A More Credible Approach to Parallel Trends^{*}

Ashesh Rambachan[†]

Jonathan Roth[‡]

April 1, 2022

Abstract

This paper proposes tools for robust inference in difference-in-differences and eventstudy designs where the parallel trends assumption may be violated. Instead of requiring that parallel trends holds exactly, we impose restrictions on how different the post-treatment violations of parallel trends can be from the pre-treatment differences in trends ("pre-trends"). The causal parameter of interest is partially identified under these restrictions. We introduce two approaches that guarantee uniformly valid inference under the imposed restrictions, and we derive novel results showing that they have desirable power properties in our context. We illustrate how economic knowledge can inform the restrictions on the possible violations of parallel trends in two economic applications. We also highlight how our approach can be used to conduct sensitivity analyses showing what causal conclusions can be drawn under various restrictions on the possible violations of the parallel trends assumption.

Keywords: Difference-in-differences, event-study, parallel trends, robust inference, sensitivity analysis, partial identification.

[†]Harvard University. Email: <u>asheshr@g.harvard.edu</u> [‡]Brown University. Email: <u>jonathanroth@brown.edu</u>

^{*}This paper was previously titled, "An Honest Approach to Parallel Trends." We are grateful to Isaiah Andrews, Elie Tamer, and Larry Katz for their invaluable advice and encouragement. We also thank Clément de Chaisemartin, Gary Chamberlain, Raj Chetty, Peter Ganong, Ed Glaeser, Nathan Hendren, Ryan Hill, Ariella Kahn-Lang, Jens Ludwig, Sendhil Mullainathan, Claudia Noack, Frank Pinter, Adrienne Sabety, Pedro Sant'Anna, Jesse Shapiro, Neil Shephard, Jann Spiess, Jim Stock, and seminar participants at Brown, Chicago Booth, Dartmouth, Harvard, Michigan, Microsoft, Princeton, Rochester, UCL, Yale, SEA 2021, and ASSA 2022 for helpful comments, and Dorian Carloni for kindly sharing data. We gratefully acknowledge financial support from the NSF Graduate Research Fellowship under Grant DGE1745303 (Rambachan) and Grant DGE1144152 (Roth).

1 Introduction

Researchers using difference-in-differences (DiD) and related methods are often unsure about the validity of the parallel trends assumption needed for point identification of the causal parameter of interest. It has therefore become common practice to assess the plausibility of the parallel trends assumption by testing for pre-treatment differences in trends ("pretrends"). Although pre-trends tests are intuitive, recent research has shown that they may suffer from low power (Freyaldenhoven, Hansen and Shapiro, 2019; Kahn-Lang and Lang, 2020; Bilinski and Hatfield, 2020; Roth, Forthcoming), and that conditioning the analysis on passing pre-trends tests introduces statistical issues related to pre-testing (Roth, Forthcoming). How then should researchers proceed when they are unsure about the validity of the parallel trends assumption?

This paper proposes methods for robust inference and sensitivity analysis in empirical settings where the parallel trends assumption may not hold. Building on work by Manski and Pepper (2018), we show that the causal parameter of interest can be (partially) identified under a large class of restrictions that impose that the post-treatment violations of parallel trends cannot be "too different" from the pre-trends. We then introduce methods that yield uniformly valid inference for the treatment effect under the imposed restrictions. Intuitively, our inference methods account for both statistical uncertainty (we can only noisily estimate the true pre-trend) as well as "identification uncertainty" (even if the true pre-trend were known, we may not know exactly how to extrapolate it). Our approach thus formalizes the intuition motivating pre-trends testing while avoiding the statistical issues described above.

More concretely, we consider a setting in which the researcher estimates a vector of "eventstudy" coefficients $\hat{\beta} = (\hat{\beta}'_{pre}, \hat{\beta}'_{post})' \in \mathbb{R}^{T+\bar{T}}$, where $\hat{\beta}_{pre}$ and $\hat{\beta}_{post}$ respectively correspond with estimates for \bar{T} pre-treatment periods and \bar{T} post-treatment periods. We assume that $\hat{\beta}$ is consistent for the reduced-form parameter β , which can be decomposed as

$$\beta = \underbrace{\begin{pmatrix} 0 \\ \tau_{post} \end{pmatrix}}_{=:\tau} + \underbrace{\begin{pmatrix} \delta_{pre} \\ \delta_{post} \end{pmatrix}}_{=:\delta}, \tag{1}$$

where τ is a causal parameter of interest that is assumed to be 0 in the pre-treatment period and δ is a bias from a difference in trends. For instance, in the canonical (nonstaggered) difference-in-differences framework, $\hat{\beta}$ may be the coefficients from an "eventstudy regression" specification, τ the vector of period-specific average treatment effects on the treated (ATT) for some policy of interest, and δ the difference in trends of untreated potential outcomes between the treated and comparison groups. As we discuss in Section 2, this framework also applies to more complicated empirical settings, such as those with staggered treatment timing (e.g. Callaway and Sant'Anna, 2020; Sun and Abraham, 2020). The usual parallel trends assumption used to point identify τ_{post} is that $\delta_{post} = 0$, and researchers frequently assess the plausibility of this assumption by testing the null hypothesis $\delta_{pre} = 0$ (a "pre-trends" test).

Instead of imposing that the parallel trends assumption holds exactly, we place restrictions on the possible post-treatment difference in trends δ_{post} given the point identified pretrend δ_{pre} . Such restrictions formalize the intuition motivating pre-trends tests, namely that pre-trends are informative about counterfactual post-treatment differences in trends. Formally, we assume that $\delta \in \Delta$ for some researcher-specified set Δ , and show that the causal parameter τ_{post} is partially identified under such restrictions.

Restrictions of this form can be used to formalize a wide variety of intuitions about possible violations of the parallel trends assumption that are commonly expressed in applied work. For example, as discussed in Manski and Pepper (2018), researchers may be willing to assume that the confounding factors that create post-treatment violations of parallel trends are similar in magnitude to those in the pre-treatment period. This intuition can be formalized by specifying a Δ that bounds the maximal post-treatment violation of parallel trends by a parameter \overline{M} times the maximal pre-treatment violation of parallel trends. In other contexts, researchers are concerned about violations of parallel trends from secular trends that are assumed to evolve smoothly over time. This intuition can be formalized by bounding the extent to which the slope of the violation of parallel trends can change over time. We adopt a flexible framework that allows researchers to capture these intuitions, as well as many other restrictions that are implied by context-specific knowledge about possible confounding factors.

We then introduce methods to conduct uniformly valid inference on a scalar parameter of the form $\theta = l' \tau_{post}$ under the restriction $\delta \in \Delta$. As emphasized in the recent literature on pre-trends testing, the pre-treatment coefficients $\hat{\beta}_{pre}$ are often imprecise estimates of δ_{pre} . It is therefore important to introduce inference methods that account for the statistical uncertainty in the estimation of the event-study coefficients. We introduce two main inference approaches, with different desirable properties depending on the shape of Δ .

We first introduce a general inference approach that can accommodate a wide variety of restrictions of Δ . This approach is based on the observation that conducting inference on θ can be cast as the problem of testing a system of moment inequalities, allowing us to leverage a large econometrics literature on moment inequality testing (Canay and Shaikh (2017) provide a recent review). The moments have a potentially large number of nuisance parameters that enter linearly, and we therefore consider an implementation of this approach based on the conditional and hybrid approaches of Andrews, Roth and Pakes (2021, henceforth ARP), who considered moment inequalities with this structure. Uniform size control for these tests follows nearly immediately from results in ARP.

We then prove that the tests proposed by ARP have some desirable power properties in our context. First, the conditional and hybrid tests are consistent, in the sense that they have power approaching 1 against fixed alternatives outside of the identified set. Second, we prove that the conditional test has optimal local asymptotic power under a linear independence constraint qualification (LICQ) assumption. As described in Kaido, Molinari and Stoye (2021), LICQ and related constraint qualifications have been used widely in the partial identification literature, and are often imposed to ensure size control. By contrast, we show that ARP conditional test is asymptotically valid even when LICQ fails, but has optimal local asymptotic power envelope when LICQ is satisfied. Intuitively, this result implies that the conditional test will perform well when the binding and non-binding moments are "far apart" relative to the sampling variation in the data. We provide several intuitive examples to illustrate when this result will and will not be applicable. Our result also implies that the ARP hybrid test will have near-optimal local asymptotic power under LICQ. These power results are new, and exploit additional structure in our context not contained in ARP.

Our second approach to inference is based on fixed length confidence intervals (FLCIs) (Donoho, 1994). FLCIs have desirable finite-sample guarantees for particular Δ s of interest. In particular, results from Armstrong and Kolesár (2018, 2020b) imply that when Δ is convex and centrosymmetric, FLCIs have near-optimal expected length in the finite-sample normal model. These results are applicable for one of our leading examples, Δ^{SD} , which restricts the smoothness of the difference in trends. In Monte Carlo simulations, we find that the use of such FLCIs can lead to substantial power gains over the conditional/hybrid approaches for Δ^{SD} when the length of the identified set is short relative to the sampling variation in the data. This is intuitive since the asymptotic power guarantees for the conditional/hybrid approaches are in the asymptotic regime where sampling uncertainty is small relative to the length of the identified set, in contrast to the finite-sample guarantees for FLCIs. On the other hand, FLCIs are applicable for a much smaller range of Δ s: indeed, we show that for many other choices of Δ , they will be inconsistent in the strong sense that power against fixed points outside the identified set need not converge to one asymptotically.

Based on our theoretical results and Monte Carlo simulations, we recommend the ARP hybrid approach for general forms of Δ , but prefer the FLCI approach in special cases (such as for Δ^{SD}) where the conditions for consistency and finite-sample near-optimality are met.

We recommend that applied researchers use our methods to construct robust confidence sets under economically-motivated restrictions on how the pre-trends relate to the posttreatment violations of parallel trends. Our tools can also be used to to conduct sensitivity analyses in which the researcher reports confidence sets under varying restrictions on the possible differences in trends. For example, if the researcher suspects that the confounding factors in the post-treatment periods are similar in magnitude to those in the pre-treatment period, then it may be reasonable to impose that the post-treatment violations of parallel trends are no larger than the maximum pre-treatment violation of parallel trends. As a sensitivity analysis, the researcher might also report confidence sets that allow the post-treatment maximum violations of parallel trends to be up to \overline{M} times larger than the maximum pretreatment violation for different values of \overline{M} . Performing such sensitivity analyses makes clear what must be assumed about the possible differences in trends in order to draw specific causal conclusions. We provide an R package, HonestDiD, that implements our recommended methods.¹ We illustrate our recommended approach with applications to two recently published papers, in which we show how the choice of the restrictions Δ can be tailored to the economic context.

Related literature: The approach in this paper builds on the foundational partial identification analysis for DiD in Manski and Pepper (2018). Manski and Pepper consider identification under researcher-specified bounds on the magnitude of δ_{post} (what they call "bounded DiD variation"), and calibrate these bounds using the maximal pre-treatment violation of parallel trends in their empirical application on the effects of right-to-carry gun laws.² One of our leading classes of restrictions, Δ^{RM} , formalizes this calibration approach by bounding the magnitude of post-treatment violations of parallel trends by \overline{M} times the maximal pre-treatment violation. Our framework also allows for many other intuitive restrictions such as bounds on how far δ can deviate from linearity — and it can be applied to a variety of difference-in-differences estimators, including recent proposals for settings with staggered treatment timing. Most importantly, while Manski and Pepper (2018) provide a framework for identification, we provide inference methods to construct uniformly valid confidence sets for the treatment effect of interest. This allows applied researchers to account for statistical uncertainty in their analyses, which can be important since event-study coefficients are often imprecisely estimated in practice.

Several other recent papers consider various relaxations of the parallel trends assumption. Keele, Small, Hsu and Fogarty (2019) develop techniques for testing the sensitivity

¹The latest version of the R package can be downloaded by visiting http://github.com/ asheshrambachan/HonestDiD.

²Manski and Pepper (2018) also consider "bounded time" and "bounded state" restrictions that bound how much the mean of Y(0) can differ either across treatment groups or within-groups over time. Such restrictions could also be incorporated into our framework by augmenting the vector $\hat{\beta}$ to include group-specific sample averages.

of DiD designs to violations of the parallel trends assumption, but they do not incorporate information from the observed pre-trends in their sensitivity analysis. Empirical researchers commonly adjust for the extrapolation of a linear trend from the pre-treatment periods when there are concerns about violations of the parallel trends assumption, which is valid if the difference in trends is exactly linear (e.g., Dobkin, Finkelstein, Kluender and Notowidigdo, 2018; Goodman-Bacon, 2018, 2021; Bhuller, Havnes, Leuven and Mogstad, 2013). Our methods nest this approach as a special case, but allow for valid inference under less restrictive assumptions about the class of possible differences in trends (such as when δ is only approximately linear). Freyaldenhoven et al. (2019) propose a method that allows for violations of the parallel trends assumption but requires an additional covariate that is affected by the same confounding factors as the outcome but not by the treatment of interest. Ye, Keele, Hasegawa and Small (2020) consider partial identification of treatment effects when there exist two control groups whose outcomes have a bracketing relationship with the outcome of the treated group. Leavitt (2020) proposes an empirical Bayes approach calibrated to pre-treatment differences in trends, and Bilinski and Hatfield (2020) and Dette and Schumann (2020) propose approaches based on pre-tests for the magnitude of the pre-treatment violations of parallel trends.

Our methods address several concerns related to current empirical practice in differencein-differences and event-study designs. First, pre-trends tests may be underpowered against meaningful violations of parallel trends, potentially leading to severe undercoverage of conventional confidential intervals (Freyaldenhoven et al., 2019; Bilinski and Hatfield, 2020; Kahn-Lang and Lang, 2020; Roth, Forthcoming). Second, statistical distortions from pretrends tests may further undermine the performance of conventional inference procedures (Roth, Forthcoming). Third, parametric approaches to controlling for pre-existing trends may be sensitive to functional form assumptions (Wolfers, 2006; Lee and Solon, 2011). We address these issues by providing tools for inference that do not rely on an exact parallel trends assumption, incorporate statistical uncertainty about the estimated event-study coefficients, and make clear the mapping between the researcher's assumptions about the potential differences in trends and the strength of their causal conclusions.

Our work complements a growing literature on the causal interpretation of event-study coefficients in two-way fixed effects models in the presence of staggered treatment timing and heterogeneous treatment effects (Borusyak and Jaravel, 2016; Athey and Imbens, 2021; Goodman-Bacon, 2021; Callaway and Sant'Anna, 2020; de Chaisemartin and D'Haultfœuille, 2020; Sun and Abraham, 2020). A key finding is that regression coefficients from conventional approaches may not produce convex weighted averages of treatment effects even if parallel trends holds. Several alternative estimators have been proposed that consistently estimate

interpretable causal estimands under a suitable parallel trends assumption. Our methodology can be used in conjunction with these alternative estimators to assess their sensitivity to violations of the corresponding parallel trends assumption; see Section 2.1 for additional details.

More broadly, our work contributes to a larger econometric literature that uses partial identification to provide empirical researchers with tractable tools to conduct inference under assumptions that may be more credible in empirical practice; see, for example, Manski (2003, 2007, 2013), Tamer (2010), Ho and Rosen (2017), and Molinari (2020) for reviews.

2 Model set-up

2.1 Event-study Coefficients

We suppose that the researcher has estimated a vector of "event-study coefficients" $\hat{\beta}_n \in \mathbb{R}^{T+\bar{T}}$, which can be partitioned into vectors of coefficients corresponding with the pretreatment and post-treatment periods, $\hat{\beta}_n = (\hat{\beta}'_{n,pre}, \hat{\beta}'_{n,post})$, where $\hat{\beta}_{n,pre} \in \mathbb{R}^T$ and $\hat{\beta}_{n,post} \in \mathbb{R}^{\bar{T}}$. Event-study estimates of this form arise from non-staggered DiD as well as a variety of related estimators, as we illustrate with several examples

Example 1 (Non-staggered DiD). Consider the canonical DiD setting in which we have a balanced panel of units from period $t = -\underline{T}, ..., \overline{T}$, and units with $D_i = 1$ receive a treatment beginning in period t = 1, while units with $D_i = 0$ never receive the treatment. It is common to report difference-in-differences estimates of the form

$$\hat{\beta}_s = (\bar{Y}_{s1} - \bar{Y}_{s0}) - (\bar{Y}_{01} - \bar{Y}_{00}).$$

where \bar{Y}_{sd} is the sample mean of the outcome for units with $D_i = d$ in period t = s. Intuitively, $\hat{\beta}_s$ compares the change in the mean outcome between period 0 and period s for the treated and comparison units. In this setting, the estimates $\hat{\beta}_s$ are numerically equivalent to the OLS coefficients from the regression

$$Y_{it} = \lambda_i + \phi_t + \sum_{s \neq 0} \beta_s \times \mathbb{1}[t=s] \times D_i + \epsilon_{it}.$$
(2)

In this case, $\hat{\beta}_{post}$ collects the estimated coefficients corresponding with treated periods, $(\hat{\beta}_1, ..., \hat{\beta}_{\bar{T}})$, while $\hat{\beta}_{pre}$ collects the estimated coefficients corresponding with periods before treatment $(\hat{\beta}_{-\bar{T}}, ..., \hat{\beta}_{-1})$.

Example 2 (Staggered DiD). Event-study coefficients can also be obtained from more complicated DiD procedures. For example, in settings with staggered treatment timing, Callaway and Sant'Anna (2020) propose event-study estimates of the form

$$\hat{\beta}_r = \sum_g w_g \widehat{ATT}(g, g+r),$$

where $\widehat{ATT}(g,t)$ is a difference-in-differences estimate that compares the evolution of the outcome for units first treated at period g to units first-treated after period t between time periods g-1 and t, and the w_g are weights that sum to one (e.g. proportional to sample size). In this case, $\hat{\beta}_{post}$ collects the values of $\hat{\beta}_r$ for $r \ge 0$ (i.e. estimates where one of the groups is treated), and $\hat{\beta}_{pre}$ collects the values of $\hat{\beta}_r$ for values of r < 0. Several other related procedures have been proposed for constructing event-study coefficients in contexts with staggered treatment timing; see de Chaisemartin and D'Haultfœuille (2021) and Roth, Sant'Anna, Bilinski and Poe (2022) for reviews.

Example 3 (Other related estimators). Other examples of estimators that can be used to produce event-studies coefficients of the form considered here include the GMM procedure proposed by Freyaldenhoven et al. (2019), instrumental variables event-studies (Hudson, Hull and Liebersohn, 2017), as well estimators that flexibly control for differences in covariates between treated and comparison groups (e.g., Heckman, Ichimura, Smith and Todd, 1998; Abadie, 2005; Sant'Anna and Zhao, 2020).

2.2 Causal Decomposition

Under mild regularity conditions, all of the estimators described above will be asymptotically normally distributed, satisfying $\sqrt{n}(\hat{\beta}_n - \beta) \rightarrow \mathcal{N}(0, \Sigma^*)$ for some parameter vector β . We assume the parameter vector β satisfies the following causal decomposition.

Assumption 1. The parameter vector β can be decomposed as

$$\beta = \underbrace{\begin{pmatrix} \tau_{pre} \\ \tau_{post} \end{pmatrix}}_{=: \tau} + \underbrace{\begin{pmatrix} \delta_{pre} \\ \delta_{post} \end{pmatrix}}_{=: \delta} \quad with \ \tau_{pre} = 0.$$
(3)

The first term, τ , represents the treatment effects of interest. We assume the treatment has no causal effect prior to its implementation, so $\tau_{pre} = 0$. The second term, δ , represents the difference in trends between the treated and comparison groups that would have occurred absent treatment. The parallel trends assumption imposes that $\delta_{post} = 0$, and therefore $\beta_{post} = \tau_{post}$ under parallel trends.

Example: Non-staggered DiD (continued) Suppose the observed outcome satisfies $Y_{it} = D_i Y_{it}(1) + (1 - D_i) Y_{it}(0)$, where $Y_{it}(1)$ and $Y_{it}(0)$ are respectively the potential outcomes when unit *i* is ultimately treated / not treated. Assume further that there is no anticipation of treatment, so that $Y_{it}(1) = Y_{it}(0)$ for all t < 1. Then, for any *s*, under mild regularity conditions $\hat{\beta}_s$ will be consistent for

$$\beta_s = \tau_{ATT,s} + \underbrace{\mathbb{E}\left[Y_{is}(0) - Y_{i0}(0) \mid D_i = 1\right] - \mathbb{E}\left[Y_{is}(0) - Y_{i0}(0) \mid D_i = 0\right]}_{\text{Differential trend} =: \delta_s},$$

where $\tau_{ATT,s} = \mathbb{E}[Y_{is}(1) - Y_{is}(0)|D_i = 1]$ is the average treatment effect on the treated in period s, and δ_s is the difference in trends in potential outcomes between period 0 and period s.³ Since the no-anticipation assumption implies that $\tau_{ATT,s} = 0$ for s < 0, this yields the decomposition (3).

Example: Staggered DiD (continued) Likewise, in the staggered DiD context, define $Y_{it}(g)$ to be the potential outcome for unit i in period t if they are first treated at period g and $Y_{it}(\infty)$ to be the never-treated potential outcome. Then $\hat{\beta}_r$ will be consistent for the parameter $\beta_r = \tau_r + \delta_r$, where $\tau_r = \sum_g w_g ATT(g, g + r)$ and $ATT(g, g + r) = \mathbb{E}[Y_{i,g+r}(g) - Y_{i,g+r}(\infty)|G_i = g]$ is the ATT in period g + r for units first treated at period g. Likewise $\delta_r = \sum w_g \delta_{g,g+r}$, where $\delta_{g,g+r} = \mathbb{E}[Y_{i,g+r}(\infty) - Y_{i,g-1}(\infty)|G_i = g] - \mathbb{E}[Y_{i,g+r}(\infty) - Y_{i,g-1}(\infty)|G_i > g + r]$ is the difference in trends in never-treated potential outcomes between units first treated at period g and units first treated after period g+r. Under a no-anticipation assumption, $\tau_r = 0$ for r < 0, which again yields the decomposition (3).

Example: Other related estimators (continued) We can decompose β as in (3) for other estimators as well. For example, for event-study IVs (with non-staggered timing), τ_{post} is a vector containing the local average treatment effect for each period, and δ represents the violation of the exclusion restriction at each period. For methods that flexibly control for covariate differences between treated and comparison groups, τ_{post} is again a vector of ATTs, and δ_{post} represents a weighted average (across covariates) of the violation of the conditional parallel trends assumption.

³We focus on the ATT as the target parameter, as in most of the DiD literature. If one is interested in the population-wide average treatment effect (ATE), one could obtain bounds on the ATE under restrictions on treatment effect heterogeneity, or other assumptions that allow one to bound the treatment effects for untreated units; see Manski and Pepper (2013) for an insightful discussion.

2.3 Target parameter and identification

We suppose the target parameter is a linear combination of the post-treatment causal effects, $\theta := l' \tau_{post}$ for some known \overline{T} -vector l. For example, θ equals the t-th period causal effect τ_t when the vector l equals the t-th standard basis vector. Similarly, θ equals the average causal effect across all post-treatment periods when $l = (\frac{1}{\overline{T}}, ..., \frac{1}{\overline{T}})'$.

We obtain partial identification of θ by assuming that δ lies in a researcher-specified set of possible differences in trends $\Delta \subseteq \mathbb{R}^{T+\bar{T}}$. This nests the usual parallel trends assumption as the special case with $\Delta = \{\delta : \delta_{post} = 0\}$. Since $\delta_{pre} = \beta_{pre}$ is identified, the assumption that $\delta = (\delta'_{pre}, \delta'_{post})' \in \Delta$ restricts the possible values of δ_{post} given the identified value of the pre-treatment difference in trends δ_{pre} .

It is natural to place restrictions on the relationship between δ_{pre} and δ_{post} , since applied researchers frequently test the null hypothesis $\delta_{pre} = 0$ in order to assess the plausibility of the assumption that $\delta_{post} = 0$. Our identification framework, which generalizes the partial identification framework in Manski and Pepper (2018), thus helps formalize the intuition motivating pre-trends testing.

Under the assumption that $\delta \in \Delta \neq \{\delta : \delta_{post} = 0\}$, the parameter θ will typically be set-identified. The *identified set* is the set of values for θ that are consistent with a given value of β under the restriction $\delta \in \Delta$,

$$\mathcal{S}(\beta, \Delta) := \left\{ \theta : \exists \delta \in \Delta, \tau_{post} \in \mathbb{R}^{\bar{T}} \text{ s.t. } l' \tau_{post} = \theta, \beta = \delta + \begin{pmatrix} 0 \\ \tau_{post} \end{pmatrix} \right\}.$$
(4)

When Δ is a closed and convex set, the identified set has a simple characterization.

Lemma 2.1. If Δ is closed and convex, then $S(\beta, \Delta)$ is an interval in \mathbb{R} , $S(\beta, \Delta) = [\theta^{lb}(\beta, \Delta), \theta^{ub}(\beta, \Delta)]$, where

$$\theta^{lb}(\beta, \Delta) := l'\beta_{post} - \underbrace{\left(\max_{\delta} l'\delta_{post}, s.t. \ \delta \in \Delta, \delta_{pre} = \beta_{pre}\right)}_{=:b^{max}(\beta_{pre}, \Delta)},\tag{5}$$

$$\theta^{ub}(\beta, \Delta) := l'\beta_{post} - \underbrace{\left(\min_{\delta} l'\delta_{post}, s.t. \ \delta \in \Delta, \delta_{pre} = \beta_{pre}\right)}_{=:b^{min}(\beta_{pre}, \Delta)}.$$
(6)

Proof. Re-arranging terms in (4), the identified set can be equivalently written as $S(\beta, \Delta) = \{\theta : \exists \delta \in \Delta \text{ s.t. } \delta_{pre} = \beta_{pre}, \theta = l' \beta_{post} - l' \delta_{post} \}$. The result is then immediate.

Example: Non-staggered DiD (continued) In the three-period DiD model ($\underline{T} = \overline{T} = 1$), the ATT in period 1 is point identified if we assume that the counterfactual post-treatment

difference in trends δ_1 is exactly zero (parallel trends). Instead, we assume $\delta = (\delta_{-1}, \delta_1)' \in \Delta$ for some set Δ . When Δ is closed and convex, the identified set for the ATT in period 1 is $[\beta_1 - b^{max}, \beta_1 - b^{min}]$, where $b^{max} = \max_{\delta} \delta_1$ s.t $(\delta_{-1}, \delta_1)' \in \Delta$ is the maximum possible bias of $\hat{\beta}_1$ given $\delta_{-1} = \beta_{-1}$ and b^{min} is defined analogously.

Additionally, it is immediate from the definition of the identified set in (4) that if Δ is the finite union of sets, $\Delta = \bigcup_{k=1}^{K} \Delta_k$, then its identified set is the union of the identified sets for its subcomponents,

$$\mathcal{S}(\beta, \Delta) = \bigcup_{k=1}^{K} \mathcal{S}(\beta, \Delta_k).$$
(7)

This fact will be useful, since several empirically relevant choices of Δ can be written as the finite union of convex sets as we will see below.

2.4 Possible choices of Δ

The class of possible differences in trends Δ must be specified by the researcher, and the choice of Δ will depend on the economic context. We highlight several possible choices of Δ that may be reasonable in empirical applications and formalize intuitive arguments that are commonly made by applied researchers regarding possible violations of parallel trends. Throughout our discussion, we write $\delta_{pre} = (\delta_{-T}, ... \delta_{-1})'$ and $\delta_{post} = (\delta_1, ... \delta_{\bar{T}})'$, with δ_0 normalized to zero. This aligns the notation with Example 1, where δ corresponds to the difference in trends between treated and comparison groups, and δ_0 is normalized to zero.

2.4.1 Bounding Relative Magnitudes

In empirical applications, researchers may be willing to assume that the confounding factors which produce non-parallel trends in the post-treatment periods are not too much larger in magnitude than the confounding factors in the pre-treatment periods. In their empirical application to right-to-carry gun laws, Manski and Pepper (2018) operationalize this intuition by calibrating bounds on $|\delta_1|$ to the largest violations of parallel trends in the pre-treatment period (see their Table 3).⁴ Such a restriction can be formalized in our framework by imposing that $\delta \in \Delta^{RM}(\bar{M})$ for $\bar{M} \ge 0$, where

$$\Delta^{RM}(\bar{M}) = \{ \delta : \forall t \ge 0, \ |\delta_{t+1} - \delta_t| \le \bar{M} \cdot \max_{s < 0} |\delta_{s+1} - \delta_s| \}.$$

⁴In their application, Manski and Pepper (2018) observe the outcome for the entire population of interest, and thus their observed pre-treatment data corresponds with δ_{pre} rather than $\hat{\beta}_{pre}$.

 $\Delta^{RM}(\bar{M})$ bounds the maximum post-treatment violation of parallel trends between consecutive periods by \bar{M} times the maximum pre-treatment violation of parallel trends. We use the abbreviation RM for "relative magnitudes". The choice $\Delta^{RM}(\bar{M})$ may be reasonable if the researcher suspects that possible violations of parallel trends are driven by confounding economic shocks that are of a similar magnitude to confounding economics shocks in the pre-period. When the number of pre-treatment and post-treatment periods is similar, a natural benchmark may be $\bar{M} = 1$, which bounds the worst-case post-treatment difference in trends by the equivalent maximum in the pre-treatment period.

Example: Non-staggered DiD (continued) In the three-period DiD model ($\underline{T} = \overline{T} = 1$), assuming $\delta \in \Delta^{RM}(\overline{M}) = \{(\delta_{-1}, \delta_1)' : |\delta_1| \leq \overline{M} |\delta_{-1}|\}$ bounds the magnitude of δ_1 based on the magnitude of δ_{-1} . The larger the magnitude of the pre-treatment violation in parallel trends, $|\delta_{-1}|$, the wider the range of possible post-treatment violations of parallel trends.

