum_{j \in \Pi_n^*} M'_j \mu_j}_{T_n \in \mathcal{R}(M'_{\Pi_n^*})} \quad (124)$$

Notice that, if  $\sum_{j \in V^*} M'_j \notin \mathcal{R}(M'_{\Pi_n^*})$ , then the RHS of (124) has unbounded norm:

$$\begin{aligned} & \left\| w_n (I - P_{\Pi_n^*}) \sum_{j \in V_n^*} M'_j + T_n \right\|^2 = \\ & = w_n^2 \left\| (I - P_{\Pi_n^*}) \sum_{j \in V_n^*} M'_j \right\|^2 + \|T_n\|^2 \end{aligned} \quad (125)$$

Since the square norm of the LHS of (124) is bounded from above by  $\|p\|^2 +$

$O_p(\frac{w_n^2}{n}) = \|p\|^2 + o_p(1)$ , (125) will contradict the equality in (124) w.p.a.1. Suppose, alternatively, that  $\exists v : \sum_{j \in V_n^*} M'_j = M'_{\Pi_n^*} v$ . Equation (124) rewrites:

$$p - O_p\left(\frac{w_n}{\sqrt{n}}\right) = M'_{\Pi_n^*}(\mu_{\Pi_n^*} + w_n v),$$

which implies, for example, that:

$$(I - P_{\Pi_n^*})p + P_{\Pi_n^*}p - M'_{\Pi_n^*}(\mu_{\Pi_n^*} + w_n v) = O_p\left(\frac{w_n}{\sqrt{n}}\right) \quad (126)$$

The norm of the LHS of (126) must go to 0, however, if  $p \notin \mathcal{R}(M'_{\Pi_n^*})$ , we have, by orthogonality:

$$\|(I - P_{\Pi_n^*})p\|^2 + \|P_{\Pi_n^*}p - M'_{\Pi_n^*}(\mu_{\Pi_n^*} + w_n v)\|^2 \geq \|(I - P_{\Pi_n^*})p\|^2 > 0,$$

which will also yield a contradiction w.p.a.1. To complete the proof, one applies the same probabilistic arguments as used in step 5.a above, which we omit here. Thus,  $\mathbb{P}[F_n] \rightarrow 0$ .

6. We define the *correct set of vertices*,  $\mathcal{G}$ , as follows:

$$\mathcal{G} \equiv \{A \subseteq [q] : \exists x \in \mathcal{A}(\theta_0) \text{ s.t. } M_A x = c_A, p \in \mathcal{R}(M'_A)\}$$

In line with previous notation, let  $G_n \equiv \{\Pi_n^* \in \mathcal{G}\}$ . The results of point 5 imply that  $\mathbb{P}[E'_n \cap F'_n] = \mathbb{P}[G_n] \rightarrow 1$ .

Consider any  $A \in \mathcal{G}$ . Suppose  $p = M'_A v$  for some  $v \in \mathbb{R}^{|A|}$ . Further, fix any  $x \in \mathcal{A}(\theta_0) : M_A x = c_A$ , then:

$$B(\theta_0) = p'x = v'M_A x = v'c_A \quad (127)$$

The conclusion then follows from the following chain of equalities:

$$G_n \implies p'x_n^* - B(\theta_0) = v'M_{\Pi_n^*}x_n^* - v'c_{\Pi_n^*} = \quad (128)$$

$$= v'\hat{M}_{\Pi_n^*}x_n^* - v'c_{\Pi_n^*} + v'(M_{\Pi_n^*} - \hat{M}_{\Pi_n^*})x_n^* = \quad (129)$$

$$= v'(\hat{c}_{\Pi_n^*} - c_{\Pi_n^*}) + v'(M_{\Pi_n^*} - \hat{M}_{\Pi_n^*})x_n^* \quad (130)$$

Finally, from (130), applying the triangle and Cauchy-Shwartz inequalities as well as



noting that over the event  $G_n$  one has  $\Pi_n^* \in \mathcal{G}$  by definition, it follows that:

$$G_n \implies |p'x_n^* - B(\theta_0)| \leq \varpi_n \equiv \max_{A \in \mathcal{G}} \left\{ \left( \|\hat{c}_A - c_A\| + \|x\|_\infty \|M_A - \hat{M}_A\| \right) \cdot \min_{v \in \mathbb{R}^{|A|}: M'_A v = p} \|v\| \right\} \quad (131)$$

One concludes by noting that the RHS is clearly  $O_p(1/\sqrt{n})$ , as  $\mathcal{G}$  is finite and  $\hat{\theta}_n - \theta_0 = O_p(1/\sqrt{n})$  by assumption. Formally, for any  $\varepsilon > 0$ :

$$\mathbb{P}[r_n |p'x_n^* - B(\theta_0)| > \varepsilon] = \mathbb{P}[r_n |p'x_n^* - B(\theta_0)| > \varepsilon, G_n] + o(1) \leq \quad (132)$$

$$\mathbb{P}[r_n \varpi_n > \varepsilon, G_n] + o(1) \leq \mathbb{P}[r_n \varpi_n > \varepsilon] + o(1) \quad (133)$$

and  $r_n \varpi_n = O_p(\frac{r_n}{\sqrt{n}})$  for any  $r_n \rightarrow \infty$ , where we used the fact that  $\mathbb{P}[G'_n \cap O_n] \leq \mathbb{P}[G'_n] = o(1)$  for any measurable  $O_n$ . Recalling that the choice of  $x_n^* \in \mathcal{V}_x^*(\hat{\theta}_n)$  was arbitrary and that neither  $\varpi_n$ , nor the  $o(1)$  depend on  $x_n^*$ , one gets:

$$\sup_{x \in \mathcal{V}_x^*(\hat{\theta}_n)} |p'x - B(\theta_0)| = O_p(1/\sqrt{n}) \quad (134)$$

But because any  $x \in \mathcal{A}(\hat{\theta}_n; w_n)$  can be represented as a convex combination of vertices,  $\{x_j\}_{j=1}^K \subseteq \mathcal{V}_x^*(\hat{\theta}_n)$ , as:  $x = \sum_j \omega_j x_j$ , where  $\omega_j \in [0; 1]$  and  $\sum_j \omega_j = 1$ . Using that, applying the triangle inequality and taking maximum, one gets, for any  $x \in \mathcal{A}(\hat{\theta}_n; w_n)$ :

$$|p'x - B(\theta_0)| = \left| \sum_j \omega_j (p'x_j - B(\theta_0)) \right| \leq \max_j |p'x_j - B(\theta_0)| \leq \sup_{x \in \mathcal{V}_x^*(\hat{\theta}_n)} |p'x - B(\theta_0)| = O_p(1/\sqrt{n})$$

taking supremum on the left hand side establishes the claim of the theorem. ■

## 6.6. Proof of Theorem 2.3

*Proof.* For  $x \in \mathbb{R}^{qd}$ , define the inverse-vectorization operator as

$$\text{vec}_{q \times d}^{-1}(x) \equiv (\text{vec}(I_d)' \otimes I_q) (I_d \otimes x).$$

Further, define selector matrices  $C_c$  and  $C_M$  that select the  $c$  and  $M$  components of  $\theta$  respectively:

$$C_c \theta = c, \quad C_M \theta = \text{vec}(M).$$

Moreover, for a subset of rows  $A \subseteq \{1, 2, \dots, q\}$ , define the row-selector  $C(A)$  as

$$C(A)M = M_A, \quad C(A)c = c_A.$$

We first work on  $\mathcal{D}_n^{(1)}$ . From Step 5.b in the proof of Theorem 2.2, it follows that solving the penalized problem is w.p.a.1 equivalent to solving a relaxed LP, i.e., w.p.a.1,

$$\tilde{\mathcal{A}}(\hat{\theta}_1; w_{n_1}) = \min_{x \in \mathbb{R}^d, a \in \mathbb{R}^q} p'x + w_{n_1}t'a, \quad \text{s.t.: } a \geq \hat{c}^{(1)} - \hat{M}^{(1)}x, \quad a \geq 0. \quad (135)$$

Denote the set of vertex-solutions of (135) by  $\hat{\mathcal{V}}_x$  and define

$$\hat{x} \in \arg \min_{x \in \hat{\mathcal{V}}_x} p'x, \quad \hat{A} = J(\hat{x}; \hat{\theta}^{(1)}).$$

From Step 6 of the proof of Theorem 2.2 it follows that  $\hat{A} \in \mathbb{A}$  w.p.a.1.

For a nonempty  $A \in 2^{[q]}$  and any  $\tilde{M} \in \mathbb{R}^{q \times d}$ , define

$$\mathbb{S}(A, \tilde{M}) \equiv \arg \min_{v \in \mathbb{R}^{|A|}; \|v\| \leq \check{v}} \|p - \tilde{M}'_A v\|^2$$

The optimization problem above is continuous in  $\tilde{M}$ , the constraint correspondence is constant and compact. Hence, for any nonempty  $A \in 2^{[q]}$ , by Berge's Maximum Theorem,  $\mathbb{S}(A, \cdot)$  is compact, nonempty, and upper-hemicontinuous (see Theorem 17.31 in Aliprantis and Border (2007)).

Because  $\hat{M}^{(1)} \xrightarrow{p} \hat{M}^{(1)}$ , for any nonempty  $A \subseteq 2^{[q]}$  it follows by the usual M-estimation argument<sup>34</sup>, there exists a deterministic  $s_n(A) \downarrow 0$ , such that

$$\mathbb{S}(A, \hat{M}^{(1)}) \subseteq \mathbb{S}(A, M)^{s_n(A)}, \quad \text{w.p.a.1,} \quad (136)$$

so that also

$$\mathbb{S}(\hat{A}, \hat{M}^{(1)}) \subseteq \mathbb{S}(\hat{A}, M)^{s_n(\hat{A})} \quad \text{w.p.a.1} \quad (137)$$

Observing that the objective function is convex,  $\mathbb{S}(A, \tilde{M})$  is also convex-valued for any nonempty  $A \in 2^{[q]}$  and  $\tilde{M} \in \mathbb{R}^{q \times d}$ .

Define some measurable  $\check{v} \in \mathbb{S}(\hat{A}, \hat{M}^{(1)})$  and denote its projection onto  $\mathbb{S}(\hat{A}, M)$  by  $\tilde{v}_n$ . Both  $\check{v}$  is a random sequence, but we suppress the dependence on  $n$  for simplicity.  $\tilde{v}_n$  is well-defined by the Hilbert Projection Theorem.  $\check{v}$  is well-defined, and  $\tilde{v}_n$  is measurable by

---

<sup>34</sup>Theorem 18.19 in Aliprantis and Border (2007) establishes measurability of  $\mathbb{S}(\hat{A}, \hat{M}^{(1)})$  and  $\mathbb{S}(\hat{A}, M)$

Theorem 18.19 in Aliprantis and Border (2007). From (137) it follows that

$$\|\check{v} - \tilde{v}_n\| = o_p(1) \quad (138)$$

By definition,  $A \in \mathbb{A}$  implies that  $\exists v^* \in \mathcal{S}_A$  such that  $\|v^*\| \leq \max_{\hat{A} \in \mathbb{A}} \min_{v \in \mathcal{S}_{\hat{A}}} \|v\| \leq \bar{v}$ , where the last inequality is by Assumption B3. This implies that, for any  $A \in \mathbb{A}$ ,

$$\inf_{v \in \mathbb{R}^{|A^*|}: \|v\| \leq \bar{v}} \|p - M'_{A^*} v\|^2 = \min_{v \in \mathbb{R}^{|A^*|}: \|v\| \leq \bar{v}} \|p - M'_{A^*} v\|^2 = 0. \quad (139)$$

By (139), we have  $\mathbb{S}(A, M) \subseteq \mathcal{S}_A$  if  $A \in \mathbb{A}$ . Therefore,  $\tilde{v}_n \in \mathcal{S}_{\hat{A}}$  w.p.a.1.

The following Lemma justifies our construction:

**Lemma 6.4.** Suppose  $\hat{A} \in \mathbb{A}$ . Then,

$$\tilde{v}'_n c_{\hat{A}} = B(\theta_0), \quad (140)$$

$$\tilde{v}'_n M_{\hat{A}} \hat{x} = p' \hat{x}. \quad (141)$$

*Proof.* If  $\hat{A} \in \mathbb{A}$ , condition (10) holds for some  $x \in \mathcal{A}(\theta_0)$  such that  $M_{\hat{A}} x = c_{\hat{A}}$ . Since such  $x$  is a minimizer, it follows that  $p' x = B(\theta_0)$ . As  $\tilde{v}_n \in \mathcal{S}_{\hat{A}}$ , we have  $p = M'_{\hat{A}} \tilde{v}_n$ . Taking transpose and multiplying by  $\hat{x}$  yields (141). To show (140), write:

$$p' x = \tilde{v}'_n M_{\hat{A}} x = \tilde{v}'_n c_{\hat{A}}$$

This concludes the proof of the Lemma. ■

To avoid dealing with changing dimension, we let  $\hat{v} \in \mathbb{R}^q$  be such that  $\hat{v}_{\hat{A}} = \check{v}$  and  $\hat{v}_j = 0$  if  $j \notin \hat{A}$ . Similarly, define  $\hat{v}_n: (\hat{v}_n)_{\hat{A}} = \tilde{v}_n$  and  $(\hat{v}_n)_j = 0$  if  $j \notin \hat{A}$ . Note that  $\|\hat{v} - \hat{v}_n\| = \|\check{v} - \tilde{v}_n\|$ .

