

Inference in Models Defined by Infinitely Many Inequalities: A Survey

Xiaoxia Shi*

August 14, 2025

Abstract

This paper complements the surveys by Molinari (2020) and Canay and Shaikh (2017) by reviewing the extensive body of work on inference based on an infinite number of inequalities. In addition to synthesizing existing results, I aim to identify gaps in the literature and outline promising avenues for future research.

Keywords: Moment Inequalities, Hypothesis Testing, Test Inversion, Gaussian Approximation, Power Regret

*Department of Economics, University of Wisconsin-Madison (xshi@ssc.wisc.edu). I thank Francesca Molinari, Zhipeng Liao, Haitian Xie, Zhenting Sun, Qingliang Fan, and Anders Kock for helpful discussion. The paper is written while I am visiting Guanghai Management School at Peking University.

The study of partially identified models originated from the recognition that economic models often fail to yield point identification. Foundational contributions by Manski (1989, 1990, 1993) developed the use of bounds and incomplete models to draw inference under minimal assumptions, particularly in treatment effect and selection settings.¹ Building on these early insights, a rich literature has developed on partially identified models. The applications of such models span all major areas of empirical economics including labor economics (e.g. Blundell et al. (2007)), industrial organization (e.g. Ciliberto and Tamer (2009)), trade (e.g. Morales et al. (2019), Kalouptside et al. (2020)), market design (e.g. He (2017) and Fack et al. (2019)), macroeconomics (e.g. Giacomini and Kitagawa (2021)), network formation and interaction (e.g. Sheng (2020)), and political economics (e.g. Iaryczower et al. (2018)). A more thorough overview of the applications can be found in the survey papers Ho and Rosen (2017), Pakes et al. (2015), Molinari (2020), and Kline et al. (2021). Canay et al. (2023) provide a user’s guide for doing inference in such models.

In many applications, partially identified models are characterized by infinitely many moment (in)equalities. The combination of partial identification and the sheer number of inequality constraints poses distinct challenges for statistical inference. Over the past two decades, a rich set of methods has been developed, including diverse testing procedures and new asymptotic tools. Nevertheless, important questions remain open. This review seeks to synthesize key contributions from the existing literature and to highlight unresolved issues and promising avenues for future

¹In parallel, Phillips (1989) formalized partial identification as a consequence of singularity in the information matrix, highlighting how weak instruments and local identification failure can lead to nonstandard asymptotic behavior.

research.

This paper complements the comprehensive review of other aspects of partially identified models, by Molinari (2020) and Canay and Shaikh (2017). The former provides a high-level overview of the key issues including consistent set estimation, uniform coverage, and the effect of misspecification, with more detailed treatment on random set theory and computation. The latter focuses on inference methods for models defined by a finite number of moment inequalities.

1 Setup and Examples

A generic form of models defined by inequalities is as follows:

$$g_P(\theta, t) \leq 0 \text{ for all } t \in \mathcal{T}, \quad (1)$$

for an index set \mathcal{T} , where $g_P(\cdot, \cdot)$ is an unknown function that is determined by the unknown data generating process (DGP) P , and θ is an unknown parameter living in the parameter space Θ . This setup allows equalities: when there is an equality restriction, we can simply write it as a pair of opposing inequalities. For much of the subsequent discussion, it is not necessary to write the equalities out, but it should be noted that some test statistics used in the literature may benefit from an explicit equality/inequality notation.

This model is defined by an infinite number of inequalities when \mathcal{T} contains an infinite number of points. There are a few reasons that \mathcal{T} may be infinite, which I illustrate with examples:

Example 1 (Conditional Moment Inequalities). *Many models used in structural estimation are conditional ones, where the model specifies the conditional generating process of the endogenous variables given exogenous variables. When model incompleteness and/or data imperfection lead to moment inequalities, they are conditional moment inequalities given the exogeneous variables. Mathematically, they are*

$$\mathbb{E}_P[m(W, \theta)|Z] \leq 0 \text{ almost surely,} \quad (2)$$

where $m(\cdot, \cdot)$ is a known R^d -valued moment function for a $d > 0$, W is the vector of observables which may include exogenous and endogenous variables, and Z is the vector of exogenous variables. When Z is a continuous variables, (2) stands for a continuum of inequalities.

There are two ways to write (2) in the form of (1). The first is the non-parametric conditional mean approach, where

$$g_P(\theta, t) = \mathbb{E}_P[m(W, \theta)|Z = t], \text{ for all } t \in \mathcal{T}, \quad (3)$$

where $\mathcal{T} = \text{Supp}(Z)$, where $\text{Supp}(Z)$ stands for the support of Z . The second is the instrumental function approach, where

$$g_P(\theta, t) = \mathbb{E}_P[m(W, \theta)t(Z)], \text{ for all } t \in \mathcal{T}, \quad (4)$$

where \mathcal{T} is a sufficiently rich class of functions mapping the support of Z to $[0, \infty)$. The classes of functions that are sufficiently rich are discussed in Andrews and Shi

(2013). These two representations are equivalent, but they motivate different inference procedures as I discuss later.

Example 2 (Sharp Identified Set). *A common feature of many incomplete structural models is that the correspondence between observables and unobservables are multi-valued. Let Y stand for the vector of observables and ε stand for the vector of unobservables. The model imposes the following restrictions:*

$$\varepsilon \in \mathcal{E}(\theta, Y), \quad (5)$$

where $\mathcal{E}(\theta, y)$ is a known closed set for each value of θ and y , and $\varepsilon \sim G_\varepsilon(\cdot|\theta)$. Then the sharp identified set is the set of $\theta \in \Theta$ such that the following inequalities hold:

$$P(\mathcal{E}(\theta, Y) \subseteq A) \leq G_\varepsilon(A|\theta). \quad (6)$$

for all measurable subsets A of the support of ε . See e.g. Chesher and Rosen (2017). Clearly, (6) stands for an infinite number of inequalities. The inequalities in (6) can be written in the form of (1) as follows:

$$g_P(\theta, t) = \mathbb{E}[1\{\mathcal{E}(\theta, Y) \subseteq t\}] - G_\varepsilon(t|\theta), \text{ for all } t \in \mathcal{T}, \quad (7)$$

where \mathcal{T} is the set of all measurable subsets of the support of Y . For specific models, such as discrete choice models with instrumental variables, selectively observed data, and some auction models, Galichon and Henry (2011), Chesher et al. (2013) and Chesher and Rosen (2017) develop methods to reduce the class of all measurable

subsets to a much smaller core-determining class, but even that class can contain many elements.

Example 3 (Support Function Characterization of Sharp Identified Set). *When certain moment conditions are imposed and some mild conditions are satisfied, Beresteanu et al. (2011) show that (6) can be equivalently expressed in terms of support functions. For example, $\mathbb{E}[\varepsilon] = \mathbf{0}$ and (6) together can be written equivalently as:*

$$\mathbb{E}h_{\mathcal{E}(\theta, Y)}(u) \geq 0 \text{ for all } u \in \mathbb{S}^{d_u-1}, \quad (8)$$

where $h_A(u) = \sup_{e \in A} u'e$ is the support function of the set A , $\mathbb{S}^{d_u-1} = \{u \in \mathbb{R}^{d_u} : \|u\| = 1\}$ and d_u is the dimension of ε .

Example 4 (Infinite Number of Conditional Moment Inequalities). *In Examples 2 and 3, I abstracted away from exogenous variables for simplicity. However, in practice, there often are exogenous variables, and the model is often regarding conditional distributions of Y given exogenous variables. To be precise, let X be the vector of exogenous variables. Instead of (5), we have*

$$\varepsilon \in \mathcal{E}(\theta, Y, X), \quad (9)$$

where $\varepsilon \sim G_\varepsilon(\cdot|X, \theta)$. *Instead of (6), we have*

$$P(\mathcal{E}(\theta, Y) \subseteq A|X) \leq G_\varepsilon(A|X, \theta), \text{ a.s.} \quad (10)$$

for all measurable subsets A of the support of ε . *Instead of $\mathbb{E}_P[\varepsilon] = 0$, we have*

$\mathbb{E}_P[\varepsilon|X] = 0$, and instead of (8), we have

$$\mathbb{E}[h_{\mathcal{E}(\theta, Y, X)}(u)|X] \geq 0 \text{ for all } u \in \mathcal{S}^{d_u-1}. \quad (11)$$

When X contains a continuous variable, both (10) and (11) involve infinite number of conditional moment inequalities. We can write them in the form of (1) either by taking the non-parametric conditional mean approach as in (3) or by taking the instrumental function approach as in (4).

1.1 Inference by Test Inversion

The inequality model (1) often does not point identify θ . Instead, it defines an identified set for θ :

$$\Theta_0(P) = \{\theta \in \Theta : g_P(\theta, t) \leq 0 \text{ for all } t \in \mathcal{T}\}. \quad (12)$$

If $g_P(\cdot, \cdot)$ is known, one can calculate this set using numerical or analytical tools. This is the exercise that, for example, Chesher et al. (2013) do in the numerical part of their paper, where they design a DGP and calculate $\Theta_0(P)$ under this artificial DGP. Such exercises are useful for studying the identification power of various model assumptions under a designed P , but not applicable in an empirical environment.

In an empirical environment, the researcher has a dataset drawn from P . The goal is to infer about $\Theta_0(P)$ based on the dataset. In standard point-identified models, one often calculates a consistent parameter estimator and builds a confidence interval around it. However, in the literature of partially identified models, the estimation

of $\Theta_0(P)$ is out-shadowed by confidence set construction.² Part of the reason may be that the sample analogue estimator $\hat{\Theta}_n = \{\theta \in \Theta : \hat{g}(\theta, t) \leq 0 \text{ for all } t \in \mathcal{T}\}$ is inward biased unless strong assumptions are imposed, and removing the inward bias requires tuning parameters that the resulting estimator is sensitive to. The other part may be that confidence sets based on a consistent set estimator are difficult to develop except in special cases. I do not discuss the estimation problem further but refer interested readers to Section 4.2 of Molinari (2020) for a thorough review of existing methods.

I focus on confidence set construction. In particular, I mainly discuss confidence sets, denoted $CS_n(1 - \alpha)$, that cover the true value of the parameter with a given probability (asymptotically):

$$\inf_{P \in \mathcal{P}} \inf_{\theta_0 \in \Theta_0(P)} \Pr_P(\theta_0 \in CS_n(1 - \alpha)) \geq 1 - \alpha + o(1), \quad (13)$$

where \mathcal{P} is a set of DGPs allowed by the model and $\alpha \in (0, 1)$ is a nominal significance level. Since the researcher does not know the true DGP and does not know which point in $\Theta_0(P)$ is the true value even given P , we would like the minimum coverage probability under all possible combinations of (P, θ_0) allowed by the model to be bounded from below. Confidence sets that satisfy (13) are said to have uniform asymptotic coverage for the true value of the parameter.³

²A confidence set generalizes the concept of a confidence interval. It is a subset of the parameter space that has certain coverage probability guarantee.