2.4.2 Smoothness restrictions

In other empirical settings, researchers may be worried about confounding from secular trends (e.g. long-run changes in labor supply) that they suspect evolve smoothly over time. In such settings, it is common for empirical researchers to control for a linear group-specific time trend.⁵ This approach is valid if the difference in trends is linear, i.e. $\Delta = \{\delta : \delta_t = \gamma \cdot t, \gamma \in \mathbb{R}\}$. There are often concerns, however, that the linear specification is not exactly correct (Wolfers, 2006; Lee and Solon, 2011). A natural relaxation is therefore to impose only that the differential trends evolve smoothly over time by bounding the extent to which its slope may change across consecutive periods. Such a restriction can be formalized in our framework by imposing that $\delta \in \Delta^{SD}(M)$ for $M \ge 0$, where

$$\Delta^{SD}(M) := \{ \delta : |(\delta_{t+1} - \delta_t) - (\delta_t - \delta_{t-1})| \leq M, \forall t \}.$$

$$\tag{8}$$

The parameter $M \ge 0$ governs the amount by which the slope of δ can change between consecutive periods, and thus bounds the discrete analog of the second derivative. We use the abbreviation SD for "second differences" or "second derivative."⁶ In the special case where $M = 0, \Delta^{SD}(0)$ requires that the difference in trends be exactly linear, which corresponds

⁵Specifically, researchers often augment specification (2) with group-specific linear trends, an approach Dobkin et al. (2018) refer to as a "parametric event-study." An analogous approach is to estimate a linear trend using only observations prior to treatment, and then subtract out the estimated linear trend from the observations after treatment (Bhuller et al., 2013; Goodman-Bacon, 2018, 2021).

⁶Restrictions on the second derivative of the conditional expectation function or density have been used in regression discontinuity settings (Kolesár and Rothe, 2018; Frandsen, 2016; Noack and Rothe, 2020). Smoothness restrictions are also used to obtain partial identification in Kim, Kwon, Kwon and Lee (2018).

with the assumption underlying the parametric linear specification common in applied work.

Example: Non-staggered DiD (continued). In the three-period DiD model, assuming the differential trend is exactly linear is equivalent to assuming $\Delta = \{\delta : \delta_1 = -\delta_{-1}\}$. Assuming $\delta \in \Delta^{SD}(M)$ requires only that the linear extrapolation be *approximately* correct, $\delta_1 \in [-\delta_{-1} - M, -\delta_{-1} + M]$.

2.4.3 Combining smoothness and relative magnitudes bounds

In some contexts, researchers may be willing to assume that the difference in trends evolves relatively smoothly over time but may be unsure about the smoothness bound $M \ge 0$ introduced above. In such cases, it may be reasonable to assume that the possible non-linearities in the post-treatment difference in trends are bounded by the observed non-linearities in the pre-treatment difference in trends. This can be formalized with the restriction

$$\Delta^{SDRM}(\bar{M}) = \{ \delta : \forall t \ge 0, \, |(\delta_{t+1} - \delta_t) - (\delta_t - \delta_{t-1})| \le \bar{M} \cdot \max_{s < 0} |(\delta_{s+1} - \delta_s) - (\delta_s - \delta_{s-1})| \},\$$

which bounds the maximum deviation from a linear trend in the post-treatment period by $\overline{M} \ge 0$ times the equivalent maximum in the pre-treatment period. The set $\Delta^{SDRM}(\overline{M})$ is thus similar to $\Delta^{SD}(M)$ introduced above, except it allows the magnitude of the possible non-linearity to explicitly depend on the observed pre-trends.

2.4.4 Sign and monotonicity restrictions

Context-specific knowledge may sometimes also suggest sign or monotonicity restrictions on the differential trend. For instance, if the policy of interest occurs at the same time as a confounding policy change that we expect to have a positive effect on the outcome, we might restrict the post-treatment bias to be positive, $\delta \in \Delta^{PB} := \{\delta : \delta_t \ge 0 \ \forall t \ge 0\}$. Likewise, there may be secular pre-existing trends that we expect would have continued following the treatment date.⁷ We may then wish to impose that the differential trend be increasing, $\delta \in \Delta^I := \{\delta : \delta_t \ge \delta_{t-1} \ \forall t\}$, or monotone with unknown sign, $\delta \in \Delta^{Mon} := \Delta^I \cup (-\Delta^I)$. Sign and monotonicity restrictions may be combined with the previously discussed restrictions, such as $\Delta^{SDPB}(M) := \Delta^{SD}(M) \cap \Delta^{PB}$, $\Delta^{SDI}(M) := \Delta^{SD}(M) \cap \Delta^I$, and $\Delta^{RMI}(\bar{M}) :=$ $\Delta^{RM}(\bar{M}) \cap \Delta^I$.

⁷Monotone violations of parallel trends are often discussed in applied work. For example, Lovenheim and Willen (2019) argue that violations of parallel trends cannot explain their results because "pre-[treatment] trends are either zero or in the wrong direction (i.e., opposite to the direction of the treatment effect)." Greenstone and Hanna (2014) estimate upward-sloping pre-existing trends and argue that "if the pre-trends had continued" their estimates would be upward biased.

2.4.5 Polyhedral restrictions

Although the restrictions described above will be sensible in many empirical contexts, researchers will often have context-specific knowledge that motivates alternative restrictions than what we introduced above. To accommodate such cases, we consider the broad class of Δ s that can be written as polyhedra (sets defined by linear inequalities), or the finite union of polyhedra.

Definition 1 (Polyhedral restriction). The class Δ is polyhedral if it takes the form $\Delta = \{\delta : A\delta \leq d\}$ for some known matrix A and vector d.

All of the examples described above can be written either as polyhedral restrictions or finite unions of such restrictions. For instance, $\Delta^{SD}(M)$ and $\Delta^{SDPB}(M)$ can be written directly as polyhedra.⁸ Likewise, $\Delta^{RM}(\bar{M})$ or $\Delta^{SDRM}(\bar{M})$ can be written as the finite union of polyhedra, where each polyhedron corresponds with a different location for the maximum pre-treatment violation.⁹

The class of (finite unions of) polyhedra is quite broad, and allows for a variety of other restrictions that may be relevant in empirical work. For example researchers studying labor market training and related programs may be concerned about Ashenfelter's dip (Ashenfelter, 1978), in which earnings for the treated group trend downwards (relative to control) before treatment and upwards afterwards. In this type of setting, researchers might naturally use a polyhedral Δ to impose i) restrictions on the signs of the pre-treatment and post-treatment biases, as well as ii) restrictions on the magnitude of the rebound effect relative to the pre-treatment shock.

2.5 Inferential Goal

As discussed above, the event study coefficients $\hat{\beta}_n$ will satisfy $\sqrt{n}(\hat{\beta}_n - \beta) \rightarrow_d \mathcal{N}(0, \Sigma^*)$ for a wide variety of commonly-used estimators. This suggests the finite-sample normal approximation

$$\hat{\beta}_n \approx_d \mathcal{N}\left(\beta, \Sigma_n\right),\tag{9}$$

where \approx_d denotes approximate equality in distribution and $\Sigma_n = \Sigma^*/n$. We will construct confidence sets that are uniformly valid for all parameter values θ in the identified set when

⁸In our ongoing three-period difference-in-differences example, $\Delta^{SD}(M) = \{\delta : A^{SD}\delta \leq d^{SD}\}$ for $A^{SD} = \begin{pmatrix} -1 & 1 \\ 1 & -1 \end{pmatrix}$ and $d^{SD} = (M, M)'$. This generalizes naturally when there are multiple pre-periods and multiple post-periods.

⁹For example, define the polyhedra $\Delta_{s,+}^{RM}(\bar{M}) = \{\delta : \forall t \ge 0, |\delta_{t+1} - \delta_t| \le \bar{M}(\delta_{s+1} - \delta_s)\}$ and $\Delta_{s,-}^{RM} = \{\delta : \forall t \ge 0, |\delta_{t+1} - \delta_t| \le -\bar{M}(\delta_{s+1} - \delta_s)\}$. Then $\Delta^{RM}(\bar{M}) = \bigcup_{s<0} (\Delta_{s,+}^{RM}(\bar{M}) \cup \Delta_{s,-}^{RM}(\bar{M}))$.

the approximation in (9) holds exactly with Σ_n known. That is, we construct confidence sets $C_n(\hat{\beta}_n, \Sigma_n)$ satisfying

$$\inf_{\delta \in \Delta, \tau} \inf_{\theta \in \mathcal{S}(\delta + \tau, \Delta)} \mathbb{P}_{\hat{\beta}_n \sim \mathcal{N}(\delta + \tau, \Sigma_n)} \left(\theta \in \mathcal{C}_n(\hat{\beta}_n, \Sigma_n) \right) \ge 1 - \alpha.$$
(10)

In Section 3.3, we will show that finite-sample size control in the normal model in the sense of (10) translates to uniform asymptotic size control over a large class of data-generating processes when Σ_n is replaced by a consistent estimate $\hat{\Sigma}_n$. That is, we will show that the constructed confidence sets further satisfy

$$\liminf_{n \to \infty} \inf_{P \in \mathcal{P}} \inf_{\theta \in \mathcal{S}(\delta_P + \tau_P, \Delta)} \mathbb{P}_P\left(\theta \in \mathcal{C}_n(\hat{\beta}_n, \hat{\Sigma}_n)\right) \ge 1 - \alpha.$$
(11)

for a large class of distributions \mathcal{P} such that $\delta_P \in \Delta$ for all $P \in \mathcal{P}$.

We will focus on constructing confidence sets for the case where Δ is a polyhedron. For the case where Δ is the finite union of polyhedra, a valid confidence set can be constructed by taking the union of the confidence sets for each of its components.

Lemma 2.2. Suppose that for each k = 1, ..., K, the confidence set $C_{n,k}(\hat{\beta}_n, \Sigma_n)$ satisfies (10) with $\Delta = \Delta_k$. Then the confidence set $C_n(\hat{\beta}_n, \Sigma_n) = \bigcup_{k=1}^K C_{n,k}(\hat{\beta}_n, \Sigma_n)$ satisfies (10) with $\Delta = \bigcup_{k=1}^K \Delta_k$.

In the next two sections, we introduce two approaches to obtain confidence sets satisfying (10), with different desirable properties depending on the form of Δ . The first approach, based on moment inequalities, accommodates a wide range of restrictions Δ and has some desirable asymptotic power guarantees. The second approach, based on fixed length confidence intervals, can potentially offer finite-sample power improvements for certain special classes of Δ of interest, such as $\Delta^{SD}(M)$.

3 Inference using Moment Inequalities

In this section, we introduce a general approach for inference that has good asymptotic properties over a large class of possible restrictions Δ . We show that inference on the partially identified parameter $\theta = l' \tau_{post}$ in this setting is equivalent to testing a system of moment inequalities with a potentially large number of nuisance parameters that enter the moments linearly. We consider an implementation based on the conditional approach developed in ARP, which allows us to obtain computationally tractable confidence sets with desirable power properties for many parameter configurations.

Representation as a moment inequality problem with linear nui-3.1sance parameters

Consider the problem of conducting inference on $\theta = l' \tau_{post}$ when Δ takes the polyhedral form $\Delta = \{\delta : A\delta \leq d\}$. We will develop tests that have exact size under the null hypothesis $H_0: \theta = \overline{\theta}, \delta \in \Delta$ when the normal approximation (9) holds exactly with known variance matrix Σ_n . In Section 3.3, we will provide conditions under which size control in the finite sample normal model translates to uniform asymptotic size control over a large class of data-generating processes.

As a first step, we show that testing H_0 is equivalent to testing a system of moment inequalities with linear nuisance parameters in the normal model. Observe that if $\hat{\beta}_n \sim$ $\mathcal{N}(\beta, \Sigma_n)$ for β satisfying (3), then $\mathbb{E}_{\hat{\beta}_n \sim \mathcal{N}(\delta + \tau, \Sigma_n)} \left[\hat{\beta}_n - \tau \right] = \delta$. It follows that $\delta \in \Delta =$ $\{\delta : A\delta \leq d\}$ if and only if $\mathbb{E}_{\hat{\beta}_n \sim \mathcal{N}(\delta + \tau, \Sigma_n)} \left[A\hat{\beta}_n - A\tau \right] \leq d$. Defining $Y_n = A\hat{\beta}_n - d$ and $L_{post} = [0, I]'$ to be the matrix such that $\tau = L_{post} \tau_{post}$, it is immediate that the null hypothesis H_0 is equivalent to the composite null

$$H_0: \exists \tau_{post} \in \mathbb{R}^{\bar{T}} \text{ s.t. } l' \tau_{post} = \bar{\theta} \text{ and } \mathbb{E}_{\hat{\beta}_n \sim \mathcal{N}(\delta + \tau, \Sigma_n)} \left[Y_n - AL_{post} \tau_{post} \right] \leq 0.$$
(12)

Testing the null hypothesis H_0 is therefore equivalent to testing that the moment inequalities $\mathbb{E}_{\hat{\beta}_n \sim \mathcal{N}(\delta + \tau, \Sigma_n)} \left[Y_n - AL_{post} \tau_{post} \right] \leq 0 \text{ hold for some value of } \tau_{post} \text{ satisfying } l' \tau_{post} = \theta.$

For the purposes of developing tests, it will be useful to re-cast this null hypothesis in terms of moments involving an unrestricted nuisance parameter $\tilde{\tau}$ of dimension $\bar{T} - 1$. By applying a change of basis to the matrix AL_{post} , we can re-write the expression $AL_{post}\tau_{post}$ as $\tilde{A}(\theta, \tilde{\tau}')'$ for $\tilde{\tau} \in \bar{\tau}^{-1}$.¹⁰ The null H_0 is then equivalent to

$$H_0: \exists \tilde{\tau} \in \mathbb{R}^{\bar{T}-1} \text{ s.t. } \mathbb{E}\left[\tilde{Y}_n(\bar{\theta}) - \tilde{X}\tilde{\tau}\right] \leqslant 0,$$
(13)

where $\tilde{Y}(\bar{\theta}) = Y_n - \tilde{A}_{(\cdot,1)}\bar{\theta}$ and $\tilde{X} = \tilde{A}_{(\cdot,-1)}$.¹¹ Since $\tilde{Y}_n(\bar{\theta})$ is normally distributed with covariance matrix $\tilde{\Sigma}_n = A \Sigma_n A'$ under the finite-sample normal model, testing $H_0: \theta = \bar{\theta}, \delta \in$ Δ is equivalent to testing a system of moment inequalities with linear nuisance parameters.

The testing problem (13) is a special case of the problem studied in ARP, which focuses on

¹⁰Let Γ be a square matrix with the vector l' in the first row and remaining rows chosen so that Γ has full rank. Define $\tilde{A} := AL_{post}\Gamma^{-1}$. Then $AL_{post}\tau_{post} = \tilde{A}\Gamma\tau_{post} = \tilde{A}\left(\underbrace{\begin{array}{c}\theta\\ \Gamma_{(-1,\cdot)}\tau_{post}\end{array}}_{\tilde{L}}\right)$. If $\bar{T} = 1$, then $\tilde{\tau}$ is

⁰⁻dimensional and should be interpreted as 0.

¹¹We use the notation $V_{(\cdot,1)}$ to denote the first column of a matrix V, and $V_{(\cdot,-1)}$ to denote the matrix containing all but the first column of V.

testing null hypotheses of the form $H_0: \exists \tau \text{ s.t. } \mathbb{E}[Y(\theta) - X\tau \mid X] \leq 0$ (a.s.). Our setting is a special case of this framework in which: i) the variable X takes the degenerate distribution $X = \tilde{X}$, and ii) $Y(\theta) = \tilde{Y}(\theta)$ is linear in θ . This additional structure will play an important role in the development of our asymptotic power results below.

3.2 Constructing conditional and hybrid confidence sets

We consider tests for the system of moment inequalities describe above using the conditional and hybrid methods proposed by ARP. This is for both computational and efficiency reasons. From the computational perspective, a practical challenge to testing the hypothesis (13) in our setting is that the dimension of the nuisance parameter $\tilde{\tau}$ is $\bar{T} - 1$, and thus will be large if there are many post-treatment periods. For example, 5 of the 12 recent event-study papers reviewed in Roth (Forthcoming) have $\bar{T} > 10$. This renders many moment inequality methods, especially those which rely on test inversion over a grid for the full parameter vector, computationally infeasible. To tractably deal with the nuisance parameter, we consider tests based on the conditional and hybrid approaches of ARP, which directly exploit the linear structure of the hypothesis (13) to deliver computationally tractable tests even when the number of post-treatment periods \bar{T} is large.¹² From the perspective of power, we will show that the tests proposed by ARP have (near-)optimal local asymptotic power in our setting when an LICQ condition is satisfied.

We briefly sketch the construction of the conditional testing approach in our setting, and refer the reader to ARP for full details. These tests are implemented in the R package, HonestDiD, that accompanies the paper.

Conditional confidence sets. Suppose we wish to test (13) for some fixed $\bar{\theta}$. The conditional testing approach considers tests based on the profiled test statistic

$$\hat{\eta} := \min_{\eta, \tilde{\tau}} \eta \text{ s.t. } \tilde{Y}_n(\bar{\theta}) - \tilde{X}\tilde{\tau} \leqslant \tilde{\sigma}_n \cdot \eta,$$
(14)

where $\tilde{\sigma}_n = \sqrt{\operatorname{diag}(\tilde{\Sigma}_n)}$. This linear program selects the value of the nuisance parameters $\tilde{\tau} \in \mathbb{R}^{\bar{T}-1}$ that minimizes the maximum studentized moment. Duality results from linear

¹²Other moment inequality methods have been proposed for subvector inference, but typically do not exploit the linear structure of our setting — see, e.g., Romano and Shaikh (2008); Chernozhukov, Newey and Santos (2015); Bugni, Canay and Shi (2017); Chen, Christensen and Tamer (2018); Kaido, Molinari and Stoye (2019). Cho and Russell (2019), Gafarov (2019), and Flynn (2019) also provide methods for subvector inference with linear moment inequalities, but in contrast to our approach require a linear independence constraint qualification (LICQ) assumption for size control. More recently, Cox and Shi (2022) introduced new tests for the linear moment inequality setting in ARP; see Section 3.5 below for further discussion.

programming (e.g. Schrijver (1986), Section 7.4) imply that the value $\hat{\eta}$ obtained from the primal program (14) equals the optimal value of the dual program¹³

$$\hat{\eta} = \max_{\gamma} \gamma' \tilde{Y}_n(\bar{\theta}) \text{ s.t. } \gamma' \tilde{X} = 0, \gamma' \tilde{\sigma}_n = 1, \gamma \ge 0.$$
(15)

If a vector γ_* is optimal in the dual problem above, then it is a vector of Lagrange multipliers for the primal problem. Standard results in linear programming imply that the optimum is always obtained at one of the finite set of vertices, $V(\Sigma_n)$ (also known as the set of basic feasible solutions). We denote by $\hat{V}_n \subset V(\Sigma_n)$ the set of optimal vertices of the dual program.¹⁴

To construct critical values, ARP use the fact that the distribution of $\hat{\eta}$ has a truncated normal distribution conditional on the event that γ_* is optimal in the dual problem. Specifically,

$$\hat{\eta} \mid \{\gamma_* \in \hat{V}_n, S_n = s\} \sim \xi \mid \xi \in [v^{lo}, v^{up}],$$

where $\xi \sim \mathcal{N}\left(\gamma'_*\tilde{\mu}(\bar{\theta}), \gamma'_*\tilde{\Sigma}_n\gamma_*\right), \tilde{\mu}(\bar{\theta}) = \mathbb{E}\left[\tilde{Y}_n(\bar{\theta})\right], S_n = (I - \frac{\tilde{\Sigma}_n\gamma_*}{\gamma'_*\tilde{\Sigma}_n\gamma_*}\gamma'_*)\tilde{Y}_n(\bar{\theta}), \text{ and } v^{lo}, v^{up}$ are known functions of $\tilde{\Sigma}_n, s, \gamma_*$ (see Lemma 2 in ARP).¹⁵ Intuitively, the distribution of $\hat{\eta}$ depends on the vector $\tilde{\mu}(\bar{\theta})$, and so to eliminate the dependence on the components of $\tilde{\mu}(\bar{\theta})$ other than $\gamma'\tilde{\mu}(\bar{\theta})$, we condition on S_n , which is a sufficient statistic for the components of $\tilde{\mu}(\bar{\theta})$ that are orthogonal to $\gamma'_*\tilde{\mu}(\bar{\theta})$.

ARP show that all quantiles of the conditional distribution of $\hat{\eta}$ in the previous display are increasing in $\gamma'_* \tilde{\mu}(\bar{\theta})$. Moreover, the null hypothesis (13) implies $\gamma'_* \tilde{\mu}(\bar{\theta}) \leq 0$. To see why this is the case, note that the definition of the dual problem (15) implies that $\gamma_* \geq 0$ and $\gamma'_* \tilde{X} = 0$, whereas the null hypothesis implies that there exists $\tilde{\tau}$ such that $\tilde{\mu}(\bar{\theta}) - \tilde{X}\tilde{\tau} \leq 0$. It follows that $\gamma'_* \tilde{\mu}(\bar{\theta}) = \gamma'_* (\tilde{\mu}(\bar{\theta}) - \tilde{X}\tilde{\tau}) \leq 0$ under the null. The ARP conditional test therefore uses the critical value max $\{0, c_{C,\alpha}\}$, where $c_{C,\alpha}$ is the $1 - \alpha$ quantile of the truncated normal distribution $\xi | \xi \in [v^{lo}, v^{up}]$ under the worst-case assumption that $\gamma'_* \tilde{\mu}(\bar{\theta}) = 0$.¹⁶ We denote

¹³Technically, the duality results require that $\hat{\eta}$ be finite. However, one can show that $\hat{\eta}$ is finite with probability 1, unless the span of \tilde{X} contains a vector with all negative entries, in which case the identified set for θ is the real line. We therefore trivially define our test never to reject if $\hat{\eta} = -\infty$.

¹⁴In general, there may not be a unique solution to the dual program. ARP show that in the context of the finite sample normal, conditional on any one vertex of the dual program's feasible set being optimal, every other vertex is optimal with either probability 0 or 1. In the finite sample normal model it thus suffices to condition on the event that a vector $\gamma_* \in \hat{V}$. Our conditions for asymptotic validity of the conditional test below, however, ensure that the optimal vertex will be unique w.p.a. 1.

¹⁵The cutoffs v^{lo} and v^{up} are the maximum and minimum of the set $\{x : x = \max_{\gamma \in F_n} \gamma'(s + \frac{\Sigma_n \gamma_*}{\gamma'_* \tilde{\Sigma}_n \gamma_*} x)\}$ when $\gamma'_* \tilde{\Sigma}_n \gamma_* \neq 0$, where F_n is the feasible set of the dual program (15). When $\gamma'_* \tilde{\Sigma}_n \gamma_* = 0$, we define $v^{lo} = -\infty$ and $v^{up} = \infty$, so the conditional test rejects if and only if $\hat{\eta} > 0$.

¹⁶As noted in ARP, the truncation at 0 is not necessary for the conditional test to control size in the finite sample normal model, but it simplifies asymptotic arguments. It also prevents the test from rejecting when

by $\psi_{\alpha}^{C}(\hat{\beta}_{n}, A, d, \bar{\theta}, \Sigma_{n})$ an indicator for whether the conditional test rejects the null that $\theta = \bar{\theta}$ for $\Delta = \{\delta : A\delta \leq d\}$.

We can then form a confidence set for θ by test inversion, $\mathcal{C}^{C}_{\alpha,n}(\hat{\beta}_{n}, \Sigma_{n}) := \{\bar{\theta} : \psi^{C}_{\alpha}(\hat{\beta}_{n}, A, d, \bar{\theta}, \Sigma_{n}) = 0\}$. The construction of the conditional test implies that $\mathbb{E}_{\hat{\beta}_{n}\sim\mathcal{N}(\delta+\tau,\Sigma_{n})}\left[\psi^{C}_{\alpha}(\hat{\beta}_{n}, A, d, l'\tau_{post}, \Sigma_{n})\right] \leq \alpha$ for any $\delta \in \Delta$. It therefore follows that $\mathcal{C}^{C}_{\alpha,n}(\hat{\beta}_{n}, \Sigma_{n})$ satisfies the finite-sample coverage requirement (10). In Section 3.3 below, we show that coverage in the normal model translates to uniform asymptotic coverage over a large class of DGPs.

Example 4. An instructive example is when $\overline{T} = 1$ (so that there are no nuisance parameters), and $\tilde{\Sigma}_n = I$. Then $\hat{\eta} = \max_j \tilde{Y}_{n,j}$ is the maximum component of \tilde{Y}_n , $v^{lo} = \max_{j \neq \hat{j}} \tilde{Y}_{n,j}$ is the second-largest element of \tilde{Y}_n (where \hat{j} denotes the index of the max), and $v^{up} = \infty$. Thus, the conditional test rejects when $\hat{\eta}$ exceeds the $1 - \alpha$ quantile of the standard normal distribution truncated to $[v^{lo}, \infty)$. Intuitively, this means that the conditional test will tend to reject when the maximum sample moment is far enough away from the second-largest sample moment. Two special cases are worth special consideration. First, consider the case where in population one moment is violated and the remaining moments are very slack, e.g. $\tilde{\mu}_1 > 0$ while $\tilde{\mu}_j \ll 0$ for $j \neq 1$. Then with high probability $\hat{\eta}$ will equal $\tilde{Y}_{n,1}$ and v^{lo} will be very negative. Thus, the conditional test will behave similarly to a one-sided t-test using $Y_{n,1}$, which can be shown to be the most powerful test in the finite-sample normal model in this example. On the other hand, if $\mu_1 \approx \mu_2 > 0$, then the maximum and second-largest sample moments (i.e. $\hat{\eta}$ and v^{lo}) will be close together with high probability, so the conditional test may not reject with substantial probability even if both μ_1 and μ_2 are large, and thus the conditional test may have poor power. To improve power in these settings where the binding and non-binding moments are close together (relative to sampling variation), ARP introduce a "hybrid" test, which we describe next.

Hybrid confidence sets. ARP propose a "hybrid" test that combines the conditioning approach above with a test based on the "least-favorable" assumption that $\tilde{\mu}(\bar{\theta}) = 0$. In particular, ARP show that the distribution of $\hat{\eta}$ under the null is bounded above (in the sense of first-order stochastic dominance) by the distribution of $\hat{\eta}$ when $\tilde{\mu}(\bar{\theta}) = 0$ (see Section 3.2 of ARP). One can therefore construct a size- κ least-favorable (LF) test in the finite-sample normal model that rejects whenever $\hat{\eta}$ exceeds the $1 - \kappa$ quantile of $\max_{\gamma \in V(\Sigma)} \gamma' \xi$, where $\xi \sim \mathcal{N}\left(0, \tilde{\Sigma}_n\right)$. This critical value, which we will denote by $c_{LF,\kappa}$ can easily be calculated by simulation. For $0 < \kappa < \alpha$, the ARP conditional-LF hybrid test is defined to reject if a first-stage, size- κ LF test rejects. If this first-stage test does not reject, then in the second

all moments are satisfied in sample.

stage the hybrid test conducts a modified version of the size- $\left(\frac{\alpha-\kappa}{1-\kappa}\right)$ conditional test that also conditions on the event that the first-stage LF test did not reject. In particular, by similar logic as for the conditional test, we have that

$$\hat{\eta} \mid \{\gamma_* \in \hat{V}_n, S_n = s, \hat{\eta} \leqslant c_{LF,\kappa}\} \sim \xi \mid \xi \in [v^{lo}, v_H^{up}],$$

where $v_H^{up} = \min\{v^{lo}, c_{LF,\kappa}\}$ (see Section 3.4 of ARP). The second-stage of the hybrid test rejects if $\hat{\eta}$ exceeds the critical value for the size- $\left(\frac{\alpha-\kappa}{1-\kappa}\right)$ conditional test that uses v_H^{up} instead of v^{up} . We will denote by $\psi_{\kappa,\alpha}^{C-LF}(\hat{\beta}_n, A, d, \bar{\theta}, \Sigma_n)$ an indicator for whether the hybrid test rejects at a particular value $\bar{\theta}$, and denote by $\mathcal{C}_{\kappa,\alpha,n}^{C-LF}(\hat{\beta}_n, \Sigma_n)$ the confidence set that collects the values of $\bar{\theta}$ for which the hybrid test does not reject. As with the conditional test, by construction the hybrid confidence set satisfies that coverage criterion (10) in the finite-sample normal model. In our implementation below, we use $\kappa = \alpha/10$, following ARP.

3.3 Uniform asymptotic size control

We now provide conditions under which size control in the finite sample normal model translates to uniform asymptotic size control over a large class of data-generating processes \mathcal{P} under which $\hat{\beta}_n$ is asymptotically normally distributed and Σ_n is replaced with a consistent estimator $\hat{\Sigma}_n$. In particular, we provide sufficient conditions on $\hat{\beta}_n$, $\hat{\Sigma}_n$, and Δ such that the higher-level conditions for size control in ARP are satisfied (Proposition 2 in ARP).

Throughout this section, we fix $\Delta = \{A\delta \leq d\}$ for some A with all non-zero rows, and assume that Δ is non-empty. We consider a class of data-generating processes, indexed by $P \in \mathcal{P}$, under which $\sqrt{n}(\hat{\beta}_n - \beta_P)$ is asymptotically normal, where β_P satisfies the causal decomposition in (3), i.e. $\beta_P = \delta_P + L_{post}\tau_{P,post}$ for $\delta_P \in \Delta$ and $\tau_{P,post} \in \mathbb{R}^T$. The parameter of interest is $\theta_P := l'\tau_{P,post}$, for some fixed $l \neq 0$. Our first assumption imposes uniform asymptotic normality of $\hat{\beta}_n$.

Assumption 2. Let BL_1 denote the set of Lipschitz functions which are bounded by 1 in absolute value and have Lipschitz constant bounded by 1. We assume

$$\lim_{n \to \infty} \sup_{P \in \mathcal{P}} \sup_{f \in BL_1} \left| \mathbb{E}_P \left[f(\sqrt{n}(\hat{\beta}_n - \beta_P)) \right] - \mathbb{E} \left[f(\xi_P) \right] \right| = 0,$$

where $\xi_P \sim \mathcal{N}(0, \Sigma_P)$, and $\beta_P = \delta_P + L_{post} \tau_{P,post}$ for $\delta_P \in \Delta$ and $\tau_{P,post} \in \mathbb{R}^{\overline{T}}$.

Convergence in distribution is equivalent to convergence in bounded Lipschitz metric (see Theorem 1.12.4 in van der Vaart and Wellner (1996)), so Assumption 2 formalizes the notion of uniform convergence in distribution of $\sqrt{n}(\hat{\beta}_n - \beta_P)$ to a $\mathcal{N}(0, \Sigma_P)$ variable under P. Our next two assumptions require that the eigenvalues of the asymptotic variance of $\hat{\beta}_n$ be bounded above and away from zero, and that there exists a uniformly consistent estimator for the variance of $\hat{\beta}_n$.

Assumption 3. Let **S** denote the set of matrices with eigenvalues bounded below by $\underline{\lambda} > 0$ and above by $\overline{\lambda} \ge \underline{\lambda}$. For all $P \in \mathcal{P}, \Sigma_P \in \mathbf{S}$.