Equipped with  $\hat{v}$ ,  $\hat{A}$  and  $\hat{x}$ , we can now move onto the second fold. For  $(A, v, x) \in 2^{\overline{1,q}} \setminus \{\emptyset\} \times \mathbb{R}^q \times \mathcal{X}$ , define

$$H_n(A, x, v) \equiv \frac{\sqrt{n_2}}{\sigma_n(A, x, v)} v'_A \left( \hat{c}_A^{(2)} - c_A - (\hat{M}_A^{(2)} - M_A)x \right), \quad H_n \equiv H_n(\hat{A}, \hat{x}, \hat{v}).$$

Let  $Z_n^{(2)} \equiv \sqrt{n_2}(\hat{\theta}^{(2)} - \theta_0)$ . One can rewrite

$$\sqrt{n_2} v'_A \left( \hat{c}_A^{(2)} - c_A - (\hat{M}_A^{(2)} - M_A)x \right) = v'_A C(A) \left( C_c Z_n^{(2)} - \text{vec}_{q \times d}^{-1}(C_M Z_n^{(2)})x \right). \quad (142)$$

Applying the definition of  $\text{vec}_{q \times d}^{-1}$  and using bilinearity of Kronecker product, one notes that (142) is linear in  $Z_n^{(2)}$  and therefore, under Assumption B1, for any  $(A, v, x) \in 2^{\overline{1,q}} \setminus \{\emptyset\} \times \mathbb{R}^q \times \mathcal{X}$ ,

we have

$$\sqrt{n_2}v'_A \left( \hat{c}_A^{(2)} - c_A - (\hat{M}_A^{(2)} - M_A)x \right) \xrightarrow{d} \mathcal{N}(0, \sigma^2(A, x, v, \Sigma)),$$

where  $\sigma^2(\cdot)$  is given in Lemma (6.6).

By assumption B4 we then have, for a fixed optimal triplet  $A, x, v$  and CMT,

$$\frac{\sqrt{n_2}}{\sigma(A, x, v, \Sigma)} v'_A \left( \hat{c}_A^{(2)} - c_A - (\hat{M}_A^{(2)} - M_A)x \right) \xrightarrow{d} \mathcal{N}(0, 1) \quad (143)$$

We begin by taking the infeasible  $\hat{\sigma}_n(A, v, x) = \sigma(A, v, x, \Sigma)$ . Consider the set:

$$\mathfrak{N}(\underline{v}, \underline{\sigma}) \equiv \{(A, v, x) \in 2^{\overline{1,q}} \setminus \{\emptyset\} \times \mathbb{R}^q \times \mathcal{X} : \underline{v} \leq \|v_n\| \leq \bar{v}, \sigma(A, v, x, \Sigma) \geq \underline{\sigma}\} \quad (144)$$

We now fix an arbitrary deterministic sequence  $(A_n, v_n, x_n) \in \mathfrak{N}(\underline{v}, \underline{\sigma})$  for all  $n \in \mathbb{N}$  for some small  $\underline{v} > 0$  and  $\underline{\sigma} > 0$  that we pick below. Consider the limit (integration is with respect to  $\mathcal{D}_n^2$  only):

$$\lim_{n \rightarrow \infty} \mathbb{P}[H_n(A_n, v_n, x_n) \leq z_{1-\alpha}]$$

The space  $2^{\overline{1,q}} \setminus \{\emptyset\}$ , to which  $A_n$  belongs, is endowed with a discrete metric, and we consider the space  $\mathfrak{N}(\underline{v}, \underline{\sigma})$  as endowed with the maximum product metric  $\rho_\infty$ . It is straightforward to notice that  $\sigma(\cdot)$  is continuous in its first three arguments with respect to  $\rho_\infty$  even on the unrestricted space  $2^{\overline{1,q}} \setminus \{\emptyset\} \times \mathbb{R}^q \times \mathcal{X}$ , and thus  $\mathfrak{N}(\underline{v}, \underline{\sigma})$  is a compact space for any  $\underline{v} > 0, \underline{\sigma} > 0$ . It is also non-empty for some small enough  $\underline{v} > 0, \underline{\sigma} > 0$  by Assumption B4. Suppose  $\underline{v} > 0, \underline{\sigma} > 0$  are small enough and pick any convergent subsequence  $(A_{n_k}, v_{n_k}, x_{n_k}) \rightarrow (A, v, x)$ . Recall that:

$$H_n(A_n, v_n, x_n) = g(\sqrt{n_2}(\hat{\theta}^{(2)} - \theta_0), A_n, v_n, x_n) \quad (145)$$

for a continuous function  $g$  and:

$$\begin{pmatrix} \sqrt{(n_2)_k}(\hat{\theta}^{(2)} - \theta_0) \\ A_{n_k} \\ v_{n_k} \\ x_{n_k} \end{pmatrix} \xrightarrow{d} \begin{pmatrix} \mathcal{N}(0, \Sigma) \\ A \\ v \\ x \end{pmatrix} \quad (146)$$

we conclude that, by continuous mapping theorem, as  $k \rightarrow \infty$ :

$$g(\sqrt{(n_2)_k}(\hat{\theta}_{n_k}^{(2)} - \theta_0), A_{n_k}, v_{n_k}, x_{n_k}) = H_{n_k}(A_{n_k}, v_{n_k}, x_{n_k}) \xrightarrow{d} g(Z, A, v, x), \quad (147)$$

where  $Z \sim \mathcal{N}(0, \Sigma)$ . By (143), this implies:

$$\lim_{k \rightarrow \infty} \mathbb{P}[H_{n_k}(A_{n_k}, v_{n_k}, x_{n_k}) \leq z_{1-\alpha}] = 1 - \alpha \quad (148)$$

We claim that this further implies that:

$$\lim_{n \rightarrow \infty} \mathbb{P}[H_n(A_n, v_n, x_n) \leq z_{1-\alpha}] = 1 - \alpha \quad (149)$$

Suppose, by contradiction,  $\lim_{n \rightarrow \infty} \mathbb{P}[H_n(A_n, v_n, x_n) \leq z_{1-\alpha}] \neq 1 - \alpha$ . It means that  $\exists \varepsilon > 0$  such that  $\forall N \in \mathbb{N} \exists n \geq N$  such that:

$$|\mathbb{P}[H_n(A_n, v_n, x_n) \leq z_{1-\alpha}] - (1 - \alpha)| > \varepsilon \quad (150)$$

Thus, we can construct a subsequence  $n_k$  such that:

$$|\mathbb{P}[H_{n_k}(A_{n_k}, v_{n_k}, x_{n_k}) \leq z_{1-\alpha}] - (1 - \alpha)| > \varepsilon \quad (151)$$

for all  $k \in \mathbb{N}$ . Noting that  $A_{n_k}, v_{n_k}, x_{n_k}$  still belongs to a compact metric space, we can find a further subsequence  $n_{k_j}$  such that  $A_{n_{k_j}}, v_{n_{k_j}}, x_{n_{k_j}}$  is convergent. But for this subsequence our previous result, (148), yields that:

$$\mathbb{P}[H_{n_{k_j}}(A_{n_{k_j}}, v_{n_{k_j}}, x_{n_{k_j}}) \leq z_{1-\alpha}] \rightarrow (1 - \alpha), \quad (152)$$

which yields a contradiction. Thus, for any  $(A_n, v_n, x_n)$  satisfying  $x_n \in \mathcal{X}$ ,  $\underline{v} < \|v_n\| \leq \bar{v}$  and  $\sigma(A_n, v_n, x_n, \Sigma) \geq \underline{\sigma}$  for all  $n \in \mathbb{N}$ :

$$\lim_{n \rightarrow \infty} \mathbb{P}[H_n(A_n, v_n, x_n) \leq z_{1-\alpha}] = 1 - \alpha \quad (153)$$

**Lemma 6.5.** There exists a measurable sequence  $\check{x}$ , such that  $(\hat{A}, v_n, \check{x})$  is an optimal triplet w.p.a.1 and  $\rho_\infty((\hat{A}, \hat{v}, \hat{x}), (\hat{A}, \hat{v}_n, \check{x})) = o_p(1)$ .

*Proof.* For  $A \in 2^{[q]} \setminus \{\emptyset\}$  and  $\varepsilon \geq 0$ , define

$$\mathbb{X}(A, \varepsilon) \equiv \{x \in \mathbb{R}^d : p'x = B(\theta_0), Mx \geq c, \|M_A x - c_A\| \leq \varepsilon\},$$

From the assumption of Theorem (2.3), that  $\mathcal{A}(\theta_0) \subseteq \text{Int}(\mathcal{X})$ , it follows that  $\mathbb{X}(A, \varepsilon) \subseteq \mathcal{A}(\theta_0)$  if  $A \in \mathbb{A}$ . Further, define

$$\bar{\mathbb{X}}(A, \varepsilon) \equiv \{x \in \mathbb{R}^d : p'x = B(\theta_0), Mx \geq c, M_A x - c_A \leq \varepsilon \iota_{|A|}\}.$$

Observe that for any  $A \in \mathbb{A}$  and  $\varepsilon \geq 0$ ,  $\mathbb{X}(A, \varepsilon), \overline{\mathbb{X}}(A, \varepsilon)$  are nonempty, and

$$\mathbb{X}(A, \varepsilon) \subseteq \overline{\mathbb{X}}(A, \varepsilon),$$

so

$$\begin{aligned} d_H(\mathbb{X}(A, \varepsilon), \mathbb{X}(A, 0)) &= \max\left\{ \sup_{x \in \mathbb{X}(A, \varepsilon)} d(x, \mathbb{X}(A, 0)), \sup_{x \in \mathbb{X}(A, 0)} d(x, \mathbb{X}(A, \varepsilon)) \right\} = \\ &\sup_{x \in \mathbb{X}(A, \varepsilon)} d(x, \mathbb{X}(A, 0)) \leq \sup_{x \in \overline{\mathbb{X}}(A, \varepsilon)} d(x, \mathbb{X}(A, 0)) = d_H(\overline{\mathbb{X}}(A, \varepsilon), \mathbb{X}(A, 0)) \leq C|A|\varepsilon, \end{aligned}$$

where the last inequality, for some  $C > 0$ , follows from Lipschitz-continuity of polytopes with respect to the RHS perturbations, see Li (1993). We conclude that

$$d_H(\mathbb{X}(A, \varepsilon), \mathbb{X}(A, 0)) \leq C|A|\varepsilon \tag{154}$$

From the proof of Theorem (2.2) it follows that the projection  $\tilde{x}$  of  $\hat{x}$  onto  $\mathcal{A}(\theta_0)$  is such that

$$\|\tilde{x} - \hat{x}\| = o_p(1)$$

By triangle and Cauchy-Schwartz inequalities,

$$\begin{aligned} \|M_{\hat{A}}\tilde{x} - c_{\hat{A}}\| &\leq \|M_{\hat{A}}\| \cdot \|\tilde{x} - \hat{x}\| + \|M_{\hat{A}}\hat{x} - c_{\hat{A}}\| = \\ &\|M_{\hat{A}}\| \cdot \|\tilde{x} - \hat{x}\| + \|M_{\hat{A}}\hat{x} - \hat{M}_{\hat{A}}^{(1)}\hat{x} + \hat{c}_{\hat{A}}^{(1)} - c_{\hat{A}}\| \leq \\ &\|M_{\hat{A}}\| \cdot \|\tilde{x} - \hat{x}\| + \|M_{\hat{A}} - \hat{M}_{\hat{A}}^{(1)}\| \cdot \|x\|_{\infty} + \|\hat{c}_{\hat{A}}^{(1)} - c_{\hat{A}}\| \leq \\ &\|M\| \cdot \|\tilde{x} - \hat{x}\| + \|M - \hat{M}^{(1)}\| \cdot \|x\|_{\infty} + \|\hat{c}^{(1)} - c\|, \end{aligned}$$

because the right-hand side vanishes in probability, it follows that for any  $\varepsilon > 0$ ,  $\tilde{x} \in \mathbb{X}(\hat{A}, \varepsilon)$  w.p.a.1. Denote the projection of  $\tilde{x}$  onto  $\mathbb{X}(\hat{A}, 0)$  by  $\check{x}$ . It is measurable by the usual arguments, and, by (154),

$$\|\tilde{x} - \check{x}\| \leq C|A| \cdot \|M_{\hat{A}}\tilde{x} - c_{\hat{A}}\| = o_p(1)$$

Finally, by triangle inequality,

$$\|\hat{x} - \check{x}\| \leq \|\hat{x} - \tilde{x}\| + \|\tilde{x} - \check{x}\| = o_p(1). \tag{155}$$

Observe that  $\mathbb{X}(A, 0)$  for  $A \in \mathbb{A}$  is the set of  $x \in \mathcal{A}(\theta_0)$  that satisfy the respective requirements of an optimal triplet jointly with  $A$ . Thus, whenever  $\hat{A} \in \mathbb{A}$ ,  $\check{x}, \hat{A}, v_n$  form an optimal triplet,

which occurs w.p.a.1. Combining (155) and (138),

$$\rho_\infty((\hat{A}, \hat{v}, \hat{x}), (\hat{A}, v_n, \check{x})) = o_p(1)$$

This concludes the proof of the Lemma. ■

We now show that we can pick  $\underline{\sigma}$  and  $\underline{v}$  such that the event

$$E_n \equiv \{\sigma(\hat{A}, \hat{v}, \hat{x}, \Sigma) < \bar{\sigma}\} \cup \{\|\hat{v}\| < \underline{v}\}$$

vanishes asymptotically.