³The literature has also defined a different notion of coverage:

$$\inf_{P \in \mathcal{P}} \Pr_P(\Theta_0(P) \subseteq CS_n(1 - \alpha)) \geq 1 - \alpha + o(1). \quad (14)$$

Confidence sets that satisfy this coverage guarantee are said to have uniform asymptotic coverage

Confidence sets satisfying (13) are constructed by test inversion. Specifically, one constructs a family of tests $\{\varphi_{n,\alpha}(\theta) : \theta \in \Theta\}$, where for each $\theta \in \Theta$, $\varphi_{n,\alpha}(\theta)$ is a test for the hypothesis

$$H_0 : g_P(\theta, t) \leq 0 \text{ for all } t \in \mathcal{T}. \quad (15)$$

Based on the tests, one defines the confidence set to be

$$CS_n(1 - \alpha) = \{\theta \in \Theta : \varphi_{n,\alpha}(\theta) = 0\}, \quad (16)$$

that is, the set of θ values at which the test does not reject. Since $\mathbb{E}_P[\varphi_{n,\alpha}(\theta)] = 1 - \Pr_P(\varphi_{n,\alpha}(\theta) = 0)$, the coverage guarantee (13) is satisfied if the test has the following uniform asymptotic level control:

$$\sup_{P \in \mathcal{P}} \sup_{\theta \in \Theta_0(P)} \mathbb{E}_P[\varphi_{n,\alpha}(\theta)] \leq \alpha + o(1). \quad (17)$$

In practice, $CS_n(1 - \alpha)$ is often computed by conducting the test $\varphi_{n,\alpha}(\theta)$ for a grid of θ values on Θ , and by inferring the boundary of $CS_n(1 - \alpha)$ from the acceptance and the rejection regions on the grid. When the projection of $CS_n(1 - \alpha)$ on a scalar parameter, say $\lambda(\theta)$, is desired, one can also compute the lower and upper end points of the projection via the following constrained optimization problems:

$$\min \setminus \max_{\theta \in \Theta : \varphi_{n,\alpha}(\theta) = 0} \lambda(\theta).$$

for the identified set. This notion of coverage is stronger than (13) and typically requires the confidence set to be wider. It is worth noting that some papers do not add $\inf_{P \in \mathcal{P}}$ in (14). See Chernozhukov et al. (2007). That results in the pointwise-in P asymptotic coverage of the identified set, which is not stronger (or weaker) than (13).

As we can see, constructing the confidence set (16) is naturally related to the literature of jointly testing an infinite number of inequalities. The hypotheses tested in this literature include stochastic dominance (e.g. Barrett and Donald (2003), Donald and Hsu (2016), Chetverikov et al. (2021)), conditional stochastic dominance (e.g. Delgado and Escanciano (2013)), stochastic monotonicity (e.g. Lee et al. (2009), Seo (2018)), regression monotonicity (e.g. Ghosal et al. (2000), Hsu et al. (2019)), density ratio ordering (Carolan and Tebbs (2005), Beare and Moon (2015), Beare and Shi (2019)), conditional predictive superiority (e.g. Li et al. (2022)) and so on. These hypotheses can be viewed as a special case of (15) where a parameter θ is not there. They can be tested using the tests developed for (15), although in some of the aforementioned papers, specialized tests that are not applicable to (15) are developed. Here I focus on the generic tests and do not discuss the specialized procedures.

2 Existing Tests

Now I discuss the existing tests for (15). Since the tests are constructed for each given θ , we omit θ and henceforce write $g_P(\theta, t)$ as $g_P(t)$ for notational simplicity.

There are two scenarios in which the literature has considered (15), depending on the large sample property of the estimator of $g_P(t) : t \in \mathcal{T}$. Let $\hat{g}_n(t) : t \in \mathcal{T}$ denote the estimator. The two scenarios are as follows:

1. The sequence of stochastic processes $\{r_n(\hat{g}_n(t) - g_P(t)) : t \in \mathcal{T}\}_{n=1}^{\infty}$ converges weakly in $\ell^\infty(\mathcal{T}, \mathbb{R}^d)$ to a tight Gaussian process $G(t) : t \in \mathcal{T}$, for a normalizing

sequence $\{r_n\}_{n=1}^\infty$.⁴

2. $\{r_n(\hat{g}_n(t) - g_P(t)) : t \in \mathcal{T}\}_{n=1}^\infty$ does not converge weakly to a tight Gaussian process for any normalizing sequence $\{r_n\}$.

2.1 Scenario 1

In the first scenario, standard empirical process techniques, in e.g. van der Vaart and Wellner (1996), can be applied. This allows one to aggregate the information contained in each dimension of $\hat{g}_n(t)$ for each t in a variety of ways to form the test statistic. In particular, it allows us to derive asymptotic distributions for tests statistics of the form

$$T_n^{\text{CvM}} = \int_{\mathcal{T}} S(r_n \hat{g}_n(t), \bar{\Sigma}_n(t)) d\mu(t), \quad (18)$$

where $\bar{\Sigma}_n(t)$ is an estimator of the variance-covariance matrix of $r_n(\hat{g}_n(t) - g_P(t))$ and $S(\cdot, \cdot)$ is a user-chosen function to aggregate different dimensions of $\hat{g}_n(t)$. This is called the Cramér-von-Mises type statistic in Andrews and Shi (2013).

Andrews and Shi (2013) define the CvM type statistic for conditional moment inequality models. Andrews and Shi (2014) define it for a non-parametric conditional moment inequality model where conditional moment inequalities hold at a given value of some of the conditional variables. Andrews and Shi (2017) define it for a model defined by infinitely many conditional moment inequalities. Hsu et al. (2019) extend it to testing generalized regression monotonicity. In these papers, a large

⁴Here $\ell^\infty(\mathcal{T}, \mathbb{R}^d)$ consists of all bounded functions $f : \mathcal{T} \rightarrow \mathbb{R}^d$.

sample distributional approximation of T_n is derived instead of a limit distribution. However, under a mild assumption on the variance-covariance estimator, it is possible to derive a limit distribution for T_n . I do so in Proposition 1.

Proposition 1. *Let $\{P_n\}$ be a sequence of distributions such that H_0 in (15) holds for each n . Suppose:*

- (i) $r_n(\widehat{g}_n(t) - g_{P_n}(t)) \Rightarrow G(t)$ in $\ell^\infty(\mathcal{T}, \mathbb{R}^d)$, where G is a tight Gaussian process;
- (ii) $\sup_{t \in \mathcal{T}} \|\bar{\Sigma}_n(t) - \bar{\Sigma}(t)\| \rightarrow_p 0$, for some $\bar{\Sigma} \in \ell^\infty(\mathcal{T}, \mathbb{S}_+^d)$ with eigenvalues uniformly bounded below by $\varepsilon > 0$, where \mathbb{S}_+^d is the set of positive definite matrices of size d ;
- (iii) $-r_n g_{P_n}(t) \rightarrow h(t)$ pointwise, for some $h : \mathcal{T} \rightarrow [0, \infty]^d$;
- (iv) $S : [-\infty, \infty]^d \times \mathbb{S}_+^d \rightarrow \mathbb{R}_+$ is continuous, non-decreasing in the first argument, and satisfies $S(m, \Sigma + \Sigma_1) \leq S(m, \Sigma)$ for all m and $\Sigma, \Sigma_1 \in \mathbb{S}_+^d$.

Then

$$T_n := \int_{\mathcal{T}} S(r_n \widehat{g}_n(t), \bar{\Sigma}_n(t)) d\mu(t) \longrightarrow_d T_\infty := \int_{\mathcal{T}} S(G(t) - h(t), \bar{\Sigma}(t)) d\mu(t).$$

Proof. Conditions (i) and (ii) imply, via the almost sure representation theorem, that there exists a version of $(r_n(\widehat{g}_n(t) - g_{P_n}(t)), \bar{\Sigma}_n(t)) : t \in \mathcal{T}$ with the same distribution that converges almost surely to a version of $(G(t), \bar{\Sigma}(t)) : t \in \mathcal{T}$. Without loss of generality, assume we are working with this version. Then, almost surely,

$$(r_n(\widehat{g}_n(t) - g_{P_n}(t)), \bar{\Sigma}_n(t)) \rightarrow (G(t), \bar{\Sigma}(t)) \quad \text{uniformly over } t \in \mathcal{T}. \quad (19)$$

Fix a sample path along which the convergence in (19) holds and $\sup_{t \in \mathcal{T}} G(t) < \infty$. Consider convergence along this path. Then, by the continuity of S (Condition (iv)) and Condition (iii),

$$S(r_n \hat{g}_n(t), \bar{\Sigma}_n(t)) \rightarrow S(G(t) - h(t), \bar{\Sigma}(t)) \quad \text{for all } t \in \mathcal{T}. \quad (20)$$

Moreover, the uniform convergence of $\bar{\Sigma}_n(\cdot)$ to $\bar{\Sigma}(\cdot)$ and condition (ii) imply that the minimum eigenvalue of $\bar{\Sigma}_n(t)$ is eventually uniformly bounded below by $\varepsilon/2$. This, combined with $g_{P_n}(t) \leq 0$, the non-decreasing property of S in its first argument and the non-increasing property in its second argument (condition (iv)), implies that $S(r_n \hat{g}_n(t), \bar{\Sigma}_n(t)) \leq S(r_n(\hat{g}_n(t) - g_{P_n}(t)), \frac{\varepsilon}{2}I)$. By the continuity of $S(\cdot, \frac{\varepsilon}{2}I)$ (via condition (iv)) and $\sup_{t \in \mathcal{T}} G(t) < \infty$, the Heine–Cantor theorem implies that $S(r_n(\hat{g}_n(t) - g_{P_n}(t)), \frac{\varepsilon}{2}I)$ converges uniformly to $S(G(t), \frac{\varepsilon}{2}I)$. Thus,

$$\sup_{t \in \mathcal{T}} S(r_n \hat{g}_n(t), \bar{\Sigma}_n(t)) < \sup_{t \in \mathcal{T}} S(G(t), \frac{\varepsilon}{2}I) + \varepsilon, \quad \text{eventually.} \quad (21)$$

Hence, the bounded convergence theorem applies and yields

$$\int_{\mathcal{T}} S(r_n \hat{g}_n(t), \bar{\Sigma}_n(t)) d\mu(t) \rightarrow \int_{\mathcal{T}} S(G(t) - h(t), \bar{\Sigma}(t)) d\mu(t). \quad (22)$$

This holds for all sample paths where (19) holds, and $\sup_{t \in \mathcal{T}} G(t) < \infty$. Thus, it holds almost surely and therefore also in distribution. This concludes the proof. \square

The proposition is a new result. It strengthens the distributional approximation result in Theorem 1 of Andrews and Shi (2013) to a limit distribution one. Based

on these results, one can obtain a simulation-based critical value after estimating a lower bound for $h(\cdot)$.

We can only bound $h(\cdot)$ instead of consistently estimating it because it is a limit of $r_n g_{P_n}(t)$ and $g_{P_n}(t)$ can at best be estimated r_n -consistently. On the other hand, a lower-bound is sufficient for constructing a valid test because replacing $h(\cdot)$ by a lower bound enlarges the integral $\int_{\mathcal{T}} S(G(t) - h(t), \bar{\Sigma}(t)) d\mu(t)$ due to the monotonicity of $S(g, \Sigma)$ in its first argument. Therefore, the resulting simulated critical value is asymptotically valid (i.e. not leading to excessive over-rejection under H_0).

Typically, there are two ways to bound $h(t)$. The first is sometimes called “least favorable” and sometimes called “plug-in asymptotics (PA)” in the literature. It is to bound $h(t)$ by 0 for all t , and is justified by the fact that $g_{P_n}(t) \geq 0$ for all P_n satisfying H_0 in (15).

The second is called “generalized moment selection (GMS)” in the literature. The idea is to approximate $h_j(t)$ by $+\infty$ or something that diverges to $+\infty$ if there is strong evidence that $h_j(t) = \infty$ and to replace it by zero or something that converges to zero otherwise. Andrews and Shi (2013) recommend the following GMS bound for $h(\cdot)$ (when $r_n = n^{1/2}$):

$$\underline{h}_{n,j}(t) = B_n 1\{\kappa_n^{-1} n^{1/2} \hat{g}_{n,j}(t) / \bar{\sigma}_{n,j}(t) > 1\}, \quad (23)$$

where κ_n and B_n (e.g. $\kappa_n = (0.3 \log(n))^{1/2}$, $B_n = (0.4 \log(n) / \log(\log(n)))^{1/2}$) are user-chosen positive constants such that $\kappa_n \rightarrow \infty$ and $B_n / \kappa_n \rightarrow 0$ as $n \rightarrow \infty$, and $\bar{\sigma}_{n,j}(t)$ is the j th diagonal element of $\bar{\Sigma}_n(t)$.

