

# Valid Post-Selection and Post-Regularization Inference: An Elementary, General Approach

Victor Chernozhukov,<sup>1</sup> Christian Hansen,<sup>2</sup>  
and Martin Spindler<sup>3</sup>

<sup>1</sup>Department of Economics, Massachusetts Institute of Technology, Cambridge, Massachusetts 02142; email: vchern@mit.edu

<sup>2</sup>University of Chicago Booth School of Business, Chicago, Illinois 60637; email: chistian.hansen@chicagobooth.edu

<sup>3</sup>Munich Center for the Economics of Aging, 80799 Munich, Germany; email: spindler@mea.mpisoc.mpg.de

Annu. Rev. Econ. 2015. 7:649–88

The *Annual Review of Economics* is online at [economics.annualreviews.org](http://economics.annualreviews.org)

This article's doi:  
10.1146/annurev-economics-012315-015826

Copyright © 2015 by Annual Reviews.  
All rights reserved

JEL codes: C18, C55, C26

## Keywords

Neyman, orthogonalization,  $C(\alpha)$  statistics, optimal instrument, optimal score, optimal moment, efficiency, optimality

## Abstract

We present an expository, general analysis of valid post-selection or post-regularization inference about a low-dimensional target parameter in the presence of a very high-dimensional nuisance parameter that is estimated using selection or regularization methods. Our analysis provides a set of high-level conditions under which inference for the low-dimensional parameter based on testing or point estimation methods will be regular despite selection or regularization biases occurring in the estimation of the high-dimensional nuisance parameter. A key element is the use of so-called immunized or orthogonal estimating equations that are locally insensitive to small mistakes in the estimation of the high-dimensional nuisance parameter. As an illustration, we analyze affine-quadratic models and specialize these results to a linear instrumental variables model with many regressors and many instruments. We conclude with a review of other developments in post-selection inference and note that many can be viewed as special cases of the general encompassing framework of orthogonal estimating equations provided in this article.

## 1. INTRODUCTION

Analysis of high-dimensional models, models in which the number of parameters to be estimated is large relative to the sample size, is becoming increasingly important. Such models arise naturally in readily available high-dimensional data, which have many measured characteristics available per individual observation, as in, for example, large survey data sets, scanner data, and text data. Such models also arise naturally even in data with a small number of measured characteristics in situations where the exact functional form with which the observed variables enter the model is unknown. Examples of this scenario include semi-parametric models with nonparametric nuisance functions. More generally, models with many parameters relative to the sample size often arise when attempting to model complex phenomena.

The key concept underlying the analysis of high-dimensional models is that regularization, such as model selection or shrinkage of model parameters, is necessary if one is to draw meaningful conclusions from the data. For example, the need for regularization is obvious in a linear regression model with the number of right-hand-side variables greater than the sample size but arises far more generally in any setting in which the number of parameters is not small relative to the sample size. Given the importance of the use of regularization in analyzing high-dimensional models, it is then important to explicitly account for the impact of this regularization on the behavior of estimators if one wishes to accurately characterize their finite-sample behavior. The use of such regularization techniques may easily invalidate conventional approaches to inference about model parameters and other interesting target parameters. A major goal of this article is to present a general, formal framework that provides guidance about setting up estimating equations and making appropriate use of regularization devices so that inference about parameters of interest will remain valid in the presence of data-dependent model selection or other approaches to regularization.

It is important to note that understanding estimators' behavior in high-dimensional settings is also useful in conventional low-dimensional settings. As noted above, dealing formally with high-dimensional models requires that one explicitly account for model selection or other forms of regularization. Providing results that explicitly account for this regularization then allows us to accommodate and coherently account for the fact that low-dimensional models estimated in practice are often the result of specification searches. As in the high-dimensional setting, failure to account for this variable selection will invalidate the usual inference procedures, whereas the approach that we outline will remain valid and can easily be applied in conventional low-dimensional settings.