By Lemma 6.5, continuity of  $\sigma(\cdot)$  in the first three arguments with respect to the  $\rho_\infty$  metric, and Assumption B4 combined with the fact that the set of optimal triplets with the additional requirement that  $\|v\| \leq \bar{v}$  in Assumption B3 is compact, if we consider

$$\underline{\sigma} = 0.5 \min_{A, x, v \text{--optimal triplet, } \|v\| \leq \bar{v}} \sigma(A, x, v, \Sigma),$$

then

$$\mathbb{P}[\sigma(\hat{A}, \hat{v}, \hat{x}, \Sigma) < \bar{\sigma}] \rightarrow 0.$$

For the second part of  $E_n$ , majorize

$$\|p\| = \|M' \dot{v}_n\| = \|\dot{v}_n\| \cdot \|M' \frac{\dot{v}_n}{\|\dot{v}_n\|}\| \leq \sigma_1(M) \|\dot{v}_n\|,$$

so that

$$\|\dot{v}_n\| \geq \frac{\|p\|}{\sigma_1(M)}.$$

Set  $\underline{v} \equiv 0.5 \frac{\|p\|}{\sigma_1(M)}$ . Using (138), triangle inequality and recalling that  $\|\dot{v}_n\| = \|\tilde{v}\|$ , and  $\|\hat{v}\| = \|\tilde{v}\|$  establishes

$$\mathbb{P}[\|\hat{v}\| < \underline{v}] \rightarrow 0,$$

so a union bound yields

$$\mathbb{P}[E_n] \rightarrow 0.$$

Note that

$$\mathbb{P}[H_n \leq z_{1-\alpha} | \mathcal{D}_n^{(1)}] = \mathbb{P}[H_n \leq z_{1-\alpha} | \hat{A}, \hat{v}, \hat{x}], \quad (156)$$

because the data in  $\mathcal{D}_n^{(1)}$  is independent from  $\mathcal{D}_n^{(2)}$  and all dependencies of  $H_n$  on  $\mathcal{D}_n^{(1)}$  can be described as measurable functions of  $\hat{A}, \hat{v}, \hat{x}$ .

Observe that:

$$\mathbb{1}_{E'_n} \inf_{A,v,x \in \mathcal{N}(\underline{v}, \underline{\sigma})} \mathbb{P}[H_n(A, v, x) \leq z_{1-\alpha}] \leq \mathbb{P}[H_n \leq z_{1-\alpha} | \hat{A}, \hat{v}, \hat{x}] \leq \quad (157)$$

$$\leq \sup_{A,v,x \in \mathcal{N}(\underline{v}, \underline{\sigma})} \mathbb{P}[H_n(A, v, x) \leq z_{1-\alpha}] + \mathbb{1}_{E_n} \quad (158)$$

It follows that:

$$\mathbb{P}[H_n \leq z_{1-\alpha} | \hat{A}, \hat{v}, \hat{x}] = 1 - \alpha + o_p(1) \quad (159)$$

Therefore:

$$\mathbb{P}[H_n \leq z_{1-\alpha} | \hat{A}, \hat{v}, \hat{x}] = 1 - \alpha + o_p(1) \quad (160)$$

By Portmanteau and because probability is bounded, we can integrate (160) over  $\mathcal{D}_n^{(1)}$ ,

$$\mathbb{P}[H_n \leq z_{1-\alpha}] = 1 - \alpha + o(1). \quad (161)$$

Finally, define

$$G_n \equiv \frac{\sqrt{n_2}}{\hat{\sigma}_n(\hat{A}, \hat{v}, \hat{x})} (\check{v} - \tilde{v}_n)' (c_{\hat{A}} - M_{\hat{A}} \hat{x}).$$

Using  $\|c_{\hat{A}} - M_{\hat{A}} \hat{x}\| = O_p(\frac{1}{\sqrt{n}})$  (see Proof of Lemma (6.5)), (138) and CMT, one concludes that  $G_n = o_p(1)$ . Applying Lemma 6.4 yields

$$\frac{\sqrt{n_2}}{\hat{\sigma}_n(\hat{A}, \hat{v}, \hat{x})} \left( \check{v}'(\hat{c}_{\hat{A}}^{(2)} - \hat{M}_{\hat{A}}^{(2)} \hat{x}) + p' \hat{x} - B(\theta_0) \right) = H_n - G_n$$

Finally, because  $G_n = o_p(1)$ , we have, for any  $\varepsilon > 0$ ,

$$o(1) + \mathbb{P}[H_n \leq z_{1-\alpha} - \varepsilon] \leq \mathbb{P}[H_n - G_n \leq z_{1-\alpha}] \leq \mathbb{P}[H_n \leq z_{1-\alpha} + \varepsilon] + o(1). \quad (162)$$

Letting  $\alpha^+(\varepsilon) \equiv 1 - \Phi(z_{1-\alpha} - \varepsilon)$  and  $\alpha^-(\varepsilon) \equiv 1 - \Phi(z_{1-\alpha} + \varepsilon)$ , applying (161), one obtains:

$$o(1) + 1 - \alpha^+(\varepsilon) \leq \mathbb{P}[H_n - G_n \leq z_{1-\alpha}] \leq o(1) + 1 - \alpha^-(\varepsilon) \quad (163)$$

Taking  $\varepsilon \rightarrow 0$  and using continuity of the normal's cdf, we obtain:

$$\mathbb{P}[H_n - G_n \leq z_{1-\alpha}] = 1 - \alpha + o(1) \quad (164)$$



To extend the proof to other consistent  $\hat{\sigma}_n$ , refer to CMT. This concludes the proof of the Theorem. ■

## 6.7. Asymptotic variance

**Lemma 6.6.** At fixed  $A, x, v$ ,

$$\sigma^2(A, x, v, \Sigma) = J_1 \Sigma J_1' - 2J_2 (I_d \otimes C_M \Sigma J_1') x + J_2 (x x' \otimes C_M \Sigma C_M') J_2',$$

where

$$J_1 \equiv \check{v}' C(\hat{A}) C_c, \quad J_2 \equiv \check{v}' C(\hat{A}) (\text{vec}(I_d)' \otimes I_q).$$

*Proof.*

$$\text{Var} \left( \check{v}' C(\hat{A}) \left( C_c Z - \text{vec}_{q \times d}^{-1}(C_M Z) \hat{x} \right) \right) = \quad (165)$$

$$= \text{Var} \left( \check{v}' C(\hat{A}) C_c Z \right) - 2 \text{Cov} \left( \check{v}' C(\hat{A}) C_c Z, \check{v}' C(\hat{A}) \text{vec}_{q \times d}^{-1}(C_M Z) \hat{x} \right) + \quad (166)$$

$$+ \text{Var} \left( \check{v}' C(\hat{A}) \text{vec}_{q \times d}^{-1}(C_M Z) \hat{x} \right) \quad (167)$$

where  $Z \sim \mathcal{N}(0, \Sigma)$  has the asymptotic distribution of  $Z_n^{(2)}$ . The first term rewrites as:

$$\text{Var} \left( \check{v}' C(\hat{A}) C_c Z \right) = J_1 \Sigma J_1' \quad (168)$$

To deal with the last term, rewrite:

$$\text{Var} \left( \check{v}' C(\hat{A}) \text{vec}_{q \times d}^{-1}(C_M Z) \hat{x} \right) = J_2 \text{Var} \left( (I_d \otimes C_M Z) \hat{x} \right) J_2' \quad (169)$$

Direct computation yields:

$$(I_d \otimes C_M Z) \hat{x} = \begin{pmatrix} C_M Z \hat{x}_1 \\ C_M Z \hat{x}_2 \\ \dots \\ C_M Z \hat{x}_d \end{pmatrix} \quad (170)$$

So:

$$\text{Var} \left( (I_d \otimes C_M Z) \hat{x} \right) = \hat{x} \hat{x}' \otimes C_M \Sigma C_M' \quad (171)$$

Consider:

$$\text{Cov}(\check{v}'C(\hat{A})C_cZ, \check{v}'C(\hat{A})\text{vec}_{q \times d}^{-1}(C_MZ)\hat{x}) = \mathbb{E}[J_1ZJ_2(I_d \otimes C_MZ)\hat{x}] = \quad (172)$$

$$= J_2\mathbb{E}[(I_d \otimes C_MZZ'J_1')]\hat{x} = J_2(I_d \otimes C_M\Sigma J_1')\hat{x} \quad (173)$$

Combining everything, we get:

$$\sigma(\hat{A}, \hat{x}, \hat{v}, \Sigma) = J_1\Sigma J_1' - 2J_2(I_d \otimes C_M\Sigma J_1')\hat{x} + J_2(\hat{x}\hat{x}' \otimes C_M\Sigma C_M')J_2' \quad (174)$$

We thus have, for fixed  $\hat{A}, \hat{v}, \hat{x}$  with  $\hat{v} \neq 0$ . ■

### 6.8. Proof of Lemma 2.4

*Proof.* Let  $\delta > 0$  be a jump at  $\mathbb{P}_0$ . Construct a sequence  $\{\mathbb{P}_n\} \subset \mathcal{P}$  such that for some  $0 < \vartheta < 1$ :

$$\|\mathbb{P}_0 - \mathbb{P}_n\|_{TV} < \vartheta n^{-1} \quad (175)$$

While  $\|V(\mathbb{P}_0) - V(\mathbb{P}_n)\| > \delta$ . Recall that:

$$\|\mathbb{P}_0^n - \mathbb{P}_n^n\|_{TV} \leq n\|\mathbb{P}_0 - \mathbb{P}_n\|_{TV} \quad (176)$$

It follows that:

$$\|\mathbb{P}_0^n - \mathbb{P}_n^n\|_{TV} \leq \vartheta \quad (177)$$

Using the binary Le Cam's method<sup>35</sup>, one obtains  $\forall n$ :

$$\inf_{\hat{V}_n} \sup_{\mathbb{P} \in \mathcal{P}} \mathbb{E}_{\mathbb{P}}[\|V(\mathbb{P}) - \hat{V}_n(X(\mathbb{P}^n))\|] \geq \frac{\delta(1 - \vartheta)}{2} \quad (178)$$

Recalling that  $0 < \vartheta < 1$  and  $\delta$  were chosen arbitrarily and taking supremum over  $\delta$  as well as sending  $\vartheta \rightarrow 0$  yields the result. ■

### 6.9. Proof of Proposition 2.5

*Proof.* Consider the problem and its associated Lagrangean:

$$(P) : \min_x p'x \quad \text{s.t.} : Mx \geq c, \quad \mathcal{L} \equiv p'x + \lambda'(c - Mx)$$

---

<sup>35</sup>See Chapter 15 of [Wainwright \(2019\)](#)

FOCs:

$$\begin{aligned} [x] : p - M'\lambda &= 0 \\ [\lambda] : c - Mx &\leq 0 \\ [\text{CS}] : \lambda'(c - Mx) &= 0 \\ [\text{POS}] : \lambda &\geq 0 \end{aligned}$$

Because  $\mathcal{X}$  is a compact, whenever the problem has a solution, it must be that there is also a solution  $\lambda^*, x^*$  at which  $\exists J \subseteq \{1, 2, \dots, q\}$  with  $|J| = k \geq d$ :

$$M_J x^* = c_J,$$

where  $M_J \in \mathbb{R}^{k \times d}$  is a matrix of full column rank:  $\text{rk}(M_J) = d$ . Define the set of inactive constraints  $I \equiv \{1, 2, \dots, q\} \setminus J$  where:

$$M_I x^* > c_I$$

It follows that  $\lambda_J^* = 0$ . Notice that the KKT condition that:

$$p = M_J' \lambda_J$$

for some  $\lambda_J \geq 0$  means that  $p \in \text{Cone}(M_J')$ . By the conical hull version of Caratheodory's Theorem, it follows that  $\exists J^* \subseteq J$  such that  $|J^*| = r \leq d$  and  $p \in \text{Cone}(M_{J^*}')$  and, moreover, the columns of  $M_{J^*}'$  are linearly independent. If the Caratheodory number  $r$  is strictly smaller than the dimension of  $x$ , i.e.  $r < d$ , then we shall complement  $J^*$  with  $d - r$  vectors from  $M_J'$  such that we obtain  $\text{rk}(M_{J^*}') = d$ , setting the appropriate  $\lambda_i^*$  to 0. By necessity and sufficiency of KKT for LP problems, this constitutes a solution.  $\blacksquare$

## 6.10. Proof of Theorem 2.5

*Proof.* We first establish a well-known Lemma.