Assumption 4. We have an estimator $\hat{\Sigma}_n$ that is uniformly consistent for Σ_P ,

$$\lim_{n \to \infty} \sup_{P \in \mathcal{P}} \mathbb{P}_P\left(\left| \left| \hat{\Sigma}_n - \Sigma_P \right| \right| > \epsilon \right) = 0,$$

for all $\epsilon > 0$.

Finally, we impose some regularity conditions on the matrix A.

Assumption 5. At least one of the following holds.

- (A) For $k_1 + k_2 = \dim(\delta)$, the matrix A can be written as TQ, where Q has full row-rank and $T = \begin{pmatrix} I_{k_1} & 0 \\ -I_{k_1} & 0 \\ 0 & I_{k_2} \end{pmatrix}$. (We allow for the case where one of k_1 or k_2 is 0, in which case the zero-dimensional blocks can be ignored).
- (B) Let $\bar{\gamma}_1, ..., \bar{\gamma}_K$ be the elements of V(I). Then for all k, either $\bar{\gamma}'_k A = 0$ or $\inf_{a \ge 0} \inf_{j \ne k} ||(\bar{\gamma}_k a\bar{\gamma}_j)'A|| > 0$.

Part (A) of Assumption 5 imposes that the only source of degeneracy in the rows of A is matching inequalities of opposite signs. This is the case for many restrictions of interest, such as $\Delta^{SD}(M)$ and the polyhedra that form $\Delta^{RM}(\bar{M})$. Part (B) provides an alternative, higher-level condition that ensures that for distinct vertices $\bar{\gamma}_k, \bar{\gamma}_j$, the random variables $\bar{\gamma}'_j \tilde{Y}$ and $\bar{\gamma}'_k \tilde{Y}$ are not perfectly positively correlated with each other. Assumption 5 is used to guarantee that degeneracy in the asymptotic distribution of $\gamma' \tilde{Y}$ arises only from known degeneracies in A. We note, however, that Assumption 5 does not rule out settings where the solutions to the bounds of the identified set given in equation (5) and (6) are non-unique or degenerate (i.e., where the extreme points for θ occur at "flat faces" of the identified set). If A is full-rank, for example, then Assumption 5(A) holds trivially with T = I, and thus the mean of the moments $\tilde{\mu}(\bar{\theta})$ is completely unrestricted.¹⁷

¹⁷Some values of A satisfying Assumption 5 may imply that certain pairs of moments cannot simultaneously be binding. For example, the restriction that $|\delta_1| \leq 1$ can be represented as $\delta_1 \leq 1$ and $-\delta_1 \leq 1$, which satisfies Assumption 5(A), but clearly both moments cannot simultaneously bind. Nevertheless Assumption 5(A) is compatible with "flat faces" even when A is not full rank. For example, if Δ corresponds with the restrictions $|\delta_1| \leq 1$ and $|\delta_2| \leq 1$, then the extreme points for $\theta = \tau_2$ occur at flat faces of the identified set for (τ_1, τ_2) .

The assumptions stated above are sufficient for the conditions in Proposition 2 in ARP, which establishes uniform size control for the conditional and hybrid tests.

Proposition 3.1. Suppose Assumptions 2 to 5 hold. Then the conditional and LF-hybrid tests uniformly control size. That is, for any $\alpha < 0.5$,

$$\limsup_{n \to \infty} \sup_{P \in \mathcal{P}} \mathbb{E}_{P} \left[\psi_{\alpha}^{C}(\hat{\beta}_{n}, A, d, \theta_{P}, \frac{1}{n} \hat{\Sigma}_{n}) \right] \leq \alpha.$$
$$\limsup_{n \to \infty} \sup_{P \in \mathcal{P}} \mathbb{E}_{P} \left[\psi_{\kappa, \alpha}^{C-LF}(\hat{\beta}_{n}, A, d, \theta_{P}, \frac{1}{n} \hat{\Sigma}_{n}) \right] \leq \alpha.$$

3.4 Uniform asymptotic consistency

We next provide conditions under which the conditional and hybrid tests are uniformly asymptotically consistent, in the sense that power against fixed alternatives outside the identified set converges uniformly to 1. To establish uniform consistency of the conditional and hybrid tests, we strengthen Assumptions 2 and 3 as follows.

Assumption 6. Let $W_n = ((\hat{\beta}_n - \beta_P)', (vec(\hat{\Sigma}_n) - vec(\Sigma_P))')'$, where $vec(\Sigma)$ is the vector of the elements of the matrix Σ . We assume

$$\lim_{n \to \infty} \sup_{P \in \mathcal{P}} \sup_{f \in BL_1} \left\| \mathbb{E}_P \left[f(\sqrt{n}W_n) \right] - \mathbb{E} \left[f(\xi_P^+) \right] \right\| = 0$$

where $\xi_P^+ \sim \mathcal{N}(0, V_P), V_P = \begin{pmatrix} \Sigma_P & V_{P,\beta\Sigma} \\ V_{P,\Sigma\beta} & V_{P,\Sigma} \end{pmatrix}.$

Assumption 7. For all $P \in \mathcal{P}$, $\Sigma_P \in \mathbf{S}$ and the matrix V_P defined in Assumption 6 lies in a compact set \mathbf{V} . Additionally, $(\Sigma_P - V_{P,\beta\Sigma}V_{P,\Sigma\beta}^{\dagger})$ has eigenvalues bounded below by $\tilde{\lambda} > 0$, where \dagger denotes the Moore-Penrose inverse.

Assumption 6 strengthens Assumption 2 to require that $\hat{\beta}_n$ and $\hat{\Sigma}_n$ have a joint normal asymptotic distribution. Although somewhat more restrictive, event-study estimates are often estimated via OLS, and standard covariance estimators for OLS, including cluster-robust variance estimators, produce asymptotically normal estimates as the number of clusters grows large (Hansen, 2007; Stock and Watson, 2008; Hansen and Lee, 2019). We do not impose that the asymptotic distributions of $\hat{\beta}_n$ and $\hat{\Sigma}_n$ are independent, as would occur in linear models if the linear model is properly specified. Assumption 7 strengthens Assumption 3 to require that the asymptotic distribution of $\hat{\beta}_n$ is not perfectly asymptotically colinear with $\hat{\Sigma}_n$. Under the imposed assumptions, we obtain uniform consistency of the conditional and hybrid tests.

Proposition 3.2. Suppose Assumptions 4 to 7 hold. Then for any x > 0 and $\alpha < 0.5$,

$$\lim_{n \to \infty} \inf_{P \in \mathcal{P}} \mathbb{E}_P \left[\psi_{\alpha}^C(\hat{\beta}_n, A, d, \theta_P^{ub} + x, \frac{1}{n} \hat{\Sigma}_n) \right] = 1$$
$$\lim_{n \to \infty} \inf_{P \in \mathcal{P}} \mathbb{E}_P \left[\psi_{\kappa, \alpha}^{C\text{-}LF}(\hat{\beta}_n, A, d, \theta_P^{ub} + x, \frac{1}{n} \hat{\Sigma}_n) \right] = 1,$$

where $\theta_P^{ub} = \sup \mathcal{S}(\beta_P, \Delta)$ is the upper bound of the identified set. The analogous result holds replacing $\theta_P^{ub} + x$ with $\theta_P^{lb} - x$ for $\theta_P^{lb} = \inf \mathcal{S}(\beta_P, \Delta)$.

3.5 Optimal local asymptotic power

We next provide conditions under which the conditional test has optimal local asymptotic power. We first state the conditions and our formal results, and then provide several examples highlighting when the assumptions will and will not hold.

3.5.1 Main results

We begin by defining the linear independence constraint qualification (LICQ). Recall that the upper bound of the identified set is given by

$$\theta^{ub}(\beta, \Delta) = l'\beta_{post} - \left(\min_{\delta} l'\delta_{post}, \text{ s.t. } A\delta \leqslant d, \delta_{pre} = \beta_{pre}\right).$$

Since $\delta_{post} = \beta_{post} - \tau_{post}$, we can re-write the upper bound as a maximization over τ_{post} ,

$$\theta^{ub}(\beta, \Delta) = \max_{\tau_{post}} l' \tau_{post}, \text{ s.t. } -A_{(\cdot, post)} \tau_{post} \leqslant d - A\beta,$$
(16)

where $A_{(\cdot,post)}$ contains the columns of A corresponding with δ_{post} . Let τ_{post}^* denote a solution to the optimization for $\theta^{ub}(\beta, \Delta)$ in (16), and let B^* denote the indices of the binding constraints, so that $-A_{(B^*,post)}\tau_{post}^* = d_{B^*} - A_{(B^*,\cdot)}\beta$ and $-A_{(-B^*,post)}\tau_{post}^* < d_{-B^*} - A_{(-B^*,\cdot)}\beta$.

Definition 2 (LICQ). We say that the linear independence constraint qualification (LICQ) holds in direction l if there exists a solution τ_{post}^* to (16) such that the gradient of the binding constraints with respect to τ_{post} , $-A_{(B^*,post)}$, has full row rank.¹⁸ We define LICQ in the direction -l analogously for the optimization that replaces max with min in (16).

¹⁸The definition of LICQ in Kaido et al. (2021) would require that this condition holds for all solutions τ_{post}^* . For our results, however, it is sufficient for the condition to hold for some solution τ_{post}^* .

For $\epsilon > 0$, we define \mathcal{P}_{ϵ} to be the set of distributions $P \in \mathcal{P}$ such that LICQ holds in the direction l and the non-binding constraints are slack by at least ϵ , i.e. $-A_{(-B^*,post)}\tau_{post}^* < d - A\beta_P - \epsilon$.

Our next result states that for $P \in \mathcal{P}_{\epsilon}$, the local power of the conditional test converges to the power envelope for tests that control size in the finite sample normal model. To state this result formally, we define $\mathcal{I}_{\alpha}(\Delta, \Sigma_n)$ to be the collection of confidence sets that control size in the finite sample normal model, i.e. confidence sets satisfying (10).

Proposition 3.3. Suppose Assumptions 2 to 4 hold. Let $\theta_P^{ub} = \sup \mathcal{S}(\beta_P, \Delta)$. Then for any $\epsilon > 0, x > 0, and \alpha < 0.5,$

$$\lim_{n \to \infty} \sup_{P \in \mathcal{P}_{\epsilon}} \left| \mathbb{E}_{P} \left[\psi_{\alpha}^{C}(\hat{\beta}_{n}, A, d, \theta_{P}^{ub} + \frac{1}{\sqrt{n}}x, \frac{1}{n}\hat{\Sigma}_{n}) \right] - \rho_{\alpha}^{*}(P, x) \right| = 0,$$

where

$$\rho_{\alpha}^{*}(P,x) = \lim_{n \to \infty} \sup_{\mathcal{C}_{\alpha,n} \in \mathcal{I}_{\alpha}(\Delta,\frac{1}{n}\Sigma_{P})} \mathbb{P}_{\hat{\beta}_{n} \sim N(\beta_{P},\frac{1}{n}\Sigma_{P})} \left(\left(\theta_{P}^{ub} + \frac{1}{\sqrt{n}}x \right) \notin \mathcal{C}_{\alpha,n} \right)$$

is the optimal local asymptotic power of a size- α test in the finite sample normal model. An analogous result holds for the lower bound under the class of distributions where LICQ holds in direction -l.

Since the LF-hybrid test rejects whenever the conditional test with size $\frac{\alpha-\kappa}{1-\kappa}$ rejects, it is immediate that the local asymptotic power of the LF-hybrid test is at least as good as the power of the optimal size- $\left(\frac{\alpha-\kappa}{1-\kappa}\right)$ test.

Corollary 3.1. Under the conditions of Proposition 3.3,

$$\liminf_{n \to \infty} \inf_{P \in \mathcal{P}_{\epsilon}} \left(\mathbb{E}_{P} \left[\psi_{\kappa, \alpha}^{C-LF}(\hat{\beta}_{n}, A, d, \theta_{P}^{ub} + \frac{1}{\sqrt{n}}x, \frac{1}{n}\hat{\Sigma}_{n}) \right] - \rho_{\frac{\alpha-\kappa}{1-\kappa}}^{*}(P, x) \right) \ge 0.$$

We emphasize that Proposition 3.3 and Corollary 3.1 are new, and exploit structure in our context not contained in the more general setting considered in ARP.

3.5.2 Discussion and Examples

As discussed in Kaido et al. (2021), LICQ and related constraint qualifications have been used frequently in the partial identification literature. Intuitively, LICQ ensures that the bounds of the identified set are differentiable with respect to the means of the moments $(\tilde{\mu}(\bar{\theta}))$, and thus avoids challenges related to estimation and inference on non-differentiable parameters (Hirano and Porter, 2012). Uniform LICQ conditions have been invoked recently by Gafarov (2019) and Cho and Russell (2019), and a related Slater constraint qualification is used in Kaido and Santos (2014). One important distinction between our results and previous results using LICQ is that we do not require LICQ for our size control results (Proposition 3.1). Thus, our tests control size even when LICQ fails (and so the bounds may be non-differentiable), but Proposition 3.3 shows that this does not come at the cost of power asymptotically when indeed LICQ holds.¹⁹

Figure 1: Diagram illustrating when LICQ (Assumption 2) will and will not hold in the case where $\overline{T} = 2$.



Note: In each panel, we assume that the rows associated with binding moments are ordered first in the matrix A for ease of notation. The blue shading denotes the identified set for (τ_1, τ_2) and the dashed red arrow points in the direction $l = (\frac{1}{2}, \frac{1}{2})'$. In panel (a), LICQ is satisfied since there is a unique τ_{post}^* (colored in red) at which two linearly independent moments are binding. In panel (b), even though τ_{post}^* is not unique, LICQ is satisfied as there is either one or two linearly independent binding moments at the values of τ_{post}^* colored in red. In panel (c), there are three binding moments at τ_{post}^* (colored in red), and so LICQ is violated.

Figure 1 provides geometric intuition for when LICQ will and will not hold in the case where $\overline{T} = 2$ and the target parameter is the average of the post-treatment effects, $\theta = \frac{1}{2}(\tau_1 + \tau_2)$. In panel (a), there is a unique τ_{post}^* (colored in red) at which two linearly independent moments bind, so LICQ is satisfied. LICQ is likewise satisfied in panel (b), where the optimal τ_{post}^* is not unique (a so-called "flat-face" problem). This is because at the indicated values τ_{post}^* (colored in red), there is either one or two linearly independent binding moments. A failure of LICQ is shown in panel (c). In this example, there are three binding moments at τ_{post}^* (colored in red), so the binding constraints cannot be linearly independent in \mathbb{R}^2 . Such a situation may arise when there are both smoothness restrictions and sign or shape restrictions that are simultaneously binding at the boundary of the identified set.

¹⁹We view this result as loosely parallel to results in the weak identification literature showing that certain procedures control size under weak identification but are efficient under strong identification (e.g. Moreira, 2003).

In the three period DiD model (where there are no nuisance parameters, since $\overline{T} = 1$), LICQ is satisfied when the bounds of the identified set are each determined by one moment. This holds everywhere for $\Delta^{SD}(M)$ when M > 0. It holds almost everywhere for $\Delta^{SDPB}(M)$ when M > 0, although it fails when both the sign restrictions and smoothness restrictions are simultaneously binding. (For LICQ to hold with non-binding moments slack by at least ϵ , i.e. $P \in \mathcal{P}_{\epsilon}, \delta_P$ must not be local to a point at which LICQ fails.) When M = 0, both the upper and lower bounds for $\Delta^{SD}(M)$ and $\Delta^{SDPB}(M)$ are binding, so LICQ fails.

More generally, the result in Proposition 3.3 is under the asymptotic regime where the sampling variation grows small relative to the length of the identified set, and thus the binding and non-binding moments are "far" apart relative to sampling variation. Importantly, it can be shown that the LICQ condition rules out settings where θ is point identified. Thus, the asymptotics considered in Proposition 3.3 may not provide a good approximation to the finite-sample performance of the conditional test in settings where θ is point-identified, or when the length of the identified set is "small" relative to sampling variation.

We are not aware of results analogous to Proposition 3.3 for any test that controls size in the finite-sample normal model. Kaido and Santos (2014) provide an efficiency result under a related Slater constraint qualification condition, but their test does not control size when the constraint qualification fails. It is worth highlighting that if LICQ holds for a particular set of moments, then it also holds if one adds moments that are slack at the optimal τ_{post}^* . Proposition 3.3 thus requires that the asymptotic power of the test is not affected by the inclusion of slack moments. The only other tests that we are aware of that control size in the finite-sample normal model and have this form of insensitivity to slack moments are the tests proposed by Cox and Shi (2022). An interesting open question is whether the tests proposed by Cox and Shi (2022) also converge to the power envelope under LICQ.²⁰

3.5.3 Extensions

Proposition 3.3 is stated for the case when Δ is a single polyhedron. An immediate corollary, however, is that when $\Delta = \bigcup_{k=1}^{K} \Delta_k$, the conditional test based on the union of confidence sets has optimal local asymptotic power when the Δ_k that determines the identified set bounds is unique and satisfies the conditions of Proposition 3.3. This implies, for example, that when $\Delta = \Delta^{RM}(\bar{M})$, the power of the conditional test converges to the power envelope when there is a unique (non-zero) pre-treatment maximum violation, i.e. when $\max_{s<0} |\delta_{s+1} - \delta_s| > 0$ and has a unique solution.²¹ Likewise, the conditional test has opti-

 $^{^{20}}$ Extending the results to the Cox and Shi tests is non-trivial given that they use a different test statistic and form critical values in a different way.

²¹For this convergence to hold uniformly, the non-binding moments must be slack by ϵ , so we would need that $\max_{s<0} |\delta_{s+1} - \delta_s|$ is at least ϵ greater than the second largest difference.

mal local asymptotic power for $\Delta^{SDRM}(\bar{M})$ when there is a unique maximum non-linearity in the pre-treatment period. Intuitively, this is because the upper bound of the identified set is determined by a single Δ_{k*} satisfying LICQ, so the conditional test for this Δ_{k*} has optimal local asymptotic power, whereas our consistency results imply that the tests for the remaining Δ_k that do not determine the identified set bound reject with probability approaching 1. See Corollary 4.1 in the working paper version of this paper for a formal derivation (Rambachan and Roth, 2021).

Proposition 3.3 shows that under LICQ the local asymptotic power of the conditional test converges to the power envelope for tests controlling size in the finite-sample normal model. In the working paper version of this paper, we showed that the power envelope from the finite-sample normal model corresponds with the power-envelope among tests that control size asymptotically and have certain invariance properties using results in Müller (2011) (Proposition E.4 in Rambachan and Roth (2021)).

4 Inference using Fixed Length Confidence Intervals

We next consider fixed length confidence intervals (FLCIs) based on affine estimators. While the conditional and hybrid confidence sets offer attractive asymptotic power guarantees under asymptotics in which sampling variation grows small relative to the length of the identified set, FLCIs offer finite-sample power guarantees (in the normal model) for certain classes Δ of interest. In certain special cases, FLCIs may thus outperform the ARP tests when sampling variation is large relative to the length of the identified set. For brevity of exposition, we focus on the properties of FLCIs in the case where the finite-sample normal approximation (9) holds exactly with Σ_n known; Armstrong and Kolesár (2020b) provide uniform asymptotic results for FLCIs under conditions similar to those in Section 3.3.

4.1 Constructing FLCIs

Following Donoho (1994) and Armstrong and Kolesár (2018, 2020a), we consider fixed length confidence intervals based on an affine estimator for θ , denoted by $\mathcal{C}_{\alpha,n}(a, v, \chi) := (a + v'\hat{\beta}_n) \pm \chi$, where a and χ are scalars and $v \in \mathbb{R}^{T+\bar{T}}$. We minimize the half-length of the confidence interval, χ , subject to the constraint that $\mathcal{C}_{\alpha,n}(a, v, \chi)$ satisfies the coverage requirement (10) in the finite-sample normal model.

Observe that if $\hat{\beta}_n \sim \mathcal{N}(\beta, \Sigma_n)$, then $a + v'\hat{\beta}_n \sim \mathcal{N}(a + v'\beta, v'\Sigma_n v)$, and hence $|a + v'\hat{\beta}_n - \theta| \sim |\mathcal{N}(b, v'\Sigma_n v)|$, where $b = a + v'\beta - \theta$ is the affine estimator's bias for θ . Observe further that $\theta \in \mathcal{C}_{\alpha,n}(a, v, \chi)$ if and only if $|a + v'\hat{\beta}_n - \theta| \leq \chi$. For fixed values a and v, the

smallest value of χ that satisfies (10) is therefore the $1 - \alpha$ quantile of the $|\mathcal{N}(\bar{b}, v'\Sigma_n v)|$ distribution, where \bar{b} is the affine estimator's worst-case bias

$$\bar{b}(a,v) := \sup_{\delta \in \Delta, \tau_{post} \in \mathbb{R}^{\bar{T}}} |a + v' \left(\delta + L_{post} \tau_{post}\right) - l' \tau_{post}|.$$
(17)

Let $cv_{\alpha}(t)$ denote the $1 - \alpha$ quantile of the folded normal distribution $|\mathcal{N}(t, 1)|^{22}$ For fixed a and v, the smallest value of χ satisfying the coverage requirement (10) is thus

$$\chi_n(a,v;\alpha) = \sigma_{v,n} \cdot cv_\alpha(\bar{b}(a,v)/\sigma_{v,n}), \tag{18}$$

where $\sigma_{v,n} := \sqrt{v' \Sigma_n v}$. The optimal (i.e., minimum-length) FLCI is constructed by choosing the values of a and v to minimize (18). When Δ is convex, this minimization can be solved as a nested optimization problem, where both the inner and outer minimizations are convex (Low, 1995; Armstrong and Kolesár, 2018, 2020a). We denote the $1 - \alpha$ level, optimal FLCI by $C_{\alpha,n}^{FLCI}(\hat{\beta}_n, \Sigma_n) := (a_n + v'_n \hat{\beta}_n) \pm \chi_n$, where $\chi_n := \inf_{a,v} \chi_n(a, v; \alpha)$ and a_n , v_n are the optimal values in the minimization.

Example: $\Delta^{SD}(M)$. Suppose $\theta = \tau_1$. For $\Delta^{SD}(M)$, the affine estimator used by the optimal FLCI takes the form $a + v'\hat{\beta}_n = \hat{\beta}_{n,1} - \sum_{s=-T+1}^0 w_s \left(\hat{\beta}_{n,s} - \hat{\beta}_{n,s-1}\right)$, where the weights w_s sum to one (but may be negative). This estimator adjusts the event-study coefficient for t = 1 by an estimate of the differential trend between t = 0 and t = 1 formed by taking a weighted average of the differential trends in periods prior to treatment. The worst-case bias will be smaller if more weight is placed on pre-treatment periods closer to the treatment date, but it may reduce variance to place more weight on earlier pre-periods. The weights w_s are optimally chosen to balance this tradeoff.

4.2 Finite-sample near optimality

In particular cases of interest, such as when $\Delta = \Delta^{SD}(M)$, the optimal FLCIs introduced above have near-optimal expected length in the finite-sample normal model. The following result, which is an immediate consequence of results in Armstrong and Kolesár (2018, 2020a), bounds the ratio of the expected length of the shortest possible confidence interval that controls size relative to the length of the optimal FLCI.

Assumption 8. Assume i) Δ is convex and centrosymmetric (i.e. $\tilde{\delta} \in \Delta$ implies $-\tilde{\delta} \in \Delta$), and ii) $\delta \in \Delta$ is such that $(\tilde{\delta} - \delta) \in \Delta$ for all $\tilde{\delta} \in \Delta$.

²²If $t = \infty$, we define $cv_{\alpha} = \infty$.

Proposition 4.1. Suppose δ and Δ satisfy Assumption 8. Let $\mathcal{I}_{\alpha}(\Delta, \Sigma_n)$ denote the class of confidence sets that satisfy the coverage criterion (10) at the $1 - \alpha$ level. Then, for any τ with $\tau_{pre} = 0$ and Σ_n positive definite,

$$\frac{\inf_{\mathcal{C}_{\alpha,n}\in\mathcal{I}_{\alpha}(\Delta,\Sigma_{n})}\mathbb{E}_{\hat{\beta}_{n}\sim\mathcal{N}(\delta+\tau,\Sigma_{n})}\left[\lambda(\mathcal{C}_{\alpha,n})\right]}{2\chi_{n}} \geqslant \frac{z_{1-\alpha}(1-\alpha)-\tilde{z}_{\alpha}\Phi(\tilde{z}_{\alpha})+\phi(z_{1-\alpha})-\phi(\tilde{z}_{\alpha})}{z_{1-\alpha/2}},$$

where $\lambda(\cdot)$ denotes the length (Lebesgue measure) of a set and $\tilde{z}_{\alpha} = z_{1-\alpha} - z_{1-\alpha/2}$.

Part i) of Assumption 8 is satisfied for $\Delta^{SD}(M)$ but not for our other ongoing examples. For example, $\Delta^{SDPB}(M)$ is convex but not centrosymmetric, and $\Delta^{RM}(\bar{M})$ is neither convex nor centrosymmetric. Part ii) of Assumption 8 is satisfied whenever parallel trends holds in both the pre-treatment and post-treatment periods ($\delta = 0$) and whenever δ is a linear trend for the case of $\Delta^{SD}(M)$.

FLCIs thus offer attractive guarantees for the case of $\Delta^{SD}(M)$. When $\alpha = 0.05$, the lower bound in Proposition 4.1 evaluates to 0.72, meaning that the expected length of the shortest possible confidence set that satisfies the coverage requirement (10) is at most 28% shorter than the length of the optimal FLCI when the conditions of the proposition hold.

4.3 (In)Consistency of FLCIs

As discussed above, these finite-sample guarantees do not apply for several types of restrictions Δ of importance, including those that construct bounds using the maximum pretreatment violation or incorporate sign and shape restrictions. We now show that the FLCIs can perform poorly under such restrictions. We first provide two illustrative examples, and then state a formal inconsistency result.

Example: $\Delta^{SDPB}(M)$ and $\Delta^{SDI}(M)$. Suppose $\theta = \tau_1$. It can be shown that the worstcase bias of an affine estimator over $\Delta^{SDPB}(M)$ or $\Delta^{SDI}(M)$ is the same as the worst-case bias for that estimator over $\Delta^{SD}(M)$.²³ Since the construction of the optimal FLCI depends only on the worst-case bias and variance of the affine estimator, it follows that the optimal FLCI constructed using $\Delta^{SDPB}(M)$ or $\Delta^{SDI}(M)$ is the same as the one constructed using $\Delta^{SD}(M)$. Therefore, the optimal FLCI does not adapt to additional sign or monotonicity restrictions.

²³Suppose the vector $\overline{\delta}$ maximizes the bias for an affine estimator (a, v) over $\Delta^{SD}(M)$. The vector that adds a constant slope to $\overline{\delta}$, say $\widetilde{\delta}_c = \overline{\delta} + c \cdot (-\overline{T}, ..., \overline{T})'$, also lies in $\Delta^{SD}(M)$, and for c sufficiently large, $\widetilde{\delta}_c$ will lie in $\Delta^{SDPB}(M)$. Moreover, the worse-case bias will be the same for δ and $\widetilde{\delta}_c$, since if (a, v) has finite worst-case bias it must subtract out a weighted average of the pre-treatment slopes.

Example: $\Delta^{RM}(\bar{M})$. Suppose $\theta = \tau_1$. If $\Delta = \Delta^{RM}(\bar{M})$ and $\bar{M} > 0$, then all affine estimators for τ_1 have infinite worst-case bias, since $\delta \in \Delta^{RM}(\bar{M})$ can have $|\delta_1|$ arbitrarily large if $|\delta_{-1}|$ is also sufficiently large. Thus, the only valid FLCI is the entire real line.

We next provide a formal result on the (in)consistency of the FLCIs. Specifically, we will show that even as the sampling variation Σ_n converges to 0, the optimal FLCI will include fixed points outside of the identified set with positive probability unless certain special conditions are met.²⁴ Recall from Lemma 2.1 that the identified set $S(\beta, \Delta)$ is an interval when Δ is convex, with length equal to $\theta^{ub}(\beta, \Delta) - \theta^{lb}(\beta, \Delta) = b^{max}(\beta_{pre}, \Delta) - b^{min}(\beta_{pre}, \Delta)$. Since the length of the identified set only depends on β_{pre} and Δ , denote it by $LID(\beta_{pre}, \Delta)$. Our next result shows that $C_{\alpha,n}^{FLCI}(\hat{\beta}_n, \Sigma_n)$ is consistent if and only if $LID(\beta_{pre}, \Delta)$ is its maximum possible value, provided that the identified set is not the entire real line (in which case any procedure is trivially consistent).

Assumption 9 (Identified set maximal length and finite). Suppose $\delta \in \Delta$ is such that $LID(\delta_{pre}, \Delta) = \sup_{\tilde{\delta}_{pre} \in \Delta_{pre}} LID(\tilde{\delta}_{pre}, \Delta) < \infty$, where $\Delta_{pre} = \{\delta_{pre} \in \mathbb{R}^T : \exists \delta_{post} \ s.t. \ (\delta'_{pre}, \delta'_{post})' \in \Delta\}$ is the set of possible values for δ_{pre} .

Proposition 4.2. Suppose Δ is convex and $\alpha < 0.5$. Fix $\delta \in \Delta$ and τ with $\tau_{pre} = 0$, and suppose $\mathcal{S}(\delta + \tau, \Delta) \neq \mathbb{R}$. Then (δ, Δ) satisfy Assumption 9 if and only if $\mathcal{C}_{\alpha,n}^{FLCI}(\hat{\beta}_n, \Sigma_n)$ is consistent, meaning that for $\Sigma_n = \Sigma^*/n$,

$$\lim_{n \to \infty} \mathbb{P}_{\hat{\beta}_n \sim \mathcal{N}(\delta + \tau, \Sigma_n)} \left(\theta^{out} \in \mathcal{C}_{\alpha, n}^{FLCI}(\hat{\beta}_n, \Sigma_n) \right) = 0 \text{ for all } \theta^{out} \notin \mathcal{S}(\delta + \tau, \Delta).$$

Thus, if Assumption 9 fails, then $C_{\alpha,n}^{FLCI}(\hat{\beta}_n, \Sigma_n)$ is inconsistent in the strong sense that it includes fixed points outside of the identified set with non-vanishing probability. It follows that there will be some $\delta \in \Delta$ such that the FLCI is inconsistent under δ unless the identified set is always the same length. Proposition 4.2 is new, and may be relevant for other settings in which FLCIs are used.