Once a feasible bound $\underline{h}_n(t)$ is chosen, one can define the critical value, $cv_n(\alpha)$,

to be the simulated $1 - \alpha$ quantile of:

$$T_n^{cv} = \int_{\mathcal{T}} S(G_n^*(t) + \underline{h}_n(t), \bar{\Sigma}_n(t)) d\mu(t), \quad (24)$$

where $G_n^*(t)$ is the random component (conditional on data) to be simulated. It can be the bootstrap empirical process $r_n(\hat{g}_n^*(t) - \hat{g}_n(t)) : t \in \mathcal{T}$ where $\hat{g}_n^*(t)$ is $\hat{g}_n(t)$ calculated from a bootstrap sample. It can also be a Gaussian process with variance covariance kernel $\hat{\Sigma}_n(t_1, t_2) : t_1, t_2 \in \mathcal{T}$ which for each (t_1, t_2) is a consistent estimator of $\text{Cov}(r_n(\hat{g}_n(t_1) - g_P(t_1)), r_n(\hat{g}_n(t_2) - g_P(t_2)))$. Finally, the test is defined to be

$$\varphi_{n,\alpha} = 1\{T_n > cv_n(\alpha)\}. \quad (25)$$

A few things are left out in the foregoing discussion. First, the eigenvalues of the matrix $\bar{\Sigma}(t)$ is required to be uniformly bounded away from zero. This is typically not satisfied if $\bar{\Sigma}_n(t)$ is a *consistent* estimator of the variance-covariance matrix of $r_n(\hat{g}_n(t) - g_P(t))$. This is because often the latter matrix can be arbitrarily close to singularity when \mathcal{T} has infinitely many elements. Thus, in order to take advantage of the empirical process result as one does in Proposition 1 and in results like Theorem 1 of Andrews and Shi (2013), the variance-covariance matrix needs to be regularized. In this literature, $\bar{\Sigma}_n(t)$ is often taken as $\hat{\Sigma}_n(t) + \varepsilon \hat{\Sigma}_n$ where $\hat{\Sigma}_n$ is $\hat{\Sigma}_n(t)$ at a particular t where the variance-covariance matrix is not degenerate. The regularization parameter ε is not allowed to converge to zero as $n \rightarrow \infty$. It is interesting to ask what happens to the limit distribution of T_n if a sequence $\varepsilon_n \rightarrow 0$ is used instead of

a fixed ε . For the CvM statistic, this is an open question.⁵

Second, the critical value that I define above is not exactly the same as that proposed in Andrews and Shi (2013) and the subsequent papers. Specifically, the critical values used in that literature are η plus the simulated $1 - \alpha + \eta$ quantile, where η is the so-called infinitesimal constant. The constant is needed to prove the asymptotic level control of the test because the limit distribution of T_n is not derived in that literature, but instead, an approximate distribution that still depends on $r_n g_{P_n}(t)$ is derived. That alone is not sufficient to justify the use of the $1 - \alpha$ quantile of the approximating distribution as the critical value.

Proposition 1 establishes the limit distribution of T_n^{CvM} , and Lemma B.3 in the supplemental appendix of Andrews and Shi (2013) proves the continuity and strict monotonicity of the limit distribution when the following S functions are used

$$S^{\max}(g, \Sigma) = \max_j \frac{\max\{g_j, 0\}^2}{\sigma_j^2}, \text{ and } S^{\text{sum}}(g, \Sigma) = \sum_j \frac{\max\{g_j, 0\}^2}{\sigma_j^2}. \quad (26)$$

In light of Lemma 5 of Andrews and Guggenberger (2010), these results together should obviate the need for the infinitesimal constant. The validity of this conjecture and the breadth of its applicability are open questions.

Third, there is a natural alternative to the CvM test statistic, the Kolmogorov-

⁵A truly studentized Kolmogorov-Smirnov (KS) type statistic is discussed extensively in Scenario 2 where penultimate distributional approximations are derived but not limit distributions. A limit distribution result for that KS statistic is derived in Armstrong (2015) for the conditional moment inequality hypothesis under the data generating processes that the conditional moments are binding ($= 0$) on a measure-zero set. The derivation does not apply when the binding set is not of measure-zero.

Smirnov Statistic:

$$T_n^{KS} = \sup_{t \in \mathcal{T}} S(r_n \hat{g}_n(t), \bar{\Sigma}_n(t)) \quad (27)$$

However, the arguments in the proof of Proposition 1 does not work for the KS statistic because there is no bounded convergence theorem for the supremum operator. Nevertheless, under a fixed P satisfying H_0 in (15) (as opposed to a drifting sequence $\{P_n\}$), Barrett and Donald (2003) establish the limit distribution for an identity weighted KS statistic⁶ for the hypothesis of stochastic dominance. The arguments do not appear to be specific to stochastic dominance or to using identity weighting, but do seem to depend on the fixed P . It is an open question how it extends to drifting P . The drifting P result is necessary to prove an asymptotic size result that is uniform over DGPs, that is, a result like (17) with the $\sup_{P \in \mathcal{P}}$ (and $\sup_{\theta \in \Theta}$ if there is a θ).

On the other hand, a limit distribution result might not be necessary for certain KS-type statistics to justify a uniformly asymptotically valid testing procedure, as shown in the literature on Scenario 2, which I move on to now.

2.2 Scenario 2

In the second scenario, weak convergence fails to hold. This arises, for instance, when $g_P(t)$ includes nonparametrically estimated conditional expectations or if \mathcal{T} is a discrete set that does not have a particular structure.

⁶Barrett and Donald's KS statistic equals the square root of T_n^{KS} with $\bar{\Sigma}_n = I$ and S being S^{\max} defined above

In the conditional moment inequality models described in Examples 1 and 4, if one writes down $g_P(t)$ using the instrumental function approach illustrated in (3), and the set of instrumental functions is appropriately chosen to be rich enough but not too large (ref. Andrews and Shi (2013)), then one can form a $\widehat{g}_n(t)$ to satisfy weak convergence and hence to use the techniques developed for Scenario 1. On the other hand, if one writes down $g_P(t)$ using the nonparametric conditional mean approach described in (4), weak convergence will not hold and one is in Scenario 2.

In Example 3, and in Example 2 when the core determining class is a finite union of half spaces, one can be in Scenario 1, as proved in Sections 7-9 in Andrews and Shi (2017). Otherwise, one is in Scenario 2.

In the literature that works under Scenario 2, the variety of test statistics considered has been much more limited. The predominant choice is a studentized KS statistic, or a supremum (SUP) statistics, with the exception of Lee et al. (2013). I discuss Lee et al. (2013) at the end of this section.

Suppose that $g_P(t)$ is scalar-valued.⁷ The SUP statistic is

$$T_n^{\text{sup}} = \sup_{t \in \mathcal{T}} \frac{r_n \widehat{g}_n(t)}{\widehat{\sigma}_n(t)}, \quad (28)$$

where $\widehat{\sigma}_n(t)$ is a consistent estimator of the asymptotic standard deviation of $r_n(\widehat{g}_n(t) - g_P(t))$.⁸

Let $\widehat{Z}_n(t) = r_n(\widehat{g}_n(t) - g_{P_n}(t))/\widehat{\sigma}_n(t)$ for $t \in \mathcal{T}$. Note that this process may not

⁷It is without loss of generality to assume that $g_P(t)$ is scalar-valued when the SUP statistic is used. This is because each component of $g_P(t)$ enters the statistic separately. Thus, one can simply lump the component index with t .

⁸Note that this statistic is the square root of T_n^{KS} above with S being S^{\max} and with the diagonals of $\bar{\Sigma}_n(t)$ being $\widehat{\sigma}_n(t)^2$.

weakly converge to a tight Gaussian process even when $r_n(\widehat{g}_n(t) - g_{P_n}(t)) : t \in \mathcal{T}$ does if $\widehat{\sigma}_n(t)$ is not bounded away from zero, let alone when $r_n(\widehat{g}_n(t) - g_{P_n}(t)) : t \in \mathcal{T}$ does not converge weakly. Thus, different tools are needed to analyze the distributional behavior of T_n^{sup} .

One tool used in this literature is the strong approximation of $\widehat{Z}_n(t) : t \in \mathcal{T}$ by a penultimate Gaussian process $Z_n^*(t) : t \in \mathcal{T}$:

$$\sup_{t \in \mathcal{T}} |\widehat{Z}_n(t) - Z_n^*(t)| = o_p(\delta_n), \quad (29)$$

where δ_n is a sequence such that $\delta_n \rightarrow 0$. See Chernozhukov et al. (2013). Strong approximation of this type can be established using coupling arguments such as the Yurinskii coupling for sums of independent random variables (Chapter 10 of Pollard (2002)) and for partial sums of mixingales (Li and Liao (2020)). I describe a distributional approximation result for T_n^{sup} based on this tool.

A related tool is the Gaussian approximation of the supreme of $\widehat{Z}_n(t)$, such as

$$\left| \sup_{t \in \mathcal{T}} \widehat{Z}_n(t) - \sup_{t \in \mathcal{T}} Z_n^*(t) \right| = o_p(\delta_n), \quad (30)$$

The verification of such approximation and results based on this tool are derived in Chernozhukov et al. (2014b). Also in this genre, Li et al. (2022) use a Gaussian approximation for the supreme over many projection directions of $\widehat{Z}_n(t)$ to derive a test for superior predictivity. They prove the validity of the Gaussian approximation for the nonparametric series estimator for dependent data. Not requiring full approximation of the entire vector of $\widehat{Z}_n(t)$ allows the dimension of the vector to grow

at a faster rate.

A different type of tools is central limit theorem (CLT) in high dimensions (e.g. Chernozhukov et al. (2017) and Fang and Koike (2024)). Such results are stated with a finite but increasing \mathcal{T} . Specifically, denote the \mathcal{T} used for a given sample size n as \mathcal{T}_n . The set \mathcal{T}_n consists of p_n elements.⁹ Let \vec{Z} denote $(Z(t_1), \dots, Z(t_{p_n}))'$. A CLT in high dimensions is of the form:

$$\sup_{A \in \mathcal{A}_n} |\Pr(\vec{Z}_n \in A) - \Pr(\vec{Z}_n^* \in A)| = o(\xi_n), \quad (31)$$

where \mathcal{A}_n is a collection of measurable sets on \mathbb{R}^{p_n} , such as the collection of hyper-rectangles or of Euclidean balls, and $\xi_n \rightarrow 0$. Such results are used to derive tests for increasingly many moment inequalities, for example, in Chetverikov (2018) and Bai et al. (2022).

One last tool used is the analytical bounds for the tail probability of $\sup_{t \in \mathcal{T}_n} \widehat{Z}_n(t)$, derived using the subadditivity of probability measures and tail bounds for standard normal random variables (ref. Chernozhukov et al. (2019)). These bounds are not as accurate as the CLT approximations, but they do not require bootstrap to implement and can allow p_n to be much larger than n .

Now I describe a result from Chernozhukov et al. (2013). Consider a sequence of data generating processes $\{P_n\}$ satisfying H_0 . Let $\sigma_n(t)$ be the population counterpart of $\widehat{\sigma}_n(t)$. Let $T_n^* = \sup_{t \in \mathcal{T}_n^*} Z_n^*(t)$ where $\mathcal{T}_n^* = \{t \in \mathcal{T} : r_n g_{P_n}(t) \geq -\sigma_n(t) k_n(\gamma_n)\}$, $k_n(\gamma)$ is the $100(1 - \gamma)\%$ quantile of $\sup_{t \in \mathcal{T}} Z_n^*(t)$, and γ_n is a sequence of positive numbers such that $\gamma_n \rightarrow 0$.

⁹The finite set \mathcal{T}_n can be viewed as an increasingly fine discretization of an uncountable \mathcal{T} .

Proposition 2. *Let \mathcal{T} be a compact subset of \mathbb{R}^d and suppose that \mathcal{T}_n^* is compact for all n . Moreover, consider a sequence of DGPs $\{P_n\}$ satisfying H_0 . Suppose that*

- (i) Strong Approximation: *under the sequence $\{P_n\}$ (29) holds,*
- (ii) Anti-Concentration: $\sup_{x \in R} \Pr(|\sup_{t \in \mathcal{T}_n} Z_n^*(t) - x| \leq \delta_n) \rightarrow 0$ *for any compact subset $\mathcal{T}_n \subseteq \mathcal{T}$.*
- (iii) Variance Convergence: $\sup_{t \in \mathcal{T}} \left| \frac{\sigma_n(t)}{\widehat{\sigma}_n(t)} - 1 \right| = o_p(\delta_n k_n(\gamma_n)^{-1})$.