**Lemma 6.7.** For any  $A \in \mathbb{R}^{l \times m}$  and  $b \in \mathbb{R}^m$  the following inequality holds:

$$\|Ab\|_\infty \leq \|A\|_2 \|b\| = \sigma_1(A) \|b\|$$

*Proof.* Suppose  $a_i$ ,  $i \in [l]$  are rows of  $A$ . Then,

$$\|Ab\|_\infty = \max_i \{|(Ab)_i|\} = \max_i \{|a_i' b|\} \leq \|b\| \max_i \|a_i\| \quad (179)$$

Recall that the operator norm is transpose-invariant, and can be written as:

$$\|A\|_2 = \|A'\|_2 = \sup_{\|y\| \leq 1} \|A'y\| \geq \max_i \|A'e_i\| = \max_i \|a_i\| \quad (180)$$

Combining (179) and (180) yields the result.  $\blacksquare$

We now prove the Theorem. We write  $M(\mathbb{P}), c(\mathbb{P})$  for components of  $\theta_0(\mathbb{P})$ . Fix  $\delta > 0$ . By definition of  $\mathcal{P}^\delta$  and using Proposition 2.5, for any  $\mathbb{P} \in \mathcal{P}^\delta$  there exists  $J^* = J^*(\mathbb{P}, \delta) \subseteq [q]$  and the associated KKT vector  $\lambda^* = \lambda^*(\mathbb{P}, \delta) \in \Lambda(\theta_0(\mathbb{P}))$ , such that  $M_{J^*} = M(\mathbb{P})_{J^*(\mathbb{P}, \delta)}$  is invertible, and

$$\lambda_{J^*}^* = M_{J^*}^{-1'} p, \quad \sigma_1(M_{J^*}^{-1'}) = \sigma_d^{-1}(M_{J^*}) < \delta^{-1}.$$

Using Lemma 10, one observes that

$$\|\lambda^*\|_\infty \leq \delta^{-1} \|p\|.$$

One concludes that for any  $\mathbb{P} \in \mathcal{P}^\delta$ ,

$$E_n \equiv \{w_n > \delta^{-1} \|p\|\} \subseteq \{\tilde{B}(\theta_0(\mathbb{P}), w_n) = B(\theta_0(\mathbb{P}))\}. \quad (181)$$

In what follows, we denote  $B = B(\theta_0(\mathbb{P}))$ ,  $\tilde{B} = \tilde{B}(\theta_0(\mathbb{P}); w_m(\mathbb{P}))$  and  $\tilde{B}_m = \tilde{B}(\hat{\theta}_m(\mathbb{P}); w_m(\mathbb{P}))$ , and  $\tilde{r}_m \equiv \frac{r_m}{w_m}$ . Furthermore, let  $F_n \equiv \{\inf_{m \geq n} \mathbf{1}_{E_m} = 1\}$ . Consider

$$\begin{aligned} \mathbb{P} \left[ \sup_{m \geq n} \tilde{r}_m |\tilde{B}_m - B| > \varepsilon \right] &= \mathbb{P} \left[ \sup_{m \geq n} \tilde{r}_m \left( \mathbf{1}_{E_m} |\tilde{B}_m - \tilde{B}| + \mathbf{1}_{E'_m} |\tilde{B}_m - B| \right) > \varepsilon \right] = \quad (182) \\ &\quad \mathbb{P} \left[ \sup_{m \geq n} \tilde{r}_m |\tilde{B}_m - \tilde{B}| > \varepsilon, F_n \right] + \\ &\quad \mathbb{P} \left[ \sup_{m \geq n} \tilde{r}_m \left( \mathbf{1}_{E_m} |\tilde{B}_m - \tilde{B}| + \mathbf{1}_{E'_m} |\tilde{B}_m - B| \right) > \varepsilon, F'_n \right] \leq \\ &\quad \mathbb{P} \left[ \sup_{m \geq n} \tilde{r}_m |\tilde{B}_m - \tilde{B}| > \varepsilon \right] + \mathbb{P} [F'_n]. \end{aligned}$$

Using (182), the fact that for any sequences  $g_k, h_k$ ,  $\sup_k g_k + h_k \leq \sup_k g_k + \sup_k h_k$  and the fact that  $\sup A \leq \sup B$  whenever  $A \subseteq B$ , we get

$$\sup_{\mathbb{P} \in \mathcal{P}^\delta} \mathbb{P} \left[ \sup_{m \geq n} \tilde{r}_m |\tilde{B}_m - B| > \varepsilon \right] \leq \sup_{\mathbb{P} \in \mathcal{P}} \mathbb{P} \left[ \sup_{m \geq n} \tilde{r}_m |\tilde{B}_m - \tilde{B}| > \varepsilon \right] + 1 - \inf_{\mathbb{P} \in \mathcal{P}} \mathbb{P} [F_n] \quad (183)$$

Observe that

$$\mathbb{P}[F_n] = \mathbb{P}[\cap_{m \geq n} E_m] \geq \mathbb{P}[\inf_{m \geq n} w_m > \delta^{-1} \|p\|].$$

Using this and condition ii), taking limits in (183) yields

$$\lim_{n \rightarrow \infty} \sup_{\mathbb{P} \in \mathcal{P}^\delta} \mathbb{P} \left[ \sup_{m \geq n} \tilde{r}_m \left| \tilde{B}_m - B \right| > \varepsilon \right] \leq \lim_{n \rightarrow \infty} \sup_{\mathbb{P} \in \mathcal{P}} \mathbb{P} \left[ \sup_{m \geq n} \tilde{r}_m \left| \tilde{B}_m - \tilde{B} \right| > \varepsilon \right]. \quad (184)$$

From the bound in (47) and condition i) it follows that

$$\lim_{n \rightarrow \infty} \sup_{\mathbb{P} \in \mathcal{P}} \mathbb{P} \left[ \sup_{m \geq n} \tilde{r}_m \left| \tilde{B}_m - \tilde{B} \right| > \varepsilon \right] = 0 \quad (185)$$

Combining (184), (185) and recalling that  $\delta > 0$  was arbitrary, one can take suprema on both sides, as

$$\sup_{\delta > 0} \lim_{n \rightarrow \infty} \sup_{\mathbb{P} \in \mathcal{P}^\delta} \mathbb{P} \left[ \sup_{m \geq n} \tilde{r}_m \left| \tilde{B}_m - B \right| > \varepsilon \right] = 0$$

This concludes the proof of the Theorem. ■

### 6.11. Proof of Lemma 2.5

*Proof.* If  $x \in \Theta$ , the inequality holds trivially. Consider  $x$  such that  $d(x, \Theta) = \varepsilon > 0$ . We construct a projection of  $x$  onto the polytope. It must be a solution of the following program:

$$\min_{y \in \Theta} \frac{1}{2} (y - x)' (y - x) \quad (186)$$

Construct Lagrangean:

$$\mathcal{L} = (y - x)' (y - x) + \lambda' (c - My) \quad (187)$$

FOCs:

$$y - x - M' \lambda = 0 \quad (188)$$

$$\lambda_j (M'_j y - c_j) = 0 \quad (189)$$

$$My \geq c \quad (190)$$

This problem is convex and thus has a global minimum characterized by the KKT conditions. Let that minimum be  $y^*$ . Denote the subset of binding equalities:

$$J \equiv \{j \in \overline{1, q} \mid M'_j y^* = c_j\} \quad (191)$$

Suppose  $y^*$  belongs to at least  $k^*$ -face  $f^*$ , meaning that face  $f^*$  is given by:

$$f^* = \bigcap_{f\text{-face of } \Theta_I: y \in f} f, \quad (192)$$

with the associated set of binding equalities  $J$  such that  $|J| \geq d - k^*$  and  $\text{rk}(M_J) = d - k^*$ . By construction:

$$y - x \in \text{Cone}(M'_J), \quad (193)$$

Therefore, by Caratheodory's Conical Hull theorem, there exists a subset  $J^* \subseteq J$  such that  $|J^*| = r \leq d - k^*$  and a corresponding  $\lambda_{J^*}^* > 0$ :

$$y - x = M'_{J^*} \lambda_{J^*}^* \quad (194)$$

Forming  $\lambda^*$  as  $(\lambda^*)_{J^*} \equiv \lambda_{J^*}^*$  and setting  $\lambda_j^* = 0$  for  $j \notin J^*$ , one can observe that  $y^*, \lambda^*$  solve the above of KKT conditions. Moreover, if  $r < d - k^*$ , we can complement  $J^*$  with  $d - k^* - r$  linearly independent constraints from  $J \setminus J^*$  to obtain  $J^{**} \supseteq J^*$ , such that:  $|J^{**}| = \text{rk}(M_{J^{**}}) = d - k^*$ . Finally, setting  $\lambda^{**} \equiv \lambda_{J^{**}}^*$ , we get:

$$y^* - x = M'_{J^{**}} \lambda^{**} \quad (195)$$

From where it follows that:

$$\lambda^{**} = (M_{J^{**}} M'_{J^{**}})^{-1} M_{J^{**}} (y^* - x) \quad (196)$$

Recall that  $\|y^* - x\| = \varepsilon > 0$ , and note that, because  $(M_{J^{**}} M'_{J^{**}})^{-1} M_{J^{**}}$  is the left inverse of  $M'_{J^{**}}$ :

$$\|(M_{J^{**}} M'_{J^{**}})^{-1} M_{J^{**}}\| \leq \sigma_{d-k^*}^{-1}(M_{J^{**}}) \leq \kappa^{-1}(\Theta_I) \quad (197)$$

By Cauchy-Schwarz, we then obtain:

$$\|\lambda^{**}\| \leq \varepsilon \kappa^{-1}(\Theta_I) \quad (198)$$

Since  $\|\lambda^{**}\|_\infty \leq \|\lambda^{**}\|$ , it also follows that:

$$\|\lambda^{**}\|_\infty \leq \varepsilon \kappa^{-1}(\Theta_I) \quad (199)$$

Further, since  $M_{J^{**}}y^* = c_{J^{**}}$  by construction, multiplying both sides of (195) by  $M_{J^{**}}$  yields:

$$c_{J^{**}} - M_{J^{**}}x = M_{J^{**}}M'_{J^{**}}\lambda^{**} \quad (200)$$

And plugging (195) into the value function, one gets:

$$\varepsilon^2 = \lambda'_{J^{**}}M_{J^{**}}M'_{J^{**}}\lambda^{**} = \lambda'_{J^{**}}(c_{J^{**}} - M_{J^{**}}x) \quad (201)$$

Combining (201), the fact that at least one of the components of  $\lambda^{**}$  is positive from  $y^* - x \neq 0$  and (195), as well as  $0 \leq \lambda^{**}_j \leq \varepsilon \kappa^{-1}(\Theta_I)$  from the definition and the bound on  $\|\lambda^*\|_\infty$ , one observes that:

$$\exists j \in J^{**} : (c - M_{J^{**}}x)_j \geq \frac{\kappa(\Theta_I)\varepsilon}{d - k^*} \quad (202)$$

It then follows that:

$$l'(c - Mx)^+ \geq \frac{\kappa(\Theta_I)\varepsilon}{d - k^*} \quad (203)$$

Taking the minimum over  $k^*$  yields the claim of the proposition. ■

## 6.12. Proof of Theorem 2.6

*Proof.* Note that:

$$\begin{aligned} p'x_n^* - B(\theta_0) &\geq \min_{x \in \Theta_I^{d(x_n^*, \Theta_I)}} p'x - \min_{x \in \Theta_I} p'x = \\ &\|p\| \left( s \left( \frac{p}{\|p\|}, \Theta_I^{d(x_n^*, \Theta_I)} \right) - s \left( \frac{p}{\|p\|}, \Theta_I \right) \right) \geq \\ &- \|p\| \max_{\|y\| \leq 1} \left| s(y, \Theta_I^{d(x_n^*, \Theta_I)}) - s(y, \Theta_I) \right| = \\ &- \|p\| d_H \left( \Theta_I^{d(x_n^*, \Theta_I)}, \Theta_I \right) \geq - \frac{\|p\| d l'(c - Mx_n^*)^+}{\kappa(\Theta_I)} \end{aligned}$$

Thus:

$$\begin{aligned} O_p(1) &= \frac{\sqrt{n}}{w_n} \left( p'x_n^* - B(\theta_0) + w_n t'(\hat{c}_n - \hat{M}_n x_n^*) \right) = \\ &= \frac{\sqrt{n}}{w_n} \left( p'x_n^* - B(\theta_0) + w_n t'(c - Mx_n^*) \right) + O_p(1) \geq \\ &= \sqrt{n} \left( 1 - \frac{1}{w_n} \frac{\|p\|d}{\kappa(\Theta_I)} \right) t'(c - Mx_n^*)^+ + O_p(1) \end{aligned}$$

From where it follows that:

$$t'(c - Mx_n^*)^+ = O_p\left(\frac{1}{\sqrt{n}}\right)$$

Using:

$$\begin{aligned} \frac{\sqrt{n}}{w_n} \left( p'x_n^* - B(\theta_0) + w_n t'(\hat{c}_n - \hat{M}_n x_n^*) \right) &\geq \frac{\sqrt{n}}{w_n} (p'x_n^* - B(\theta_0)) \geq \\ &= \frac{-1}{w_n} \frac{\sqrt{n} \|p\| d t'(c - Mx_n^*)^+}{\kappa(\Theta_I)} \end{aligned}$$

One deduces that  $p'x_n^* - B(\theta_0)$  is  $O_p\left(\frac{w_n}{\sqrt{n}}\right)$ . All arguments above are uniform if the convergence of  $\hat{\theta}_n$  is uniform and as  $\kappa(\Theta_I) \geq \delta$  for some  $\delta > 0$ . ■

### 6.13. Penalty parameter selection

To develop an intuition for the tradeoff involved in selecting  $\delta > 0$  and therefore the  $w_n$  penalizing sequence in Theorem 6, let us return to the example in Proposition 5:

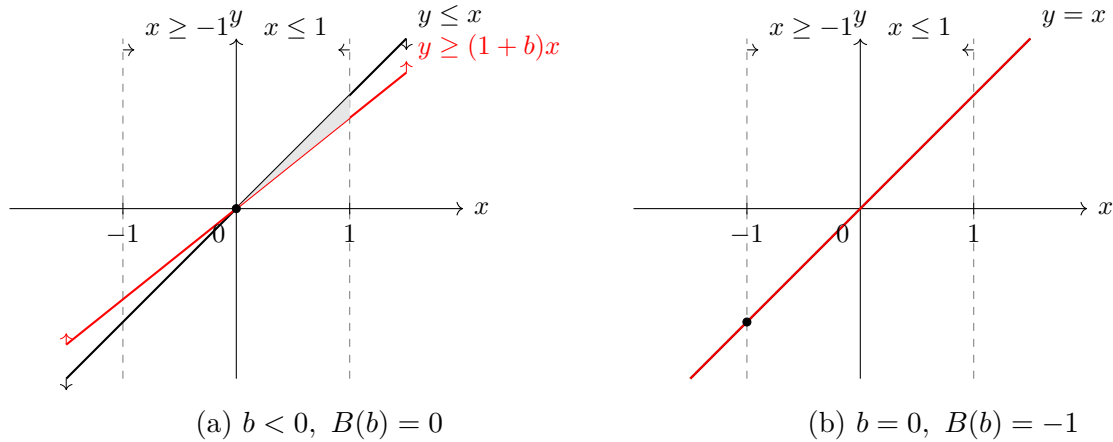


Figure 12:  $B(b) = \min_{x,y} x$  s.t. :  $y \geq (1+b)x, y \leq x, x \in [-1; 1]$

In this case, the smallest singular value at the binding constraint for  $b < 0$  is simply  $|b|$ .