The intuition for the possible inconsistency of FLCIs is as follows: to ensure that an FLCI satisfies the coverage requirement (10), its length must be at least $\sup_{\tilde{\delta}_{pre} \in \Delta_{pre}} LID(\tilde{\delta}_{pre}, \Delta)$. However, this implies that if in fact $LID(\delta_{pre}, \Delta) < \sup_{\tilde{\delta}_{pre} \in \Delta_{pre}} LID(\tilde{\delta}_{pre}, \Delta)$, then the FLCI is strictly longer than the length of the identified set, regardless of the value of Σ_n , and thus some points outside of the identified set must be covered with non-vanishing probability. This reflects the fact that FLCIs are by construction *fixed length*, and thus their length does not adapt to information in the data about the length of the identified set. By contrast, the

²⁴For ease of exposition, we present a result using "small- Σ " asymptotics in the normal model, as in e.g., Kadane (1971) and Moreira and Ridder (2019).

length of the conditional/hybrid confidence sets can depend on $\hat{\beta}_{pre}$ and thus "adapts" to the length of the identified set.

In the three-period difference-in-differences example, Assumption 9 holds everywhere for $\Delta^{SD}(M)$ (since the identified set is always the same length, 2M), for values of δ where the sign restrictions do not bind for $\Delta^{SDPB}(M)$, and nowhere for the polyhedra that form $\Delta^{RM}(\bar{M})$. The restrictiveness of Assumption 9 thus depends greatly on Δ .

The results in this section establish that when certain conditions on Δ are satisfied, optimal FLCIs are consistent and have desirable finite-sample guarantees in terms of expected length. FLCIs are thus attractive for our baseline smoothness class $\Delta^{SD}(M)$, since they are guaranteed to be consistent and offer attractive finite-sample guarantees. Our inconsistency result shows, however, that FLCIs may perform poorly for other choices of Δ that may be of interest in empirical applications, such as those that construct bounds using a pre-treatment maximum or incorporate sign and monotonicity restrictions.

5 Simulation study

In this section, we conduct a simulation study to investigate the performance of the discussed confidence sets across a range of relevant data-generating processes. We find good size control for all of the procedures, and therefore focus in the main text on a comparison of power to provide concrete recommendations on the best approach in practice. In the supplementary material, we present results on size control and other additional simulation results.

5.1 Simulation Design

Our simulations are calibrated using the estimated covariance matrix from the 12 recentlypublished papers surveyed in Roth (Forthcoming). For any given paper in the survey, we denote by $\hat{\Sigma}$ the estimated variance-covariance matrix from the event-study in the paper, calculated using the clustering scheme specified by the authors. For a chosen mean vector β , we simulate event-study coefficients $\hat{\beta}_s$ from a normal model, $\hat{\beta}_s \sim \mathcal{N}\left(\beta, \hat{\Sigma}\right)$.²⁵ In simulation s, we construct nominal 95% confidence sets for the parameter of interest θ using the pair $(\hat{\beta}_s, \hat{\Sigma})$ for each proposed procedure. The parameter of interest is the causal effect in the first post-treatment period ($\theta = \tau_1$); in the supplementary material, we present simulation results in which the parameter of interest is the average causal effect in the post-treatment periods ($\theta = \bar{\tau}_{post}$), with qualitatively similar results.

 $^{^{25}}$ We focus on the normal simulations in the main text since it allows for a tractable computation of the optimal excess length of procedures that control size. In the supplementary material, we show that our procedures perform similarly in simulations based on the empirical distribution in the original paper.

For a given choice of Δ , we compute the identified set $S(\beta, \Delta)$ and calculate the expected excess length for each of the proposed confidence sets. We benchmark the expected excess length of our proposed confidence sets relative to an efficiency bound for confidence sets that satisfy the uniform coverage requirement.²⁶ We report the efficiency ratio of each procedure, which is defined as the ratio of the optimal benchmark relative to the average excess length for the procedure. All results are calculated over 1000 simulations per paper.

	Parallel Trends		Pulse Pre-Trend	
	$\Delta^{SD}(M)$	$\Delta^{SDPB}(M)$	$\Delta^{SDRM}(\bar{M})$	$\Delta^{RM}(\bar{M})$
Conditional and Hybrid				
Consistent	\checkmark	\checkmark	\checkmark	\checkmark
Asymptotically (near-)optimal	\checkmark	\checkmark	\checkmark	×
FLCI				
Consistent	\checkmark	×	×	×
Finite-sample near-optimal	\checkmark	×	×	×

Table 1: Summary of expected properties for each simulation design

We consider four choices of Δ to highlight the performance of our proposed confidence sets across a range of conditions: $\Delta^{SD}(M)$, $\Delta^{SDPB}(M)$, $\Delta^{RM}(\bar{M})$, and $\Delta^{SDRM}(\bar{M})$. We consider simulations under the assumption of zero treatment effects, so that $\tau = 0$ and thus $\beta = \delta$. We consider two forms for δ . First, we consider the baseline case of parallel trends $(\delta = 0)$. Second, we consider a "pulse" pre-trend in which δ_{-1} is non-zero and the remaining elements of δ are zero. Such a pre-trend might arise in practice if there are confounding policy changes or other events close to the time of treatment. These different choices of δ allow us to highlight the relative strengths of the proposed inference procedures. For example, FLCIs have near-optimal expected length when $\delta = 0$ and $\Delta = \Delta^{SD}(M)$, whereas the conditional test has optimal local asymptotic power under the pulse design when $\Delta = \Delta^{SDPB}(M)$. Table 1 summarizes which of our theoretical results hold for each of the simulation designs when M and \bar{M} are non-zero.

In practice, we find that for $\Delta^{SD}(M)$ and $\Delta^{SDPB}(M)$, the results depend on M but are qualitatively similar across values of δ . By contrast, for $\Delta^{SDRM}(\bar{M})$ and $\Delta^{RM}(\bar{M})$, the choice of δ is more important than the choice of \bar{M} . Therefore, to highlight the most important

²⁶For choices of Δ that are convex (e.g., $\Delta^{SD}(M)$ and $\Delta^{SDPB}(M)$), we benchmark the expected excess length of our proposed confidence sets against a sharp optimal bound over confidence sets that satisfy the finite-sample coverage requirement (10). This optimal bound is provided in the supplementary materials, and follows as a corollary from results in Armstrong and Kolesár (2018) on the optimal expected length of a confidence set satisfying the uniform coverage requirement (10). For choices of Δ that can be written as the union of convex sets (e.g., $\Delta^{RM}(\bar{M})$ and $\Delta^{SDRM}(\bar{M})$), we compare the expected excess length of our proposed confidence sets against the maximal optimal bound over each set in the union, which is a potentially non-sharp bound for any confidence set with correct coverage.

dimensions for each of the simulation designs, in the main text of the paper we report results for $\Delta^{SD}(M)$ and $\Delta^{SDPB}(M)$ under different values of M and $\delta = 0$ (parallel trends), whereas for $\Delta^{RM}(\bar{M})$ and $\Delta^{SDRM}(\bar{M})$ we vary the magnitude of the pre-treatment pulse δ_{-1} , holding $\bar{M} = 1$ constant. In the supplementary materials, we report results for additional choices of these parameters.

We report results for three methods for constructing confidence sets: FLCIs, conditional confidence sets, and conditional-least favorable hybrid confidence sets.²⁷ For $\Delta^{RM}(\bar{M})$ and $\Delta^{SDRM}(\bar{M})$, we omit results for the FLCI since the FLCIs have infinite length.

5.2 Simulation Results

To compare results easily across the 12 papers in the simulation study, we normalize the units of δ_{-1} and M by the standard deviation of $\hat{\beta}_1$ (denoted σ_1). Large normalized values of M or δ_{-1} correspond with the case where the identified set is large relative to sampling variation, mimicking our asymptotic power results in which sampling variation grows small relative to the identified set. In the graphs below, we report the median value of excess length efficiency across the papers in the survey. The normalization described above implies that the units of the *x*-axis correspond with the worst-case bias of the naive estimator $\hat{\beta}_1$ divided by its standard error.²⁸

Results for $\Delta^{SD}(M)$: The left panel of Figure 2 plots the efficiency ratio for each procedure as a function of M/σ_1 when $\Delta = \Delta^{SD}(M)$. All procedures perform well as M/σ_1 grows large with efficiency ratios approaching 1, illustrating our asymptotic (near-)optimality results for this design. However, the FLCIs perform best for smaller values of M/σ_1 , including the point-identified case where M = 0, illustrating the finite-sample near-optimality results for the FLCIs when Assumption 8 holds. Although the conditional and hybrid confidence sets have efficiency approaching the optimal bound for M/σ_1 large, their efficiency is only about 50% when $M/\sigma_1 = 0$, in which case θ is point identified and thus LICQ does not hold. The conditional and hybrid confidence sets perform similarly.

Results for $\Delta^{SDPB}(M)$: The right panel of Figure 2 plots the efficiency ratio for each procedure as a function of M/σ_1 when $\Delta = \Delta^{SDPB}(M)$. The efficiency ratios for the conditional and hybrid confidence sets are again (near-)optimal as M/σ_1 grows large, highlighting

²⁷For the conditional-least favorable hybrid confidence sets, we use a first-stage least-favorable test of size $\kappa = \alpha/10$, following ARP and Romano, Shaikh and Wolf (2014).

²⁸For $\hat{\beta}_1$ normally distributed, the worst-case coverage of a conventional 95% confidence interval as a function of the normalized worst-case bias b is $\Phi(1.96 + b) - \Phi(-1.96 + b)$, which is 0.95 for b = 0, 0.83 for b = 1, 0.48 for b = 2, etc.

Figure 2: Simulation results for $\Delta^{SD}(M)$ and $\Delta^{SDPB}(M)$: Median efficiency ratios for proposed procedures.



Note: Median efficiency ratios for our proposed confidence sets over $\Delta^{SD}(M)$ and $\Delta^{SDPB}(M)$ under the assumption of parallel trends and zero treatment effects (i.e., $\beta = 0$). The efficiency ratio for a procedure is defined as the efficiency bound divided by the procedure's expected excess length. The results for the FLCI are plotted in purple, conditional-LF ("C-LF Hybrid") hybrid in blue, and conditional confidence set in green. Results are averaged over 1000 simulations for each of the 12 papers surveyed, and the median across papers is reported here.

Figure 3: Simulation results for $\Delta^{SDRM}(\bar{M})$ and $\Delta^{RM}(\bar{M})$: Median efficiency ratios for proposed procedures.



Note: Median efficiency ratios for our proposed confidence sets over $\Delta^{SDRM}(\bar{M})$ and $\Delta^{RM}(\bar{M})$ with $\bar{M} = 1$ under the assumption of zero treatment effects and a "pulse" pre-trend (i.e., $\beta_{-1} = \delta_{-1}$ and $\beta_t = 0$ for all $t \neq -1$). The efficiency ratio for a procedure is defined as the efficiency bound divided by the procedure's expected excess length. The results for the conditional-least favorable ("C-LF") hybrid are plotted in blue, and conditional confidence set in green. Results are averaged over 1000 simulations for each of the 12 papers surveyed, and the median across papers is reported here.

our asymptotic (near-)optimality results for these procedures in this simulation design. By contrast, the efficiency ratios for the FLCIs steadily decrease as M/σ_1 increases, reflecting

that the FLCIs are not consistent in this simulation design when M > 0. The conditional-LF hybrid confidence sets slightly improve efficiency relative to the conditional when M/σ_1 is small and retain near-optimal performance as M/σ_1 grows large.

Results for $\Delta^{SDRM}(\bar{M})$: The left panel of Figure 3 plots the efficiency ratios for the conditional and conditional-least favorable hybrid confidence sets as a function of δ_{-1}/σ_1 when $\Delta = \Delta^{SDRM}(\bar{M})$. We omit results for the optimal FLCI since the optimal FLCI has infinite length for this design. Both procedures perform well as δ_{-1}/σ_1 grows large with efficiency ratios approaching 1, illustrating our asymptotic (near-) optimality result for this design. Both procedures also have similar power curves, with slightly higher power for the conditional.

Results for $\Delta^{RM}(\bar{M})$: The right panel of Figure 3 plots the efficiency ratio for the conditional and conditional-least favorable hybrid confidence sets as a function of δ_{-1}/σ_1 when $\Delta = \Delta^{RM}(\bar{M})$. We again omit results for the optimal FLCI since the optimal FLCI has infinite length for this design. The conditions for our asymptotic (near-) optimality result for unions of convex sets do not hold in this simulation design (as the maximum pre-period violation is not unique). Nonetheless, we find that the conditional-least favorable hybrid confidence set and the conditional confidence set perform quite well for large values of δ_{-1}/σ_1 , with efficiency ratios approaching about 83%. This is encouraging as it shows that these procedures may perform well even in cases where LICQ fails. Once again, we also find that the conditional and conditional-least favorable hybrid have similar power.

Takeaways from Simulations: Two clear patterns emerge from our simulations. First, the conditional and hybrid confidence sets perform well across a wide range of specifications, with particularly good power when the length of the identified set is large relative to sampling variation. Second, the FLCIs have the best performance for $\Delta^{SD}(M)$, particularly when M is small, which aligns with the finite-sample near-optimality results in Section 4. However, FLCIs can perform quite poorly for other classes of Δ .

Overall, we therefore recommend to use the conditional-LF hybrid confidence sets for generic forms of Δ , and optimal FLCIs for the special case of $\Delta^{SD}(M)$ (or other special cases where the consistency/finite-sample near-optimality of FLCIs is guaranteed). Although the conditional and hybrid approaches perform similarly in our simulations, we recommend the hybrid approach in general based on the guidance provided in ARP. We implement these recommendations in our applications in the next section.

6 Practical Guidance and Empirical Illustrations

6.1 Practical Guidance

We recommend that researchers use our methods to construct robust confidence intervals under restrictions on the possible violations of parallel trends Δ that are motivated by domain knowledge in their empirical setting. We also suggest that researchers report sensitivity analyses to illustrate the sensitivity of their causal conclusions to alternative assumptions on the possible violations of parallel trends.

Choice of Δ . The choice of Δ should be motivated by economic knowledge about the types of possible confounding factors that would produce non-parallel trends. We now provide some guidance on how the choice of Δ can be motivated by domain knowledge, highlighting some cases where our leading examples, $\Delta^{RM}(\bar{M})$ and $\Delta^{SD}(M)$, would be sensible choices.

In some empirical settings, researchers may be concerned about differential economic shocks to the treated and control groups that generate violations of parallel trends. If the researcher believes that the magnitude of these differential shocks in the post-treatment period is not too different from the magnitude in the pre-treatment period, then it may be reasonable to assume $\delta \in \Delta^{RM}(\bar{M})$, which explicitly bounds the relative magnitudes of violations of parallel trends in the post-treatment based on observed violations in the pretreatment period. In other settings, researchers may be worried about violations of parallel trends that arise due to differences in smoothly evolving secular trends that differentially affect treated and comparison groups. In this case, it may be reasonable to assume $\delta \in \Delta^{SD}(M)$, which explicitly bounds the extent to which the slope of the difference in trends can vary across consecutive periods. Economic knowledge may imply additional restrictions as well. For example, if the researcher knows of a confounding policy change that would have a positive effect on the outcome, then it is reasonable to further assume that post-treatment difference in trends must be positive (i.e., $\delta_t \ge 0$ for t > 0).

In our empirical applications below, we illustrate how domain knowledge about the types of possible violations of parallel trends can inform the choice of Δ . We encourage applied researchers to use such domain knowledge to inform the restrictions they impose on the possible choices of parallel trends in their context.

Choice of inference procedure. Based on our theoretical results and Monte Carlo simulations, we recommend the ARP hybrid confidence sets for generic, polyhedral forms of Δ . For the special case of $\Delta^{SD}(M)$ — or other choices of Δ for which the consistency and finite-sample-near-optimality of FLCIs are guaranteed — we recommend FLCIs. Our
recommended choice of inference procedure is implemented in the R package, HonestDiD, that accompanies the paper.²⁹ Furthermore, these confidence sets are quick to compute. Each sensitivity analysis plot in the empirical applications below took less than 9 minutes to compute on a 2012 Macbook Pro.

Sensitivity analyses. Once the researcher has chosen a baseline class of restrictions on the possible violations of trends (e.g. relative magnitudes bounds $\Delta^{RM}(\bar{M})$ or smoothness bounds $\Delta^{SD}(M)$), we recommend conducting sensitivity analysis over the associated parameter ($\bar{M} \ge 0$ or $M \ge 0$, respectively) that governs how different the post-treatment violations of parallel trends can be from the pre-trends. It is natural to report both the sensitivity of the researcher's causal conclusion to the choice of this parameter and the "breakdown" parameter value at which particular hypotheses of interest can no longer be rejected; similar "breakdown" concepts appear in the partial identification settings of Horowitz and Manski (1995); Kline and Santos (2013); Manski and Pepper (2018); Masten and Poirier (2020).³⁰ We illustrate how one can interpret the magnitudes of the breakdown points in our two empirical illustrations below.

6.2 Estimating the incidence of a value-added tax cut

Benzarti and Carloni (2019, henceforth, BC) study the incidence of a decrease in the valueadded tax (VAT) on restaurants in France. France reduced its VAT on sit-down restaurants from 19.6 to 5.5 percent in July of 2009. BC analyze the impact of this change using a dynamic difference-in-differences design that compares restaurants to a control group of other market services firms that were not affected by the VAT change, estimating

$$Y_{it} = \sum_{s \neq 2008} \beta_s \times 1[t=s] \times D_i + \phi_i + \lambda_t + \epsilon_{it}, \tag{19}$$

where Y_{it} is the log of (before-tax) profits for firm *i* in year *t*; D_i is an indicator for whether firm *i* is a restaurant; ϕ_i and λ_t are firm and year fixed effects; and standard errors are

 $^{^{29}{\}rm The}$ latest version of the R package can be downloaded by visiting http://github.com/asheshrambachan/HonestDiD.

³⁰Our main focus in this paper is on constructing robust confidence sets given a particular restriction $\Delta(M)$, rather than inference on the identification breakdown point or breakdown frontier as in e.g. Masten and Poirier (2020). Note, however, that if we define $M^* = \min M$ s.t. $0 \in \mathcal{S}(\beta, \Delta(M))$ to be the identification breakdown point for a null effect, and $\hat{M}^* = \min M$ s.t. $0 \in \mathcal{C}(\hat{\beta}_n, \hat{\Sigma}_n; \Delta(M))$ to be the sample breakdown point, then $P(\hat{M}^* \ge M^*) \ge P(0 \in \mathcal{C}(\hat{\beta}_n, \hat{\Sigma}_n; \Delta(M^*)))$. It follows that $(-\infty, \hat{M}^*]$ is a valid $(1 - \alpha)$ -level confidence interval for M^* provided that our conditions for size control are satisfied for $\Delta(M^*)$. We suspect that our results could be extended to allow for uniform coverage of the breakdown frontier under additional regularity conditions, but leave this to future work.

clustered at the regional level. BC's main finding is that the VAT reduction had a large, positive effect on restaurant profits. Figure 4 shows the estimated event-study coefficients $\{\hat{\beta}_s\}$ from specification (19). We can formally reject the hypothesis that $\beta_{pre} = 0$ (p < 0.01), as there appears to have been a difference in trends between 2006 and 2007. Nevertheless, the changes in profits after the policy change appear to be larger in magnitude than any of the pre-trends.

Figure 4: Event-study coefficients $\{\beta_s\}$ for log profits, estimated using the event-study specification in (19).



A key concern in this empirical setting is that there may be unobserved, industry-specific or macroeconomic shocks that would have affected restaurants differently from other marketservices firms even in the absence of a change in VAT. It seems reasonable to impose that the industry-specific shocks to restaurants in the post-treatment period are not too much larger than those in the pre-treatment period — whereas imposing that industry-specific shocks follow a smooth trend seems unreasonable — and so we base our analysis on bounds on relative magnitudes $\Delta^{RM}(\bar{M})$.

The left panel of Figure 5 shows robust confidence sets for the treatment effect in 2009 for $\Delta^{RM}(\bar{M})$ using different values of \bar{M} . The figure shows that if we impose $\bar{M} = 1$, meaning that we restrict the post-treatment violations of parallel trends to be no larger than the maximal pre-treatment violation of parallel trends, then we obtain a robust confidence set of [0.07, 0.31] for the causal effect on restaurant profits in 2009. This is wider than the original OLS confidence interval which is only valid if parallel trends holds exactly, but nevertheless rules out a null effect on restaurant profits in 2009. Looking further to the right, we see that the "breakdown value" for a null effect is around $\bar{M} = 2$. Thus, our conclusion of a significant effect on restaurant profits depends on whether we are willing to restrict that

the post-treatment violations of parallel trends can be no more than twice as large as the maximal pre-treatment violation. Given that the first year after the treatment coincided with a large recession in France (2009), it may be plausible that the differential factors affecting restaurants were larger in that year than in the pre-treatment period. Our approach helps formalize how much larger they would need to be to reject the conclusion of a null effect (or other hypotheses).





The right panel of Figure 5 shows analogous results when the estimand is the average causal effect on restaurant profits across all four post-treatment periods $(\bar{\tau})$. When $\bar{M} = 1$, our robust confidence set now includes zero, and is about twice as large as for the first-period effect (notice the difference in scale across the panels). The intuition for why the confidence sets are larger when looking at $\bar{\tau}$ than τ_{2009} is that $\Delta^{RM}(\bar{M})$ bounds the violation of parallel trends across *consecutive periods* by \bar{M} times the max in the pre-treatment period. Thus, the identified set will be larger for later periods, since the treatment and control groups have more time to diverge (e.g., the identified set for the second period will be twice as larger as for the first period). If we are willing to bound the magnitude of economic shocks by the max in the pre-treatment period, we will thus typically obtain wider confidence sets for parameters involving later periods.

6.3 The effect of duty-to-bargain laws on long-run student outcomes

Lovenheim and Willen (2019, henceforth LW) study the impact of state-level public sector

duty-to-bargain (DTB) laws, which mandated that school districts bargain in good faith with teachers' unions. LW examine the impacts of these laws on the adult labor market outcomes of people who were students around the time that these laws were passed, comparing individuals across different states and different birth cohorts to exploit the differential timing of the passage of DTB laws across states. The authors estimate the following regression specification separately for men and women, using data from the American Community Survey (ACS),

$$Y_{sct} = \sum_{r=-11}^{21} D_{scr}\beta_r + X'_{sct}\gamma + \lambda_{ct} + \phi_s + \epsilon_{sct}.$$
 (20)

 Y_{sct} is an average outcome for the cohort of students born in state s in cohort c in ACS calendar year t. D_{scr} is an indicator for whether state s passed a DTB law r years before cohort c turned age 18.³¹ The event-study coefficients $\{\hat{\beta}_r\}$ estimate the dynamic treatment effects (or placebo effects) r years after DTB passage.³² The remaining terms include time-varying controls, birth-cohort-by-ACS-year fixed effects, and state fixed effects. We normalize the event-study coefficient β_{-2} to 0.³³ We focus on the results where the outcome is employment.

Figure 6 plots the estimated event-study coefficients $\{\hat{\beta}_r\}$ from specification (20). In the event-study for men (left panel), the pre-period coefficients are relatively close to zero, whereas the longer-run post-period coefficients are negative. By contrast, the results for women (right panel) suggest a downward-sloping pre-existing trend.

LW write that, the "primary concern in our identification strategy is the existence of secular trends that differ systematically with treatment" (p. 318), such as confounding changes in labor supply or educational attainment. Given that the concern is long-run trends that are likely to evolve smoothly over time, smoothness restrictions of the form Δ^{SD} seem natural in this context. Indeed, in some of their robustness checks, LW estimate models with group-specific linear trends, which roughly corresponds with the case $\Delta^{SD}(0)$.³⁴ It thus

 $[\]overline{{}^{31}D_{sc,-11}}$ is set to 1 if state *s* passed a law 11 years or more after cohort *c* turned 18. Likewise, $D_{sc,21}$ is set to 1 if state *s* passed a law 21 or more years before cohort *c* turned 18.

³²Treatment timing in LW is staggered, and therefore the results in Sun and Abraham (2020) imply that β_r can be interpreted as a sensible weighted average of causal effects under parallel trends only if treatment effects are homogeneous across adoption cohorts. For simplicity, we focus on the robustness of the results to violations of parallel trends using the original specification in LW, which is valid under the assumption of homogeneous treatment effects. As discussed in Section 2.1, our sensitivity analysis can also be applied to estimators that are robust to treatment effect heterogeneity.

³³LW normalize event time -1 to 0, but discuss how cohorts at event time -1 may have been partially treated since LW impute the year that a student starts school with error. Since our robust confidence sets assume that there is no causal effect in the pre-period ($\tau_{pre} = 0$), we instead treat event-time -2 as the reference period in our analysis.

 $^{^{34}}$ The two are not exactly equivalent, however, because LW include parametric trends into what they call a "parametric event-study" model (see their specification (2)), which imposes that treatment effects are linear in time since treatment, rather than the flexible dynamic event-study specification (20).

seems natural to consider relaxations of the form $\Delta^{SD}(M)$, which allows for deviations from non-linearity of no more than M between consecutive periods.



Figure 6: Event-study coefficient $\{\beta_r\}$ for employment, estimated using the event-study specification in (20).

Figure 7: Sensitivity analysis for $\theta = \tau_{15}$ using $\Delta = \Delta^{SD}(M)$



Figure 7 reports results for the treatment effect on employment for the cohort 15 years after the passage of a DTB law (as in Table 2 of LW), constructing robust confidence sets about how non-linear the difference in trends can be. In blue, we plot the original OLS confidence intervals for $\hat{\beta}_{15}$ from specification (20). In red, we plot FLCIs when $\Delta = \Delta^{SD}(M)$ for different values of M; recall that M = 0 corresponds with allowing only for linear violations of parallel trends, and larger values of M allow for larger deviations from linearity. In the analysis for men (left panel), the FLCIs are similar to those from OLS when allowing

for violations of parallel trends that are approximately linear $(M \approx 0)$, but become wider as we allow for more non-linearity; the breakdown value for a significant effect is $M \approx 0.01$. For women (right panel), the original OLS estimates are negative and the confidence interval rules out 0. When we allow for linear violations of parallel trends (M = 0), however, the picture changes substantially owing to the pre-existing downward trend that is visible in Figure 6. Indeed, for M < 0.01 the robust confidence set contains only positive values. Intuitively, this is because the point estimate for t = 15 lies above a linear extrapolation of the negative pre-trend. Thus, if we were to impose the same smoothness restrictions for men as for women, we would either have to reconcile significant effects of opposite signs by gender (if M < 0.01) or we would not be able to rule out null effects for both genders $(M \ge 0.01)$.

How can we interpret the magnitudes of M in this example? We consider a calibration exercise based on the magnitudes of possible possible confounds: if violations of parallel trends were driven by confounding changes in education quality, what would a given value of M imply about the evolution of those confounds? Chetty, Friedman and Rockoff (2014) estimate that a 1 standard deviation increase in teacher value-added (VA) corresponds with a 0.4 percentage point increase in adult employment. Hence, a value of M = 0.01 would correspond with allowing the slope of the differential trend to change by the equivalent of a one-fourtieth of a standard deviation of teacher VA across consecutive periods. Since the robust confidence sets for both men and women begin to include zero around this value of M, the strength with which we can rule out a null effect depends on our assessment of the economic plausibility of such non-linearities.

7 Conclusion

This paper considers the problem of conducting inference in difference-in-differences and related designs that is robust to violations of the parallel trends assumption. We introduce a variety of restrictions on the class of possible differences in trends that formalize commonly made arguments in empirical work, generalizing the framework for partial identification in Manski and Pepper (2018). We provide inference procedures that are uniformly valid so long as the difference in trends satisfies these restrictions, and derive novel results on the power of these procedures. We recommend that applied researchers report robust confidence sets under economically-motivated restrictions on parallel trends. We also recommend that researchers conduct formal sensitivity analyses, in which they report confidence sets for the causal effect of interest under a variety of possible restrictions on the underlying trends. Such sensitivity analyses make transparent what assumptions are needed in order to draw particular conclusions.

References

- Abadie, Alberto, "Semiparametric Difference-in-Differences Estimators," The Review of Economic Studies, 2005, 72 (1), 1–19.
- Andrews, Isaiah, Jonathan Roth, and Ariel Pakes, "Inference for Linear Conditional Moment Inequalities," arXiv:1909.10062 [econ], December 2021. arXiv: 1909.10062.
- Armstrong, Timothy and Michal Kolesár, "Optimal Inference in a Class of Regression Models," *Econometrica*, 2018, 86, 655–683.
- and _ , "Sensitivity Analysis using Approximate Moment Condition Models," Quantitative Economics, 2020. Forthcoming.
- and _, "Simple and Honest Confidence Intervals in Nonparametric Regression," Quantitative Economics, 2020, 11 (1), 1–39. arXiv: 1606.01200.
- Ashenfelter, Orley, "Estimating the Effect of Training Programs on Earnings," *The Review* of Economics and Statistics, 1978, 60 (1), 47–57. Publisher: The MIT Press.
- Athey, Susan and Guido W. Imbens, "Design-based analysis in Difference-In-Differences settings with staggered adoption," *Journal of Econometrics*, April 2021.
- Benzarti, Youssef and Dorian Carloni, "Who Really Benefits from Consumption Tax Cuts? Evidence from a Large VAT Reform in France," American Economic Journal: Economic Policy, February 2019, 11 (1), 38–63.
- Bhuller, Manudeep, Tarjei Havnes, Edwin Leuven, and Magne Mogstad, "Broadband Internet: An Information Superhighway to Sex Crime?," The Review of Economics and Statistics, 2013, 80, 1237—1266.
- Bilinski, Alyssa and Laura A. Hatfield, "Nothing to see here? Non-inferiority approaches to parallel trends and other model assumptions," *arXiv:1805.03273 [stat.ME]*, 2020.
- **Borusyak, Kirill and Xavier Jaravel**, "Revisiting Event Study Designs," SSRN Scholarly Paper ID 2826228, Social Science Research Network, Rochester, NY August 2016.
- Bugni, Federico A., Ivan A. Canay, and Xiaoxia Shi, "Inference for subvectors and other functions of partially identified parameters in moment inequality models," *Quantitative Economics*, 2017, 8 (1), 1–38.