The chief goal of this article is to offer a general framework that encompasses many existing results regarding inference on model parameters in high-dimensional models. The encompassing framework we present and the key theoretical results are new, although they are clearly heavily influenced and foreshadowed by previous, more specialized results. As an application of the framework, we also present new results on inference in a reasonably broad class of models, termed affine-quadratic models, that includes the usual linear model and linear instrumental variables (IV) model and then apply these results to provide new ones regarding post-regularization inference on the parameters on endogenous variables in a linear IV model with very many instruments and controls (and also allowing for some misspecification). We also provide a discussion of previous research that aims to highlight that many existing results fall within the general framework.

Formally, we present a series of results for obtaining valid inferential statements about a low-dimensional parameter of interest,  $\alpha$ , in the presence of a high-dimensional nuisance parameter,  $\eta$ .

The general approach we offer relies on two fundamental elements. First, it is important that estimating equations used to draw inferences about  $\alpha$  satisfy a key orthogonality or immunization condition.<sup>1</sup> For example, when estimation and inference for  $\alpha$  are based on the empirical analog of a theoretical system of equations

$$M(\alpha, \eta) = 0,$$

we show that setting up the equations in a manner such that the orthogonality or immunization condition

$$\partial_{\eta} M(\alpha, \eta) = 0$$

holds is an important element in providing an inferential procedure for  $\alpha$  that remains valid when  $\eta$  is estimated using regularization. We note that this condition can generally be established. For example, we can apply Neyman's classic orthogonalized score in likelihood settings (see, e.g., Neyman 1959, 1979). We also describe an extension of this classic approach to the generalized method of moments (GMM) setting. In general, applying this orthogonalization will introduce additional nuisance parameters that will be treated as part of  $\eta$ .

Second, it is important to use high-quality, structured estimators of  $\eta$ . Crucially, additional structure on  $\eta$  is needed for informative inference to proceed, and it is thus important to use estimation strategies that leverage and perform well under the desired structure. An example of a structure that has been usefully employed in the recent literature is approximate sparsity (e.g., Belloni et al. 2012). Within this framework,  $\eta$  is well approximated by a sparse vector, which suggests the use of a sparse estimator such as the Lasso (Frank & Friedman 1993, Tibshirani 1996). The Lasso estimator solves the general problem

$$\hat{\eta}_L = \arg \min_{\eta} \ell(\text{data}, \eta) + \lambda \sum_{j=1}^p |\psi_j \eta_j|,$$

where  $\ell(\text{data}, \eta)$  is some general loss function that depends on the data and the parameter  $\eta$ ,  $\lambda$  is a penalty level, and  $\psi_j$ 's are penalty loadings. The choice of the regularization parameter  $\lambda$  is an important issue. We provide some discussion of this issue in the context of the linear model in Appendix A (see also, e.g., Belloni & Chernozhukov 2011 for additional detailed discussion). The leading example is the usual linear model in which  $\ell(\text{data}, \eta) = \sum_{i=1}^n (y_i - x_i' \eta)^2$  is the usual least-squares loss, with  $y_i$  denoting the outcome of interest for observation  $i$  and  $x_i$  denoting predictor variables, and we provide further discussion of this example in Appendix A. Other examples of  $\ell(\text{data}, \eta)$  include suitable loss functions corresponding to well-known M-estimators, the negative of the log-likelihood, and GMM criterion functions. This estimator and related methods, such as those in Candès & Tao (2007), Meinshausen & Yu (2009), Bickel et al. (2009), Belloni & Chernozhukov (2013), and Belloni et al. (2011), are computationally efficient and have been shown to have good estimation properties even when perfect variable selection is not feasible under approximate sparsity. These good estimation properties then translate into providing good-enough estimates of  $\eta$  to result in valid inference about  $\alpha$  when coupled with orthogonal estimating equations, as discussed above. Finally, it is important to note that the general results we present do not require or leverage approximate sparsity or sparsity-based estimation strategies. We provide

<sup>1</sup>We refer to the condition as an orthogonality or immunization condition, as orthogonality is a much-used term and our usage differs from some other usage in defining orthogonality conditions used in econometrics.

this discussion here simply as an example and because the structure offers one concrete setting in which the general results we establish may be applied.