Then, uniformly over $x \in [0, \infty)$,

$$\Pr_{P_n}(T_n^{\sup} \leq x) \geq \Pr(T_n^* \leq x) - o(1). \quad (32)$$

Proof. Consider the derivation

$$\begin{aligned} T_n^{\sup} &= \sup_{t \in \mathcal{T}} [\widehat{Z}_n(t) + r_n g_{P_n}(t) / \widehat{\sigma}_n(t)] \\ &\leq \sup_{t \in \mathcal{T}} [Z_n^*(t) + r_n g_{P_n}(t) / \widehat{\sigma}_n(t)] + \sup_{t \in \mathcal{T}} |Z_n^*(t) - \widehat{Z}_n(t)| \\ &\leq \max \left\{ \sup_{t \in \mathcal{T}_n} \left[Z_n^*(t) + \frac{r_n g_{P_n}(t)}{\widehat{\sigma}_n(t)} \right], \sup_{t \in \mathcal{T} / \mathcal{T}_n^*} \left[Z_n^*(t) - k_n(\gamma_n) \frac{\sigma_n(t)}{\widehat{\sigma}_n(t)} \right] \right\} + o_p(\delta_n) \\ &\leq \max \left\{ \sup_{t \in \mathcal{T}_n} Z_n^*(t), \sup_{t \in \mathcal{T} / \mathcal{T}_n^*} \left[Z_n^*(t) - k_n(\gamma_n) \frac{\sigma_n(t)}{\widehat{\sigma}_n(t)} \right] \right\} + o_p(\delta_n) \\ &\leq \max \left\{ T_n^*, \sup_{t \in \mathcal{T}} Z_n^*(t) - k_n(\gamma_n) + k_n(\gamma_n) \sup_{t \in \mathcal{T}} \left| \frac{\sigma_n(t)}{\widehat{\sigma}_n(t)} - 1 \right| \right\} + o_p(\delta_n), \quad (33) \end{aligned}$$

where the first inequality holds because $\widehat{Z}_n(t) \leq Z_n^*(t) + \sup_{t \in \mathcal{T}_n} |Z_n^*(t) - \widehat{Z}_n(t)|$, the second inequality holds by $\sup_{t \in \mathcal{T}} [\cdot] = \max\{\sup_{t \in \mathcal{T}_n^*} [\cdot], \sup_{t \in \mathcal{T} / \mathcal{T}_n^*} [\cdot]\}$, the definition of $k_n(\gamma_n)$, and Condition (i), the third inequality holds because $g_{P_n}(t) \leq 0$ for all $t \in \mathcal{T}$

since H_0 holds under P_n , and the fourth inequality uses the triangular inequality and the fact that $\mathcal{T}/\mathcal{T}_n^* \subseteq \mathcal{T}$. Therefore, for any $x \geq 0$,

$$\begin{aligned} & \Pr(T_n^{\sup} \leq x) \\ & \geq \Pr(T_n^* + o_p(\delta_n) \leq x) - \Pr\left(\sup_{t \in \mathcal{T}} Z_n^*(t) > k_n(\gamma_n) - o_p(\delta_n) - k_n(\gamma_n) \sup_{t \in \mathcal{T}} \left| \frac{\sigma_n(t)}{\widehat{\sigma}_n(t)} - 1 \right| \right). \end{aligned} \quad (34)$$

Note that $\Pr(T_n^* + o_p(\delta_n) \leq x) \geq \Pr(T_n^* \leq x) - \Pr(|T_n^* - x| \leq o_p(\delta_n))$ and that $\Pr(|\sup_{t \in \mathcal{T}_n^*} Z_n^*(t) - x| \leq o_p(\delta_n)) \rightarrow 0$ uniformly over $x \in R$ by Condition (ii). Similarly the subtracted term in the above display is bounded above by

$$\begin{aligned} & \Pr\left(\sup_{t \in \mathcal{T}} Z_n^* > k_n(\gamma_n)\right) + \Pr\left(\left|\sup_{t \in \mathcal{T}} Z_n^* - k_n(\gamma_n)\right| < o_p(\delta_n) + k_n(\gamma_n) \sup_{t \in \mathcal{T}} \left| \frac{\sigma_n(t)}{\widehat{\sigma}_n(t)} - 1 \right| \right) \\ & = \Pr\left(\sup_{t \in \mathcal{T}} Z_n^* > k_n(\gamma_n)\right) + o(1), \end{aligned} \quad (35)$$

where the equality holds by Condition (ii) and $k_n(\gamma_n) \sup_{t \in \mathcal{T}} \left| \frac{\sigma_n(t)}{\widehat{\sigma}_n(t)} - 1 \right| = o_p(\delta_n)$ (Condition (iii)). By the definition of $k_n(\gamma_n)$, $\Pr(\sup_{t \in \mathcal{T}} Z_n^* > k_n(\gamma_n)) \leq \gamma_n$. Thus it is also $o(1)$. Therefore, $\Pr(T_n^{\sup} \leq x) \geq \Pr(T_n^* \leq x) + o(1)$ uniformly over $x \in [0, \infty)$ \square

Based on a distributional approximation result as that in Proposition 2, we can define a simulated critical value after bounding the contact set \mathcal{T}_n^* , that is, finding an index set $\widehat{\mathcal{T}}_n$ such that

$$\Pr(\mathcal{T}_n^* \subseteq \widehat{\mathcal{T}}_n) \rightarrow 1. \quad (36)$$

Clearly, the least-favorable option, $\widehat{\mathcal{T}}_n = \mathcal{T}$, satisfies the requirement. Chernozhukov et al. (2013) consider a more sophisticated $\widehat{\mathcal{T}}_n$ that can be a much smaller set than \mathcal{T} while still satisfying (36) when the inequalities $g_P(t) \leq 0$ are slack on most of \mathcal{T} . The smaller $\widehat{\mathcal{T}}_n$ serves a similar role as moment selection and can greatly reduce under-rejection comparing to the least-favorable test.

Once $\widehat{\mathcal{T}}_n$ is constructed, one defines the critical value $cv_n^{\text{sup}}(\alpha)$ to be the $100(1 - \alpha)\%$ quantile of $\sup_{t \in \widehat{\mathcal{T}}_n} \tilde{Z}_n(t)$, where $\tilde{Z}_n(t)$ is the random component to be simulated. Its conditional distribution (given data) approximates that of $\sup_{t \in \mathcal{T}_n^*} Z_n^*(t)$. The process $\tilde{Z}_n(t) : t \in \mathcal{T}$ typically is a Gaussian multiplier bootstrap process because it is easier to establish strong approximation results for such processes. With the critical value, the test is

$$\varphi_{n,\alpha}^{\text{sup}} = 1\{T_n^{\text{sup}} > cv_n^{\text{sup}}(\alpha)\}. \quad (37)$$

It is worth discussing the difference between the conditions for Propositions 1 and 2. Unlike Proposition 1, Proposition 2 does not require the weak convergence of the empirical process and it allows unregularized studentization, that is, the $\widehat{\sigma}_n(t)$ does not need to be bounded away from zero. However, Proposition 2 requires an anti-concentration property for the distribution of $\sup_{t \in \mathcal{T}} Z_n^*$ (Condition (ii)). It also needs $k_n(\gamma_n)$ to not grow too fast because otherwise Condition (iii) is violated. This is typically verified through a concentration property, or in other words, a tail probability bound, of the supremum of Gaussian processes. The anti-concentration and concentration properties have been established under general conditions in Chernozhukov et al. (2014a) and Chernozhukov et al. (2015) for $\sup_{t \in \mathcal{T}} Z_n^*(t)$.

Anti-concentration and concentration results are not available for more general functionals of Gaussian processes. This may be the reason that test statistics other than T_n^{sup} are rarely discussed in Scenario 2.

Nevertheless, if CLT results with a convergence rate, such as those in Chernozhukov et al. (2017), are used instead of strong approximation, one may not need to separately establish anti-concentration or concentration of the approximating Gaussian process, which may be a viable direction to develop tests based on other test statistics in Scenario 2.

A different test statistic is considered in Lee et al. (2013) in the setting of conditional moment inequality hypotheses of the form (4), though not using Gaussian approximation tools described above. They propose a test statistic that is a weighted integral of the one-sided norm of the numerator of the Nadaraya-Watson estimator of the conditional mean:

$$T_n^{LSW} = \sum_{j=1}^d \int_{\mathcal{Z}} \max \left\{ 0, \sum_{i=1}^n m(W_i) K_h(Z_i - z) \right\}^p w(z) dz, \quad (38)$$

where $K_h(x) = \frac{1}{h} K(x/h)$, h is a bandwidth that converges to zero as $n \rightarrow \infty$, $K(\cdot)$ is a kernel function, and $w(z)$ is a nonnegative and square-integrable weight function. They use a Poissonization technique to show that $\sigma_n^{-1}(T_n^{LSW} - a_n) \rightarrow_d \mathcal{N}(0, 1)$ where a_n and σ_n are quantities that need to be estimated. They show that the convergence still holds when estimators of σ_n and a_n are used instead. Thus, their test rejects H_0 if $\hat{\sigma}_n^{-1}(T_n^{LSW} - \hat{a}_n) > z_\alpha$ where z_α is the $100(1 - \alpha)\%$ quantile of $\mathcal{N}(0, 1)$.

3 Comparison and Hybridation of Tests

As reviewed in the previous section, for some inequality testing problems, such as conditional moment inequalities, there are quite a few uniformly asymptotically valid tests available in the literature. In some settings, there currently is only one test statistic proposed, that is T_n^{sup} , but new asymptotic theory developed in the literature may allow new tests to be introduced. Then, it is natural to ask what one should do with these options. How should one select a test to use? Should one select just one test or is there a way to combine different tests for better performance?

These questions are intriguing because the hypothesis in (15) is a composite null hypothesis and as such, there is no uniformly most powerful (UMP) test. Moreover, unlike in the context of equalities, the inequality hypothesis is not rotation invariant, and hence it is not obvious that one could restrict attention to a reasonable class of invariant tests and find a UMP one in the class.

In the absence of a UMP test, the relative power of competing tests necessarily depends on the specific alternative under which the data are generated. Since the data generating process is unobserved, selecting among tests becomes a classic decision-under-uncertainty problem. Two decision criteria has been suggested in the literature of inequality testing: a weighted average power (WAP) maximization criterion and a Maximin power criterion. I now review some existing work on this topic.

3.1 Existing Optimality Discussion

The WAP maximization criterion is considered in Chiburis (2008) for moment inequality hypotheses and in Elliott et al. (2015) to develop a general framework for designing a nearly optimal test when testing a null hypothesis in the presence of a nuisance parameter.^{10 11} These two papers do not compare tests in the literature. Instead, they try to design a test that maximizes WAP over all tests that control size.

Assume again that $g_P(t)$ is a scalar for each t . Chiburis’s approach discretizes \mathcal{T} into $\{t_1, \dots, t_K\}$ and the space of \mathbb{R}^K into S regions, and defines a test to be a vector of rejection probabilities on the S regions. He then numerically calculates the WAP of such a test for a given weight over the alternative space, as well as the null rejection probabilities (NRP) for points on a grid for the null space. Finally, he calculates the optimal test by maximizing the WAP subject to the constraint that the NRPs are all less than or equal to a nominal significance level. The procedure can be computationally intensive or even infeasible when K is larger than a handful because many regions (large S) and many grid points need to be considered in order to get a reasonably good approximation.

Elliott et al. (2015) improves on Chiburis (2008) by making use of the Neyman-Pearson (NP) Lemma. When applied in our context, their approach also discretizes \mathcal{T} and $\{\mu_P \in \mathbb{R}^K : \mu_P \leq 0\}$, but instead of solving a linear programming problem,

¹⁰In inequality testing problems, the slackness parameter $\sqrt{n}g_P(t) : t \in \mathcal{T}$ are the nuisance parameters because their values are not uniquely determined by H_0 .