- Callaway, Brantly and Pedro H. C. Sant'Anna, "Difference-in-Differences with multiple time periods," *Journal of Econometrics*, December 2020.
- Canay, Ivan and Azeem Shaikh, "Practical and Theoretical Advances in Inference for Partially Identified Models," in Bo Honoré, Ariel Pakes, Monika Piazzesi, and Larry Samuelson, eds., Advances in Economics and Econometrics, 2017, pp. 271–306.
- Chen, Xiaohong, Timothy M. Christensen, and Elie Tamer, "Monte Carlo Confidence Sets for Identified Sets," *Econometrica*, 2018, *86* (6), 1965–2018.
- Chernozhukov, Victor, Whitney K. Newey, and Andres Santos, "Constrained Conditional Moment Restriction Models," arXiv:1509.06311 [math, stat], September 2015. arXiv: 1509.06311.
- Chetty, Raj, John N. Friedman, and Jonah E. Rockoff, "Measuring the Impacts of Teachers II: Teacher Value-Added and Student Outcomes in Adulthood," American Economic Review, September 2014, 104 (9), 2633–2679.
- Cho, JoonHwan and Thomas M. Russell, "Simple Inference on Functionals of Set-Identified Parameters Defined by Linear Moments," arXiv:1810.03180 [econ.EM], 2019.
- Cox, Gregory and Xiaoxia Shi, "Simple Adaptive Size-Exact Testing for Full-Vector and Subvector Inference in Moment Inequality Models," *The Review of Economic Studies*, March 2022.
- de Chaisemartin, Clément and Xavier D'Haultfœuille, "Two-Way Fixed Effects Estimators with Heterogeneous Treatment Effects," American Economic Review, September 2020, 110 (9), 2964–2996.
- and _ , "Two-Way Fixed Effects and Differences-in-Differences with Heterogeneous Treatment Effects: A Survey," SSRN Scholarly Paper ID 3980758, Social Science Research Network, Rochester, NY December 2021.
- **Dette, Holger and Martin Schumann**, "Difference-in-Differences Estimation Under Non-Parallel Trends," *Working paper*, 2020.
- Dobkin, Carlos, Amy Finkelstein, Raymond Kluender, and Matthew J. Notowidigdo, "The Economic Consequences of Hospital Admissions," American Economic Review, 2018, 108 (2), 308–352.
- **Donoho, David L.**, "Statistical Estimation and Optimal Recovery," The Annals of Statistics, 1994, 22 (1), 238–270.

- Flynn, Zach, "Inference Based on Continuous Linear Inequalities via Semi-Infinite Programming," SSRN Scholarly Paper ID 3390788, Social Science Research Network, Rochester, NY May 2019.
- Frandsen, Brigham R, "The Effects of Collective Bargaining Rights on Public Employee Compensation: Evidence from Teachers, Firefighters, and Police," *ILR Review*, 2016, 69 (1), 84–112.
- Freyaldenhoven, Simon, Christian Hansen, and Jesse Shapiro, "Pre-event Trends in the Panel Event-study Design," *American Economic Review*, 2019, 109 (9), 3307–3338.
- Gafarov, Bulat, "Inference in high-dimensional set-identified affine models," arXiv:1904.00111 [econ.EM], 2019.
- Goodman-Bacon, Andrew, "Public Insurance and Mortality: Evidence from Medicaid Implementation," *Journal of Public Economics*, 2018, 126 (1), 216–262.
- _ , "Difference-in-differences with variation in treatment timing," *Journal of Econometrics*, June 2021.
- Greenstone, Michael and Rema Hanna, "Environmental Regulations, Air and Water Pollution, and Infant Mortality in India," *American Economic Review*, October 2014, 104 (10), 3038–3072.
- Hansen, Bruce E. and Seojeong Lee, "Asymptotic theory for clustered samples," *Journal* of *Econometrics*, June 2019, *210* (2), 268–290.
- Hansen, Christian B., "Asymptotic properties of a robust variance matrix estimator for panel data when T is large," *Journal of Econometrics*, December 2007, 141 (2), 597–620.
- Heckman, James, Hidehiko Ichimura, Jeffrey Smith, and Petra Todd, "Characterizing Selection Bias Using Experimental Data," *Econometrica*, 1998, *66* (5), 1017–1098.
- Hirano, Keisuke and Jack R. Porter, "Impossibility Results for Nondifferentiable Functionals," *Econometrica*, 2012, 80 (4), 1769–1790.
- Ho, Kate and Adam Rosen, "Partial Identification in Applied Research: Benefits and Challenges," in Bo Honoré, Ariel Pakes, Monika Piazzesi, and Larry Samuelson, eds., Advances in Economics and Econometrics, 2017, pp. 307–359.
- Horowitz, Joel L. and Charles F. Manski, "Identification and Robustness with Contaminated and Corrupted Data," *Econometrica*, 1995, 63 (2), 281–302.

- Hudson, Sally, Peter Hull, and Jack Liebersohn, "Interpreting Instrumented Difference-in-Differences," Working Paper 2017.
- Kadane, Joseph B., "Comparison of k-Class Estimators When the Disturbances Are Small," *Econometrica*, 1971, *39* (5), 723–737.
- Kahn-Lang, Ariella and Kevin Lang, "The Promise and Pitfalls of Differences-in-Differences: Reflections on 16 and Pregnant and Other Applications," *Journal of Business* and Economic Statistics, 2020, 38 (3), 613–620.
- Kaido, Hiroaki and Andres Santos, "Asymptotically Efficient Estimation of Models Defined by Convex Moment Inequalities," *Econometrica*, 2014, 82 (1), 387–413.
- _, Francesca Molinari, and Jörg Stoye, "Confidence Intervals for Projections of Partially Identified Parameters," *Econometrica*, 2019, 87, 1397–1432.
- _ , _ , and _ , "CONSTRAINT QUALIFICATIONS IN PARTIAL IDENTIFICATION," Econometric Theory, June 2021, pp. 1–24. Publisher: Cambridge University Press.
- Keele, Luke J., Dylan S. Small, Jesse Y. Hsu, and Colin B. Fogarty, "Patterns of Effects and Sensitivity Analysis for Differences-in-Differences," arXiv:1901.01869 [stat.AP], 2019.
- Kim, Wooyoung, Koohyun Kwon, Soonwoo Kwon, and Sokbae Lee, "The identification power of smoothness assumptions in models with counterfactual outcomes," *Quantitative Economics*, 2018, 9 (2), 617–642.
- Kline, Patrick and Andres Santos, "Sensitivity to missing data assumptions: Theory and an evaluation of the U.S. wage structure," *Quantitative Economics*, 2013, 4 (2), 231–267.
- Kolesár, Michal and Christoph Rothe, "Inference in Regression Discontinuity Designs with a Discrete Running Variable," *American Economic Review*, 2018, 108 (8), 2277–2304.
- Leavitt, Thomas, "Beyond Parallel Trends: Improvements on Estimation and Inference in the Difference-in-Differences Design," *Working paper*, 2020.
- Lee, Jin Young and Gary Solon, "The Fragility of Estimated Effects of Unilateral Divorce Laws on Divorce Rates," The B.E. Journal of Economic Analysis & Policy, 2011, 11 (1).
- Lovenheim, Michael F. and Alexander Willen, "The Long-Run Effects of Teacher Collective Bargaining," American Economic Journal: Economic Policy, 2019, 11 (3), 292– 324.

- Low, Mark G., "Bias-Variance Tradeoffs in Functional Estimation Problems," *The Annals of Statistics*, 1995, 23 (3), 824–835.
- Manski, Charles F., Partial Identification of Probability Distributions, Springer, 2003.
- _, Identification for Prediction and Decision, Harvard University Press, 2007.
- _, Public Policy in an Uncertain World: Analysis and Decisions, Harvard University Press, 2013.
- and John V. Pepper, "Deterrence and the Death Penalty: Partial Identification Analysis Using Repeated Cross Sections," *Journal of Quantitative Criminology*, March 2013, 29 (1), 123–141.
- and _ , "How Do Right-to-Carry Laws Affect Crime Rates? Coping with Ambiguity Using Bounded-Variation Assumptions," *Review of Economics and Statistics*, 2018, 100 (2), 232–244.
- Masten, Matthew A. and Alexandre Poirier, "Inference on Breakdown Frontiers," *Quantitative Economics*, 2020, 11 (1), 41–111.
- Molinari, Francesca, "Chapter 5 Microeconometrics with partial identification," in Steven N. Durlauf, Lars Peter Hansen, James J. Heckman, and Rosa L. Matzkin, eds., Handbook of Econometrics, Vol. 7 of Handbook of Econometrics, Volume 7A, Elsevier, January 2020, pp. 355–486.
- Moreira, Marcelo J., "A Conditional Likelihood Ratio Test for Structural Models," *Econo*metrica, 2003, 71 (4), 1027–1048.
- and Geert Ridder, "Efficiency Loss of Asymptotically Efficient Tests in an Instrumental Variables Regression," SSRN Scholarly Paper ID 3348716, Social Science Research Network, Rochester, NY March 2019.
- Müller, Ulrich K, "Efficient tests under a weak convergence assumption," *Econometrica*, 2011, 79 (2), 395–435.
- Noack, Claudia and Christoph Rothe, "Bias-Aware Inference in Fuzzy Regression Discontinuity Designs," arXiv:1906.04631 [econ.EM], 2020.
- Rambachan, Ashesh and Jonathan Roth, "An Honest Approach to Parallel Trends," Working paper, 2021, p. 126.

- Romano, Joseph P. and Azeem M. Shaikh, "Inference for identifiable parameters in partially identified econometric models," *Journal of Statistical Planning and Inference*, September 2008, 138 (9), 2786–2807.
- _ , _ , and Michael Wolf, "A Practical Two-Step Method for Testing Moment Inequalities," *Econometrica*, 2014, 82 (5), 1979–2002.
- Roth, Jonathan, "Pre-test with Caution: Event-study Estimates After Testing for Parallel Trends," *American Economic Review: Insights*, Forthcoming.
- _, Pedro H. C. Sant'Anna, Alyssa Bilinski, and John Poe, "What's Trending in Difference-in-Differences? A Synthesis of the Recent Econometrics Literature," arXiv:2201.01194 [econ, stat], January 2022. arXiv: 2201.01194.
- Sant'Anna, Pedro H. C. and Jun B. Zhao, "Doubly Robust Difference-in-Differences Estimators," *Journal of Econometrics*, 2020, 219 (1), 101–122.
- Schrijver, Alexander, Theory of Linear and Integer Programming, Wiley-Interscience, 1986.
- Stock, James and Mark Watson, "Heteroskedasticity-Robust Standard Errors for Fixed Effects Panel Data Regression," *Econometrica*, 2008, 76 (1), 155–174.
- Sun, Liyan and Sarah Abraham, "Estimating Dynamic Treatment Effects in Event Studies with Heterogeneous Treatment Effects," *Journal of Econometrics*, 2020.
- Tamer, Elie, "Partial Identification in Econometrics," Annual Review of Economics, 2010, 2, 167–195.
- van der Vaart, Aad W and Jon A Wellner, "Weak Convergence and Empirical Processes: With Applications to Statistics," 1996.
- Wolfers, Justin, "Did Unilateral Divorce Laws Raise Divorce Rates? A Reconciliation and New Results," American Economic Review, 2006, 96, 1802–1820.
- Ye, Ting, Luke Keele, Raiden Hasegawa, and Dylan S. Small, "A Negative Correlation Strategy for Bracketing in Difference-in-Differences with Application to the Effect of Voter Identification Laws on Voter Turnout," arXiv:2006.02423 [stat.ME], 2020.

A More Credible Approach to Parallel Trends

Online Appendix

Ashesh Rambachan Jonathan Roth

April 1, 2022

This online appendix contains proofs and additional results for the paper "A More Credible Approach to Parallel Trends" by Ashesh Rambachan and Jonathan Roth. Section A contains proofs and auxilliary lemmas for results stated in the main text. Section B contains additional details and results from our simulations.

A Proofs of Results in Main Text

Proof of Lemma 2.2

Proof. By equation (7), we can write the coverage requirement as

$$\inf_{\delta \in \Delta, \tau} \inf_{k} \inf_{\theta \in \mathcal{S}(\delta + \tau, \Delta_{k})} \mathbb{P}_{\hat{\beta}_{n} \sim \mathcal{N}(\delta + \tau, \Sigma_{n})} \left(\theta \in \bigcup_{k'} \mathcal{C}_{n,k'}(\hat{\beta}_{n}, \Sigma_{n}) \right) \ge 1 - \alpha.$$

The left-hand side is bounded below by

$$\inf_{\delta \in \Delta, \tau} \inf_{k} \inf_{\theta \in \mathcal{S}(\delta + \tau, \Delta_{k})} \mathbb{P}_{\hat{\beta}_{n} \sim \mathcal{N}(\delta + \tau, \Sigma_{n})} \left(\theta \in \mathcal{C}_{n,k}(\hat{\beta}_{n}, \Sigma_{n}) \right),$$

which is at least $1 - \alpha$ since $C_{n,k}(\hat{\beta}_n, \Sigma_n)$ satisfies (10) for $\Delta = \Delta_k$ for all k.

Proof of Proposition 3.1

Proof. We verify that the conditions of the proposition are sufficient for the conditions for size control for the conditional and hybrid tests given in Proposition 2 of ARP. Note that in our setting, the non-stochastic variable \tilde{X} plays the role of the instruments Z in ARP, so all statements in ARP conditional on Z can be interpreted as unconditional in our context.

First, suppose that Assumption 5(A) holds. Then we can write $\tilde{Y}_n(\theta) = A\hat{\beta}_n - d - \tilde{A}_{(\cdot,1)}\theta = TU_n(\theta) - \zeta(\theta)$, where $U_n(\theta) = Q\hat{\beta}_n$ and $\zeta(\theta) = d + \tilde{A}_{(\cdot,1)}\theta$ is non-stochastic,

which is the structure required by the first part of Assumption 1 of ARP.³⁵ Note that $\Omega_P := Var_P(U_n(\theta)) = Q\Sigma_P Q'$. Since Q is full-rank by assumption and Σ_P has eigenvalues bounded away from zero by Assumption 3, so too does $\Omega_P = Q\Sigma_P Q'$, as required by the latter part of Assumption 1 in ARP. Next, note that our estimate of the variance of $\tilde{Y}_n(\theta)$, $A\hat{\Sigma}_n A'$, can be expressed as $T\hat{\Omega}_n T$, for $\hat{\Omega}_n = Q\hat{\Sigma}_n Q'$. It is immediate from Assumption 4 that $\hat{\Omega}_n$ is uniformly consistent for Ω_P , as required in Assumption 2 in ARP. Next, note that if $f \in BL_1$, then $g(x) = ||G||_{op}^{-1} f(Gx)$ is also in BL_1 , where $||\cdot||_{op}$ is the operator norm. This implies that

$$\sup_{f \in BL_1} \left| \mathbb{E}_P \left[f(\sqrt{n}Q(\hat{\beta} - \beta_P)) \right] - \mathbb{E} \left[f(Q\xi_P) \right] \right| \leq ||Q||_{op} \sup_{f \in BL_1} \left| \mathbb{E}_P \left[f(\sqrt{n}(\hat{\beta} - \beta_P)) \right] - \mathbb{E} \left[f(\xi_P) \right] \right|.$$

Since $U_n(\theta) = Q\hat{\beta}_n$, Assumption 2 together with the previous argument implies that

$$\lim_{n \to \infty} \sup_{P \in \mathcal{P}} \sup_{f \in BL_1} \left| \mathbb{E}_P \left[f(\sqrt{n}(U_n(\theta) - Q\beta_P)) \right] - \mathbb{E} \left[f(\tilde{\xi}_P) \right] \right| = 0,$$

where $\tilde{\xi}_P \sim \mathcal{N}(0, \Omega_P)$. This verifies Assumption 3 in ARP. Note that Assumption 5(A) implies that Assumption C.1 in ARP is satisfied, and Assumption C.2 in ARP is trivially satisfied for $\mathcal{X} = {\tilde{X}}$. Hence, Proposition C.1 in ARP implies that Assumption 4 in ARP is satisfied. We have thus verified the conditions for size control in Proposition 2 of ARP.

Second, consider the case where Assumption 5(B) holds. In this case, we can write $\tilde{Y}_n(\theta) = TU_n(\theta) - \zeta(\theta)$, where now T = A, $U_n(\theta) = \hat{\beta}_n$, and $\zeta(\theta) = d + \tilde{A}_{(\cdot,1)}\theta$. Assumptions 1-3 in ARP can be verified analogously to the arguments above for the case where T is as given in Assumption 5(A). To verify Assumption 4 in ARP, we must show that

$$\sup_{\Sigma_P \in \mathbf{S}} \min_{\gamma, \tilde{\gamma} \in V_{\dagger}(\Sigma_P), \gamma \neq \tilde{\gamma}, a \ge 0} (\gamma - a \tilde{\gamma})' A \Sigma_P A'(\gamma - c \tilde{\gamma}) > 0,$$

where $V_{\dagger}(\Sigma)$ is the subset of vertices in $V(\Sigma)$ that can be optimal when $\hat{\eta} > 0$ (see Lemma 5 in ARP). By Lemma A.1 below, each $\gamma \in V(\Sigma_P)$ can be written as $c_j(\Sigma_P)\bar{\gamma}_j$ for some element $\bar{\gamma}_j \in V(I)$. Moreover, $c_j(\Sigma_P) = (\bar{\gamma}'_j \tilde{\sigma}(\Sigma_P))^{-1}$, where $\tilde{\sigma}(\Sigma_P)$ is the square root of the diagonal elements of $\Omega_P = A\Sigma_P A'$. However, the *j*th diagonal element of Ω_P is $A_{(j,\cdot)}\Sigma_P A'_{(j,\cdot)}$, where $A_{(j,\cdot)}$ is the *j*th row of A. Since the eigenvalues of Σ_P are bounded above by $\bar{\lambda}$, it follows that $A_{(j,\cdot)}\Sigma_P A'_{(j,\cdot)}$ is bounded above by $\bar{\lambda}||A_{(j,\cdot)}||^2$. The elements of $\tilde{\sigma}(\Sigma_P)$ are thus

³⁵Assumption 1 of ARP imposes the structure $Y_i = TU_i + \zeta_i$, where the index *i* corresponds with individual observations and the sample moments are formed by averaging across *i*. However, this structure is only used in the proofs of size control to show that the scaled sample moments, denoted $Y_{n,0}$ in ARP, have the structure $Y_{n,0} = TU_{n,0} + \zeta_{n,0}(\theta)$, where $U_{n,0}$ and $\zeta_{n,0}$ are sample averages of U_i and ζ_i . In our setting \tilde{Y}_n is analogous to $\frac{1}{\sqrt{n}}Y_{n,0}$ in ARP, and we thus verify this structure directly.

bounded above, and hence $c_j(\Sigma_P)$ is bounded away from zero. Thus, there exists a <u>c</u> such that $c_j(\Sigma_P) \ge c$ for all $\Sigma_P \in \mathcal{S}$. Hence,

$$\sup_{\Sigma_{P}\in\mathbf{S}}\min_{\gamma,\tilde{\gamma}\in V_{\dagger}(\Sigma_{P}),\gamma\neq\tilde{\gamma},a\geqslant 0}(\gamma-a\tilde{\gamma})'A\Sigma_{P}A'(\gamma-a\tilde{\gamma})\geqslant\underline{c}^{2}\left(\sup_{\Sigma_{P}\in\mathbf{S}}\min_{\gamma,\tilde{\gamma}\in V(I),\gamma\neq\tilde{\gamma},a\geqslant 0}(\gamma-a\tilde{\gamma})'A\Sigma_{P}A'(\gamma-a\tilde{\gamma})\right)$$
$$\geq\underline{c}^{2}\left(\min_{\gamma,\tilde{\gamma}\in V_{\dagger}(I),\gamma\neq\tilde{\gamma},a\geqslant 0}||(\gamma-a\tilde{\gamma})'A||^{2}\underline{\lambda}\right),$$

where the second inequality uses the fact that the minimal eigenvalue of Σ_P is at least $\underline{\lambda}$. To complete the proof, it thus suffices to show that $V_{\dagger}(I)$ contains only vertices such that $\bar{\gamma}'_j A \neq 0$, so that the lower bound obtained in the previous display is strictly positive by Assumption 5(B). To show this, note that if $\bar{\gamma}'_j A = 0$, then $\bar{\gamma}'_j \tilde{Y}_n(\bar{\theta}) = \bar{\gamma}'_j (A\hat{\beta}_n - d - \tilde{A}_{(\cdot,1)}\bar{\theta}) = -\bar{\gamma}'_j d$. Since Δ is non-empty, there exists some δ such that $A\delta - d \leq 0$, which implies that $-\bar{\gamma}'_j d = \bar{\gamma}'_j (A\delta - d) \leq 0$ since $\bar{\gamma}_j \geq 0$ by construction. We have thus established that $\bar{\gamma}'_j \tilde{Y}(\bar{\theta}) \leq 0$, and hence $\bar{\gamma}_j$ can never be optimal when $\hat{\eta} > 0$, so $\bar{\gamma}_j \notin V_{\dagger}(I)$. We have thus verified that Assumption 4 in ARP holds, as needed.

A.1 Proof and auxiliary lemmas for uniform consistency

Proof of Proposition 3.2

Proof. Towards contradiction, suppose that the conditional test is not consistent. Then there exists an increasing sequence of sample sizes and distributions (n_m, P_m) , x > 0, and $\omega > 0$ such that

$$\limsup_{m \to \infty} \mathbb{E}_{P_m} \left[\psi_{\alpha}^C(\hat{\beta}_{n_m}, A, d, \theta_{P_m}^{ub} + x, \frac{1}{n} \hat{\Sigma}_{n_m}) \right] \leq 1 - \omega.$$

It is straightforward to verify that the conditional test is invariant to a re-scaling of the units of $\hat{\beta}$, so that $\psi_{\alpha}^{C}(\hat{\beta}_{n_{m}}, A, d, \theta_{P_{m}}^{ub} + x, \frac{1}{n_{m}}\hat{\Sigma}_{n_{m}}) = \psi_{\alpha}^{C}(\sqrt{n_{m}}\hat{\beta}_{n_{m}}, A, \sqrt{n_{m}}d, \sqrt{n_{m}}(\theta_{P_{m}}^{ub} + x), \hat{\Sigma}_{n_{m}}).$ Thus, along this sequence,

$$\limsup_{m \to \infty} \mathbb{E}_{P_m} \left[\psi_{\alpha}^C(\sqrt{n_m} \hat{\beta}_{n_m}, A, \sqrt{n_m} d, \sqrt{n_m} (\theta_{P_m}^{ub} + x), \hat{\Sigma}_{n_m}) \right] \leq 1 - \omega.$$

Since **V** is compact, we can extract a further subsequence m_1 under which $V_{P_{m_1}} \to V^*$ for $V^* \in \mathbf{V}$. Denote the top left block of V^* by Σ^* .

To obtain a contradiction, we will construct a further subsequence such that the conditions of Lemma A.2 hold asymptotically with probability at least $1-\omega/2$. From Lemma A.1, each element $\gamma_{P_{m_1}} \in V(\Sigma_{P_{m_1}})$ can be written as $c_j(\Sigma_{P_{m_1}})\bar{\gamma}_j$, where $\bar{\gamma}_1, ..., \bar{\gamma}_J$ are the elements of V(I). We argued in the proof to Proposition 3.1 that there exists a constant \underline{c} such that $c_j(\Sigma_P) \ge \underline{c}$ for all j whenever Σ_P has eigenvalues bounded above by $\overline{\lambda}$. By an analogous argument, we can show that there exists a constant \overline{c} such that $c_j(\Sigma_P) \le \overline{c}$ whenever Σ_P has eigenvalues bounded below by $\underline{\lambda}$. Thus, $\underline{c} \le c_j(\Sigma_P) \le \overline{c}$ for $\Sigma_P \in \mathcal{S}$. For $\gamma \in V(\Sigma_P)$, $\gamma' A \Sigma_P A' \gamma = c_j(\Sigma_P)^2 \overline{\gamma}'_j A \Sigma_P A' \overline{\gamma}_j$ for some j, and thus for $\Sigma_P \in \mathbf{S}$, we have that

$$\underline{c}^2 ||\bar{\gamma}_j'A||^2 \underline{\lambda} \leqslant \gamma' A \Sigma_P A' \gamma \leqslant \overline{c}^2 ||\bar{\gamma}_j'A||^2 \overline{\lambda}.$$

Thus, either $\gamma' A \Sigma_P A' \gamma = 0$ (if $\bar{\gamma}'_j A = 0$), or

$$\underline{c}^{2} \min_{j:\bar{\gamma}'_{j}A\neq 0} ||\bar{\gamma}'_{j}A||^{2} \underline{\lambda} \leqslant \gamma' A \Sigma_{P} A' \gamma \leqslant \bar{c}^{2} \max_{j:\bar{\gamma}'_{j}A\neq 0} ||\bar{\gamma}'_{j}A||^{2} \overline{\lambda},$$

where the upper and lower bounds are finite and positive since V(I) is finite. Now consider the vertex $\hat{\gamma}_{m_1,j} = c_j(\hat{\Sigma}_{n_{m_1}})\bar{\gamma}_j$. By the continuous mapping theorem, $\hat{\gamma}'_{m_1,j}A\hat{\Sigma}_P A'\hat{\gamma}_{m_1,j} \rightarrow_p c_j(\Sigma^*)^2 \bar{\gamma}'_j A \Sigma^* A' \bar{\gamma}_j$. From this convergence and the inequalities in the previous display, it follows that there exist constants $\underline{\sigma}^2$ and $\bar{\sigma}^2$ such that condition (i) of Lemma A.2 is satisfied w.p.a. 1.

Next, define

$$\eta(\beta, A, d, \bar{\theta}, \Sigma) := \min_{\eta, \tilde{\tau}} \eta \text{ s.t. } A\beta - d - \tilde{A}_{(\cdot, 1)} \bar{\theta} - \tilde{A}_{(\cdot, -1)} \tilde{\tau} \leqslant \eta \tilde{\sigma},$$
(21)

where $\tilde{\sigma}$ is the square root of the diagonal elements of $A\Sigma A'$. Since $\theta_P^{ub} \in S(\beta_P, \Delta)$, $\eta(\beta_P, A, d, \theta_P^{ub}, \Sigma_P) \leq 0$. By duality, we can write $\eta(\beta_P, A, d, \theta_P^{ub}, \Sigma_P) = \max_{\gamma \in V(\Sigma_P)} \gamma'(A\beta_P - d - \tilde{A}_{(\cdot,1)}\theta^{ub})$. It follows that there exists some $\tilde{\gamma}_P \in V(\Sigma_P)$ such that $\tilde{\gamma}'_P(A\beta_P - d - \tilde{A}_{(\cdot,1)}\theta^{ub}) = 0$ and $-\tilde{\gamma}'_P\tilde{A}_{(\cdot,1)} > 0$, since otherwise for $\epsilon > 0$ sufficiently small we would have that $\eta(\beta_P, A, d, \theta_P^{ub} + \epsilon, \Sigma_P) = \max_{\gamma \in V(\Sigma_P)} \gamma'(A\beta_P - d - \tilde{A}_{(\cdot,1)}(\theta^{ub} + \epsilon)) \leq 0$, which would imply that $\theta^{ub} + \epsilon \in S(\beta_P, \Delta)$, which is a contradiction. From Lemma A.1, $\tilde{\gamma}_{Pm_1} \in V(\Sigma_{Pm_1})$ can be written as $c_j(\Sigma_{Pm_1})\bar{\gamma}_j$, where $c_j(\Sigma) \geq c > 0$ for all $\Sigma \in \mathbf{S}$ and $\bar{\gamma}_1, ..., \bar{\gamma}_J$ are the elements of V(I). Since V(I) is finite, we can extract a further subsequence (n_l, P_l) such that $\tilde{\gamma}_{Pl} = c_{j^*}(\Sigma_P)\tilde{\gamma}_{j^*}$ for fixed j^* . For ease of notation, without loss of generality we assume $j^* = 1$. It follows that

$$\begin{split} \eta(\sqrt{n_l}\hat{\beta}_{n_l}, A, \sqrt{n_l}d, \sqrt{n_l}(\theta_{P_l}^{ub} + x), \hat{\Sigma}_{n_l}) &= \max_{\gamma \in V(\hat{\Sigma}_{n_l})} \gamma' \sqrt{n_l} (A\hat{\beta}_{n_l} - d - \tilde{A}_{(\cdot,1)}(\theta_{P_l}^{ub} + x)) \\ &\geqslant \sqrt{n_l} c_1(\hat{\Sigma}_{n_l}) \bar{\gamma}_1' (A\hat{\beta}_{n_l} - d - \tilde{A}_{(\cdot,1)}(\theta_{P_l}^{ub} + x)) \\ &= c_1(\hat{\Sigma}_{n_l}) \sqrt{n_l} \bar{\gamma}_1' A(\hat{\beta}_{n_l} - \beta_{P_l}) + \sqrt{n_l} c_1(\hat{\Sigma}_{n_l}) (-\tilde{\gamma}_1' \tilde{A}_{(\cdot,1)}) x. \end{split}$$

By the continuous mapping theorem, $c_1(\hat{\Sigma}_{n_l}) \to_p c_1(\Sigma^*) > 0$. Assumption 6 and the continuous mapping theorem together imply that the first term in the previous display converges in distribution to a $\mathcal{N}(0, c_1(\Sigma^*)^2 \bar{\gamma}'_1 A \Sigma^* A' \bar{\gamma}_1)$ distribution, while the second term converges in probability to ∞ . It follows that $\eta(\sqrt{n_l}\hat{\beta}_l, A, \sqrt{n_l}d, \sqrt{n_l}(\theta_{P_l}^{ub} + x), \hat{\Sigma}_{n_l}) \to_p \infty$, and thus condition (ii) of Lemma A.2 holds w.p.a. 1 for any value of M.