Therefore, as  $b \rightarrow 0^-$ , the underlying measure belongs to a progressively smaller- $\delta$  set. For a given sample size, a higher  $w_n$  is then required to appropriately penalize the deviations, because the population Lagrange multiplier that needs to be dominated by it is equal to  $-1/b$  (see A1 and Lemma 2). On the other hand, if the true measure is the one on the right, i.e.  $b \geq 0$ , the Lagrange multiplier that needs to be dominated by  $w_n$  is fixed at  $\frac{-1+\sqrt{5}}{2}$ .

An arbitrarily large  $w_n$  will perform well in case the identified set has a 'sharp angle' ( $b \approx 0^-$ ). However, if  $b = 0$ , for example, in 50% of the cases the sample identified set will look like Figure 2.a), delivering the exploding sample Lagrange multiplier  $-\hat{b}_n^{-1}$ . If it happens to be dominated by  $w_n$  in the sense of A1, the incorrect minimum at 0, which is selected by the plug-in, is also picked by the penalty function estimator.

The aim of this section is to develop a prescription for the selection of a reasonable  $\delta$  parameter that balances finite sample performance of the estimator with sufficient robustness. As a starting point, let us note that the scale of  $\delta$  clearly depends on the scale of the singular values of  $M_{J^*}$  matrices. Any reasonable prescription for  $\delta$  parameter selection should then first normalize the constraint matrix  $M$ . More precisely, we suggest that the constraint matrix first be normalized row-wise, setting the norm of each row to 1. We further suggest rescaling it by  $s$ , where:

$$s^2 \equiv \frac{1}{Qd-1} \sum_{i,j} (\hat{M}_{ij} - \overline{\hat{M}_{..}})^2, \quad (204)$$

Once the singular values of our matrix are thus normalized,  $\delta$  may be interpreted as the degree of irregularity of the sharp identified set that one is willing to allow at the optimal solution. While uniformly and consistently estimating the sufficient  $\delta$  is infeasible (because otherwise the uniformly consistent estimator would exist), we attempt to formulate a notion of what values of  $\delta$  are *regular*. One possibility is to imagine that the population matrix of binding constraints  $M_{J^*}$ , in turn, is generated by some prior over the space of all measures. In particular, we can think of a prior such that each entry of  $M_{J^*}$  is a normalized mean zero variable, independent from other entries (but not necessarily identically distributed). In terms of the lower bound on the singular value, this prior turns out to be rather *conservative*, because it can be shown that  $\sigma_d(M_{J^*})$  goes to 0 at the rate  $\sqrt{d}$ . We therefore view this as a prudent way to characterize the irregularity of a given matrix. The random matrix theory provides the following version of the 'Central Limit Theorem' for this general prior:

**Theorem 6.1** (Tao and Vu (2010)). Let  $\Xi_n$  be a sequence of  $n \times n$  matrices with  $[\Xi_n]_{ij} \sim \xi_{ij}$ , independently across  $i, j$  where  $\xi_{ij}$  are such that  $\mathbb{E}[\xi] = 0$ ,  $Var(\xi) = 1$  and  $\mathbb{E}[|\xi|^{C_0}] < \infty$  for some sufficiently large  $C_0$ , then:

$$\sqrt{n}\sigma_n(\Xi_n) \xrightarrow{d} \Sigma \quad (205)$$

with the cdf of  $\Sigma$  given by:

$$\mathbb{P}[\Sigma \leq t] = 1 - e^{-t/2 - \sqrt{t}} \quad (206)$$

**Remark 6.3.** The distribution of mean-zero normalized  $\xi_{ij}$  in Theorem 6 is arbitrary, possibly discrete, and not necessarily identical.

This gives us the benchmark of what is 'reasonable' for a singular value of a  $d \times d$  matrix. We suggest selecting the  $0 < \alpha < 1$  quantile of this distribution, so that:

$$w_n = \|\hat{p}_n\| \delta^{-1} d_n \quad (207)$$

$$\delta = \frac{\left(\sqrt{1 - 2 \ln(1 - \alpha)} - 1\right)^2}{\sqrt{d}} \quad (208)$$

Where  $d_n \rightarrow \infty$  is some sequence that diverges slowly enough, as in Theorem 5. For example, one could set  $d_n = \ln \ln n / \ln \ln 100$  and  $\alpha = 0.15$ , seeing as the prior we selected appears rather 'conservative'. In our simulations of the example in Proposition 5, this choice of parameters delivers good uniform performance of the penalty function estimator, see Figure ??.

#### 6.14. Proof of Theorem ??

*Proof.* We first show that:

$$\Theta^* = \{\beta \in \mathbb{R} \mid \exists x \in \Theta_I : \beta = p'x + \bar{p}'\bar{x}\} \quad (209)$$

Fix  $x \in \Theta_I$ . It follows that the quantity  $m = P_m x + \bar{P}_m \bar{x}$  satisfies (19) by construction. To see that there exists at least one  $P \in \mathcal{P}$  that supports this  $m$  by generating  $m(P) = m$ , consider  $P$  under which the marginal distribution  $F_{T,Z}(\cdot)$  is as observed, and the potential outcomes have the form:

$$Y(t) = \mathbb{I}\{t \in \mathcal{O}\}(\mathbb{I}\{T \neq t\}f(t, T, Z) + \mathbb{I}\{T = t\}\eta(t)) + \mathbb{I}\{t \in \mathcal{U}\}f(t, T, Z), \quad (210)$$

where  $f : \mathcal{T}^2 \times \mathcal{Z} \rightarrow \mathbb{R}$  is a deterministic function with  $f(t, d, z)$  that maps to the component of  $x$  corresponding to the conditional moment indexed by  $t, d, z$ :  $\mathbb{E}[Y(t) \mid T = d, Z = z]$  if this

moment is counterfactual and to 0 otherwise.  $\eta(t)$  is some variable such that it aligns with  $Y(t)$  across the observed dimension:  $F_{\eta(t)|T=t,Z=z}(y) = F_{Y(t)|T=t,Z=z}(y)$ ,  $\forall y \in \mathbb{R}$  and  $\forall t \in \mathcal{O}$ ,  $\forall z \in \mathcal{Z}$ . By construction, this DGP generates  $m(P) = m$  and delivers the required identified distribution across observed dimensions,  $F_{Y|T=t,Z}(\cdot)$  for  $t \in \mathcal{O}$ . Therefore:

$$x \in \Theta_I \implies \mu^{*'}(P_m x + \bar{P}_m \bar{x}) = p'x + \bar{p}'\bar{x} \in \Theta^* \quad (211)$$

The other direction holds by construction:  $\forall \beta \in \Theta^* \exists x \in \Theta_I : p'x + \bar{p}'\bar{x} = \beta$ .

The claim of the theorem is then established by showing that the identified set is indeed an interval, a ray, or the whole line. This follows, since if  $\beta_0, \beta_1 \in \Theta^*$  with  $\beta_0 < \beta_1$ , then  $\exists x_0, x_1 \in \Theta_I$  such that  $\beta_i = p'x_i + \bar{p}'\bar{x}$  for  $i = 0, 1$ . Because  $\Theta_I$  is convex, for arbitrary  $\beta \in [\beta_0, \beta_1]$  setting  $\alpha = \frac{\beta_1 - \beta}{\beta_1 - \beta_0}$ , one obtains  $\alpha x_0 + (1 - \alpha)x_1 \in \Theta_I \implies \beta \in \Theta^*$ . ■

### 6.15. Proof of Theorem ??

**Lemma 6.8.** Fix  $K_0, \mu_v, \mu_w, K_1 \in \mathbb{R}$ :  $K_0 \leq \mu_v \leq \mu_w \leq K_1$  and  $F_w(\cdot)$  that is a valid c.d.f. with expectation  $\mu_w$ . Suppose the probability space  $(P, \Omega, \mathcal{S})$  can support a  $U[0; 1]$  random variable, and  $P[W \leq w] = F_w(w)$ . Then, there exists a random variable  $V$  s.t.  $K_0 \leq V \leq W \leq K_1$  a.s. and  $\mathbb{E}[V] = \mu_v$ .

*Proof. Proof.* Suppose  $\mu_w > K_0$  as otherwise the statement is trivial.  $W$  can be represented as:

$$W = F_w^{-1}(U) \quad (212)$$

Where  $F_w^{-1}(t) \equiv \inf\{w : F_w(w) \geq t\}$  is a generalized inverse. Consider a CDF  $G(x) \equiv \mathbb{I}\{x \geq K_0\}$  on  $[K_0; K_1]$ . Notice that by definition:

$$\int x dG(x) = K_0 \quad (213)$$

Moreover, by linearity of the Lebesgue integral  $\forall \alpha \in [0; 1]$  we have:

$$\int x d(\alpha G(x) + (1 - \alpha)F_w(x)) = \alpha K_0 + (1 - \alpha)\mu_w \quad (214)$$

Let  $F_v(x) \equiv \alpha^* G(x) + (1 - \alpha^*)F_w(x)$  where  $\alpha^* \equiv \frac{\mu_w - \mu_v}{\mu_w - K_0}$ . Then, notice that:

$$V = F_v^{-1}(U) \quad (215)$$

Yields the required random variable. ■

To prove the inverse inclusion in (??) for some  $\tilde{M}, \tilde{b}$ , note that from Theorem 1:

$$\{\beta \in \mathbb{R} | \exists P \in \mathcal{P} : \beta = \mu^{*'} m(P) \wedge b^{**} + M^{**} m(P) \geq 0\} = \{\beta \in \mathbb{R} | x : Mx \geq b : \beta = p'x + \bar{p}'\bar{x}\} \quad (216)$$

Where:

$$\bar{p} \equiv \bar{P}'_m \mu^*, \quad p \equiv P'_m \mu^* \quad (217)$$

$$M \equiv M^{**} P_m \quad b \equiv -b^{**} - M^{**} \bar{P}_m \bar{x} \quad (218)$$

Therefore proving the inclusion consists in finding such data-consistent  $\mathbb{Y}$  (or, equivalently, the measure  $P \in \mathcal{P}$ ) for any given  $x : Mx \geq b$  that it generates  $m(P) = p'x + \bar{p}'\bar{x}$  with  $M^{**} m(P) + b^{**} \geq 0$  and  $\tilde{M}\mathbb{Y} \geq \tilde{b}$   $P$  - a.s.

**1) Bounds** For any  $x : Mx \geq b$  we can once again construct the d.g.p.  $P$  from the Proof of Theorem 1:

$$Y(t) = \mathbb{I}\{t \in \mathcal{O}\}(\mathbb{I}\{T \neq t\}f(t, T, Z) + \mathbb{I}\{T = t\}\eta(t)) + \mathbb{I}\{t \in \mathcal{U}\}f(t, T, Z), \quad (219)$$

Where  $f(t, d, z), \eta(t)$  are defined as in the proof of Theorem 1 and the distribution of  $T, Z$  is as observed. Clearly,  $b^{**} + M^{**} m(P) \geq 0$  and  $P \in \mathcal{P}$  for this  $P$  holds by construction, and:  $Y(t) \in [K_0; K_1] \forall t \in \mathcal{T}$  a.s., therefore  $\tilde{M}\mathbb{Y} \geq \tilde{b}$  a.s. by construction.