¹¹Andrews and Barwick (2012) also uses the concept of WAP maximization to select tuning parameters in tests for a finite number of inequalities.

it uses an iterative procedure to find the least favorable mixed null (a mixture of the distributions in $\{P : (g_P(t_1), \dots, g_P(t_K)) \leq 0\}$). In each iteration, a NP test is constructed for the weighted mixture null with weights from the previous iteration against a given mixture alternative, and then the mixing weights for the null are updated to give more weights to P 's under which the NP test over-rejects and less weights to P 's under which the NP test under-rejects. The iteration procedure converges to a least favorable mixed null and the associated NP test is the nearly optimal test.

As we can see, the WAP consideration leads to tests different from any of the tests reviewed in the previous section, and thus is not helpful for guiding the selection or combination of those tests. In fact, if one is comfortable with choosing a weight over the alternative space and is committed to maximizing the associated WAP, and if the problem is simple enough so that the numerical algorithms in Chiburis (2008) or Elliott et al. (2015) are feasible, one probably should simply use the nearly optimal test.

Often, though, the weight on the alternative space is difficult to interpret in practice.¹² This compounded with the computational burden of the nearly optimal tests limits their practical appeal.

The maximin power rule has been considered in Chetverikov (2018), Armstrong (2018), and Chernozhukov et al. (2019). Chetverikov (2018) studies the conditional inequality moment hypothesis in (2). He shows that no test with asymptotic size

¹²It should be noted that alternatives local to the null hypothesis are the important ones for WAP consideration because most reasonable tests are consistent against fixed alternatives. Yet, data are inherently unable to provide precise information about local alternatives. Thus, we typically cannot rely on data to determine which alternatives are more relevant.

control can have asymptotic power larger than the size uniformly against alternatives belonging to a Hölder ball with smoothness parameter τ and having a sup distance from the null hypothesis at least a_n . Here $a_n = o((\log(n)/n)^{\tau/(2\tau+d_Z)})$ and d_Z is the number of continuous conditioning variables. The intuition for such a result is that there are always least-rejectable Gaussian P_n 's satisfying the sup distance requirement that is close to a P_0 in the null; the difference of the expectations of a binary statistic (that is, a test) under the P_n 's and under P_0 is bounded by the differences of P_n 's and P_0 . These differences are small when a_n is small. The least-rejectable P_n 's are typically ones such that the null hypothesis is violated on a fast shrinking neighborhood of a Z value as $n \rightarrow \infty$.

Chetverikov (2018) shows that an adaptive sup statistic-based test controls asymptotic size and is consistent against any local alternatives such that $a_n \rightarrow \infty$. In this sense, the adaptive sup statistic-based test is called maximin power rate optimal. Armstrong (2018) shows that the truly studentized KS test is also maximin power rate optimal and that a CvM test is not. Chernozhukov et al. (2019) consider testing a large number of unstructured inequalities. Using similar techniques, they also derive a maximin rate, and shows that the SUP-statistic based tests achieve this rate. As we can see, the maximin power consideration clearly recommends the truly studentized KS-type test statistic.

However, focusing only on the worst case power may be costly when there is no UMP. The cost is the power against other alternatives, in particular, alternatives under which the inequalities are violated on a non-shrinking or slowly shrinking set of \mathcal{T} albeit by a small amount at each point. Andrews and Shi (2013) show that the

tests that they propose have nontrivial power against some alternatives that converge to H_0 at $1/\sqrt{n}$ rate. The tests based on T_n^{sup} do not have nontrivial power against such alternatives.

Alternatives under which the inequalities are violated on a non-shrinking set naturally occur in moment inequality models that are close to point identification. In such models, pairs or groups of inequalities interact and restrict each other to be close to binding on a nontrivial subset of \mathcal{T} and they can be violated at all points in that subset when one moves a parameter value away from the true value.¹³ I illustrate this with a simple interval outcome regression example.

Example 5. *Consider a regression model*

$$Y = \beta_0 + \beta_1 X + \varepsilon, \quad (39)$$

where Y is missing when a missing indicator $M = 1$. Suppose that $\text{Median}(\varepsilon|X) = 0$. Then the model implies the following conditional moment inequality model:

$$\begin{aligned} m^U(X) &:= E[1\{Y \leq \beta_0 + \beta_1 X, M = 0\} + 1\{M = 1\}|X] - 0.5 \geq 0 \\ m^L(X) &:= E[1\{Y \leq \beta_0 + \beta_1 X\}|X] - 0.5 \leq 0, \end{aligned} \quad (40)$$

where the upper bound holds because $0 = E[1\{\varepsilon \leq 0\}|X] - 0.5 = E[1\{\varepsilon \leq 0, M = 0\} + 1\{\varepsilon \leq 0, M = 1\}|X] - 0.5 \leq E[1\{\varepsilon \leq 0, M = 0\} + 1\{M = 1\}|X] - 0.5$. This is a simplified version of the model considered in Blundell et al. (2007) where their Y is the wages of female individuals and their X contain covariates that may affect

¹³See Gandhi et al. (2023) for an example of point identified interval IV regression model.

wages. Missingness is caused by labor market nonparticipation.

For illustration purpose, we let $M = 1\{X+Y < C\}$, where C is the 10th percentile of $X + Y$. Let $\beta_0 = \beta_1 = 1$, and X and ε be independent standard normal random variables. Note that under this DGP, there is very little missing at large values of X . Thus, at those values of X , the two bounds are close to forming an equality. This is illustrated in Figure 1. As illustrated in the figure, as b_1 moves toward β_1 , the lower bound tilts toward the horizontal 0 line instead of parallel-shift down. Thus, the set of X values at which the bound is violated does not shrink much in the process. Andrews and Shi (2013) and Andrews and Shi (2014) contain Monte Carlo results that show the power trade-off of different tests in a model similar to this.

As the literature currently stands, there is no formal discussion of test comparison and combination that does not require one to commit to a prior distribution on the alternative space (as WAP maximization would) or to restrict attention to the most pessimistic case (as Maximin power would). In the next subsection, I suggest a direction toward such a formal discussion.

3.2 Minimax Power Regret: A Suggested Direction

In decision theory, an alternative decision criterion to maximin is minimax regret. Whereas the maximin criterion leads to the safest option by focusing on the worst-case outcome, the minimax regret criterion seeks to avoid ex post disappointment and may favor options with better upside potential if they limit worst-case regret.

In the context of hypothesis tests, we can define the regret of a test relative to a class of tests with guaranteed (asymptotic) level α . Specifically, let $\Phi = \{\varphi_{s,\alpha} : s \in$

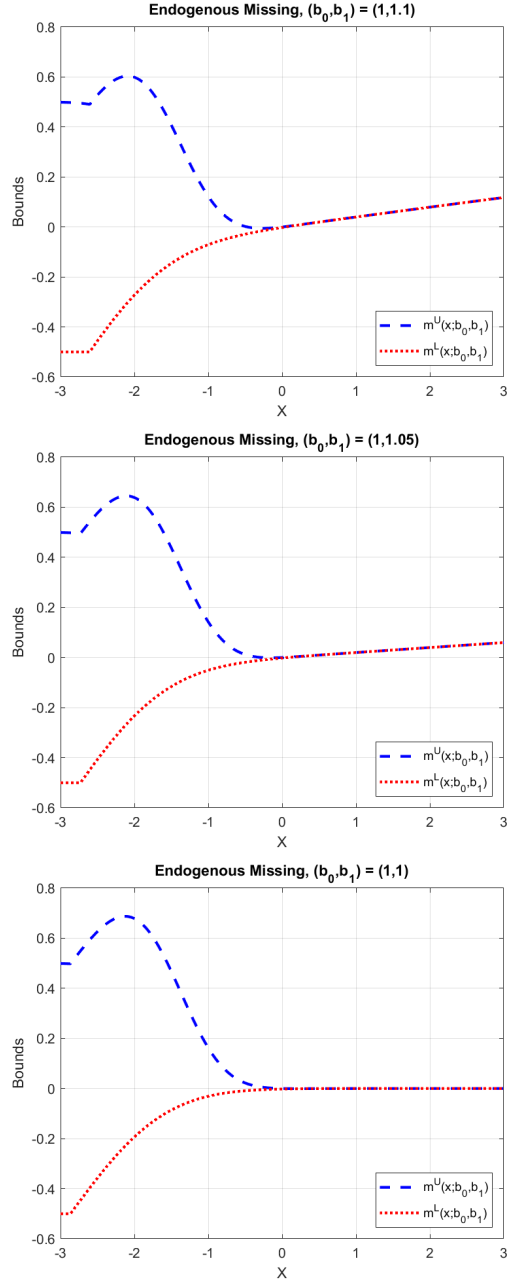


Figure 1: Bounds at Three Values of (b_0, b_1) For Example 5

\mathcal{S} be the class of tests under consideration where \mathcal{S} is a finite index set. Under a given alternative distribution P_1 , we can define the regret of test $\varphi_{s,\alpha}$ as

$$R_\Phi(\varphi_{s,\alpha}, P_1) = \max_{s' \in \mathcal{S}} \mathbb{E}_{P_1} \varphi_{s',\alpha} - \mathbb{E}_{P_1} \varphi_{s,\alpha} \text{ for } s \in \mathcal{S}. \quad (41)$$

Thus, the maximum power regret of test $\varphi_{s,\alpha}$ on the alternative space (denoted \mathcal{P}_1) is

$$\overline{R}_\Phi(\varphi_{s,\alpha}) := \max_{P_1 \in \mathcal{P}_1} R_\Phi(\varphi_{s,\alpha}, P_1). \quad (42)$$

The minimax power regret criterion recommends test s if $\overline{R}(\varphi_{s,\alpha}) \leq \overline{R}(\varphi_{s',\alpha})$ for all $s' \in \mathcal{S}$.

Let us consider an example, where K inequalities are tested:

$$g_P(t) \leq 0 \text{ for } t = 1, \dots, K. \quad (43)$$

Let $\widehat{g}(t) : t = 1, \dots, K$ be estimators of $g_P(t) : t = 1, \dots, K$. Suppose that they are independent and $(\widehat{g}(t) - g_P(t)) \sim \mathcal{N}(0, 1)$. This setting abstracts away from the many challenges in the general case, but suffices for illustrating the potential of the minimax power regret criterion.

Let us consider two test statistics:

$$\begin{aligned} T^{\text{sum}} &= \sum_{t=1}^K ([\widehat{g}(t)]_+)^2, \\ T^{\text{sup}} &= \max_{t=1, \dots, K} \widehat{g}(t), \end{aligned} \quad (44)$$

where $[x]_+ = \max\{0, x\}$. Note that since we assume mutual independence of $\{\widehat{g}(t) : t = 1, \dots, K\}$, T^{sum} is the same as the quasi-likelihood ratio statistic: $T^{\text{qlr}} = \min_{h \geq \mathbf{0}} (\widehat{g} - h)'(\widehat{g} - h)$ where $\widehat{g} = (\widehat{g}(1), \dots, \widehat{g}(K))'$.

Let us consider the least favorable critical values:

$$\begin{aligned} cv_{\alpha}^{\text{sum}}(K) &= F_{\sum_{t=1}^K [Z_t]_+^2}^{-1}(1 - \alpha) \\ cv_{\alpha}^{\text{sup}}(K) &= F_{\max_{t=1, \dots, K} Z_t}^{-1}(1 - \alpha), \end{aligned} \quad (45)$$

where $(Z_1, \dots, Z_K)' \sim N(\mathbf{0}, I_K)$ and $F_X^{-1}(1 - \alpha)$ stands for the $100(1 - \alpha)\%$ quantile of X . Let $\varphi_{\alpha, K}^{\text{sum}} = 1\{T^{\text{sum}} > cv_{\alpha}^{\text{sum}}(K)\}$ and $\varphi_{\alpha, K}^{\text{sup}} = 1\{T^{\text{sup}} > cv_{\alpha}^{\text{sup}}(K)\}$.