To complete the proof, we construct a further subsequence such that condition (iii) of Lemma A.2 holds asymptotically with probability at least $1-\omega/2$. Let $\tilde{Y}_l = A\hat{\beta}_{n_l} - d - \tilde{A}_{(\cdot,1)}(\theta_{P_l}^{ub} + x)$ and $\tilde{\mu}_l = A\beta_{P_l} - d - \tilde{A}_{(\cdot,1)}(\theta_{P_l}^{ub} + x)$. Recall that any element of $V(\hat{\Sigma}_{n_l})$, say $\gamma_{l,j}$, takes the form $\gamma_{l,j} = c_j(\hat{\Sigma}_{n_l})\bar{\gamma}_j$, and our argument above implies that $\bar{\gamma}'_j\tilde{\mu}_l \leq -\bar{\gamma}'_j\tilde{A}_{(\cdot,1)}x$. Since $c_j(\hat{\Sigma}_{n_l}) \rightarrow_p c_j(\Sigma^*) > 0$ by the continuous mapping theorem, and $\bar{\gamma}'_j\tilde{\mu}_l$ is bounded from above, we can extract a subsequence l_1 along which $\gamma'_{l_1,j}\tilde{\mu}_{l_1} \rightarrow_p \nu_j \in \mathbb{R} \cup \{-\infty\}$. The vertex set is finite, and so passing to further subsequences we obtain a subsequence indexed by ksuch that $\gamma'_{k,j}\tilde{\mu}_k \rightarrow_p \nu_j \in \mathbb{R} \cup \{-\infty\}$ for all j. Observe that for distinct vertices i and j with $\bar{\gamma}'_i A \neq 0$,

$$(c_i(\hat{\Sigma}_{n_k})\bar{\gamma}_i - c_j(\hat{\Sigma}_{n_k})\bar{\gamma}_j)'\sqrt{n_k}\tilde{Y}_k = (c_i(\hat{\Sigma}_{n_k})\bar{\gamma}_i - c_j(\hat{\Sigma}_{n_k})\bar{\gamma}_j)'\sqrt{n_k}(\tilde{Y}_k - \tilde{\mu}_k) + \sqrt{n_k}(c_i(\hat{\Sigma}_{n_k}) - c_i(\Sigma^*))\bar{\gamma}_i'\tilde{\mu}_k - \sqrt{n_k}(c_j(\hat{\Sigma}_{n_k}) - c_j(\Sigma^*))\bar{\gamma}_j'\tilde{\mu}_k + \sqrt{n_k}(c_i(\Sigma^*)\bar{\gamma}_i' - c_j(\Sigma^*)\bar{\gamma}_j')\tilde{\mu}_k$$

Consider first the case where $\gamma'_{k,i}\tilde{\mu}_k$ and $\gamma'_{k,j}\tilde{\mu}_k$ both have finite limits ν_i and ν_j . Since $\sqrt{n_k}(c_i(\Sigma^*)\bar{\gamma}'_i - c_j(\Sigma^*)\bar{\gamma}'_j)\tilde{\mu}_k$ is non-stochastic, we can extra a further subsequence k_1 such that $\sqrt{n_{k_1}}(c_i(\Sigma^*)\bar{\gamma}'_i - c_j(\Sigma^*)\bar{\gamma}'_j)\tilde{\mu}_{k_1} \rightarrow \nu^* \in \mathbb{R} \cup \{\pm \infty\}$. Assumption 6 and the continuous mapping theorem imply that $(c_i(\hat{\Sigma}_{n_{k_1}})\bar{\gamma}_i - c_j(\hat{\Sigma}_{n_{k_1}})\bar{\gamma}_j)'\sqrt{n_{k_1}}\tilde{Y}_{k_1}$ converges in distribution to

$$\zeta_{ij} = (c_i(\Sigma^*)\bar{\gamma}_i - c_j(\Sigma^*)\bar{\gamma}_j)'A\xi_\beta + \frac{\nu_i}{c_i(\Sigma^*)}Dc_i'\xi_\Sigma - \frac{\nu_j}{c_j(\Sigma^*)}Dc_j'\xi_\Sigma + \nu^*,$$

where $(\xi'_{\beta},\xi'_{\Sigma})' \sim \mathcal{N}(0, V^*)$ and Dc_i is the gradient of $c_i(\Sigma^*)$ with respect to $vec(\Sigma^*)$. The limiting distribution is normal, and limiting variance must be positive since Assumptions 5 and 7 imply that $(c_i(\Sigma^*)\bar{\gamma}_i - c_j(\Sigma^*)\bar{\gamma}_j)'A\xi_{\beta}$ has positive variance³⁶ and is not perfectly colinear with ξ_{Σ} . It follows that for any ϑ , there exists some $\epsilon > 0$ such that the probability that $\zeta_{ij} \in (-\epsilon, \epsilon)$ is less than ϑ . On the other hand, if $\bar{\gamma}'_i \tilde{\mu}_k \to -\infty$, then $c_i(\hat{\Sigma}_{n_k})\bar{\gamma}_i\sqrt{n_k}\tilde{Y}_k \to_p$ $-\infty$, so $c_i(\hat{\Sigma}_{n_k})\bar{\gamma}_i$ is optimal for $\hat{\eta}(\sqrt{n_k}\hat{\beta}_{n_k},\sqrt{n_k}d,\sqrt{n_k}(\theta^{ub}_{P_k}+x),\hat{\Sigma}_{n_k}) - c_j(\hat{\Sigma}_k)\gamma'_j\sqrt{n_k}\tilde{Y}_k \to_p \infty$. Since there

³⁶This is immediate under Assumption 5(ii). Under Assumption 5(i), the proof of Proposition C.2 in ARP shows that if there is a positive constant c such $(\bar{\gamma}_i - c\bar{\gamma}_j)'A = 0$, then $c_i(\hat{\Sigma}_{n_k})\bar{\gamma}_i$ and $c_j(\hat{\Sigma}_{n_k})\bar{\gamma}_j$ can only be optimal vertices if $\hat{\eta} \leq 0$. Since we've shown $\hat{\eta} \rightarrow_p \infty$, such vertices will be optimal w.p.a. 0, and thus can be ignored when establishing part (iii) of Lemma A.2.

are a finite number of pairs of vertices, we can choose ϑ such that the probability that $\zeta_{ij} \in (-\epsilon, \epsilon)$ for any (i, j) is bounded above by $\omega/2$, and thus condition (iii) of Lemma A.2 is satisfied with probability at least $\omega/2$, as we wished to show. The result for the hybrid test is immediate from the fact that the hybrid test rejects whenever the size- $\frac{\alpha-\kappa}{1-\kappa}$ conditional test rejects.

Lemma A.1. Let $F(\Sigma) := \{\gamma : \tilde{A}'_{(\cdot,-1)}\gamma = 0, \tilde{\sigma}(\Sigma)'\gamma = 1, \gamma \ge 0\}$ be the feasible set of the dual problem, where $\tilde{\sigma}(\Sigma)$ is the vector containing the square-roots of the diagonal elements of $A\Sigma A'$. Let $V(\Sigma)$ denote the set of vertices of $F(\Sigma)$. Then, for any Σ positive definite,

$$V(\Sigma) = \{c_1(\Sigma)\bar{\gamma}_1, ..., c_J(\Sigma)\bar{\gamma}_J\},\$$

where $\bar{\gamma}_1, ..., \bar{\gamma}_J$ are the elements of V(I) and $c_j(\Sigma) = (\bar{\gamma}'_j \tilde{\sigma}(\Sigma))^{-1}$.

Proof of Lemma A.1

Proof. Recall that v is a vertex of the polyhedron $P = \{x \in \mathbb{R}^K : Wx \leq b\}$ iff $v \in P$ and $W_{(\mathcal{J},\cdot)}x = b_{\mathcal{J}}$ for \mathcal{J} a set of indices such that $W_{(\mathcal{J},\cdot)}$ has K independent rows (see Section 8.5 of Schrijver (1986)). It follows that $v \in V(\Sigma)$ iff $v \geq 0$ and there exists \mathcal{J} such that

$$W_{\mathcal{J}} := \begin{pmatrix} \tilde{A}'_{(\cdot,-1)} \\ -I_{(\mathcal{J},\cdot)} \\ \tilde{\sigma}' \end{pmatrix}$$

has row rank equal to K, and $W_{\mathcal{J}}v = \begin{pmatrix} 0\\ 0\\ 1 \end{pmatrix}$, where K is the number of rows of A.

Now, let \mathcal{J} be the set of indices \mathcal{J} such that $\tilde{W}_{\mathcal{J}} := \begin{pmatrix} \tilde{A}'_{(\cdot,-1)} \\ -I_{(\mathcal{J},\cdot)} \end{pmatrix}$ has exactly K-1linearly independent rows and there exists a vector $v_{\mathcal{J}} \neq 0$ such that $\tilde{W}_{\mathcal{J}}v = 0$ and $v_{\mathcal{J}} \geq 0$. Since by construction $\tilde{W}_{\mathcal{J}}$ has rank K-1 and K columns, its nullspace is 1-dimensional. It is then immediate that for each $\mathcal{J} \in \mathcal{J}$, there is a unique vector $\bar{v}_{\mathcal{J}} \geq 0$ such that $\tilde{W}_{\mathcal{J}}\bar{v}_{\mathcal{J}} = 0$ and $\iota'\bar{V}_{\mathcal{J}} = 1$, where ι is the vector of ones. Moreover, \mathcal{J} is finite, since there are a finite number of possible subindices of I, and thus we can write $\{\bar{v}_{\mathcal{J}} : \mathcal{J} \in \mathcal{J}\} = \{\bar{\gamma}_1, ..., \bar{\gamma}_J\}$ for distinct vectors $\bar{\gamma}_1, ..., \bar{\gamma}_J$.

It now remains to show that $V(\Sigma) = \{c_1(\Sigma)\bar{\gamma}_1, ..., c_J(\Sigma)\bar{\gamma}_J\}$, for c_j as defined above. First, suppose that $v = c_j(\Sigma)\bar{\gamma}_j$ for some j. By construction, $\tilde{A}'_{(\cdot,-1)}v = 0, v \ge 0$, and $\tilde{\sigma}'v = (\tilde{\sigma}'v_j)^{-1}(\tilde{\sigma}'v_j) = 1$, and so $v \in F$. Additionally, there exists \mathcal{J} such that $\tilde{W}_{\mathcal{J}} = \begin{pmatrix} \tilde{A}'_{(\cdot,-1)} \\ -I_{(\mathcal{J},\cdot)} \end{pmatrix}$ has rank K-1 and $\tilde{W}_{\mathcal{J}}v = 0$. From the fact that $\tilde{W}_{\mathcal{J}}v = 0$, whereas $\tilde{\sigma}'v = 1$, we see that $\tilde{\sigma}'$ must be linearly independent from the rows of $\tilde{W}_{\mathcal{J}}$, and thus $W_{\mathcal{J}} = \begin{pmatrix} \tilde{W}_{\mathcal{J}} \\ \tilde{\sigma}' \end{pmatrix}$ has rank K. It follows that $v \in V(\Sigma)$.

Next, suppose that $v \in V(\Sigma)$. Then $v \ge 0$, and there exists \mathcal{J} such that

$$W_{\mathcal{J}} := \begin{pmatrix} \tilde{A}'_{(\cdot,-1)} \\ -I_{(\mathcal{J},\cdot)} \\ \tilde{\sigma}' \end{pmatrix}$$

has row rank equal to K, and $W_{\mathcal{J}}v = \begin{pmatrix} 0\\0\\1 \end{pmatrix}$. Let $\tilde{W}_{\mathcal{J}} = \begin{pmatrix} \tilde{A}'_{(\cdot,-1)}\\-I_{(\mathcal{J},\cdot)} \end{pmatrix}$. Note that since

 $\tilde{W}_{\mathcal{J}}v = 0$, whereas $\tilde{\sigma}'v = 1$, $\tilde{\sigma}'$ must be linearly independent of the other rows of $W_{\mathcal{J}}$, from which it follows that \tilde{W} has row rank K - 1. Thus, $\mathcal{J} \in \mathcal{J}$, and so $v = c\bar{\gamma}_j$ for some j and c > 0. Since $\tilde{\sigma}'v = 1$, we have $c\tilde{\sigma}'\bar{\gamma}_j = 1$, which implies $c = (\tilde{\sigma}'\bar{\gamma}_j)^{-1}$, which gives the desired result.

Finally, by construction, $\bar{\gamma}'_{j}\iota = 1$, and so $c_{j}(I) = 1$ for all j, so $\bar{\gamma}_{1}, ..., \bar{\gamma}_{J}$ correspond precisely with the elements of V(I).

Lemma A.2. For any positive constants $\epsilon, \underline{\sigma}^2, \overline{\sigma}^2$, there exists a finite constant \overline{C} such that the conditional test $\psi^C_{\alpha}(\hat{\beta}, A, d, \theta, \Sigma)$ rejects whenever the following conditions are satisfied

- (i) For all $\gamma \in V(\Sigma)$, either $\gamma' A \Sigma A' \gamma = 0$ or $\underline{\sigma}^2 \leqslant \gamma' A \Sigma A' \gamma \leqslant \overline{\sigma}^2$.
- (*ii*) $\hat{\eta} = \max_{\gamma \in V(\Sigma)} \gamma' \tilde{Y} > \bar{C}$, where $\tilde{Y} = A\hat{\beta} d \tilde{A}_{(\cdot,1)}\theta$.
- (iii) If the optimal vertex γ_* satisfies, $\gamma'_* A \Sigma A' \gamma_* > 0$, then for all $\tilde{\gamma} \in V(\Sigma)$ with $\tilde{\gamma} \neq \gamma_*$, we have that $|\gamma'_* \tilde{Y} \tilde{\gamma}' \tilde{Y}| > \epsilon$.

Proof. Let $\tilde{\Sigma} = A\Sigma A'$. If the optimal vertex γ_* satisfies $\gamma'_* \tilde{\Sigma} \gamma_* = 0$, then the conditional test rejects whenever $\hat{\eta} > 0$, so condition (ii) with any C > 0 suffices. For the remainder of the proof, we show that conditions (i)-(iii) are sufficient when $\gamma'_* \Sigma \gamma_* \neq 0$. Observe that the conditional test rejects if and only if $\hat{\eta} > 0$ and

$$\frac{\Phi(t) - \Phi(z^{lo})}{\Phi(z^{up}) - \Phi(z^{lo})} > 1 - \alpha,$$

where $t = \frac{\hat{\eta}}{\sigma^*}$, $z^{lo} = \frac{v^{lo}}{\sigma^*}$, $z^{up} = \frac{v^{up}}{\sigma^*}$, and $\sigma^* = \sqrt{\gamma'_* \tilde{\Sigma} \gamma_*}$. It is clear that the left-hand side of the previous display is increasing in t and decreasing in z^{up} . It is also decreasing in z^{lo} , since

the derivative with respect to z^{lo} is

$$-\frac{\phi(z_{lo})(\Phi(z^{up}) - \Phi(z^{lo}))}{(\Phi(z^{up}) - \Phi(z^{lo}))^2} < 0.$$

From Lemma A.3 below, condition (iii) implies that $\hat{\eta} - v^{lo} \ge \epsilon$, and thus $z^{lo} \le t - \tilde{\epsilon}$, for $\tilde{\epsilon} = \epsilon/\bar{\sigma}$. This, combined with the previous discussion, implies that the conditional test rejects whenever $\hat{\eta} > 0$ and

$$\frac{\Phi(t) - \Phi(t - \tilde{\epsilon})}{1 - \Phi(t - \tilde{\epsilon})} > 1 - \alpha.$$

By L'Hopitale's rule, we have that

$$\lim_{t \to \infty} \frac{\Phi(t) - \Phi(t - \tilde{\epsilon})}{1 - \Phi(t - \tilde{\epsilon})} = \lim_{t \to \infty} \frac{\phi(t - \tilde{\epsilon}) - \phi(t)}{\phi(t - \tilde{\epsilon})} = \lim_{t \to \infty} 1 - \frac{\phi(t)}{\phi(t - \tilde{\epsilon})} = 1.$$

Hence, there exists $\tilde{C} > 0$ such that the conditional test rejects whenever $t \ge \tilde{C}$. But $t = \frac{\hat{\eta}}{\sigma_*}$ and thus $t > \tilde{C}$ whenever $\hat{\eta} > \bar{C}$ for $\bar{C} = \tilde{C}\bar{\sigma}$.

Lemma A.3. Consider the conditional test $\psi_{\alpha}^{C}(\hat{\beta}, A, d, \theta, \Sigma)$. If the optimal vertex γ_{*} is such that $\gamma'_{*}A\Sigma A'\gamma_{*} > 0$, then $\hat{\eta} - v^{lo} \ge \min_{\gamma \in V(\Sigma), \gamma \neq \gamma_{*}} |\gamma'_{*}\tilde{Y} - \gamma'\tilde{Y}|$ where $\tilde{Y} = A\hat{\beta} - d - \tilde{A}_{(\cdot,1)}\theta$. Similarly, $v^{up} - \hat{\eta} \ge \frac{\gamma'_{*}A\Sigma A'\gamma_{*}}{\max_{\gamma \in V(\Sigma)}\gamma'A\Sigma A'\gamma} \min_{\gamma \in V(\Sigma), \gamma \neq \gamma_{*}} |\gamma'_{*}\tilde{Y} - \gamma'\tilde{Y}|$.

Proof. Since $\hat{\eta}$ is finite, the results hold trivially when v^{lo} and v^{up} are infinite. For the remainder of the proof, we assume that they are finite. Let $\tilde{\Sigma} = A\Sigma A'$. Lemma 2 in ARP implies that

$$v^{lo} = \min_{\gamma \in V(\Sigma): \gamma'_* \tilde{\Sigma} \gamma_* - \gamma'_* \tilde{\Sigma} \gamma > 0} \frac{\gamma'_* \Sigma \gamma_* \gamma' S}{\gamma'_* \tilde{\Sigma} \gamma_* - \gamma'_* \tilde{\Sigma} \gamma},$$

where $S = (I - \frac{\tilde{\Sigma}\gamma_*}{\gamma'_*\tilde{\Sigma}\gamma_*}\gamma'_*)\tilde{Y}$. Let $\tilde{\gamma}$ denote the vertex at which the minimum is obtained. Substituting in the definition of S and re-arranging terms, we obtain that

$$\hat{\eta} - v^{lo} = \frac{\gamma'_* \tilde{\Sigma} \gamma_*}{\gamma'_* \tilde{\Sigma} \gamma_* - \gamma'_* \tilde{\Sigma} \tilde{\gamma}} (\gamma'_* \tilde{Y} - \tilde{\gamma}' \tilde{Y}) \ge (\gamma'_* \tilde{Y} - \tilde{\gamma}' \tilde{Y}),$$

from which the result for v^{lo} is immediate. We can analogously show that

$$v^{up} - \hat{\eta} = \frac{\gamma'_* \Sigma \gamma_*}{\gamma'_* \tilde{\Sigma} \tilde{\gamma} - \gamma'_* \tilde{\Sigma} \gamma_*} (\gamma'_* \tilde{Y} - \tilde{\gamma}' \tilde{Y}),$$

for a vertex $\tilde{\gamma}$ such that $\gamma'_* \tilde{\Sigma} \tilde{\gamma} - \gamma'_* \tilde{\Sigma} \gamma_* > 0$. The result then follows from noting that

$$\frac{\gamma'_*\tilde{\Sigma}\gamma_*}{\gamma'_*\tilde{\Sigma}\tilde{\gamma} - \gamma'_*\tilde{\Sigma}\gamma_*} \ge \frac{\gamma'_*\tilde{\Sigma}\gamma_*}{\gamma'_*\tilde{\Sigma}\tilde{\gamma}} \ge \frac{\gamma'_*\tilde{\Sigma}\gamma_*}{\max_{\gamma \in V(\Sigma)}\gamma'\tilde{\Sigma}\gamma}.$$

A.2 Proof and auxiliary lemmas for uniform local asymptotic power Proof of Proposition 3.3

Proof. By an invariance to scale argument as in Proposition 3.2, it is sufficient to show that

$$\lim_{n \to \infty} \sup_{P \in \mathcal{P}_{\epsilon}} \left| \mathbb{E}_{P} \left[\psi_{\alpha}^{C}(\sqrt{n}\hat{\beta}_{n}, A, \sqrt{n}d, \sqrt{n}\theta_{P}^{ub} + x, \hat{\Sigma}_{n}) \right] - \rho_{\alpha}^{*}(P, x) \right| = 0.$$

To show this, it suffices to establish that for every subsequence (n_m, P_m) with $n_m \to \infty$, there exists a further subsequence l such that

$$\lim_{l \to \infty} \left| \mathbb{E}_{P_l} \left[\psi^C_\alpha(\sqrt{n_l} \hat{\beta}_{n_l}, A, \sqrt{n_l} d, \sqrt{n} \theta^{ub}_{P_l} + x, \hat{\Sigma}_{n_l}) \right] - \rho^*_\alpha(P_l, x) \right| = 0$$

Since $P_m \in \mathcal{P}_{\epsilon}$, for each *m* there exists a B_m^* and a value $\tilde{\tau}_m^*$ such that

$$A_{(B_m^*,\cdot)}\beta_{P_m} - d_{B_m^*} - \tilde{A}_{(B_m^*,1)}\theta_{P_m}^{ub} - \tilde{A}_{(B_m^*,-1)}\tilde{\tau}_m^* = 0$$
(22)

$$A_{(-B_{m}^{*},\cdot)}\beta_{P_{m}} - d_{-B_{m}^{*}} - \tilde{A}_{(-B_{m}^{*},1)}\theta_{P_{m}}^{ub} - \tilde{A}_{(-B_{m}^{*},-1)}\tilde{\tau}_{m}^{*} < -\epsilon.$$
(23)

Since there are a finite number of possible values of B_m^* , we can extract a subsequence m_1 along which $B_{m_1}^*$ is constant. For simplicity of notation, we'll denote the constant value $B_{m_1}^*$ by B^* . Similarly, Lemma A.4 implies that there is a unique element $\gamma_{m_1}^* \in V(\Sigma_{P_{m_1}})$ such that the elements of $\gamma_{m_1}^*$ in positions $-B^*$ are all 0. By Lemma A.1, we can write $\gamma_{m_1}^* = c_j(\Sigma_{P_{m_1}})\bar{\gamma}_j$ for $c_j(\cdot)$ a continuous function and $\bar{\gamma}_j \in V(I)$. Since V(I) is finite, we can extract a subsequence m_2 along which $\gamma_{m_2}^* = c_{j^*}(\Sigma_{P_{m_2}})\bar{\gamma}_{j^*}$ for a fixed j^* , which without loss of generality we normalize to $j^* = 1$. Moreover, since **S** is compact, we can extract a further subsequence l along which $\Sigma_{P_l} \to \Sigma^*$. By Assumption 4, $\hat{\Sigma}_{n_l} \to_p \Sigma^*$. The continuous mapping theorem then implies that $\gamma_l^* = c_1(\Sigma_{P_l})\bar{\gamma}_1 \to c_1(\Sigma^*)\bar{\gamma}_1$, and likewise, $\hat{\gamma}_l^* := c_1(\hat{\Sigma}_{n_l})\bar{\gamma}_1 \to_p c_1(\Sigma^*)\bar{\gamma}_1$. From Lemma A.9, we have that

$$\rho_{\alpha}^{*}(P_{l},x) = \Phi\left(\frac{-\gamma_{l}^{*'}\tilde{A}_{(\cdot,1)}x}{\sqrt{\gamma_{l}^{*'}A\Sigma_{P_{l}}A'\gamma_{l}^{*}}} - z_{1-\alpha}\right),$$

which combined with the convergences shown above implies that

$$\rho_{\alpha}^{*}(P_{l}, x) \to \Phi\left(\frac{-\bar{\gamma}_{1}'\tilde{A}_{(\cdot,1)}x}{\sqrt{\bar{\gamma}_{1}'A\Sigma^{*}A'\bar{\gamma}_{1}}} - z_{1-\alpha}\right).$$
(24)

Now, for the function $\eta(\cdot)$ defined in (21), let

$$\hat{\eta}_l = \eta(\sqrt{n_l}\hat{\beta}_{n_l}, A, \sqrt{n_l}d, \sqrt{n_l}\theta_{P_l}^{ub} + x, \hat{\Sigma}_{n_l}).$$

By duality, we have that

$$\hat{\eta}_{l} = \max_{\gamma \in V(\hat{\Sigma}_{n_{l}})} \gamma' \left(\sqrt{n} A \hat{\beta}_{n_{l}} - \sqrt{n_{l}} d - \sqrt{n_{l}} \tilde{A}_{(\cdot,1)} \theta_{P_{l}}^{ub} - \tilde{A}_{(\cdot,1)} x \right)$$
$$\geqslant \hat{\gamma}_{l}^{*\prime} \left(\sqrt{n} A \hat{\beta}_{n_{l}} - \sqrt{n_{l}} d - \sqrt{n_{l}} \tilde{A}_{(\cdot,1)} \theta_{P_{l}}^{ub} - \tilde{A}_{(\cdot,1)} x \right).$$

By construction, $\hat{\gamma}_l^*$ has zero elements in positions $-B^*$ and satisfies $\hat{\gamma}_l^{*'} \hat{A}_{(\cdot,-1)} = 0$. This, combined with equation (22) implies that

$$\hat{\gamma}_l^{*\prime}\left(\sqrt{n_l}A\hat{\beta}_{n_l}-\sqrt{n_l}d-\sqrt{n_l}\tilde{A}_{(\cdot,1)}\theta_{P_l}^{ub}-\tilde{A}_{(\cdot,1)}x\right)=\hat{\gamma}_l^{*\prime}A\sqrt{n_l}(\hat{\beta}_{n_l}-\beta_{P_l})-\hat{\gamma}_l^{*\prime}\tilde{A}_{(\cdot,1)}x.$$

From Assumption 2 combined with Slutsky's lemma, we have that

$$\hat{\gamma}_l^{*\prime}A\sqrt{n_l}(\hat{\beta}_{n_l}-\beta_{P_l})-\hat{\gamma}_l^{*\prime}\tilde{A}_{(\cdot,1)}x\rightarrow_d \mathcal{N}\left(-c_1(\Sigma^*)\bar{\gamma}_1'\tilde{A}_{(\cdot,1)}x,\ c_1(\Sigma^*)^2\bar{\gamma}_1'A\Sigma^*A'\bar{\gamma}_1\right).$$

Now, consider $\hat{\gamma}_{l,j} = c_j(\hat{\Sigma}_{n_l})\bar{\gamma}_j$ for $j \neq 1$. By construction $\bar{\gamma}_j \ge 0$, and Lemma A.4 implies that $\bar{\gamma}_j$ has a non-zero element in at least one component in B^* . But this, combined with equations (22) and (23) and the fact that $c_j(\hat{\Sigma}_{n_l}) \rightarrow_p c_j(\Sigma^*) > 0$, implies that

$$\hat{\gamma}_{l,j}'\left(\sqrt{n_l}A\beta_{P_l}-\sqrt{n_l}d-\sqrt{n_l}\tilde{A}_{(\cdot,1)}\theta_{P_l}^{ub}-\tilde{A}_{(\cdot,1)}x\right)\to_p-\infty,$$

and thus

$$\hat{\gamma}_{l,j}'\left(\sqrt{n_l}A\hat{\beta}_{n_l}-\sqrt{n_l}d-\sqrt{n_l}\tilde{A}_{(\cdot,1)}\theta_{P_l}^{ub}-\tilde{A}_{(\cdot,1)}x\right)\to_p-\infty,$$

as well, since as before $\hat{\gamma}'_{l,j}A\sqrt{n}(\hat{\beta}_{n_l}-\beta_{P_l})$ converges in distribution to a normal distribution with finite variance. This implies that $\hat{\gamma}^*_l$ is the optimizer of the problem for $\hat{\eta}_l$ with probability approaching 1, and thus

$$\hat{\eta}_l \to_d \mathcal{N}\left(-c_1(\Sigma^*)\bar{\gamma}_1'\tilde{A}_{(\cdot,1)}x, c_1(\Sigma^*)^2\bar{\gamma}_1'A\Sigma^*A'\bar{\gamma}_1\right).$$

This also implies that for any $j \neq 1$, $|\hat{\eta}_l - \hat{\gamma}'_{l,j}\tilde{Y}_l| \rightarrow_p \infty$, where $\tilde{Y}_l = \sqrt{n_l}A\hat{\beta}_{n_l} - \sqrt{n_l}d - \sqrt{n_l}\tilde{A}_{(\cdot,1)}\theta_{P_l}^{ub} - \tilde{A}_{(\cdot,1)}x$. Since there are a finite number of vertices, it follows that $\min_{j\neq 1} |\hat{\eta}_l - \hat{\gamma}'_{l,j}\tilde{Y}_l| \rightarrow -\infty$. This together with the result of Lemma A.3 implies that $|\hat{\eta}_l - v_l^{lo}| \rightarrow_p \infty$ and $|\hat{\eta}_l - v_l^{up}| \rightarrow_p \infty$, where v_l^{lo}, v_l^{up} are the values of v^{lo}, v^{up} associated with the $\psi^C_{\alpha}(\sqrt{n_l}\hat{\beta}_{n_l}, A, \sqrt{n_l}d, \sqrt{n_l}\theta_{P_l}^{ub} + x, \hat{\Sigma}_{n_l})$ test. Since $\hat{\eta}_l$ is stochastically bounded, and by construction $v_l^{lo} \leq \hat{\eta}_l \leq v_l^{up}$, it follows that $v_l^{lo} \rightarrow_p -\infty$ and $v_l^{up} \rightarrow_p \infty$. Let $\hat{\sigma}_l^2 = \gamma'_{*,l}A\hat{\Sigma}_{n_l}A'\gamma_{*,l}$ denote the variance at the optimal vertex used by the $\psi^C_{\alpha}(\sqrt{n_l}\hat{\beta}_{n_l}, A, \sqrt{n_l}d, \sqrt{n_l}\theta_{P_l}^{ub} + x, \hat{\Sigma}_{n_l})$ test. Since, we've shown that $\hat{\gamma}_l^*$ is optimal w.p.a. 1, we have that $\hat{\sigma}_l^2 \rightarrow_p c_1(\Sigma^*)^2 \tilde{\gamma}'_1 A \Sigma^* A' \tilde{\gamma}_1$. From another application of the continuous mapping theorem, we have that

$$\frac{\Phi(\hat{\eta}_l/\hat{\sigma}_l) - \Phi(v_l^{lo}/\hat{\sigma}_l)}{\Phi(v_l^{up}/\hat{\sigma}_l) - \Phi(v_l^{lo}/\hat{\sigma}_l)} \to_d \frac{\Phi(\xi) - \Phi(-\infty)}{\Phi(\infty) - \Phi(-\infty)} = \Phi(\xi),$$

where $\xi \sim \mathcal{N}\left(-\bar{\gamma}_1'\tilde{A}_{(\cdot,1)}x/\sqrt{\bar{\gamma}_1'A\Sigma^*A\bar{\gamma}_1}, 1\right)$. The limiting distribution is continuous, and thus

$$\mathbb{P}_{P_l}\left(\frac{\Phi(\hat{\eta}_l/\hat{\sigma}_l) - \Phi(v_l^{lo}/\hat{\sigma}_l)}{\Phi(v_l^{up}/\hat{\sigma}_l) - \Phi(v_l^{lo}/\hat{\sigma}_l)} > 1 - \alpha\right) \to \mathbb{P}\left(\Phi(\xi) > 1 - \alpha\right) = \Phi\left(\frac{-\bar{\gamma}_1\tilde{A}_{(\cdot,1)}x}{\sqrt{\bar{\gamma}_1A\Sigma^*A'\bar{\gamma}_1}} - z_{1-\alpha}\right).$$

Moreover, for $\alpha < 0.5$, z^{lo} sufficiently small, and z^{up} sufficiently large, $(\Phi(\hat{\eta}_l) - \Phi(z^{lo}))/(\Phi(z^{up}) - \Phi(z^{lo})) > 1 - \alpha$ only if $\hat{\eta}_l > 0$. It follows that

$$\mathbb{P}_{P_l}\left(\frac{\Phi(\hat{\eta}_l/\hat{\sigma}_l) - \Phi(v_l^{lo}/\hat{\sigma}_l)}{\Phi(v_l^{up}/\hat{\sigma}_l) - \Phi(v_l^{lo}/\hat{\sigma}_l)} > 1 - \alpha, \hat{\eta}_l > 0\right) \to \Phi\left(\frac{-\bar{\gamma}_1\tilde{A}_{(\cdot,1)}x}{\sqrt{\bar{\gamma}_1A\Sigma^*A'\bar{\gamma}_1}} - z_{1-\alpha}\right).$$

However, the event in the previous display is precisely the event that $\psi_{\alpha}^{C}(\sqrt{n_{l}}\hat{\beta}_{l}, A, \sqrt{n_{l}}d, \sqrt{n_{l}}\theta_{P}^{ub} + x, \hat{\Sigma}_{n_{l}}) = 1$, and thus

$$\mathbb{E}_{P_l}\left[\psi^C_{\alpha}(\sqrt{n_l}\hat{\beta}_l, A, \sqrt{n_l}d, \sqrt{n_l}\theta^{ub}_P + x, \hat{\Sigma}_{n_l})\right] \to \Phi\left(\frac{-\bar{\gamma}_1\tilde{A}_{(\cdot,1)}x}{\sqrt{\bar{\gamma}_1A\Sigma^*A'\bar{\gamma}_1}} - z_{1-\alpha}\right).$$

The result is then immediate from the previous display combined with (24).