**2) MTR** In this case it is clear that (219) fails, because it does not necessarily satisfy monotonicity almost surely. Consider:

$$\begin{aligned} \mathbb{Y} = & (\mathbb{I}\{t \in \mathcal{O}\}(\mathbb{I}\{T \neq t\}f(t, T, Z) + \mathbb{I}\{T = t\}\eta(t)) + \mathbb{I}\{t \in \mathcal{U}\}f(t, T, Z))_{t \in \mathcal{T}} + \\ & + \sum_{t \in \mathcal{O}} (\iota_{N_T} - e_t) \mathbb{I}\{T = t\}(\eta(t) - \mathbb{E}[Y(t)|T = t, Z]) \end{aligned} \quad (220)$$

Where  $e_t$  is the standard basis vector with 1 in the position of the potential outcome corresponding to  $t$  in  $\mathbb{Y}$ . Notice that the process in (220) has the same conditional means as the deterministic process of form (219), and therefore the corresponding  $m(P)$  satisfies  $M^{**} m(P) + b^{**} \geq 0$ . Furthermore, by construction of  $M^{**}$  it must be that  $\forall t \in \mathcal{O}$  and  $\forall d \in \mathcal{T} : d \neq t$ , we have:

$$\mathbb{E}[Y(d)|T = t, Z] = f(d, t, Z) \leq \mathbb{E}[Y(t)|T = t, Z] \text{ iff } d < t \quad (221)$$

and for  $d_0, d_1 \in \mathcal{T} \setminus \{t\} : d_0 < d_1$ :

$$\mathbb{E}[Y(d_0)|T = t, Z] = f(d_0, t, Z) \leq f(d_1, t, Z) = \mathbb{E}[Y(d_1)|T = t, Z] \quad (222)$$

Consider  $\mathbb{Y}$  constructed in (220) over some element of the partition of  $\Omega$  induced by  $T$ , where  $T = t$ .

i) If  $t \in \mathcal{U}$ , it is simply:

$$\mathbb{Y} = \begin{pmatrix} f(1, t, Z) \\ f(2, t, Z) \\ \dots \\ f(N_T, t, Z) \end{pmatrix} \quad (223)$$

Which satisfies  $\tilde{M}\mathbb{Y} + \tilde{b} \geq 0$  over this element of the partition a.s., by construction of  $f$ .

ii) If  $t \in \mathcal{O}$ :

$$\mathbb{Y} = \begin{pmatrix} f(1, t, z) + \eta(t) - \mathbb{E}[Y(t)|T = t, Z] \\ \dots \\ f(t-1, t, z) + \eta(t) - \mathbb{E}[Y(t)|T = t, Z] \\ \mathbb{E}[Y(t)|T = t, Z] + \eta(t) - \mathbb{E}[Y(t)|T = t, Z] \\ f(t+1, t, z) + \eta(t) - \mathbb{E}[Y(t)|T = t, Z] \\ \dots \\ f(N_T, t, z) + \eta(t) - \mathbb{E}[Y(t)|T = t, Z] \end{pmatrix} \quad (224)$$

Notice that by (221) and (222) the MTR is then satisfied, i.e.  $\tilde{M}\mathbb{Y} + \tilde{b} \geq 0$ .

**3) MTR + Bounds** It is clear that the process given in (220) does not necessarily satisfy boundedness. We therefore resort to a different constructive argument. Consider the element of the partition wrt to  $T$  corresponding to  $T = t$ . For  $t \in \mathcal{U}$  we can again set  $\mathbb{Y}$  as in (223). Because each  $f(d, t, Z)$  satisfies MTR and boundedness by construction, we have  $\tilde{M}\mathbb{Y} + \tilde{b} \geq 0$  over this element of the  $T$ -partition.

Suppose  $t \in \mathcal{O}$ . The solution of the linear programming results in some moments that are given by our map  $f(d, t, Z)$  that satisfies (221) and (222). Observe that constructing  $\mathbb{Y}$  over the considered element of partition consists in constructing the counterfactual  $Y(d)$  s.t.  $d \in \mathcal{T} : d \neq t$  such that:

$$\mathbb{E}[Y(d)|T = t, Z] = f(d, t, Z) \quad \forall d \in \mathcal{T} \setminus \{t\} \quad (225)$$

$$Y(1) \leq Y(2) \leq \dots \leq Y(t) \leq \dots \leq Y(N_T) \quad a.s. \quad (226)$$

Where the distribution of  $Y(t)$  over this element of the partition is identified. Repeated application of Lemma 7 yields this result. To construct the variables on the left, one starts from  $Y(t-1)$ , invokes Lemma 7 to construct it given the cdf of  $Y(t)$  (which is identified over this element of the partition), and proceeds to use the obtained cdf to construct  $Y(t-2)$ , etc.,

descending to  $Y(1)$ . For the variables 'above'  $Y(t)$ , the Lemma is simply applied with the negative sign. All of the variables can be constructed using the same  $U$  random variable in the proof of Lemma 7, which yields that there exists a probability space such that (225)-(226) hold jointly a.s. This concludes the proof of the Theorem. ■

### 6.16. Failure of the converse inclusion for almost sure inequalities

Consider a binary treatment  $T \in \{0, 1\}$  and suppose we estimate the sharp lower bound for  $\mathbb{E}[Y(1)|T = 0]$ . Suppose that conditional on  $T = 0$ ,  $Y(0)$  is 1 and  $-1$  with equal probability. Assume that there is the only conditional restriction that  $\mathbb{E}[Y(1)|T = 0] \geq 0$ . Further suppose that there is an almost sure restriction:

$$\begin{pmatrix} 1 & 1 \\ -2 & 1 \end{pmatrix} \begin{pmatrix} Y(0) \\ Y(1) \end{pmatrix} \geq \begin{pmatrix} 0 \\ 0 \end{pmatrix} \quad (227)$$

Note that this restriction defines the lower bound on  $Y(1)$  of 2 if  $Y(0) = 1$  and 1 if  $Y(0) = -1$ , and thus  $\mathbb{E}[Y(1)|T = 0] \geq 1.5$ . Taking the expectation of this system conditional on  $T = 0$ , however, yields that  $\mathbb{E}[Y(1)|T = 0] = 0$  is a solution. Therefore, although 0 is a lower bound, it is not sharp.

### 6.17. Identification under cMIV

Sharp identification results for cMIV conditions follow from Theorem 2. cMIV-w, however, allows for a more explicit characterization of the bounds, which may better illustrate the source of the identifying power of cMIV-w relative to MIV. This characterization is also useful in binary settings, when cMIV assumptions coincide. For didactic purposes, in this section we also show how to construct the restriction matrix  $M$  and vector  $b$  under cMIV-s and cMIV-p. While we focus on bounding potential outcomes or  $ATE$ s, other choices of  $\beta^*$  can be accommodated by applying Theorem 2.

In what follows,  $I_k$  stands for the identity matrix of dimension  $k$ , and  $\iota_k$  is the vector of ones of size  $k$ . These subscripts may be dropped in what follows without further notice. All vectors are column vectors, and  $\mathbb{R}^{n \times m}$  refers to the space of real-valued  $n \times m$  matrices. Notice that we can consider each  $t \in \mathcal{T}$  separately, because cMIV conditions do not impose any restrictions across potential outcomes.

**6.17.a. Recursive bounds under cMIV-w.** Construct the ordering on the support of  $Z$ :  $\mathcal{Z} = \{z_1, z_2, \dots, z_{N_Z}\}$ , s.t.  $z_i < z_j$  for  $i < j$ . Denote by  $l_i(t)$ ,  $u_i(t)$  the sharp lower and upper bounds for the conditional moment over the whole treatment support,  $\mathbb{E}[Y(t)|Z = z_i]$ . Similarly, let  $l_i^{-t}(t)$ ,  $u_i^{-t}(t)$  be the sharp upper and lower bounds for the counterfactual subset,  $\mathbb{E}[Y(t)|T \neq t, Z = z_i]$ . We shall suppress the dependence on  $t$  whenever it does not cause confusion.

The only bound of interest is the bound on unconditional expectation,  $l_i$ . However, it turns out to be instructive to also consider the bound for the counterfactual subset,  $l_i^{-t}$ .

**Proposition 6.4.** *If i) cMIV-w holds or ii) treatment is binary and cMIV-s or cMIV-p hold, the sharp bounds for  $\mathbb{E}[Y(t)|Z = z_j]$  are obtained through the following recursion for  $j \geq 2$ :*

$$l_j = l_{j-1} + \Delta_j \quad (228)$$

$$l_j^{-t} = l_{j-1}^{-t} + \Delta_j^{-t} \quad (229)$$

Where  $\Delta_j, \Delta_j^{-t} \geq 0$  are defined as follows:

$$\Delta_j \equiv \left( \underbrace{\frac{\Delta P[T \neq t|Z = z_j]}{P[T \neq t|Z = z_{j-1}]} (l_{j-1} - P[T = t|Z = z_{j-1}]\mathbb{E}[Y(t)|T = t, Z = z_{j-1}]) + \delta_j}_{\Delta P[T \neq t|Z = z_j]l_{j-1}^{-t}} \right)^+ \quad (230)$$

$$\Delta_j^{-t} \equiv \frac{1}{P[T \neq t|Z = z_j]} \left( -\Delta P[T \neq t|Z = z_j]l_{j-1}^{-t} - \delta_j \right)^+ \quad (231)$$

$$\delta_j \equiv \Delta(P[T = t|Z = z_j]\mathbb{E}[Y(t)|T = t, Z = z_j]) \quad (232)$$

Sharp upper bounds  $u_i, u_i^{-t}$  are obtained analogously. Moreover,

$$\sum_{i=1}^N P[Z = z_i]l_i(t) \leq \mathbb{E}[Y(t)] \leq \sum_{i=1}^N P[Z = z_i]u_i(t) \quad (233)$$

In the absence of additional information, these bounds are sharp.

*Proof.* Note that  $l_1^{-t} = K_0$  and  $u_N^{-t} = K_1$ . Moreover,  $l_1 = \mathbb{P}[T = t|Z = z_1]\mathbb{E}[Y(t)|T = t, Z = z_1] + \mathbb{P}[T \neq t|Z = z_1]K_0$ ,  $u_N = \mathbb{P}[T = t|Z = z_N]\mathbb{E}[Y(t)|T = t, Z = z_N] + \mathbb{P}[T \neq t|Z = z_N]K_1$ . First, we note that the equations above may be rearranged to yield:

$$l_j^{-t} = \max \left\{ \frac{1}{P[T \neq t|Z = z_j]} (l_{j-1} - \mathbb{E}[Y(t)|T = t, Z = z_j]P[T = t|Z = z_j]), l_{j-1}^{-t} \right\} \quad (234)$$

$$l_j = \mathbb{E}[Y(t)|T = t, Z = z_j]P[T = t|Z = z_j] + l_j^{-t}P[T \neq t|Z = z_j] \quad (235)$$

We consider the sharp lower bounds and proceed by induction on  $j$ . The proof for the sharp upper bounds is identical.

Consider  $j = 2$ . The only information about lower bounds provided by assumption cMIV-w

at  $j = 2$  is<sup>36</sup>:

$$\begin{cases} \mathbb{E}[Y(t)|Z = z_2] \geq \mathbb{E}[Y(t)|Z = z_1] \\ \mathbb{E}[Y(t)|T \neq t, Z = z_2] \geq \mathbb{E}[Y(t)|T \neq t, Z = z_1] \end{cases}$$

Which can be rewritten as a single condition on  $\mathbb{E}[Y(t)|T \neq t, Z = z_2]$ :

$$\begin{aligned} \mathbb{E}[Y(t)|T \neq t, Z = z_2] \geq \max \left\{ \mathbb{E}[Y(t)|T \neq t, Z = z_1], \right. \\ \left. P[T \neq t|Z = z_2]^{-1} \left( \mathbb{E}[Y(t)|Z = z_1] - P[T = t|Z = z_2] \mathbb{E}[Y(t)|T = t, Z = z_2] \right) \right\} \end{aligned}$$

Because  $l_1^{-t}$  is a sharp lower bound on  $\mathbb{E}[Y(t)|T \neq t, Z = z_1]$ , we get:

$$\begin{aligned} l_2^{-t} &= \max \left\{ l_1^{-t}, P[T \neq t|Z = z_2]^{-1} \left( l_1 - P[T = t|Z = z_2] \mathbb{E}[Y(t)|T = t, Z = z_2] \right) \right\} \\ l_2 &= P[T = t|Z = z_2] \mathbb{E}[Y(t)|T = t, Z = z_2] + P[T \neq t|Z = z_2] l_2^{-t} \end{aligned}$$

The base is thus proven. Now suppose that for some  $j \geq 2$ , and sharp lower bounds for  $i < j$  are defined. The information we have at  $j$  is:

$$\begin{cases} \mathbb{E}[Y(t)|Z = z_j] \geq \mathbb{E}[Y(t)|Z = z], \quad z < z_j \\ \mathbb{E}[Y(t)|T \neq t, Z = z_j] \geq \mathbb{E}[Y(t)|T \neq t, Z = z], \quad z < z_j \end{cases}$$

Or, equivalently,

$$\begin{aligned} \mathbb{E}[Y(t)|T \neq t, Z = z_j] \geq \max \left\{ \max_{i < j} \{ \mathbb{E}[Y(t)|T \neq t, Z = z_i] \}, \right. \\ \left. P[T \neq t|Z = z_j]^{-1} \left( \max_{i < j} \{ \mathbb{E}[Y(t)|Z = z_i] \} - P[T = t|Z = z_j] \mathbb{E}[Y(t)|T = t, Z = z_j] \right) \right\} \end{aligned}$$

Because  $l_i, l_i^{-t}$  are sharp and non-decreasing in  $i$  by inductive hypothesis, it follows that sharp lower bounds at  $j$  are given by:

$$\begin{aligned} l_j^{-t} &= \max \left\{ l_{j-1}^{-t}, P[T \neq t|Z = z_j]^{-1} \left( l_{j-1} - P[T = t|Z = z_j] \mathbb{E}[Y(t)|T = t, Z = z_j] \right) \right\} \\ l_j &= \mathbb{E}[Y(t)|T = t, Z = z_j] P[T = t|Z = z_j] + l_j^{-t} P[T \neq t|Z = z_j] \end{aligned}$$

The characterization in the proposition is obtained by rearranging these two equations.