Suppose that \mathcal{P}_1 is rich enough so that $\{g_P = (g_P(1), \dots, g_P(K))' : P \in \mathcal{P}_1\} = \mathbb{R}^K$. That is, we do not have prior information about $(g_P(1), \dots, g_P(K))'$ beyond that is provided by its estimator. The next proposition implies that

$$\liminf_{K \rightarrow \infty} \overline{R}_{\{\varphi_{\alpha, K}^{\text{sum}}, \varphi_{\alpha, K}^{\text{sup}}\}}(\varphi_{\alpha, K}^{\text{sum}}) \geq 1 - \alpha \text{ and } \liminf_{K \rightarrow \infty} \overline{R}_{\{\varphi_{\alpha, K}^{\text{sum}}, \varphi_{\alpha, K}^{\text{sup}}\}}(\varphi_{\alpha, K}^{\text{sup}}) \geq 1 - \alpha. \quad (46)$$

That is, both tests have large maximum power regret. Specifically, when only one inequality is violated, in order to have non-trivial power, the SUM test needs the

violation to diverge much faster than the SUP test does because the critical value of the SUM test diverges faster than the SUP test as $K \rightarrow \infty$. On the other hand, when all the inequalities are violated by a similar amount, the SUP test statistic cannot aggregate all the violations as effectively as the SUM test. As a result, in order to have nontrivial power, the SUP test needs the violation (of each inequality) to be much larger than the SUM test does.

However, we can define a hybrid test based on the two as follows

$$\varphi_{\alpha,K}^{\text{hyb}} = \max\{\varphi_{\alpha/2,K}^{\text{sum}}, \varphi_{\alpha/2,K}^{\text{sup}}\}. \quad (47)$$

That is, the level α hybrid test rejects H_0 when either the level $\alpha/2$ SUP test rejects H_0 or the level $\alpha/2$ SUM test rejects H_0 . The proposition also shows that under the sequences of P_1 that either the SUP test or the SUM test has a large power regret, the hybrid test has a zero power regret. In a way, the hybrid test adapts to the data generating process and automatically switches to the more powerful test.¹⁴

Proposition 3. (a) *Let $\{P_K\}$ be such that $g_{P_K} = a_K \mathbf{1}_K$, where $\mathbf{1}_K$ is a vector of ones, and a_K is a positive scalar sequence such that $\sqrt{\log K} a_K \rightarrow 0$ and $\sqrt{K} a_K \rightarrow \infty$ as $K \rightarrow \infty$. Then,*

$$R_{\{\varphi_{\alpha,K}^{\text{sum}}, \varphi_{\alpha,K}^{\text{sup}}\}}(\varphi_{\alpha,K}^{\text{sup}}, P_K) \rightarrow 1 - \alpha. \quad (48)$$

¹⁴It should be noted that hybridizing is not without cost. Adjusting the level from α to $\alpha/2$ guarantees level control of the hybrid test, but it reduces power across the board. The power reduction may not be made up by hybridizing under some data generating processes not described in Proposition 3. Such cost is expected because of the lack of a UMP test and is part of the reason that criteria like maximin power and minimax regret are useful in this context.

(b) Let $\{P_K\}$ be such that $g_{P_K}(1) = a_K$ and $g_{P_K}(t) = 0$ for all $t \neq 1$, where $a_K - \sqrt{2 \log(K)} \rightarrow \infty$ and $a_K = o(K^{1/4})$. Then,

$$R_{\{\varphi_{\alpha,K}^{\text{sum}}, \varphi_{\alpha,K}^{\text{sup}}\}}(\varphi_{\alpha,K}^{\text{sum}}, P_K) \rightarrow 1 - \alpha. \quad (49)$$

(c) Under each of the sequences in parts (a) and (b),

$$R_{\{\varphi_{\alpha,K}^{\text{sum}}, \varphi_{\alpha,K}^{\text{sup}}, \varphi_{\alpha,K}^{\text{hyb}}\}}(\varphi_{\alpha,K}^{\text{hyb}}, P_K) \rightarrow 0. \quad (50)$$

Proof. (a) It suffices to show that $\lim_{K \rightarrow \infty} \max\{\mathbb{E}_{P_K} \varphi_{\alpha,K}^{\text{sup}}, \mathbb{E}_{P_K} \varphi_{\alpha,K}^{\text{sum}}\} = 1$ and that $\lim_{K \rightarrow \infty} \mathbb{E}_{P_K} \varphi_{\alpha,K}^{\text{sup}} = \alpha$. Consider the derivation

$$\begin{aligned} \mathbb{E}_{P_K} \varphi_{\alpha,K}^{\text{sum}} &= \Pr_{P_K} \left(\sum_{t=1}^K ([\widehat{g}(t)]_+)^2 > cv_{\alpha}^{\text{sum}}(K) \right) \\ &= \Pr_{P_K} \left(\sum_{t=1}^K ([Z(t) + a_K]_+)^2 > cv_{\alpha}^{\text{sum}}(K) \right) \\ &= \Pr_{P_K} (W_K(a_K) + e_K > K^{-1/2}(cv_{\alpha}^{\text{sum}}(K) - K/2)), \end{aligned} \quad (51)$$

where $W_K(a) = K^{-1/2} \sum_{t=1}^K ([Z(t) + a]_+^2 - \mathbb{E}[Z(t) + a]_+^2)$ and $e_K = K^{1/2}(\mathbb{E}[Z(t) + a_K]_+^2 - 1/2)$. By Lemma 1, we have

$$W_K(a_K) \rightarrow_d N(0, 5/4). \quad (52)$$

By Lemma 2,

$$e_K = K^{1/2} \left(\frac{2}{\sqrt{2\pi}} a_K + o(a_K) \right) \quad (53)$$

Also, by the definition of $cv_\alpha^{\text{sum}}(K)$, $K^{-1/2}(cv_\alpha^{\text{sum}} - K/2)$ is the $100(1 - \alpha)\%$ quantile of $W_K(0)$. By Lemma 1, $W_K(0) \rightarrow_d N(0, 5/4)$. Thus,

$$K^{-1/2}(cv_\alpha^{\text{sum}} - K/2) \rightarrow \sqrt{5}z_\alpha/2, \quad (54)$$

where z_α is the $100(1 - \alpha)\%$ quantile of $\mathcal{N}(0, 1)$.

Therefore, the last line of (51) equals

$$\begin{aligned} & \Pr_{P_K} \left((a_K K^{1/2})^{-1} W_K(a_K) + 2/\sqrt{2\pi} + o(1) > 5(a_K K^{1/2})^{-1} z_\alpha/2 \right) \\ &= \Pr_{P_K} \left(o_p(1) + 2/\sqrt{2\pi} + o(1) > o(1) \right) \rightarrow 1, \end{aligned} \quad (55)$$

where the equality holds because $\sqrt{K}a_K \rightarrow \infty$.

For $\varphi_{\alpha,K}^{\text{sup}}$, consider the derivation

$$\begin{aligned} \mathbb{E}_{P_K} \varphi_{\alpha,K}^{\text{sup}} &= \Pr_{P_K} \left(\max_{t=1,\dots,T} \widehat{g}(t) > cv_\alpha^{\text{sup}}(K) \right) \\ &= \Pr_{P_K} \left(\max_{t=1,\dots,K} (Z(t) + g_{P_K}(t)) > cv_\alpha^{\text{sup}}(K) \right) \\ &= \Pr_{P_K} \left(a_K + \max_{t=1,\dots,K} Z(t) > cv_\alpha^{\text{sup}}(K) \right). \end{aligned} \quad (56)$$

By the Gaussian extreme value theorem, $b_K(\max_{t=1,\dots,K} Z(t) - c_K) \rightarrow_d G$, where G has the cumulative distribution function $\exp(-\exp(-x))$, $b_K = \sqrt{2 \log K}$ and

$c_K = \sqrt{2 \log(K)} - \frac{(\log \log(K) + \log(4\pi))}{2\sqrt{2 \log(K)}}$. Thus, the last line of (56) is equal to,

$$\Pr_{P_K} \left(b_K a_K + b_K \left(\max_{t=1, \dots, K} Z(t) - c_K \right) > b_K (cv_\alpha^{\sup}(K) - c_K) \right). \quad (57)$$

By design, $b_K a_K \rightarrow 0$. Thus,

$$b_K a_K + b_K \left(\max_{t=1, \dots, T} Z(t) - c_K \right) \rightarrow_d G. \quad (58)$$

Also by the definition of $cv_\alpha^{\sup}(K)$, we have that $b_K (cv_\alpha^{\sup}(K) - c_K)$ is the $100(1 - \alpha)\%$ quantile of $b_K (\max_{t=1, \dots, T} Z(t) - c_K)$. The distribution of G is continuous and strictly increasing. Thus, $b_K (cv_\alpha^{\sup}(K) - c_K)$ converges in probability to the $100(1 - \alpha)\%$ quantile of G . This and (58) together implies that

$$\mathbb{E}_{P_K} \varphi_{\alpha, K}^{\sup} \rightarrow \alpha. \quad (59)$$

(b) It is sufficient to show that $\lim_{K \rightarrow \infty} \max\{\mathbb{E}_{P_K} \varphi_{\alpha, K}^{\sup}, \mathbb{E}_{P_K} \varphi_{\alpha, K}^{\text{sum}}\} = 1$ and that

$\lim_{K \rightarrow \infty} \mathbb{E}_{P_K} \varphi_{\alpha, K}^{\text{sum}} = \alpha$. Consider the derivation:

$$\begin{aligned} & \mathbb{E}_{P_K} \varphi_{\alpha, K}^{\sup} \\ &= \Pr_{P_K} \left(\max\{Z(1) + a_K, \max_{t=2, \dots, K} Z(t)\} > cv_\alpha^{\sup}(K) \right) \\ &= \Pr_{P_K} \left(\max\{(Z(1) + a_K - c_K), (\max_{t=2, \dots, K} Z(t) - c_K)\} > (cv_\alpha^{\sup}(K) - c_K) \right) \\ &\geq \Pr_{P_K} (Z(1) + a_K - c_K > (cv_\alpha^{\sup}(K) - c_K)), \end{aligned} \quad (60)$$

where $c_K = \sqrt{2 \log(K) - \frac{(\log \log(K) + \log(4\pi))}{2\sqrt{2 \log(K)}}}$. In part (a), we have argued that $b_K(cv_\alpha^{\text{sup}}(K) - c_K)$ converges in probability to the $100(1 - \alpha)$ quantile of G , where $b_K = \sqrt{2 \log(K)}$.

Thus,

$$cv_\alpha^{\text{sup}}(K) - c_K \rightarrow_p 0. \quad (61)$$

This and $a_K - \sqrt{2 \log K} \rightarrow \infty$ together imply that the last line of (60) converges to

1. Thus,

$$\mathbb{E}_{P_K} \varphi_{\alpha, K}^{\text{sup}} \rightarrow 1. \quad (62)$$

For $\varphi_{\alpha, K}^{\text{sum}}$, consider the derivation:

$$\begin{aligned} & \mathbb{E}_{P_K} \varphi_{\alpha, K}^{\text{sum}} \\ &= \Pr_{P_K} \left([Z(1) + a_K]_+^2 + \sum_{t=2}^K [Z(t)]_+^2 > cv_\alpha^{\text{sum}}(K) \right) \\ &= \Pr_{P_K} \left(\frac{[Z(1) + a_K]_+^2}{\sqrt{K}} - \frac{1}{2\sqrt{K}} + \frac{\sqrt{K-1}}{\sqrt{K}} W_{K-1}(0) > \frac{1}{\sqrt{K}} \left(cv_\alpha^{\text{sum}}(K) - \frac{K}{2} \right) \right), \end{aligned} \quad (63)$$

Since $a_K = o(K^{1/4})$ and $Z(1) \sim \mathcal{N}(0, 1)$, we have that $\frac{[Z(1) + a_K]_+^2}{\sqrt{K}} = o_p(1)$. In the proof of part (a), we have shown that $W_{K-1}(0) \rightarrow_d N(0, 5/4)$ and $\frac{1}{\sqrt{K}} (cv_\alpha^{\text{sum}}(K) - \frac{K}{2}) \rightarrow \sqrt{5}z_\alpha/2$. Combining these facts, we deduce that the last line of (63) converges to

$$\Pr(N(0, 5/4) > \sqrt{5}z_\alpha/2) = \alpha. \quad (64)$$

Therefore, $\mathbb{E}_{P_K} \varphi_{\alpha, K}^{\text{sum}} \rightarrow \alpha$, concluding the proof of part (b).