Lemma A.4. If LICQ holds in direction l at β_P , then there exists a unique $\bar{\gamma} \in V(\Sigma_P)$ such that $\bar{\gamma}_{-B^*} = 0$, where B^* is the set of binding moments at the optimum to (16).

Proof. We first show that there is at most one such $\bar{\gamma}$. By definition, any vertex $\gamma \in V(\Sigma_P)$ satisfies $\gamma' \tilde{A}_{(\cdot,-1)} = 0$. Recall that $\tilde{A} = A_{(\cdot,post)}\Gamma^{-1}$, where Γ is full rank. LICQ implies that $A_{(B^*,post)}$ has full row rank, and thus so does $\tilde{A}_{(B^*,\cdot)}$. It follows that $\tilde{A}_{(B^*,-1)}$ has rank at

least $|B^*| - 1$. If the rank is $|B^*|$, then there are no non-zero solutions to $\gamma'_{B^*}\tilde{A}_{(B^*,-1)} = 0$, and thus there are no vertices with $\gamma_{-B^*} = 0$. If the rank is $|B^*| - 1$, then any solution to $\gamma'\tilde{A}_{(\cdot,-1)} = 0$ with $\gamma_{-B^*} = 0$ takes the form $\gamma_{B^*} = c \cdot \nu$ for some constant c and ν a vector the generates the one-dimensional nullspace of $\tilde{A}_{(B^*,-1)}$. However, any $\gamma \in V(\Sigma_P)$ also must satisfy $\gamma'\tilde{\sigma} = 1$, which uniquely pins down the constant c. Thus, there is at most one element of the feasible set with $\gamma_{-B^*} = 0$.

We next show that there exists such a $\bar{\gamma}$. Consider the optimization $\eta(\beta_P, A, d, \theta_P^{ub}, \Sigma_P)$ for $\eta(\cdot)$ defined in (21). As argued in the proof to Proposition 3.2, since θ_P^{ub} is on the boundary of the identified set, we must have $\eta(\beta_P, A, d, \theta_P^{ub}, \Sigma_P) = 0$. However, LICQ implies that there exists a value $\tilde{\tau}^*$ such that

$$A_{(B^*,\cdot)}\beta_P - d_{B^*} - \tilde{A}_{(B^*,1)}\theta_P^{ub} - \tilde{A}_{(B^*,-1)}\tilde{\tau}^* = 0$$

$$A_{(-B^*,\cdot)}\beta_P - d_{-B^*} - \tilde{A}_{(-B^*,1)}\theta_P^{ub} - \tilde{A}_{(-B^*,1)}\tilde{\tau}^* < 0.$$

In particular, this holds for $\tilde{\tau}^* = \Gamma_{(-1,\cdot)}\tau^*$. It follows that $(\eta, \tilde{\tau}) = (0, \tilde{\tau}^*)$ is a solution to $\eta(\beta_P, A, d, \theta_P^{ub}, \Sigma_P)$. By duality, there is some $\bar{\gamma} \in V(\Sigma_P)$ that is a Lagrange multiplier for this optimization problem. The complementary slackness conditions imply, however, that $\bar{\gamma}_{-B^*} = 0$, as needed.

Lemma A.5. Suppose $\hat{\beta} \sim \mathcal{N}(\beta, \Sigma)$ for Σ known. Let B_0 be a closed, convex set. Then the most-powerful size α test of $H_0 : \beta \in B_0$ against the point alternative $H_A : \beta = \beta_A$ is equivalent to the most powerful test of $H_0 : \beta = \tilde{\beta}$ against $H_A : \beta = \beta_A$, where $\tilde{\beta} =$ $\arg \min_{\beta \in B_0} ||\beta - \beta_A||_{\Sigma}$ and $||\cdot||_{\Sigma}$ is the Mahalanobis norm in Σ , $||x||_{\Sigma} = \sqrt{x'\Sigma^{-1}x}$. The most powerful test rejects for values of $(\beta_A - \tilde{\beta})'\Sigma^{-1}\hat{\beta}$ greater than $(\beta_A - \tilde{\beta})'\Sigma^{-1}\tilde{\beta} + z_{1-\alpha}||\beta_A - \tilde{\beta}||_{\Sigma}$, and has power against the alternative of $\Phi(||\beta_A - \tilde{\beta}||_{\Sigma} - z_{1-\alpha})$, for $z_{1-\alpha}$ the $1 - \alpha$ quantile of the standard normal.

Proof. Define $\langle \cdot, \cdot \rangle_{\Sigma}$ by $\langle x, y \rangle_{\Sigma} = x' \Sigma^{-1} y$, and observe that $\langle \cdot, \cdot \rangle_{\Sigma}$ is an inner product. The result then follows immediately from the discussion in Section 2.4.3 of Ingster and Suslina (2003), replacing all instances of the usual euclidean inner product with $\langle \cdot, \cdot \rangle_{\Sigma}$.

Lemma A.6. Let \mathcal{B} be a closed, convex subset of \mathbb{R}^K , and $\beta_A \notin \mathcal{B}$. Let $\tilde{\beta} = \arg \min_{\beta \in \mathcal{B}} ||\beta - \beta_A||_{\Sigma}$, where $||x||_{\Sigma}^2 = x' \Sigma^{-1} x$ for some positive definite matrix Σ . Then for any $\beta \in \mathcal{B}$, $(\tilde{\beta} - \beta_A)' \Sigma^{-1} (\beta - \tilde{\beta}) \ge 0$.

Proof. Consider any $\beta \in \mathcal{B}$. Define $\beta_{\theta} = \theta(\beta - \tilde{\beta}) + \tilde{\beta}$, and note that since \mathcal{B} is convex $\beta_{\theta} \in \mathcal{B}$

for any $\theta \in [0, 1]$. Further,

$$||\beta_{\theta} - \beta_A||_{\Sigma}^2 = \theta^2 ||\beta - \tilde{\beta}||_{\Sigma}^2 + 2\theta(\tilde{\beta} - \beta_A)' \Sigma^{-1}(\beta - \tilde{\beta}) + ||\tilde{\beta} - \beta_A||_{\Sigma}^2.$$

Differentiating with respect to θ , we have

$$\frac{\partial}{\partial \theta} ||\beta_{\theta} - \beta_A||_{\Sigma}^2 = 2\theta ||\beta - \tilde{\beta}||_{\Sigma}^2 + 2(\tilde{\beta} - \beta_A)' \Sigma^{-1} (\beta - \tilde{\beta}),$$

from which we see that the derivative evaluated at $\theta = 0$ is $2(\tilde{\beta} - \beta_A)'\Sigma^{-1}(\beta_A - \tilde{\beta})$. Since $\tilde{\beta}$ minimizes the norm, it follows that we must have $2(\tilde{\beta} - \beta_A)'\Sigma^{-1}(\beta_A - \tilde{\beta}) \ge 0$, else we could achieve a lower value of the norm at β_{θ} by choosing $\theta > 0$ sufficiently small.

Lemma A.7. Let $\mathcal{B} = \{\beta \in \mathbb{R}^K : v'\beta \leq d\}$ for some $v \in \mathbb{R}^K \setminus \{0\}$ and $d \in \mathbb{R}$. Let $\tilde{\beta} = \arg \min_{\beta \in \mathcal{B}} ||\beta - \beta_A||_{\Sigma}$ for some $\beta_A \notin \mathcal{B}$, where $||x||_{\Sigma}^2 = x'\Sigma^{-1}x$ and Σ is positive definite. Then $(\beta_A - \tilde{\beta})'\Sigma^{-1} = c \cdot v'$ for the positive constant $c = \frac{v'\beta_A - d}{v'\Sigma v}$.

Proof. Note that we can form a basis $v, \tilde{v}_2, ..., \tilde{v}_K$ such that $v'\tilde{v}_j = 0$ for j = 2, ..., K. It follows by construction that for any j = 2, ..., K and any $t \in \mathbb{R}, \tilde{\beta} + t \cdot \tilde{v}_j \in \mathcal{B}$. Hence, from Lemma A.6, $-(\beta_A - \tilde{\beta})'\Sigma^{-1}(t\tilde{v}_j) \ge 0$. Since we can choose t both positive and negative, it follows that $(\beta_A - \tilde{\beta})'\Sigma^{-1}\tilde{v}_j = 0$ for all j. Since $(\beta_A - \tilde{\beta})'\Sigma^{-1}$ is orthogonal to $\{\tilde{v}_2, ..., \tilde{v}_K\}$, and $\{v, \tilde{v}_2, ..., \tilde{v}_K\}$ form a basis, we have that $(\beta_A - \tilde{\beta})'\Sigma^{-1} = c \cdot v'$, for some $c \in \mathbb{R}$. Multiplying both sides of the equation on the right by Σv , we obtain that $(\beta_A - \tilde{\beta})'v = c \cdot v'\Sigma v$. However, since $\tilde{\beta}$ is the closest point to β_A in Mahalanobis distance, it must be on the boundary of \mathcal{B} , and so $v'\tilde{\beta} = d$. It follows that $c = (v'\beta_A - d)/(v'\Sigma v)$, which is clearly positive since $\beta_A \notin \mathcal{B}$ and thus $v'\beta_A > d$.

Lemma A.8 (Power of optimal test for linear subspace). Let $\mathcal{B} = \{\beta \in \mathbb{R}^K : v'\beta \leq d\}$ for some $v \in \mathbb{R}^K \setminus \{0\}$ and $d \in \mathbb{R}$. Suppose $\hat{\beta} \sim \mathcal{N}(\beta, \Sigma)$ for Σ positive definite known, and consider the problem of testing $H_0 : \beta \in \mathcal{B}$ against $H_A : \beta = \beta_A$ for some $\beta_A \notin \mathcal{B}$. Then the most powerful size- α test of H_0 against H_A is a one-sided t-test that rejects for large values of $v'\hat{\beta}$, and has power equal to $\Phi((v'\beta_A - d)/\sqrt{v'\Sigma v} - z_{1-\alpha})$.

Proof. From Lemma A.5, the most powerful test rejects for large values of $(\beta_A - \tilde{\beta})' \Sigma^{-1} \hat{\beta}$, where $\tilde{\beta} = \arg \min_{\beta \in \mathcal{B}} ||\beta - \beta_A||_{\Sigma}$, and has power $\Phi(||\beta_A - \tilde{\beta}||_{\Sigma} - z_{1-\alpha})$. By Lemma A.7, $(\beta_A - \tilde{\beta})' \Sigma^{-1} = cv'$, for $c = (v'\beta_A - d)/(v'\Sigma v)$. It follows that

$$||\beta_A - \beta||_{\Sigma}^2 = (\beta_A - \beta)' \Sigma^{-1} (\beta_A - \beta)$$
$$= cv'(\beta_A - \tilde{\beta})$$
$$= c(v'\beta_A - d) = (v'\beta_A - d)^2 / (v'\Sigma v),$$

where we use the fact that $v'\tilde{\beta} = d$, since $\tilde{\beta}$ must be on the boundary of \mathcal{B} , as argued in the proof to Lemma A.7. The result then follows immediately.

Lemma A.9. If $P \in \mathcal{P}_{\epsilon}$, then $\rho_{\alpha}^{*}(P, x) = \Phi\left(\frac{-\bar{\gamma}'\tilde{A}_{(\cdot,1)}x}{\sqrt{\bar{\gamma}'A\Sigma_{P}A'\bar{\gamma}}} - z_{1-\alpha}\right)$, where $\bar{\gamma} \in V(\Sigma_{P})$ is the unique element of $V(\Sigma_{P})$ with $\bar{\gamma}_{-B^{*}} = 0$ (see Lemma A.4).

Proof. Suppose $\hat{\beta}_n \sim \mathcal{N}(\beta_P, \Sigma_P/n)$. Let $\mathcal{B}_n = \{\beta : \theta_P^{ub} + x/\sqrt{n} \in \mathcal{S}(\beta, \Delta)\}$ be the set of values for β consistent with the null that $\theta = \theta_P^{ub} + x/\sqrt{n}$. Observe that $\mathcal{B}_n = \{\beta : \eta(\beta, A, d, \theta^{ub} + x/\sqrt{n}, \Sigma_P) \leq 0\}$, where $\eta(\cdot)$ is defined in (21). From Lemma A.5, the most powerful test of $H_0: \beta \in \mathcal{B}_n$ against $H_1: \beta = \beta_P$ rejects for large values of $(\beta_P - \tilde{\beta})'\Sigma_P^{-1}\hat{\beta}_n$. To derive the optimal test, it is instructive to first consider a simpler testing problem. From Lemma A.4, there exists a unique $\bar{\gamma} \in V(\Sigma_P)$ such that $\bar{\gamma}_{-B^*} = 0$, where B^* are the binding rows at the solution to (16) satisfying LICQ. Define $\mathcal{B}_n^{\bar{\gamma}} = \{\beta : \bar{\gamma}'(A\beta - d - \tilde{A}_{(\cdot,1)}(\theta_P^{ub} + x/\sqrt{n})) \leq 0\}$. We first consider testing $\tilde{H}_0: \beta \in \mathcal{B}_n^{\bar{\gamma}}$ against $H_1: \beta = \beta_P$. From Lemma A.7, the optimal test rejects for large values of $\bar{\gamma}'A\hat{\beta}_n$ and has power $\Phi(\frac{h}{\sqrt{\bar{\gamma}'A\Sigma_PA'\bar{\gamma}/n}} - z_{1-\alpha})$, where

$$h = \bar{\gamma}' (A\beta_P - d - \tilde{A}_{(\cdot,1)}(\theta_P^{ub} + x/\sqrt{n})).$$
⁽²⁵⁾

From the definition of LICQ in direction l, however, there exists a value $\tilde{\tau}^*$ such that

$$A_{(B^*,\cdot)}\beta_P - d_{B^*} - \tilde{A}_{(B,1)}\theta_P^{ub} - \tilde{A}_{(B,-1)}\tilde{\tau}^* = 0$$
(26)

$$A_{(-B^*,\cdot)}\beta_P - d_{-B^*} - \tilde{A}_{(-B,1)}\theta_P^{ub} - \tilde{A}_{(-B,1)}\tilde{\tau}^* < 0$$
(27)

By construction, $\bar{\gamma}'\tilde{A}_{(\cdot,-1)} = 0$ and $\bar{\gamma}_{-B^*} = 0$, which combined with the previous two displays implies that $h = -\bar{\gamma}'\tilde{A}_{(\cdot,1)}x/\sqrt{n}$, and hence the power of the optimal test of \tilde{H}_0 is $\Phi\left(\frac{-\bar{\gamma}'\tilde{A}_{(\cdot,1)}x}{\sqrt{\bar{\gamma}'A\Sigma_P A'\bar{\gamma}}} - z_{1-\alpha}\right)$.

To complete the proof, it thus suffices to show that the optimal test of \tilde{H}_0 against H_1 is the same as the optimal test of H_0 against H_1 for n sufficiently large. To this end, note that $\mathcal{B}_n \subseteq \mathcal{B}_n^{\tilde{\gamma}}$, since by duality,

$$\eta(\beta, A, d, \theta^{ub} + x/\sqrt{n}, \Sigma_P) = \max_{\gamma \in V(\Sigma_P)} \gamma'(A\beta - d - \tilde{A}_{(\cdot, 1)}(\theta_P^{ub} + x/\sqrt{n})) \ge \bar{\gamma}'(A\beta - d - \tilde{A}_{(\cdot, 1)}(\theta_P^{ub} + x/\sqrt{n}))$$

Thus, Lemma A.5 implies that the optimal test under H_0 coincides with the optimal test under H_1 whenever $\tilde{\beta}_n = \arg \min_{\beta \in \mathcal{B}_n^{\bar{\gamma}}} ||\beta - \beta_P||_{\Sigma_P/n}$ is in \mathcal{B}_n . From Lemma A.7, however, $\tilde{\beta}'_n = \beta'_P + \frac{h}{\sqrt{\bar{\gamma}' A \Sigma_P A' \bar{\gamma}/n}} v'(\Sigma_P/n)$, for h defined in (25). Using the equality $h = -\bar{\gamma}' \tilde{A}_{(\cdot,1)} x/\sqrt{n}$ derived above, we see that

$$\tilde{\beta}'_n = \beta'_P - \frac{1}{\sqrt{n}} \frac{\bar{\gamma}' A_{(\cdot,1)} x}{\sqrt{\bar{\gamma}' A \Sigma_P A' \bar{\gamma}}} \bar{\gamma}' A \Sigma_P,$$

and thus we can write $\tilde{\beta} = \beta_P - \nu/\sqrt{n}$ for a finite vector ν . From Lemma A.4, every $\gamma \in V(\Sigma_P)$ with $\gamma \neq \bar{\gamma}$ has $\gamma_{-B^*} \neq 0$. Since $\gamma \ge 0$ by construction, equations (26) and (27) imply that

$$\gamma'(A\beta_P - d - \tilde{A}_{(\cdot,1)}\theta_P^{ub}) < 0$$

for all $\gamma \neq \bar{\gamma}$, where we use the fact that $\gamma' \tilde{A}_{(\cdot,-1)} = 0$ by construction. We've shown, however, that $\bar{\gamma}' (A\beta_P - d - \tilde{A}_{(\cdot,1)}\theta_P^{ub}) = 0$. By continuity arguments, it follows that for *n* sufficiently large,

$$\eta(\tilde{\beta}, A, d, \theta_P^{ub} + x/\sqrt{n}, \Sigma_P) = \max_{\gamma \in V(\Sigma_P)} \gamma'(A(\beta_P - \nu/\sqrt{n}) - d - \tilde{A}_{(\cdot, 1)}(\theta_P^{ub} + x/\sqrt{n}))$$

is equal to

$$\bar{\gamma}'(A(\beta_P - \nu/\sqrt{n} - d - \tilde{A}_{(\cdot,1)}(\theta_P^{ub} + x/\sqrt{n}))),$$

and thus $\tilde{\beta}_n \in \mathcal{B}_n$, as we wished to show.

A.3 Proofs and auxiliary lemmas for FLCIs

Proof of Proposition 4.2

Proof. First, suppose Assumption 9 holds. Without loss of generality, we show $\mathbb{P}\left((\theta^{ub} + x) \in \mathcal{C}_{\alpha,n}^{FLCI}\right) \to 0$ for any x > 0. By Lemma A.11 there exists (\bar{a}, \bar{v}) such that $\bar{b}(\bar{a}, \bar{v}) = \frac{1}{2}LID(\delta_{pre}, \Delta) =: \bar{b}_{min}$ and $\mathbb{E}_{\hat{\beta}_n \sim \mathcal{N}(\delta + \tau, \Sigma_n)}\left[\bar{a} + \bar{v}'\hat{\beta}_n\right] = \frac{1}{2}(\theta^{ub} + \theta^{lb}) =: \theta^{mid}$. Let $\bar{\mathcal{C}}_n := \bar{a} + \bar{v}'\hat{\beta}_n \pm \chi_n(\bar{a}, \bar{v})$ denote the fixed length confidence interval based on (\bar{a}, \bar{v}) .

By construction, $\bar{\chi}_n := \chi_n(\bar{a}, \bar{v})$ is the $1 - \alpha$ quantile of the $|\mathcal{N}(\bar{b}_{min}, \sigma_{\bar{v},n}^2)|$ distribution. Since $\sigma_{\bar{v},n}^2 = \frac{1}{n}\sigma_{\bar{v},1}^2 \to 0$, the $|\mathcal{N}(\bar{b}_{min}, \sigma_{\bar{v},n}^2)|$ distribution collapses to a point mass at \bar{b}_{min} , and thus $\bar{\chi}_n \to \bar{b}_{min}$. By construction, the half-length of the shortest FLCI $\chi_n := \chi_n(a_n, v_n)$ must be less than or equal to $\bar{\chi}_n$, and so $\limsup_{n\to\infty}\chi_n \leq \bar{b}_{min}$. Let $b_n := \bar{b}(a_n, v_n)$ be the worst-case bias of the optimal FLCI. Since $\alpha \in (0, 0.5]$, Lemma A.12 implies that $\chi_n \geq b_n$. Additionally, Lemma A.10 implies that $b_n \geq \frac{1}{2}LID(\delta_{pre}, \Delta) = \bar{b}_{min}$, and thus $\chi_n \geq \bar{b}_{min}$. Hence, $\chi_n \to \bar{b}_{min}$ implies $b_n \to \bar{b}_{min}$. Additionally, note that for $\alpha \in (0, 0.5]$, $\chi_n(a, v)$ is increasing in both $\bar{b}(a, v)$ and $\sigma_{v,n}$. Since $\bar{b}_{min} \leq b_n$ and $\chi_n \leq \bar{\chi}_n$, it must be that $\sigma_{v_n,n} \leq \sigma_{\bar{v},n}$, from which it follows that $\sigma_{v_n,n} \to 0$.

Now, we claim that $\mu_n := \mathbb{E}_{\hat{\beta}_n \sim \mathcal{N}(\delta + \tau, \Sigma_n)} \left[a_n + v'_n \hat{\beta}_n \right]$ converges to $\theta^{mid} := \frac{1}{2} (\theta^{ub} + \theta^{lb})$. To show this, note that $\mu_n = a_n + v'_n \beta$ for $\beta = \delta + \tau$. Since $\theta^{ub}, \theta^{lb} \in \mathcal{S}(\beta, \Delta)$, by the definition of the identified set there exist $\delta^{ub}, \delta^{lb} \in \Delta$ and τ^{ub}, τ^{lb} such that $\beta = \delta^{ub} + \tau^{ub} = \delta^{lb} + \tau^{lb}, \theta^{ub} = l'\tau^{ub}_{post}$, and $\theta^{lb} = l'\tau^{lb}_{post}$. Thus, $\theta^{ub} - \mathbb{E}_{\hat{\beta}_n \sim \mathcal{N}(\beta, \Sigma_n)} \left[a_n + v'_n \hat{\beta}_n \right] = \theta^{ub} - \mu_n$ and $\mathbb{E}_{\hat{\beta}_n \sim \mathcal{N}(\beta, \Sigma_n)} \left[a_n + v'_n \hat{\beta}_n \right] - \theta^{lb} = \mu_n - \theta^{lb}$. This implies that $b_n \ge \max\{\theta^{ub} - \mu_n, \mu_n - \theta^{lb}\} = \bar{b}_{min} + |\mu_n - \theta^{mid}|$, where the equality uses the fact that $\theta^{ub} - \theta^{lb} = LID(\delta_{A,pre}, \Delta) = 2\bar{b}_{min}$. Since we've shown that $b_n \to \bar{b}_{min}$, it follows that $\mu_n \to \theta^{mid}$, as desired.

Next, note that if $\hat{\beta}_n \sim \mathcal{N}(\delta + \tau, \Sigma_n)$, then $a_n + v'_n \hat{\beta}_n \sim \mathcal{N}(\mu_n, \sigma^2_{v_n,n})$. Observe that $\bar{\theta} \in \mathcal{C}_{\alpha,n}^{FLCI}$ if and only if $a_n + v'_n \hat{\beta}_n \in [\bar{\theta} - \chi_n, \bar{\theta} + \chi_n]$. Thus,

$$\mathbb{P}_{\hat{\beta}_n \sim \mathcal{N}(\beta, \Sigma_n)} \left(\bar{\theta} \in \mathcal{C}_{\alpha, n}^{FLCI} \right) = \Phi \left(\frac{\bar{\theta} + \chi_n - \mu_n}{\sigma_{v_n, n}} \right) - \Phi \left(\frac{\bar{\theta} - \chi_n - \mu_n}{\sigma_{v_n, n}} \right)$$

Now, recalling that $\theta^{ub} = \theta^{mid} + \bar{b}_{min}$ by construction, we have $\mathbb{P}_{\hat{\beta}_n \sim \mathcal{N}(\beta, \Sigma_n)} \left((\theta^{ub} + x) \in \mathcal{C}_{\alpha, n}^{FLCI} \right)$ equals

$$\Phi\left(\frac{\theta^{mid} + \bar{b}_{min} + x + \chi_n - \mu_n}{\sigma_{v_n,n}}\right) - \Phi\left(\frac{\theta^{mid} + \bar{b}_{min} + x - \chi_n - \mu_n}{\sigma_{v_n,n}}\right).$$
 (28)

Note that the term inside the second normal CDF in the previous display equals

$$-\frac{\chi_n - b_n}{\sigma_{v_n,n}} + \frac{x + \theta^{mid} - \mu_n + \bar{b}_{min} - b_n}{\sigma_{v_n,n}}.$$

However, the first summand above is bounded between $-z_{1-\alpha/2}$ and $-z_{1-\alpha}$ by Lemma A.12. Additionally, we've shown that $\theta^{mid} - \mu_n \to 0$ and $\bar{b}_{min} - b_n \to 0$, so the numerator of the second summand converges to x > 0. Since the denominator $\sigma_{v_n,n} \to 0$, the expression in the previous display diverges to ∞ , and hence the second normal CDF term in (28) converges to 1, which implies that $\mathbb{P}\left((\theta^{ub} + x) \in \mathcal{C}_{\alpha,n}^{FLCI}\right) \to 0$, as needed.

In order to prove the other direction, we proceed via the contrapositive. Towards this, suppose Assumption 9 fails. Let $L := LID(\delta_{pre}, \Delta)$ and $\bar{L} := \sup_{\bar{\delta}_{pre} \in \Delta_{pre}} LID(\bar{\delta}_{pre}, \Delta)$. By Lemma A.10, $b_n := \bar{b}(a_n, v_n) \geq \frac{1}{2}\bar{L} =: \bar{b}_{min}$. As argued earlier in the proof, since $\alpha \in (0, .5]$, $\chi_n \geq b_n \geq \frac{1}{2}\bar{L}$. If $\bar{L} = \infty$, then $\mathcal{C}_{\alpha,n}^{FLCI}$ is the real line, and thus never rejects, so $\mathcal{C}_{\alpha,n}^{FLCI}$ is trivially inconsistent under the assumption that $\mathcal{S}(\delta + \tau, \Delta) \neq \mathbb{R}$. For the remainder of the proof, we assume $L < \bar{L} < \infty$. From Lemma 2.1, $\mathcal{S}(\delta + \tau, \Delta) = [\theta^{lb}, \theta^{ub}]$, where $\theta^{ub} - \theta^{lb} = LID(\delta_{pre}, \Delta) = L$. Let $\epsilon = \frac{1}{4}(\bar{L} - L)$, and set $\theta_1^{out} := \theta^{ub} + \epsilon$ and $\theta_2^{out} := \theta^{lb} - \epsilon$. Let $\theta^{mid} = \frac{1}{2}(\theta^{ub} + \theta^{lb})$ be the midpoint of the identified set. By construction, $\theta_1^{out} - \theta^{mid} = \theta^{mid} - \theta_2^{out} = \frac{1}{2}L + \epsilon < \frac{1}{2}\bar{L}$. Since $\mathcal{C}_{\alpha,n}^{FLCI}$ is an interval with half-length at least $\frac{1}{2}\bar{L}$, it follows that if $\theta^{mid} \in \mathcal{C}_{\alpha,n}^{FLCI}$ then at least one of $\theta_1^{out}, \theta_2^{out}$ is also in $\mathcal{C}_{\alpha,n}^{FLCI}$. Hence, $\mathbb{P}(\theta_1^{out} \in \mathcal{C}_{\alpha,n}^{FLCI}) + \mathbb{P}(\theta_2^{out} \in \mathcal{C}_{\alpha,n}^{FLCI}) \geq \mathbb{P}(\theta^{mid} \in \mathcal{C}_{\alpha,n}^{FLCI}) \geq 1 - \alpha$, where the final bound follows since $\mathcal{C}_{\alpha,n}^{FLCI}$.

satisfies the coverage requirement (10). It follows that $\limsup_{n\to\infty} \mathbb{P}\left(\theta_j^{out} \in \mathcal{C}_{\alpha,n}^{FLCI}\right) \ge \frac{1}{2}(1 - \alpha) > 0$ for at least one $j \in \{1, 2\}$.