To see that these bounds are indeed sharp, consider a process, for which  $\mathbb{E}[Y(t)|T = d, Z = z_j] = l_j^{-t}$ ,  $d \neq t, j \in \overline{1, N}$ . For such process cMIV-w will hold by construction and  $l_j$

---

<sup>36</sup>Note that we can ignore the information that  $Y(t) \geq K_0$ , as it will be implied by the bound  $l_1^{-t}$  and  $l_1$



and  $l_j^{-t}$  are both attained for all  $j$ . An example of such process is given by:

$$Y(w) = \sum_t \mathbb{I}\{t = w\} \left( \sum_j \left\{ \mathbb{I}\{Z = z_j, T = t\} \eta(t) + \sum_{d \neq t} \mathbb{I}\{Z = z_j, T = d\} l_j^{-t} \right\} \right) \quad (236)$$

Where  $\eta(t)$  is as defined in the proof of Theorem 1. ■

The intuition for Proposition 2 is that MIV bounds are obtained by 'ironing' the bounds on the population moment  $\mathbb{E}[Y(t)|Z = z]$ , which can be seen in equation (230). cMIV-w additionally 'irons' the counterfactual moments  $\mathbb{E}[Y(t)|T \neq t, Z = z]$ , as evident from (231). Figure 1 plots the derived sharp bounds as well as the benchmark MIV sharp bounds for a simulation exercise.

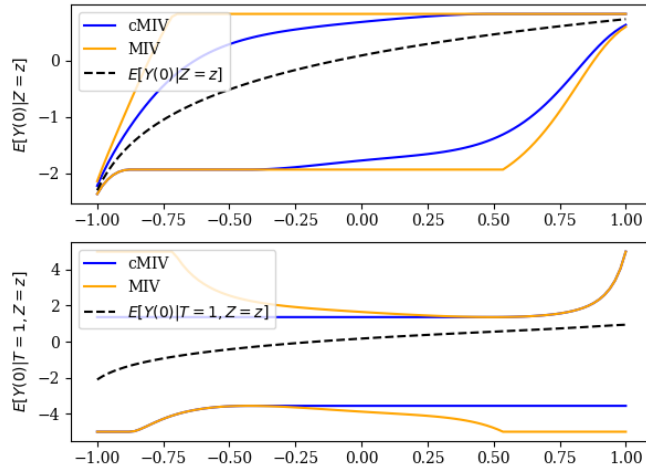


Figure 13: Bounds for the d.g.p. in Appendix 6.19.

**6.17.b. Constructing  $M$  and  $b$  for cMIV-s and cMIV-p.** Bounds given in Proposition 2 are not necessarily sharp under cMIV-s. Intuitively, cMIV-s allows us to 'iron' more moments than cMIV-w. cMIV-p, however, does not imply nor is implied by cMIV-w, so the bounds under the two conditions can compare arbitrarily. To characterize the sharp bounds under cMIV-s and cMIV-p, it is useful to introduce some notation first.

Let  $\mathcal{F} \equiv 2^{\mathcal{T}} \setminus \{\{t\}, \emptyset\}$  and its cardinality,  $Q \equiv |\mathcal{F}| = 2^{N_T} - 2$ . Fix an ordering on  $\mathcal{F}$ , so that  $\mathcal{F} = \{A^1, A^2, \dots, A^Q\}$ .

Then all information under cMIV-s can be written as:

$$\mathbb{E}[Y(t)|T \in A^k, Z = z_j] \geq \mathbb{E}[Y(t)|T \in A^k, Z = z_{j-1}], \quad k = 1, \dots, Q, \quad j = 2, \dots, N_Z \quad (237)$$

$$\mathbb{E}[Y(t)|T = d, Z = z_N] \leq K_1, \quad d \in \mathcal{T} \setminus \{t\} \quad (238)$$

$$\mathbb{E}[Y(t)|T = d, Z = z_1] \geq K_0, \quad d \in \mathcal{T} \setminus \{t\} \quad (239)$$

Where notice that the LHS of (238) is the largest marginal moment due to monotonicity in  $Z$ , while the LHS of (239) is the smallest marginal moment. Therefore, once almost sure bounds for these two moments are imposed  $\forall d \in \mathcal{T} \setminus \{t\}$ , these are also implied for all other moments through equation (237) and the law of total probability.

We now rewrite the expectations in (237) in terms of pointwise conditional moments. Let the vector of unobserved treatment responses be  $x^j \equiv (\mathbb{E}[Y(t)|T = d, Z = z_j])'_{d \neq t}$  and  $p^j \equiv (P[T = d|Z = z_j])'_{d \neq t}$  be the vector of respective probabilities at  $Z = z_j$ . Denote the element of  $x^j$  corresponding to  $T = d$  as  $x_d^j = \mathbb{E}[Y(t)|T = d, Z = z_j]$ .

For  $k = 1, \dots, Q$  and  $j = 2, \dots, N_Z$ , we can rewrite inequality (237) as follows:

$$\begin{aligned} & \sum_{d \neq t} \mathbb{I}\{d \in A^k\} \frac{P[T = d|Z = z_j]}{P[T \in A^k|Z = z_j]} x_d^j + \\ & + \mathbb{I}\{t \in A^k\} \frac{P[T = t|Z = z_j]}{P[T \in A^k|Z = z_j]} \mathbb{E}[Y(t)|T = t, Z = z_j] \geq \\ & \geq \sum_{d \neq t} \mathbb{I}\{d \in A^k\} \frac{P[T = d|Z = z_{j-1}]}{P[T \in A^k|Z = z_{j-1}]} x_d^{j-1} + \\ & + \mathbb{I}\{t \in A^k\} \frac{P[T = t|Z = z_{j-1}]}{P[T \in A^k|Z = z_{j-1}]} \mathbb{E}[Y(t)|T = t, Z = z_{j-1}] \end{aligned}$$

Inequalities (238)-(239) are just  $x_d^N \leq K_1, d \neq t$  and  $x_d^1 \geq K_0, d \neq t$ . This can be written succinctly in matrix notation. Introduce the following:

$$G_j \equiv \left( \mathbb{I}\{d \in A^k\} \frac{P[T = d|Z = z_j]}{P[T \in A^k|Z = z_j]} \right)_{k \in \overline{1, Q}, d \neq t} \in \mathbb{R}^{Q \times N_T - 1} \quad (240)$$

$$c_j \equiv \left( \mathbb{I}\{t \in A^k\} \frac{P[T = t|Z = z_j]}{P[T \in A^k|Z = z_j]} \mathbb{E}[Y(t)|T = t, Z = z_j] \right)_{k \in \overline{1, Q}} \in \mathbb{R}^Q \quad (241)$$

The whole set of information given by cMIV-s can be represented as follows:

$$G_j x^j - G_{j-1} x^{j-1} \geq -\Delta c_j, j = 2, \dots, N_Z \quad (242)$$

$$x^N \leq K_1 t \quad (243)$$

$$x^1 \geq K_0 t \quad (244)$$

The procedure for cMIV-p is similar. First, we note that all the information under it is given

by:

$$\mathbb{E}[Y(t)|Z = z_j] \geq \mathbb{E}[Y(t)|Z = z_{j-1}], \quad j = 2, \dots, N_Z \quad (245)$$

$$\mathbb{E}[Y(t)|T = d, Z = z_j] \geq \mathbb{E}[Y(t)|T = d, Z = z_{j-1}], \quad d \in \mathcal{T} \setminus \{t\}, \quad j = 2, \dots, N_Z \quad (246)$$

$$\mathbb{E}[Y(t)|T = d, Z = z_N] \leq K_1, \quad d \in \mathcal{T} \setminus \{t\} \quad (247)$$

$$\mathbb{E}[Y(t)|T = d, Z = z_1] \geq K_0, \quad d \in \mathcal{T} \setminus \{t\} \quad (248)$$

Where (245) is just MIV and (246) is the monotonicity of the pointwise conditional moments. In this case, we can once again represent all information in the matrix form (242)-(244) with the following matrices:

$$G_j \equiv \begin{pmatrix} p^{j'} \\ I_{N_T-1} \end{pmatrix} \in \mathbb{R}^{N_T \times N_T-1} \quad (249)$$

$$c_j \equiv \begin{pmatrix} P[T = t|Z = z_j] \mathbb{E}[Y(t)|T = t, Z = z_j] \\ \mathbf{0}_{N_T-1} \end{pmatrix} \in \mathbb{R}^{N_T-1} \quad (250)$$

**Corollary 2.** Under cMIV-s and cMIV-p, sharp bounds on  $\mathbb{E}[Y(t)]$  take the form:

$$\begin{aligned} & \min_{Mx \geq c} \left\{ \sum_{j=1}^N P[Z = z_j] \cdot p^{j'} x^j \right\} + \sum_{j=1}^N P[T = t, Z = z_j] \mathbb{E}[Y(t)|T = t, Z = z_j] \leq \mathbb{E}[Y(t)] \leq \\ & \leq \max_{Mx \geq c} \left\{ \sum_{j=1}^N P[Z = z_j] \cdot p^{j'} x^j \right\} + \sum_{j=1}^N P[T = t, Z = z_j] \mathbb{E}[Y(t)|T = t, Z = z_j], \end{aligned}$$

where:

$$M \equiv \begin{bmatrix} -I_{N_T-1} & \dots & 0 & 0 \\ G_N & -G_{N-1} & \dots & 0 \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \dots & G_2 & -G_1 \\ 0 & \dots & 0 & I_{N_T-1} \end{bmatrix}, \quad c \equiv \begin{pmatrix} -K_1 \cdot \iota_{N_T-1} \\ -\Delta c_N \\ \vdots \\ -\Delta c_2 \\ K_0 \cdot \iota_{N_T-1} \end{pmatrix}, \quad x = \begin{pmatrix} x^N \\ \vdots \\ x^1 \end{pmatrix}, \quad (251)$$

and  $G_j$  and  $c_j$  are given by (240) and (241) for cMIV-s and by (249) and (250) for cMIV-p respectively.

## 6.18. Proof of Proposition 4.2

Let  $\Gamma(z) \equiv \sum_{d \in \mathcal{T}} P[T = d|Z = z] \mathbb{E}[\psi(z, \eta)|Z = z]$ .

a) Let  $\tilde{g}(t) \equiv \mathbb{E}[g(t, \xi)|T = d, Z = z] = \mathbb{E}[g(t, \xi)|Z = z]$ , where we use independence of  $\xi$

and  $T, Z$ .

MIV implies:

$$\mathbb{E}[Y(t)|Z = z] = \tilde{g}(t) + h(t)\Gamma(z) \quad - \text{ increasing} \quad (252)$$

Since inequality is strict for some  $z, z'$ , it follows that  $h(t) \neq 0$  and  $h(t)/h(d) > 0$ . Note that:

$$\mathbb{E}[Y(t)|T = d, Z = z] - \tilde{g}(t) = \frac{h(t)}{h(d)} (\mathbb{E}[Y(d)|T = d, Z = z] - \tilde{g}(d)) \quad (253)$$

Therefore, cMIV-p holds iff all observed moments are monotone.

b) Let  $\tilde{g}(t, d) \equiv \mathbb{E}[g(t, \xi)|T = d, Z = z]$ , where we use independence of  $\xi$  and  $T, Z$ . We can write:

$$\mathbb{E}[Y(t)|T = d, Z = z] - \tilde{g}(t, d) = \frac{h(t)}{h(d)} (\mathbb{E}[Y(d)|T = d, Z = z] - \tilde{g}(d, d)) \quad (254)$$

Using b): ii) yields the result.