(c) It suffices to show that $\lim_{K \rightarrow \infty} \mathbb{E}_{P_K} \varphi_{\alpha, K}^{\text{hyb}} = 1$ under each of the sequences in part (a) and in part (b).

Consider a sequence satisfying the conditions in part (a). Then, by the arguments in the proof of part (a),

$$\lim_{K \rightarrow \infty} \mathbb{E}_{P_K} \varphi_{\alpha/2, K}^{\text{sum}} = 1. \quad (65)$$

No arguments in the proof there needs to be changed except that α is replaced by $\alpha/2$. By definition, $\mathbb{E}_{P_K} \varphi_{\alpha, K}^{\text{hyb}} \geq \mathbb{E}_{P_K} \varphi_{\alpha/2, K}^{\text{sum}}$. Therefore, $\mathbb{E}_{P_K} \varphi_{\alpha, K}^{\text{hyb}} \rightarrow 1$.

Consider a sequence satisfying the conditions in part (b). Then the arguments in the proof of that part show that

$$\lim_{K \rightarrow \infty} \mathbb{E}_{P_K} \varphi_{\alpha/2, K}^{\text{sup}} = 1. \quad (66)$$

No modification to the arguments is needed except that α is replaced by $\alpha/2$. Moreover, by definition, $\mathbb{E}_{P_K} \varphi_{\alpha, K}^{\text{hyb}} \geq \mathbb{E}_{P_K} \varphi_{\alpha/2, K}^{\text{sup}}$. Therefore, $\mathbb{E}_{P_K} \varphi_{\alpha, K}^{\text{hyb}} \rightarrow 1$. \square

I now investigate the implication of the proposition in a simulation exercise. I first simulate these rejection probabilities at h 's that have an Euclidean distance of 4 to the null space $\{h \in \mathbb{R}^2 : h \leq 0\}$.¹⁵ In particular, I consider three sets of points on \mathbb{R}^2 that are 4-away from the null space:

(a) $g_P(t) = 4/\sqrt{K-1}$ for $t = 2, \dots, K$, and $g_P(1)$ takes values on a grid on the

¹⁵The number 4 is somewhat arbitrarily chosen. Other distances give similar results, with smaller distance showing less dramatic contrasts between the tests.

interval $[-3, 0]$;

(b) $g_P(1)$ takes values on a grid on the interval $(0, 4]$, and $g_P(t) = \frac{\sqrt{4 - g_P(1)^2}}{\sqrt{K-1}}$ for $t = 2, \dots, K$;

(c) $g_P(1) = 4$ and $g_P(2)$ takes values on a grid on the interval $(-3/\sqrt{K-1}, 0)$.

The points can be plotted as a curve on the space of $(g_P(1), \|g_P(II)\|)'$ where $g_P(II) = (g_P(2), \dots, g_P(K))'$. This curve is shown in the upper left graph in Figure 2. The graph makes it clear that each point on the curve corresponds to the angle that the ray from $(0, 0)$ to $(-\infty, 0)$ needs to rotate (clockwise) to reach the point. Denote that angle γ . The segment on which $\gamma > \pi$ corresponds to the DGPs where only one inequality is violated and the rest are increasingly slack as we increase γ . The segment on which $\gamma \in (\pi/2, \pi)$ corresponds to the DGPs where both inequalities are violated, and the segment on which $\gamma < \pi/2$ corresponds to the DGPs where $J-1$ inequalities are violated while one inequality is increasingly slack as γ decreases.

We can plot the rejection probabilities of each test for each point on the curve against γ . These plots are shown in the three other graphs in Figure 2. Consistent with the conclusions of Proposition 3, the SUP test has large regret at small γ 's (when a large number of inequalities are violated), while the SUM test has large regret at large γ 's (when only one inequality is violated). The hybrid test seems to achieve a balance between the two and has smaller regret than both on most γ 's when $K = 5$ and on all γ 's when $K = 10$ and 50 . The advantage of the hybrid test is more and more clear as K increases.

The power lines are slices on the power surface and thus might not reflect the

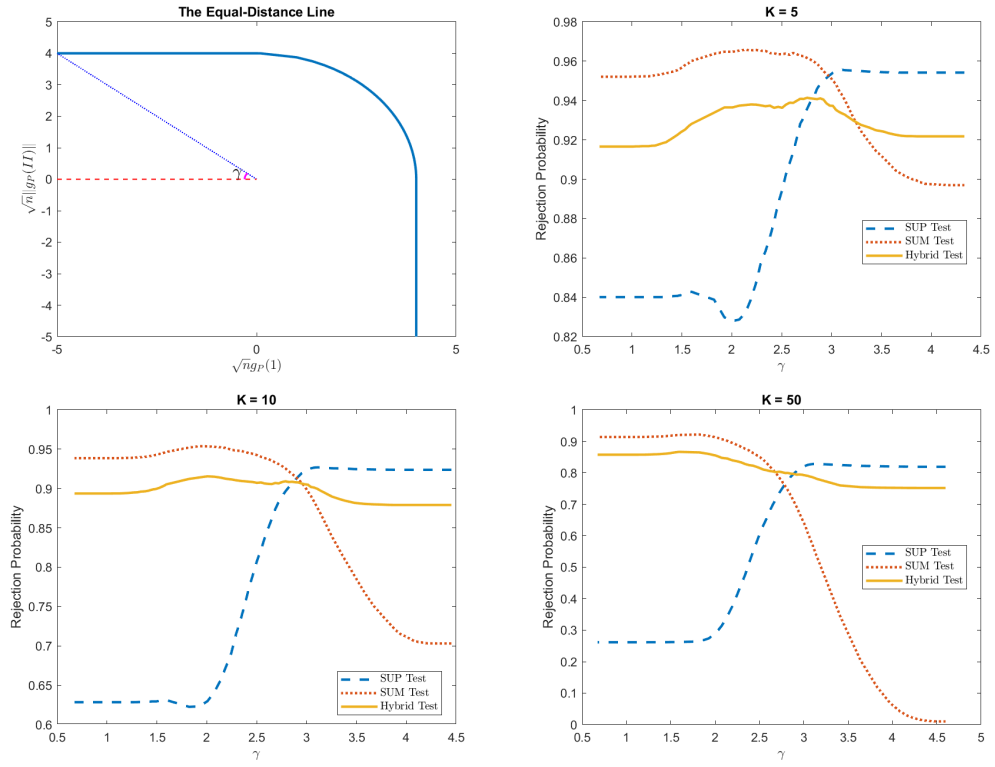


Figure 2: Power Lines for the SUP, SUM, and Hybrid Tests

regrets at all P_1 . I also construct the heat map of the regret of each test with respect to the 3-test collection. The heat maps are shown in Figure 3. As we can see from the heat bar on the right of each graph, the maximum regret of the hybrid test is much lower than the maximum regrets of the SUM and the SUP tests.

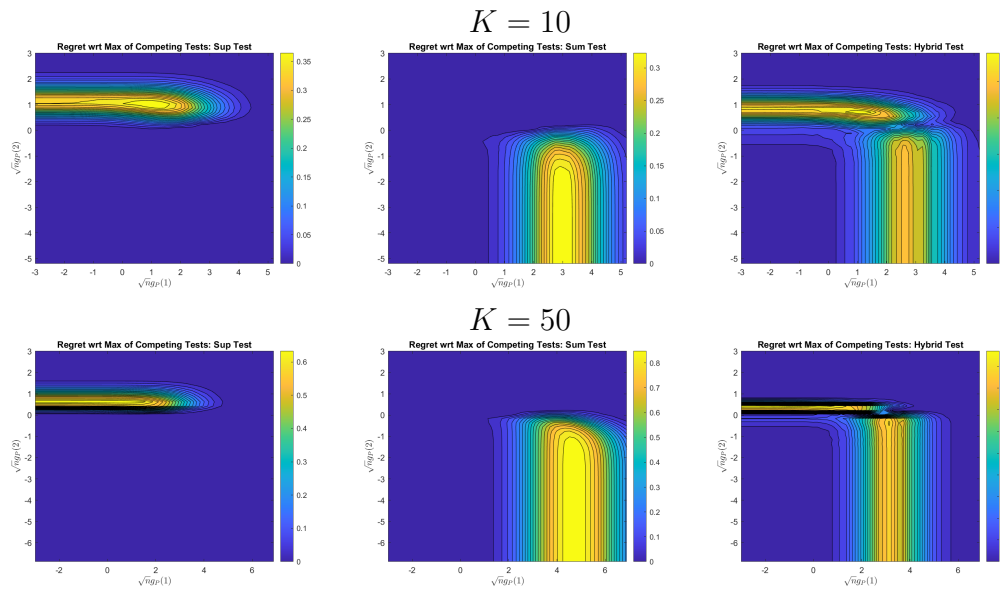


Figure 3: Heat Maps for the regret of the SUP, SUM, and Hybrid Tests against $(g_P(1), g_P(2))$ ($g_P(t) = g_P(2)$ for $t > 2$)

Proposition 3 and the simulation exercise suggest that the minimax power regret criterion may be a useful tool for test comparison. It has the potential to provide justification for various test hybridation procedures.

4 Conclusion

This paper surveys the literature on inference based on infinitely many inequalities. The emphasis is on various testing procedures for infinite dimensional inequality

hypotheses. I summarized the two main branches of the literature on testing such hypotheses and explained the differences between their theoretical foundations. I pointed out a few directions for future research. Then I reviewed the two optimality concepts used in the literature, discussed their limitations, suggested a new concept, and connected the new optimality concept with test hybridation.

Appendix: Useful Lemmas

Lemma 1. *Let $\{Z(t)\}_{t=1}^K$ be i.i.d. $\mathcal{N}(0, 1)$ random variables, and let $a_K \rightarrow 0$ as $K \rightarrow \infty$. Define the function $f_a(z) = [z + a]_+^2 = \max(z + a, 0)^2$.*

$$S_K := \frac{1}{\sqrt{K}} \sum_{t=1}^K (f_{a_K}(Z(t)) - \mathbb{E}[f_{a_K}(Z(t))]) \xrightarrow{d} N(0, 5/4).$$

Proof. Since $a_K \rightarrow 0$, for any fixed $z \in \mathbb{R}$, we have

$$f_{a_K}(z) = \max(z + a_K, 0)^2 \rightarrow \max(z, 0)^2 := f(z), \quad (67)$$

so $f_{a_K}(z) \rightarrow f(z)$ pointwise. Note that for all $a \in \mathbb{R}$,

$$f_a(z)^2 = \max(z + a, 0)^4 \leq (z + a)^4 \leq 8z^4 + 8a^2. \quad (68)$$

Thus,

$$\mathbb{E}[f_{a_K}(Z(t))^2] \leq 8\mathbb{E}[Z(t)^4] + 8a_K^4 \rightarrow 8\mathbb{E}[Z(t)^4] < \infty, \quad (69)$$

since $Z(t) \sim \mathcal{N}(0, 1)$. Therefore, the random variables $f_{a_K}(Z(t))$ have uniformly bounded second moments. This implies the Lindeberg condition holds and ensures uniform integrability.