Lemma A.10 (Bounds for worst-case bias). For any (a, v), $\bar{b}(a, v) \ge \frac{1}{2} \sup_{\delta_{pre} \in \Delta_{pre}} LID(\delta_{pre}, \Delta)$. *Proof.* Since $\beta = \delta + \tau$, we can write the bias of the affine estimator $a + v'\hat{\beta}$ as $b = a + v'\delta + (v_{post} - l)'\tau_{post}$. Since τ_{post} is unrestricted in the maximization in (17), we see that the worst-case bias will be infinite if $v_{post} \neq l$ and the lemma holds trivially. We can thus restrict attention to affine estimators with $v_{post} = l$, in which case the worst-case bias reduces to

$$\bar{b}(a,v) = \sup_{\delta \in \Delta} |a + v'\delta| = \sup_{\delta \in \Delta} |a + v'_{pre}\delta_{pre} + l'\delta_{post}|.$$
(29)

Now, pick any $\delta_{pre}^* \in \Delta_{pre}$. First, suppose that the minimum $(\min_{\delta} l' \delta_{post}, \text{ s.t. } \delta \in \Delta, \delta_{pre} = \delta_{pre}^*)$ and the maximum $(\max_{\delta} l' \delta_{post}, \text{ s.t. } \delta \in \Delta, \delta_{pre} = \delta_{pre}^*)$ are finite. Let δ^{min} and δ^{max} be the associated solutions. By construction, $\delta_{pre}^{max} = \delta_{pre}^{min} = \delta_{pre}^*$. For any v_{pre} , we can apply the triangle inequality to show that

$$\begin{aligned} \left|a + v'_{pre}\delta^{max}_{pre} + l'\delta^{max}_{post}\right| + \left|a + v'_{pre}\delta^{min}_{pre} + l'\delta^{min}_{post}\right| &\ge \left|\left(a + v'_{pre}\delta^{max}_{pre} + l'\delta^{max}_{post}\right) - \left(a + v'_{pre}\delta^{min}_{pre} + l'\delta^{min}_{post}\right)\right| \\ &= \left|l'\delta^{max}_{post} - l'\delta^{min}_{post}\right| = LID(\delta^*_{pre}, \Delta).\end{aligned}$$

Note that for any $x_1, x_2 \ge 0$, $\max\{x_1, x_2\} \ge \frac{1}{2}(x_1 + x_2)$. It then follows from the previous display that

$$\max\{\left|a+v_{pre}^{\prime}\delta_{pre}^{max}+l^{\prime}\delta_{post}^{max}\right|,\left|a+v_{pre}^{\prime}\delta_{pre}^{min}+l^{\prime}\delta_{post}^{min}\right|\} \ge \frac{1}{2}LID(\delta_{pre}^{*},\Delta).$$

Since δ_{pre}^{max} and δ_{pre}^{min} are feasible in the maximization (29), we see that $\bar{b} \ge \frac{1}{2}LID(\delta_{pre}^*, \Delta)$, as needed. To complete the proof, now suppose without loss of generality that

$$\left(\max_{\delta} l' \delta_{post}, \text{ s.t. } \delta \in \Delta, \delta_{pre} = \delta_{pre}^*\right) = \infty$$

Then, we can replay the argument above replacing δ^{max} with a sequence of values $\{\delta_j\}$ such that $l'\delta_j$ diverges, which gives that \bar{b} is infinite and the result follows.

Lemma A.11. Suppose Δ is convex, and there exists $\delta \in \Delta$ such that $LID(\delta_{pre}, \Delta) = \sup_{\tilde{\delta}_{pre} \in \Delta_{pre}} LID(\tilde{\delta}_{pre}, \Delta) < \infty$. Then there exists (a, v) such that $\bar{b}(a, v) = \frac{1}{2} \sup_{\tilde{\delta}_{pre} \in \Delta_{pre}} LID(\tilde{\delta}_{pre}, \Delta)$. Additionally, for any τ and Σ_n , $\mathbb{E}_{\hat{\beta}_n \sim \mathcal{N}(\delta + \tau, \Sigma_n)} \left[a + v' \hat{\beta}_n \right] = \frac{1}{2} (\theta^{ub} + \theta^{lb})$, where θ^{ub} and θ^{lb} are the upper and lower bounds of the identified set $\mathcal{S}(\delta + \tau, \Delta)$.

Proof. Let $b^{max}(\delta_{pre}^*) := \left(\max_{\tilde{\delta}} l' \tilde{\delta}_{post}, \text{ s.t. } \tilde{\delta} \in \Delta, \tilde{\delta}_{pre} = \delta_{pre}^*\right)$, where we define $b^{max} = -\infty$ if $\delta_{pre}^* \notin \Delta_{pre}$. Likewise, define $b^{min}(\delta_{pre}^*) := \left(\min_{\tilde{\delta}} l' \tilde{\delta}_{post}, \text{ s.t. } \tilde{\delta} \in \Delta, \tilde{\delta}_{pre} = \delta_{pre}^*\right)$, where we define $b^{min} = \infty$ if $\delta_{pre}^* \notin \Delta_{pre}$. Note that Δ convex implies that b^{max} is concave and b^{min} is convex. Thus, $-LID(\delta_{pre}^*) = b^{min}(\delta_{pre}^*) - b^{max}(\delta_{pre}^*)$ is convex (where we define $LID(\delta_{pre}^*) = -\infty$ if $\delta_{pre}^* \notin \Delta_{pre}$). The domain of $-LID(\delta_{pre}^*)$ (i.e. the set of values for which it is finite) is Δ_{pre} , since it is infinite for $\delta_{pre}^* \notin \Delta_{pre}$ by construction, and by assumption, $LID(\delta_{pre}^*)$ is finite for all $\delta_{pre}^* \in \Delta_{pre}$. Since Δ is assumed to be convex, Δ_{pre} is a non-empty convex set, and thus has non-empty relative interior, so the relative interior of the domain of -LID is non-empty.³⁷ It follows from Theorem 8.2 in Mau Nam (2019) that $\partial(-LID) = \partial(-b^{max}) + \partial(b^{min})$ where for a convex function f, ∂f is the subdifferential $\partial f(\bar{x}) := \{v : f(\bar{x}) + v'(x - \bar{x}) \leq f(x), \forall x\}$ and $\partial(-b^{max}) + \partial(b^{min})$ is the Minkowski sum of the two subdifferentials.

Additionally, if $LID(\delta_{pre}) = \sup_{\tilde{\delta}_{pre} \in \Delta_{pre}} LID(\tilde{\delta}_{pre})$, then $-LID(\delta_{pre}) = \inf_{\tilde{\delta}_{pre} \in \Delta_{pre}} -LID(\tilde{\delta}_{pre})$. Thus, standard results in convex analysis (see, e.g., Theorem 16.2 in Mau Nam (2019)) give that $0 \in \partial(-LID)(\delta_{pre}) + N(\Delta; \delta_{pre})$, where $N(\Delta; \delta_{pre}) = \{v_{pre} : v'_{pre}(\tilde{\delta}_{pre} - \delta_{pre}) \leq 0, \forall \tilde{\delta}_{pre} \in \Delta_{pre}\}$ is the normal cone to Δ_{pre} at δ_{pre} . Hence, there exist vectors $\bar{v}_{min}, \bar{v}_{max}$ such that for all $\tilde{\delta}_{pre} \in \Delta_{pre}$,

$$b^{min}(\delta_{pre}) + \bar{v}'_{min}(\tilde{\delta}_{pre} - \delta_{pre}) \leqslant b^{min}(\tilde{\delta}_{pre})$$
(30)

$$-b^{max}(\delta_{pre}) + \bar{v}'_{max}(\tilde{\delta}_{pre} - \delta_{pre}) \leqslant -b^{max}(\tilde{\delta}_{pre})$$
(31)

$$-\left(\bar{v}_{min} + \bar{v}_{max}\right)'(\tilde{\delta}_{pre} - \delta_{pre}) \leqslant 0.$$
(32)

The inequalities (31) and (32) together imply that for all $\tilde{\delta}_{pre} \in \Delta_{pre}$,

$$b^{max}(\delta_{pre}) + \bar{v}'_{min}(\tilde{\delta}_{pre} - \delta_{pre}) \ge b^{max}(\tilde{\delta}_{pre}).$$
(33)

Now, let v be the vector such that $v_{post} = l$ and $v_{pre} = -\bar{v}_{min}$. Observe that

$$\max_{\tilde{\delta}\in\Delta} a + v'_{pre}\tilde{\delta}_{pre} + l'\tilde{\delta}_{post} = \max_{\tilde{\delta}_{pre}\in\Delta_{pre}} \left(a + v'_{pre}\tilde{\delta}_{pre} + \max_{\bar{\delta}\in\Delta,\bar{\delta}_{pre}=\tilde{\delta}_{pre}} l'\bar{\delta}_{post} \right)$$
$$= \max_{\tilde{\delta}_{pre}\in\Delta_{pre}} a + v'_{pre}\tilde{\delta}_{pre} + b^{max}(\tilde{\delta}_{pre})$$
$$\leqslant a + v'_{pre}\delta_{pre} + b^{max}(\delta_{pre}), \tag{34}$$

where the first equality nests the maximization, the second equality uses the definition of

 $^{^{37}}$ The relative interior of a set is the interior of the set relative to its affine hull. See, e.g., Mau Nam (2019), Chapter 5.

 b^{max} , and the inequality follows from (33). An analogous argument using (30) yields that

$$\min_{\tilde{\delta}\in\Delta} a + v'_{pre}\tilde{\delta}_{pre} + l'\tilde{\delta}_{post} = \min_{\tilde{\delta}_{pre}\in\Delta_{pre}} a + v'_{pre}\tilde{\delta}_{pre} + b^{min}(\tilde{\delta}_{pre})$$

$$\geqslant a + v'_{pre}\delta_{pre} + b^{min}(\delta_{pre}).$$
(35)

Now, it is apparent from equation (29) that

$$\bar{b}(a,v) = \max\left\{ \left| \max_{\tilde{\delta}\in\Delta} a + v'_{pre}\tilde{\delta}_{pre} + l'\tilde{\delta}_{post} \right|, \left| \min_{\tilde{\delta}\in\Delta} a + v'_{pre}\tilde{\delta}_{pre} + l'\tilde{\delta}_{post} \right| \right\},\tag{36}$$

which is bounded above by max $\{a + v'_{pre}\delta_{pre} + b^{max}(\delta_{pre}), -(a + v'_{pre}\delta_{pre} + b^{min}(\delta_{pre}))\}$ from the results above. Setting $a = -v'_{pre}\delta_{pre} - \frac{1}{2}(b^{max}(\delta_{pre}) + b^{min}(\delta_{pre}))$, this upper bound reduces to $\frac{1}{2}(b^{max}(\delta_{pre}) - b^{min}(\delta_{pre}))$. Since $LID(\delta_{pre}, \Delta) = b^{max}(\delta_{pre}) - b^{min}(\delta_{pre})$ and $LID(\delta_{pre}, \Delta) = \sup_{\tilde{\delta}_{pre}\in\Delta_{pre}} LID(\tilde{\delta}_{pre}, \Delta)$ by assumption, it is then immediate that $\bar{b} \leq \frac{1}{2}\sup_{\tilde{\delta}_{pre}\in\Delta_{pre}} LID(\tilde{\delta}_{pre}, \Delta)$. The inequality in the opposite direction follows from Lemma A.10. Finally, substituting in the definition of a and v above and simplifying, we see that $\mathbb{E}_{\hat{\beta}_n \sim \mathcal{N}(\delta + \tau, \Sigma_n)} \left[a + v' \hat{\beta}_n \right] = l' \beta_{post} - \frac{1}{2}(b^{max}(\delta_{pre}) + b^{min}(\delta_{pre}))$, which from (5) and (6) we see is the midpoint of the identified set.

Lemma A.12. Let χ_{α} be the $1 - \alpha$ quantile of the $|\mathcal{N}(b, \sigma^2)|$ distribution for $b \ge 0$. Then $b + \sigma z_{1-\alpha} \le \chi_{\alpha} \le b + \sigma z_{1-\alpha/2}$.

Proof. Since $|\xi| \ge \xi$, we have that $q_{1-\alpha}(|\xi| | \xi \sim \mathcal{N}(b, \sigma^2)) \ge q_{1-\alpha}(\xi | \xi \sim \mathcal{N}(b, \sigma^2)) = b + \sigma z_{1-\alpha}$, which yields the first inequality. For the second inequality, observe that

$$q_{1-\alpha}(|\xi| | \xi \sim \mathcal{N}(b, \sigma^2)) = q_{1-\alpha}(|\xi+b| | \xi \sim \mathcal{N}(0, \sigma^2))$$

$$\leq b + q_{1-\alpha}(|\xi| | \xi \sim \mathcal{N}(0, \sigma^2)) = b + \sigma z_{1-\alpha/2}$$

where the first inequality uses the triangle inequality, and the final equality uses the fact that a mean-zero normal distribution is symmetric about 0.

B Additional Simulation Results

This section contains additional simulation results that complement the simulations presented in the main text. Section B.1 describes the computation of the optimal bound for expected excess length. Section B.2 contains additional results from the normal data-generating process considered in the main text. Section B.3 presents results from a non-normal datagenerating process in which the covariance matrix is estimated from the data, which show that our proposed procedures have (approximate) size control, with similar power curves to those in the normal simulations.

B.1 Optimal bounds on excess length

We now discuss the computation of optimal bounds on the excess length of confidence intervals that satisfy the uniform coverage requirement (10). In Section 5, we benchmark the performance of our proposed procedures in Monte Carlo simulations relative to these bounds.

The following result restates Theorem 3.2 of Armstrong and Kolesár (2018) in the notation of our paper, which provides a formula for the optimal expected length of a confidence set that satisfies the uniform coverage requirement.

Lemma B.1. Suppose that Δ is convex. Let \mathcal{I}_{α} denote the set of confidence sets that satisfy the coverage requirement (10). Then, for any $\delta^* \in \Delta$, $\tau^*_{post} \in \mathbb{R}^{\bar{T}}$, and Σ_n positive definite,

$$\inf_{\mathcal{C} \in \mathcal{I}_{\alpha}} \mathbb{E}_{\hat{\beta}_n \sim \mathcal{N}(\delta^* + L_{post}\tau^*, \Sigma_n)} \left[\lambda(\mathcal{C}) \right] = (1 - \alpha) \mathbb{E} \left[\bar{\omega} (z_{1-\alpha} - Z) - \bar{\omega} (z_{1-\alpha} - Z) \right] Z < z_{1-\alpha}],$$

where $Z \sim \mathcal{N}(0, 1)$, $z_{1-\alpha}$ is the $1-\alpha$ quantile of Z, and

$$\begin{split} \bar{\omega}(b) &:= \sup\{l'\tau \mid \tau \in \mathbb{R}^T, \exists \delta \in \Delta \ s.t. \ \|\delta + L_{post}\tau - \beta^*\|_{\Sigma_n}^2 \leqslant b^2\}\\ \bar{\omega}(b) &:= \inf\{l'\tau \mid \tau \in \mathbb{R}^{\bar{T}}, \exists \delta \in \Delta \ s.t. \ \|\delta + L_{post}\tau - \beta^*\|_{\Sigma_n}^2 \leqslant b^2\}, \end{split}$$

for $\beta^* := \delta^* + L_{post} \tau_{post}^*$, and $||x||_{\Sigma} = x' \Sigma^{-1} x$.

The proof of this result follows from observing that the confidence set that optimally directs power against $(\delta^*, \tau_{post}^*)$ inverts Neyman-Pearson tests of $H_0 : \delta \in \Delta, \theta = \bar{\theta}$ against $H_A : (\delta, \tau_{post}) = (\delta^*, \tau_{post}^*)$ for each value $\bar{\theta}$. The formulas above are then obtained by integrating one minus the power function of these tests over $\bar{\theta}$. By the same argument, the optimal excess length for confidence sets that control size is the integral of one minus the power function over all points $\bar{\theta}$ outside of the identified set. Additionally, for any value $\bar{\theta} \in S(\beta, \Delta)$, the null and alternative hypotheses are observationally equivalent, and so the most powerful test trivially has size α . It follows that the lowest achievable expected excess length is $(1 - \alpha) \cdot LID(\delta_{pre}^*, \Delta)$ shorter than the lowest achievable expected length, where as in Section 4, LID denotes the length of the identified set.

Corollary B.1. Under the conditions of Lemma B.1,

$$\inf_{\mathcal{C}\in\mathcal{I}_{\alpha}}\mathbb{E}_{\hat{\beta}_{n}\sim\mathcal{N}(\beta^{*},\Sigma_{n})}\left[EL(\mathcal{C};\beta^{*})\right] = \inf_{\mathcal{C}\in\mathcal{I}_{\alpha}}\mathbb{E}_{\hat{\beta}_{n}\sim\mathcal{N}(\beta^{*},\Sigma_{n})}\left[\lambda(\mathcal{C})\right] - (1-\alpha)LID(\beta^{*},\Delta),$$

where $EL(\mathcal{C};\beta) = \lambda(\mathcal{C}\setminus\mathcal{S}(\beta,\Delta))$ is the excess length of the confidence set \mathcal{C} , i.e. the length of the part of the confidence set that falls outside of the identified set.

Recall that when Δ is the union of polyhedra ($\Delta = \bigcup_{k=1}^{K} \Delta_k$), the identified set is the union of the identified sets for each of the Δ_k . Thus, any \mathcal{C}_{α} that satisfies (10) for Δ must also satisfy (10) for each Δ_k . It follows that the expected excess length for \mathcal{C} is bounded below by the optimal excess length for confidence sets satisfying (10) for Δ_k for each k. For Δ s that are unions of polyhedra, we therefore use the largest lower bound implied by the individual Δ_k , which is a potentially non-sharp lower bound on the excess length of a procedure that satisfies (10) for Δ .

B.2 Additional Results for Normal Simulations

In the main text, we report efficiency in terms of excess length for the parameter $\theta = \tau_1$ for $\Delta^{SD}(M)$, $\Delta^{SDPB}(M)$, $\Delta^{SDRM}(\bar{M})$ and $\Delta^{RM}(\bar{M})$. In this section, we provide additional simulation results.

Alternative choices of \overline{M} for $\Delta^{SDRM}(\overline{M})$ and $\Delta^{RM}(\overline{M})$. The main text reports efficiency in terms of excess length over $\Delta^{SDRM}(\overline{M})$ and $\Delta^{RM}(\overline{M})$ for $\overline{M} = 1$. We now report additional results for $\overline{M} = 1, 2, 3$. The results are qualitatively similarly, suggesting that the choice of \overline{M} does not appear to have a large effect on the performance of our proposed procedures.

Alternative choice of target parameter. The main text reports efficiency in terms of excess length for the parameter $\theta = \tau_1$. We now report additional results using the average of post-period treatment effects, $\theta = \bar{\tau}_{post}$, as the target parameter.

Figure I2 plots the efficiency results for $\theta = \bar{\tau}_{post}$ over $\Delta^{SD}(M)$ and $\Delta^{SDPB}(M)$. As in the main text, we conduct these simulations under the assumption of parallel trends and zero treatment effects (i.e., $\beta = 0$), reporting results as M/σ_1 varies.

Figure I3 plots the efficiency results for $\theta = \bar{\tau}_{post}$ over $\Delta^{SDRM}(\bar{M})$ and $\Delta^{RM}(\bar{M})$. As in the main text, we conduct these simulations under the assumption of zero treatment effects and a "pulse" pre-trend (i.e., $\beta_{-1} = \delta_{-1}$ and $\beta_t = 0$ for all $t \neq -1$), reporting results for $\bar{M} = 1$ over $\delta_{-1}/\sigma_1 = 0, 1, 2, 3$.³⁸

³⁸We note that over $\Delta^{SDRM}(\bar{M})$ the median efficiency ratio for our proposed confidence sets is larger than one for $\bar{M} = 3$. For $\bar{M} = 3$, the length of the identified set for $\theta = \bar{\tau}_{post}$ can be quite large when there are many post-treatment periods (e.g., as mentioned in the main text, 5 papers in the survey have $\bar{T} > 10$), and so this behavior occurs due to computational constraints on the grid size for the underlying test inversion.

Figure I1: $\Delta^{SDRM}(\bar{M})$ and $\Delta^{RM}(\bar{M})$: Median efficiency ratios for proposed procedures when $\theta = \tau_1$ as \bar{M} varies.



Note: This figure shows the median efficiency ratio for our proposed confidence sets for $\theta = \tau_1$ over $\Delta^{SDRM}(\bar{M})$, $\Delta^{RM}(\bar{M})$ and $\bar{M} = 1, 2, 3$. The efficiency ratio for a procedure is defined as the excess length bound divided by the procedure's expected excess length. The results for $\bar{M} = 1$ are plotted in red, $\bar{M} = 2$ are plotted in blue, and $\bar{M} = 3$ are plotted in green. The results for the conditional-least favorable hybrid confidence set ("C-LF Hybrid") are plotted in the solid line with circles. The results for the conditional confidence set are plotted in the dashed line with triangles. Results are averaged over 1000 simulations for each of the 12 papers surveyed, and the median across papers is reported here.





Note: This figure shows the median efficiency ratios for our proposed confidence sets for $\Delta^{SD}(M)$ and $\Delta^{SDPB}(M)$ when $\theta = \bar{\tau}_{post}$. The efficiency ratio for a procedure is defined as the optimal bound divided by the procedure's expected excess length. The results for the FLCI are plotted in purple, the results for the conditional-LF ("C-LF Hybrid") in blue, and the results for the conditional confidence set are in green. Results are averaged over 1000 simulations for each of the 12 papers surveyed, and the median across papers is reported here.

Figure I3: Median efficiency ratios for $\Delta^{SDRM}(\bar{M})$ and $\Delta^{RM}(\bar{M})$ when $\theta = \bar{\tau}_{post}$.



Note: This figure shows the median efficiency ratios for our proposed confidence sets for $\Delta^{SDRM}(\bar{M})$ and $\Delta^{RM}(\bar{M})$ when $\theta = \bar{\tau}_{post}$ and $\bar{M} = 1$. The efficiency ratio for a procedure is defined as the optimal bound divided by the procedure's expected excess length. The results for the conditional-least favorable ("C-LF") hybrid in blue and the results for the conditional confidence set in green. Results are averaged over 1000 simulations for each of the 12 papers surveyed, and the median across papers is reported here.

B.3 Non-normal simulation results with estimated covariance matrix

In the main text, we presented simulations results where $\hat{\beta}$ is normally distributed and its covariance matrix is treated as known. In this section, we present Monte Carlo results using a data-generating process in which $\hat{\beta}$ is not normally distributed and the covariance matrix is estimated from the data. Specifically, we consider simulations based on the empirical distribution in Bailey and Goodman-Bacon (2015). We find that all of our procedures achieve (approximate) size control, and our results on the relative power of the various procedures are quite similar to those presented in the main text.

B.3.1 Simulation design

The simulations are calibrated using the empirical distribution of the data in Bailey and Goodman-Bacon (2015).³⁹ Let $\hat{\beta}$, $\hat{\Sigma}$ denote the original, estimated event-study coefficients and variance-covariance matrix from the event-study regression in the paper. We simulate data using a clustered bootstrap sampling scheme at the county level (i.e. the level of clustering used by the authors in their event-study regression). For each bootstrap sample b, we re-estimate the event-study coefficients $\hat{\beta}_b$ and the variance-covariance matrix $\hat{\Sigma}_b$ also using the clustering scheme specified by the authors. We then re-center the bootstrapped coefficient so that under our simulated data-generating process either parallel trends holds (i.e., $\hat{\beta}_b^{centered} = \hat{\beta}_b - \hat{\beta} + \delta_{-1} * e_{-1}$ where e_{-1} is the $(T + \bar{T})$ -dimensional vector with one in t = -1 entry and zeroes everywhere else). We construct our proposed confidence sets for bootstrap draw b using the pair ($\hat{\beta}_b^{centered}, \hat{\Sigma}_b$).

As in the main text, we focus on the performance of our proposed confidence sets for $\Delta^{SD}(M)$, $\Delta^{SDPB}(M)$ under parallel trends and $\Delta^{SDRM}(\bar{M})$, $\Delta^{RM}(\bar{M})$ under the "pulse" pre-trend. The parameter of interest in these simulations is the causal effect in the first post-period ($\theta = \tau_1$). For $\Delta^{SD}(M)$ and $\Delta^{SDPB}(M)$, we report the performance of the FLCI, conditional confidence set, and conditional-least favorable confidence set. For $\Delta^{SDRM}(\bar{M})$ and $\Delta^{RM}(\bar{M})$, we report the performance of the conditional-least favorable confidence set and the conditional-least favorable confidence set and the conditional-least favorable confidence set and the conditional-least favorable confidence set. All results are averaged over 1000 bootstrap samples.

³⁹Since implementing the bootstrap in practice is logistically challenging, we do so for one paper rather than the full 12 papers in the survey. We chose the first paper alphabetically to minimize concerns about cherry-picking.
B.3.2 Size control simulations

Table 2 reports the maximum rejection rate of each procedure over a grid of parameter values θ within the identified set $S(\beta, \Delta)$ for $\Delta = \Delta^{SD}(M)$ and $\Delta = \Delta^{SDPB}(M)$ under parallel trends (i.e., $\beta = 0$). We report results for $M/\sigma_1 = 0, 1, 2, 3, 4, 5$. The table shows that all our procedures approximately control size, with null rejection rates not exceeding 0.08.

Δ	M/σ_1	Conditional	FLCI	C-LF Hybrid
$\Delta^{SD}(M)$				
	0	0.073	0.078	0.069
	1	0.046	0.061	0.044
	2	0.038	0.072	0.037
	3	0.040	0.072	0.038
	4	0.049	0.072	0.045
	5	0.059	0.072	0.051
$\Delta^{SDPB}(M)$				
	0	0.079	0.078	0.074
	1	0.052	0.047	0.048
	2	0.046	0.055	0.042
	3	0.051	0.058	0.046
	4	0.055	0.058	0.051
	5	0.059	0.058	0.057

Table 2: Maximum null rejection probability over the identified set $S(\beta, \Delta)$ for $\Delta = \Delta^{SD}(M)$ and $\Delta = \Delta^{SDPB}(M)$ under parallel trends (i.e., $\beta = 0$) using the empirical distribution from Bailey and Goodman-Bacon (2015).

Table 3 reports the maximum rejection rate of the conditional test and the conditionalleast favorable test over a grid of parameter values θ within the identified set $S(\beta, \Delta)$ for $\Delta = \Delta^{SDRM}(\bar{M})$ and $\Delta = \Delta^{RM}(\bar{M})$ under the "pulse" pre-trend (i.e., $\beta_{-1} = \delta_{-1}$ and $\beta_t = 0$ for all $t \neq -1$). We report results for $\bar{M} = 1$ and $\delta_{-1}/\sigma_1 = 1, 2, 3$. The table shows that all our procedures approximately control size, with worst-case null rejection probability of 0.058.

B.3.3 Comparison with normal simulations

We next compare results from the non-normal simulations with estimated covariance discussed above to the normal model simulations the main text, in which $\hat{\beta}$ is normal and Σ is treated as known.

Figures I4-I5 shows the rejection probabilities at different values of the parameter θ using both simulation methods for $\Delta^{SD}(M)$, $\Delta^{SDPB}(M)$ at $M/\sigma_1 = 0, 5$ respectively. The results are quite similar for all values of M/σ_1 considered, and we thus omit the intermediate values.

Δ	δ_{-1}/σ_1	Conditional	C-LF Hybrid
$\Delta^{SDRM}(\bar{M})$			
	1	0.009	0.008
	2	0.037	0.035
	3	0.058	0.054
$\Delta^{RM}(\bar{M})$			
	1	0.005	0.005
	2	0.017	0.016
	3	0.024	0.023

Table 3: Maximum null rejection probability over the identified set $S(\beta, \Delta)$ for $\Delta = \Delta^{SDRM}(\bar{M})$ and $\Delta = \Delta^{RM}(\bar{M})$ under the "pulse" pre-trend (i.e., $\beta_{-1} = \delta_{-1}$ and $\beta_t = 0$ for all $t \neq -1$) and $\bar{M} = 1$ using the empirical distribution from Bailey and Goodman-Bacon (2015). We report results for $\delta_{-1}/\sigma_1 = 1, 2, 3$.

The estimated average rejection rates of each procedure are quite similar in the non-normal simulations and the normal simulations across each choice of Δ . As a result, the relative rankings of the procedures in terms of power are the same in the non-normal simulations as in the normal simulations discussed in the main text. Similarly, Figures I6-I7 shows the rejection probabilities at different values of the parameter θ using both simulation methods for $\Delta^{SDRM}(\bar{M})$, $\Delta^{RM}(\bar{M})$ at $\delta_{-1}/\sigma_1 = 1, 2, 3$ respectively and $\bar{M} = 1$.

Figure I4: Comparison of rejection probabilities using bootstrap and normal simulations. Results are shown for $\theta = \tau_1$, and each choice of $\Delta = \Delta^{SD}(M), \Delta^{SDPB}(M)$, and $M/\sigma_1 = 0$. The average rejection rate for the non-normal simulations are in red and the average rejection rate for the normal simulations are in blue; the dashed black lines indicate the identified set bounds. Results are averaged over 1000 simulations.



Figure I5: Comparison of rejection probabilities using bootstrap and normal simulations. Results are shown for $\theta = \tau_1$, and each choice of $\Delta = \Delta^{SD}(M), \Delta^{SDPB}(M)$, and $M/\sigma_1 = 5$. The average rejection rate for the non-normal simulations are in red and the average rejection rate for the normal simulations are in blue; the dashed black lines indicate the identified set bounds. Results are averaged over 1000 simulations.



Figure I6: Comparison of rejection probabilities using bootstrap and normal simulations for $\Delta^{SDRM}(\bar{M})$ and $\Delta^{RM}(\bar{M})$. Results are shown for $\theta = \tau_1$, $\bar{M} = 1$ and $\delta_{-1}/\sigma_1 = 1$. The average rejection rate for the non-normal simulations are in red and the average rejection rate for the normal simulations are in blue; the dashed black lines indicate the identified set bounds. Results are averaged over 1000 simulations.



Figure I7: Comparison of rejection probabilities using bootstrap and normal simulations for $\Delta^{SDRM}(\bar{M})$ and $\Delta^{RM}(\bar{M})$. Results are shown for $\theta = \tau_1$, $\bar{M} = 1$ and $\delta_{-1}/\sigma_1 = 3$. The average rejection rate for the non-normal simulations are in red and the average rejection rate for the normal simulations are in blue; the dashed black lines indicate the identified set bounds. Results are averaged over 1000 simulations.



Appendix References

- Armstrong, Timothy and Michal Kolesár, "Optimal Inference in a Class of Regression Models," *Econometrica*, 2018, 86, 655–683.
- Bailey, Martha J. and Andrew Goodman-Bacon, "The War on Poverty's Experiment in Public Medicine: Community Health Centers and the Mortality of Older Americans," *American Economic Review*, March 2015, 105 (3), 1067–1104.
- Ingster, Yuri and I. A. Suslina, Nonparametric Goodness-of-Fit Testing Under Gaussian Models Lecture Notes in Statistics, New York: Springer-Verlag, 2003.
- Nam, Nguyen Mau, "Convex Analysis: An introduction to convexity and nonsmooth analysis," https://maunamn.wordpress.com/ 2019.
- Schrijver, Alexander, Theory of Linear and Integer Programming, Wiley-Interscience, 1986.