### 6.19. Simulation exercise

We now consider the following parametric example:

$$Y(t) = c + \alpha t + \beta \eta + Z \quad (255)$$

$$T = \mathbb{I}\{\varepsilon + f(Z) \geq 0\} \quad (256)$$

$$\eta = \min\{u, \max\{\varepsilon, l\}\} \quad (257)$$

$$\varepsilon \sim \mathcal{N}(0, 1) \quad (258)$$

Where  $u, l, \alpha, \beta, c \in \mathbb{R}$  and  $u > l$ . Moreover,  $\varepsilon$  is independent of all other variables. Also suppose for simplicity that  $Z \in [l; u]$  a.s. Consider:

$$\mathbb{E}[Y(t)|T = 1, Z = z] = c + \alpha t + z + \beta \mathbb{E}[\min\{u, \varepsilon\} | \varepsilon > -f(z)] = \quad (259)$$

$$= c + \alpha t + z + \beta \left( \frac{1 - \Phi(u)}{\Phi(f(z))} u + \frac{\phi(f(z)) - \phi(u)}{\Phi(f(z))} \right) \quad (260)$$

$$\mathbb{E}[Y(t)|T = 0, Z = z] = c + \alpha t + z + \beta \mathbb{E}[\max\{u, \varepsilon\} | \varepsilon \leq -f(z)] = \quad (261)$$

$$= c + \alpha t + z + \beta \left( \frac{\Phi(l)}{\Phi(-f(z))} l + \frac{\phi(l) - \phi(f(z))}{\Phi(-f(z))} \right) \quad (262)$$

For the Figure, suppose:

$$\begin{aligned}
 t &= 0 \\
 [l, u] &= [-4, 2] \\
 Z &\sim U[-1, 1] \\
 f(z) &= -2z^4 \\
 g(z) &= \ln(z + 1.1) \\
 \beta &= 0.1
 \end{aligned}$$

## 6.20. Empirical analysis

	$ATE(3, 2)$	$ATE(2, 1)$	$ATE(1, 0)$
cMIV-s	(0.059, 3.768) {0.053, 3.801}	(0.09, 3.761) {0.082, 3.81}	(0.103, 3.742) {0.094, 3.791}
cMIV-p	(0.036, 4.163) {0.033, 4.176}	(0.042, 4.185) {0.039, 4.225}	(0.053, 4.058) {0.049, 4.099}
cMIV-w	(0, 4.162) {0, 4.176}	(0, 4.072) {0, 4.102}	(0, 4.087) {0, 4.118}
MIV	(0, 4.163) {0, 4.175}	(0, 4.227) {0, 4.25}	(0, 4.108) {0, 4.134}
ETS	0.092	0.012	0.017

Table 2: Estimation results under various assumptions. CI in curly brackets are two-sided 95%, see Proposition 11.

## 6.21. Uniform rate of the debiased penalty function estimator

Our theoretical results show that under a polytope  $\delta$ -condition the debiased penalty function estimator is at least  $\sqrt{n}/w_n$  uniformly consistent. We now attempt to see if that rate is sharp uniformly, or whether the pointwise rate of  $\sqrt{n}$  is achievable. This subsection describes the design of simulations that allow us to study the uniform rate of convergence of the debiased penalty function estimator.

The proof of pointwise  $\sqrt{n}$ -consistency of the debiased penalty function estimator relies on the fact that the value  $L(x; \theta, w)$  at  $x$  outside the argmin set  $\tilde{\mathcal{A}}(\theta; w)$  is sufficiently well-separated from the optimal value  $B(\theta)$ . While at any fixed measure, including those that result in ‘flat faces’, there exists some ‘separation constant’ for a given distance from the argmin, this statement becomes problematic uniformly. In particular, around some  $\bar{\theta}$  at which there occurs a flat face, there exist sequences  $\theta_n$ , along which for any given distance of  $x$  from the argmin the difference between objective functions grows arbitrarily small.

It is worth emphasizing that the situation of an exact flat face is not problematic by itself, which is easy to see by drawing the picture of the example below at  $a = 0$ . Instead, the issue seems to occur when the measure grows arbitrarily close to a flat face. However, it seems that this is also not enough to undermine uniform  $\sqrt{n}$ -consistency: Slater's condition must also fail. Intuitively, if Slater's condition holds in the vicinity of  $\bar{\theta}$ , the estimator eventually becomes insensitive to  $w_n$  and delivers  $\sqrt{n}$ -consistency.

We consider the following linear program:

$$B(a, b, c, d) \equiv \min_{x,y} y - (1+a)x, \quad \text{s.t.:} \begin{cases} y \leq (1+b)x + d \\ y \geq (1+c)x \\ x \in [-1; 1] \end{cases}, \quad (263)$$

Where we take  $a$  to be fixed and indexing a probability measure.  $b = 0, c = 0, d = 0$  are estimated via  $b_n, c_n, d_n$  as sample averages of independent  $U[-0.5, 0.5]$  random variables. We now describe the design of our simulations:

1. We set  $w_n = \frac{\ln n}{\ln 100}(\delta/1.5)^{-1}$ , where  $\delta$  is the biggest value for which the delta condition is satisfied over  $a \in [-0.1, 0.1]$ .
2. For any fixed  $n$ , we take the grid of 9 points:

$$\mathcal{G}_n \equiv \{-0.1, 0, 0.1\} \cup \{-0.1C_1n^{-1/2}, 0.1C_1n^{-1/2}\} \cup \{-0.1C_2w_nn^{-1/2}, 0.1C_2w_nn^{-1/2}\} \cup \{-0.1C_3w_n^{-1}, 0.1C_3w_n^{-1}\},$$

where  $C_i$  are chosen so that each point is equal to  $-0.1$  at  $n = 100$ .

3. At each  $n$ , we run  $N_{sim} = 10000$  simulations, each time computing  $b_n, c_n, d_n$  and plugging in to obtain:

$$\sup_{a \in \mathcal{G}_n} |\tilde{B}(a, b_n, c_n, d_n; w_n) - B(a, 0, 0, 0)| \quad (264)$$

4. We then compute the standard deviation of (264) across simulations at each  $n$
5. We consider multiplying the resulting standard deviations by two rates:  $\sqrt{n}$  and  $\sqrt{n}/w_n$ .

In all figures below the level of the red curve is equated to the level of the blue one at the smallest  $n$  to illustrate the growth rate.

From Figure 14, it appears that standard deviations multiplied by  $\sqrt{n}$  are indeed exploding, although very slowly, while those multiplied by  $\sqrt{n}/w_n$  are stable. It may be the case that the rate of  $\sqrt{n}/w_n$  is sharp uniformly.

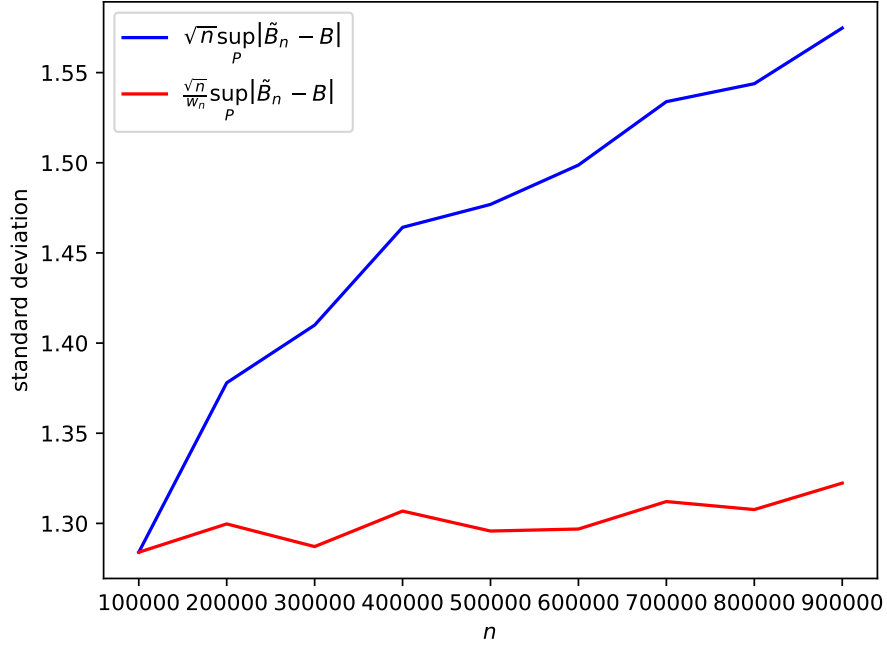


Figure 14: Uniformity of the penalized estimator: continuous vicinity of a flat face

We next consider the grid that includes the flat face itself, but restricts the measures from approaching it from the left and right. In other words, we conduct the same simulation exercise with:

$$\mathcal{G}_n \equiv \{-0.1, 0, 0.1\} \cup \{-0.05(1 + C_1 n^{-1/2}), 0.05(1 + C_1 n^{-1/2})\} \cup \{-0.05(1 + C_2 w_n n^{-1/2}), 0.05(1 + C_2 w_n n^{-1/2})\} \cup \{-0.05(1 + C_3 w_n^{-1}), 0.05(1 + C_3 w_n^{-1})\}$$

In this case, Figure 15 suggests that uniform  $\sqrt{n}$ -consistency is achieved.

Finally, we return to the original grid  $\mathcal{G}_n$ , but consider the case in which Slater's condition holds. For that reason, we take the true value of  $d = 0.5$  by sampling  $d_n$  from  $U[0, 1]$  instead.

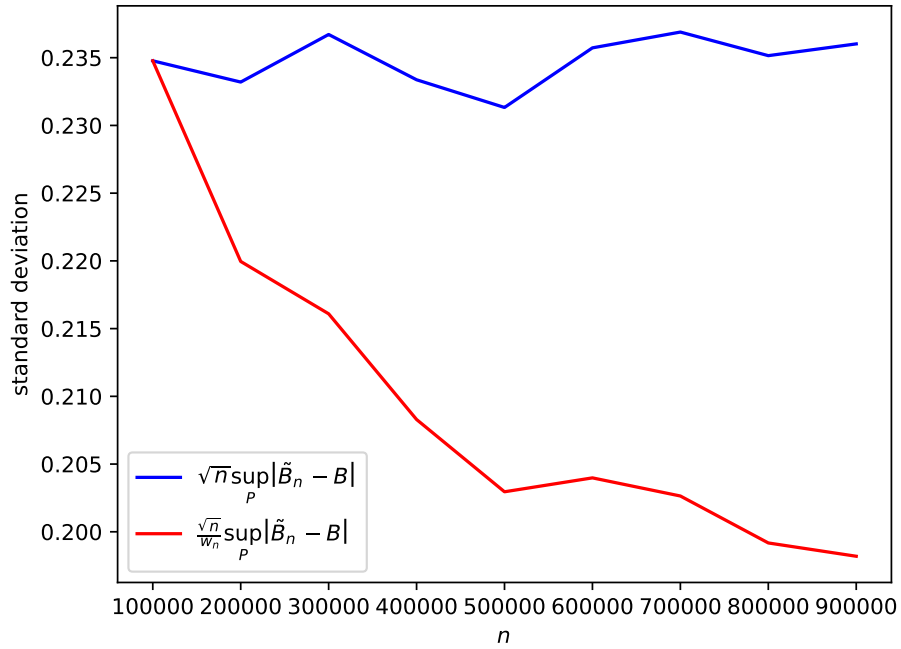


Figure 15: Uniformity of the penalized estimator: restricted vicinity of a flat face, flat face included.

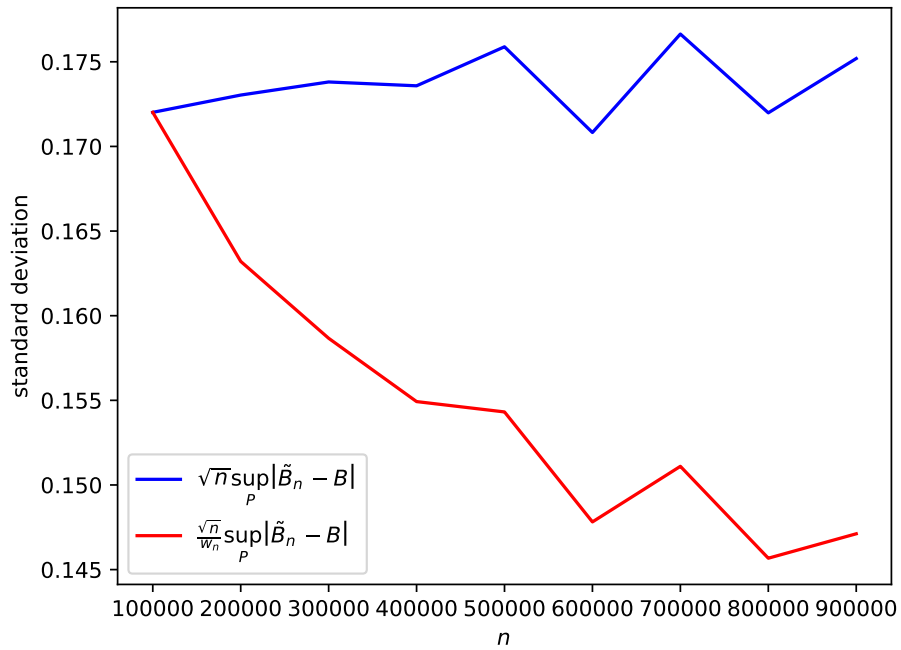


Figure 16: Uniformity of the penalized estimator: continuous vicinity of a flat face, Slater's condition holds.



Once again, it appears that we obtain uniform  $\sqrt{n}$ -consistency.

Our simulation evidence thus suggests that while our estimator is only  $\sqrt{n}/w_n$  uniformly consistent in general, it is  $\sqrt{n}$ -consistent uniformly apart from the sequences of probability measures, along which both Slater's condition fails and where a flat-face is 'approached' monotonically. It appears possible to rule out the latter scenario by considering a uniform condition similar to the  $\delta$ -condition we imposed before. This condition would restrict the set of measures under consideration to those at which the 'distance' from a flat face is either 0 or bounded away from 0 in some metric. Accordingly, it would likely cover the unrestricted set of measures in the limit. These considerations, however, are the topic of a separate exploration, and space does not permit us to include them in this paper.