For each fixed K , the terms $X_t^{(K)} := f_{a_K}(Z(t)) - \mathbb{E}[f_{a_K}(Z(t))]$ are i.i.d. mean-zero random variables with variance $\sigma_K^2 := \text{Var}(f_{a_K}(Z))$ and $a_K \rightarrow 0$. Below we show

that $\sigma_K^2 \rightarrow 5/4$. Hence, by the Lindeberg–Lévy Central Limit Theorem, we have:

$$S_K = \frac{1}{\sqrt{K}} \sum_{t=1}^K X_t^{(K)} \xrightarrow{d} N(0, 5/4).$$

Now we show that $\sigma_K^2 \rightarrow 5/4$. Since $f_{a_K}(Z) \rightarrow f(Z) := [Z]_+^2$ pointwise, and $\{f_{a_K}(Z)\}$ is uniformly integrable in L^2 , it follows that

$$\sigma_K^2 = \text{Var}(f_{a_K}(Z)) \rightarrow \text{Var}(f(Z)) = \text{Var}([Z]_+^2) = \frac{5}{4}.$$

This concludes the proof. □

Lemma 2. *Let $Z \sim \mathcal{N}(0, 1)$, and define $f(a) := \mathbb{E}[\max(Z + a, 0)^2]$. Then the function is twice-continuously differentiable and*

$$f'(0) = \frac{2}{\sqrt{2\pi}}.$$

Proof. We first write

$$f(a) = \mathbb{E}[\max(Z + a, 0)^2] = \int_{-a}^{\infty} (z + a)^2 \phi(z) dz, \tag{70}$$

where $\phi(z) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}$ is the standard normal density. Differentiate under the

integral sign using Leibniz's rule:

$$\begin{aligned}
f'(a) &= \frac{d}{da} \int_{-a}^{\infty} (z+a)^2 \phi(z) dz \\
&= -(z+a)^2 \phi(z) \Big|_{z=-a} + \int_{-a}^{\infty} \frac{\partial}{\partial a} (z+a)^2 \phi(z) dz \\
&= \int_{-a}^{\infty} 2(z+a) \phi(z) dz
\end{aligned} \tag{71}$$

Similarly, we can obtain the second derivative:

$$f''(a) = \int_{-a}^{\infty} 2\phi(z) dz = 2\Phi(a), \tag{72}$$

where $\Phi(\cdot)$ stands for the cumulative distribution function of $\mathcal{N}(0, 1)$. Therefore, the function is twice-continuously differentiable.

To evaluate $f'(0)$, change variables to $u = z + a$, so that

$$f'(a) = 2 \int_0^{\infty} u \cdot \phi(u-a) du. \tag{73}$$

Setting $a = 0$, we obtain:

$$\begin{aligned}
f'(0) &= 2 \int_0^{\infty} u \phi(u) du \\
&= 2 \cdot \mathbb{E}[Z \cdot \mathbf{1}_{Z>0}] = 2 \int_0^{\infty} z \cdot \frac{1}{\sqrt{2\pi}} e^{-z^2/2} dz = \frac{2}{\sqrt{2\pi}}.
\end{aligned} \tag{74}$$

□

References

- Andrews, D. W. and Guggenberger, P. (2010). Asymptotic size and a problem with subsampling and with the m out of n bootstrap. *Econometric Theory*, 2:426–468.
- Andrews, D. W. K. and Barwick, P. J. (2012). Inference for parameters defined by moment inequalities: A recommended moment selection procedure. *Econometrica*, 80(6):2805–2826.
- Andrews, D. W. K. and Shi, X. (2013). Inference based on conditional moment inequalities. *Econometrica*, 81(2):609–666.
- Andrews, D. W. K. and Shi, X. (2014). Nonparametric inference based on conditional moment inequalities. *Journal of Econometrics*, 179(1):31–45.
- Andrews, D. W. K. and Shi, X. (2017). Inference based on many conditional moment inequalities. *Journal of Econometrics*, 196(2):275–287.
- Armstrong, T. B. (2015). Asymptotically exact inference in conditional moment inequality models. *Journal of Econometrics*, 186(1):51–65.
- Armstrong, T. B. (2018). On the choice of test statistic for conditional moment inequalities. *Journal of Econometrics*, 203(2):241–255.
- Bai, Y., Santos, A., and Shaikh, A. M. (2022). A two-step method for testing many moment inequalities. *Journal of Business & Economic Statistics*, 40(3):1070–1080.
- Barrett, G. F. and Donald, S. G. (2003). Consistent tests for stochastic dominance. *Econometrica*, 71(1):71–104.

- Beare, B. K. and Moon, J.-M. (2015). Nonparametric tests of density ratio ordering. *Econometric Theory*, 31(3):471–492.
- Beare, B. K. and Shi, X. (2019). An improved bootstrap test of density ratio ordering. *Econometrics and Statistics*, 10:9–26.
- Beresteanu, A., Molchanov, I., and Molinari, F. (2011). Sharp identification regions in models with convex moment predictions. *Econometrica*, 79(6):1785–1821.
- Blundell, R., Gosling, A., Ichimura, H., and Meghir, C. (2007). Changes in the distribution of male and female wages accounting for employment composition using bounds. *Econometrica*, 75(2):323–363.
- Canay, I., Illanes, G., and Velez, A. (2023). A user’s guide for inference in models defined by moment inequalities. *Journal of Econometrics*, forthcoming.
- Canay, I. A. and Shaikh, A. M. (2017). Practical and theoretical advances in inference for partially identified models. In Honoré, B., Pakes, A., Piazzesi, M., and Samuelson, L., editors, *Advances in Economics and Econometrics: Eleventh World Congress*, pages 271–306. Econometric Society Monographs, Cambridge University Press.
- Carolan, C. A. and Tebbs, J. M. (2005). Nonparametric tests for and against likelihood ratio ordering in the two-sample problem. *Biometrika*, 92(1):159–171.
- Chernozhukov, V., Chetverikov, D., and Kato, K. (2014a). Anti-concentration and honest, adaptive confidence bands. *Annals of Statistics*, 42(5):1787–1818.

- Chernozhukov, V., Chetverikov, D., and Kato, K. (2014b). Gaussian approximation of suprema of empirical processes. *Annals of Statistics*, 42(4):1564–1597.
- Chernozhukov, V., Chetverikov, D., and Kato, K. (2015). Comparison and anti-concentration bounds for maxima of gaussian random vectors. *Probability Theory and Related Fields*, 162(1-2):47–70.
- Chernozhukov, V., Chetverikov, D., and Kato, K. (2017). Central Limit Theorems and Bootstrap in High Dimensions. *The Annals of Probability*, 45(4):2309–2352.
- Chernozhukov, V., Chetverikov, D., and Kato, K. (2019). Inference on causal and structural parameters using many moment inequalities. *The Review of Economic Studies*, 86(5):1867–1900.
- Chernozhukov, V., Hong, H., and Tamer, E. (2007). Estimation and confidence regions for parameter sets in econometric models. *Econometrica*, 75(5):1243–1284.
- Chernozhukov, V., Lee, S., and Rosen, A. M. (2013). Intersection bounds: estimation and inference. *Econometrica*, 81(2):667–737.
- Chesher, A. and Rosen, A. M. (2017). Generalized instrumental variable models. *Econometrica*, 85(3):959–989.
- Chesher, A., Rosen, A. M., and Smolinski, K. (2013). An instrumental variable model of multiple discrete choice. *Quantitative Economics*, 4(2):157–196.
- Chetverikov, D. (2018). Adaptive tests of conditional moment inequalities. *Econometric Theory*, 34(1):186–227.

- Chetverikov, D., Wilhelm, D., and Kim, D. (2021). An adaptive test of stochastic monotonicity. *Econometric Theory*, 37(3):495–536.
- Chiburis, R. C. (2008). Approximately most powerful tests for moment inequalities. Working paper, Princeton University, Department of Economics, Princeton, New Jersey. Unpublished manuscript.
- Ciliberto, F. and Tamer, E. (2009). Market structure and multiple equilibria in airline markets. *Econometrica*, 77(6):1791–1828.
- Delgado, M. A. and Escanciano, J. C. (2013). Conditional stochastic dominance testing. *Journal of Business & Economic Statistics*, 31(1):16–28.
- Donald, S. G. and Hsu, Y. (2016). Improving the power of tests of stochastic dominance. *Econometric Reviews*, 35(4):553–585.
- Elliott, G., Müller, U. K., and Watson, M. W. (2015). Nearly optimal tests when a nuisance parameter is present under the null hypothesis. *Econometrica*, 83(2):771–811.
- Fack, G., Grenet, J., and He, Y. (2019). Beyond truth-telling: Preference estimation with centralized school choice and college admissions. *American Economic Review*, 109(4):1486–1529.
- Fang, X. and Koike, Y. (2024). Large-dimensional central limit theorem with fourth moment error bounds on convex sets and balls. *The Annals of Applied Probability*, 34:2065–2106.

- Galichon, A. and Henry, M. (2011). Set identification in models with multiple equilibria. *Review of Economic Studies*, 78(4):1264–1298.
- Gandhi, A., Lu, Z., and Shi, X. (2023). Estimating demand for differentiated products with zeroes in market share data. *Quantitative Economics*, 14(2):381–418.
- Ghosal, S., Sen, A., and van der Vaart, A. W. (2000). Testing monotonicity of regression. *Annals of Statistics*, 28(4):1054–1082.
- Giacomini, R. and Kitagawa, T. (2021). Robust bayesian inference for set-identified models. *Econometrica*, 89(4):1519–1556.
- He, Y. (2017). Gaming the boston school choice mechanism in beijing. Technical Report 15-551, Toulouse School of Economics.
- Ho, K. and Rosen, A. M. (2017). Partial identification in applied research: Benefits and challenges. In Honoré, B., Pakes, A., Piazzesi, M., and Samuelson, L., editors, *Advances in Economics and Econometrics: Eleventh World Congress*, pages 307–359. Econometric Society Monographs, Cambridge University Press.
- Hsu, Y., Liu, C., and Shi, X. (2019). Testing generalized regression monotonicity. *Econometric Theory*, 35(6):1146–1200.
- Iaryczower, M., Shi, X., and Shum, M. (2018). Can words get in the way? the effect of deliberation in collective decision making. *Journal of Political Economy*, 126(2):688–734.

- Kalouptsi, M., Kitamura, Y., Lima, L., and Souza-Rodrigues, E. A. (2020). Partial identification and inference for dynamic models and counterfactuals. NBER Working Paper 26761, National Bureau of Economic Research.
- Kline, B., Pakes, A., and Tamer, E. (2021). Moment inequalities and partial identification in industrial organization. In Ho, K., Hortacsu, A., and Lizzeri, A., editors, *Handbook of Industrial Organization*, volume 4, pages 345–431. Elsevier.
- Lee, S., Linton, O. B., and Whang, Y. (2009). Testing for stochastic monotonicity. *Econometrica*, 77(2):585–602.
- Lee, S., Song, K., and Whang, Y. (2013). Testing functional inequalities. *Journal of Econometrics*, 172(1):14–32.
- Li, J. and Liao, Z. (2020). “uniform nonparametric inference for time series. *Journal of Econometrics*, 129:28–51.
- Li, J., Liao, Z., and Quaadvlieg, R. (2022). Conditional superior predictive ability. *Review of Economic Studies*, 89(2):843–875.
- Manski, C. F. (1989). Anatomy of the selection problem. *Journal of Human Resources*, 24(3):343–360.
- Manski, C. F. (1990). Nonparametric bounds on treatment effects. *American Economic Review, Papers and Proceedings*, 80(2):319–323.
- Manski, C. F. (1993). Identification problems in the social sciences. *Sociological Methodology*, 23:1–56.

- Molinari, F. (2020). Microeconometrics with partial identification. In Durlauf, S. N., Hansen, L. P., Heckman, J. J., and Matzkin, R. L., editors, *Handbook of Econometrics*, volume 7A, pages 355–486. Elsevier.
- Morales, E., Sheu, G., and Zahler, A. (2019). Extended gravity. *Review of Economic Studies*, 86(6):2668–2712.
- Pakes, A., Porter, J. R., Ho, K., and Ishii, J. (2015). Moment inequalities and their application. *Econometrica*, 83(1):315–334.
- Phillips, P. C. B. (1989). Partially identified econometric models. *Econometric Theory*, 5(2):181–240.
- Pollard, D. (2002). *A User’s Guide to Measure Theoretic Probability*, volume 8 of *Cambridge Series in Statistical and Probabilistic Mathematics*. Cambridge University Press.
- Seo, J. (2018). Tests of stochastic monotonicity with improved size and power. *Journal of Econometrics*, 207(1):53–70.
- Sheng, S. (2020). A structural econometric analysis of network formation games through subnetworks. *Econometrica*, 88(5):1829–1858.
- van der Vaart, A. W. and Wellner, J. A. (1996). *Weak Convergence and Empirical Processes: With Applications to Statistics*. Springer Series in Statistics. Springer, New York.