In the remainder of this article, we present the main results. In Sections 2 and 3, we provide our general set of results that may be used to establish uniform validity of inference about low-dimensional parameters of interest in the presence of high-dimensional nuisance parameters. We provide the framework in Section 2 and then discuss how to achieve the key orthogonality condition in Section 3. In Sections 4 and 5, we provide details about establishing the necessary results for the estimation quality of  $\eta$  within the approximately sparse framework. The analysis in Section 4 pertains to a reasonably general class of affine-quadratic models, and the analysis of Section 5 specializes this result to the case of estimating the parameters on a vector of endogenous variables in a linear IV model with very many potential control variables and very many potential instruments. The analysis in Section 5 thus extends results from Belloni et al. (2012, 2014a). We also provide a brief simulation example and an empirical example that looks at logit demand estimation within the linear many instrument and many control setting in Section 5. We conclude with a literature review in Section 6.

With regard to notation, we use  $\text{wp} \rightarrow 1$  to abbreviate the phrase “with probability that converges to 1,” and we use the arrows  $\rightarrow_{P_n}$  and  $\rightsquigarrow_{P_n}$  to denote convergence in probability and in distribution under the sequence of probability measures  $\{P_n\}$ . The symbol  $\sim$  means distributed as. The notation  $a \lesssim b$  means that  $a = O(b)$ , and  $a \lesssim_{P_n} b$  means that  $a = O_{P_n}(b)$ . The  $\ell_2$  and  $\ell_1$  norms are denoted by  $\|\cdot\|$  and  $\|\cdot\|_1$ , respectively, and the  $\ell_0$  norm,  $\|\cdot\|_0$ , denotes the number of nonzero components of a vector. When applied to a matrix,  $\|\cdot\|$  denotes the operator norm. We use the notation  $a \vee b = \max(a, b)$  and  $a \wedge b = \min(a, b)$ . Here and below,  $\mathbb{E}_n[\cdot]$  abbreviates the average  $n^{-1} \sum_{i=1}^n [\cdot]$  over index  $i$ . That is,  $\mathbb{E}_n[f(w_i)]$  denotes  $n^{-1} \sum_{i=1}^n [f(w_i)]$ . In what follows, we use the  $m$ -sparse norm of a matrix  $Q$  defined as

$$\|Q\|_{\text{sp}(m)} = \sup \left\{ |b' Q b| / \|b\|^2 : \|b\|_0 \leq m, \|b\| \neq 0 \right\}.$$

We also consider the pointwise norm of a square matrix  $Q$  at a point  $x \neq 0$ :

$$\|Q\|_{\text{pw}(x)} = |x' Q x| / \|x\|^2.$$

For a differentiable map  $x \mapsto f(x)$ , mapping  $\mathbb{R}^d$  to  $\mathbb{R}^k$ , we use  $\partial_x f$  to abbreviate the partial derivatives  $(\partial/\partial x')f$ , and we correspondingly use the expression  $\partial_x f(x_0)$  to mean  $\partial_x f(x)|_{x=x_0}$ , etc. We use  $x'$  to denote the transpose of a column vector  $x$ .

## 2. A TESTING AND ESTIMATION APPROACH TO VALID POST-SELECTION AND POST-REGULARIZATION INFERENCE

### 2.1. The Setting

We assume that estimation is based on the first  $n$  elements  $(w_{i,n})_{i=1}^n$  of the stationary data stream  $(w_{i,n})_{i=1}^\infty$ , which lives on the probability space  $(\Omega, \mathcal{A}, P_n)$ . The data points  $w_{i,n}$  take values in a measurable space  $\mathcal{W}$  for each  $i$  and  $n$ . Here,  $P_n$ , the probability law or data-generating process, can change with  $n$ . We allow the law to change with  $n$  to claim robustness or uniform validity of results with respect to perturbations of such laws. Thus, the data, all parameters, estimators, and other quantities are indexed by  $n$ , but we typically suppress this dependence to simplify notation.

The target parameter value  $\alpha = \alpha_0$  is assumed to solve the system of theoretical equations

$$M(\alpha, \eta_0) = 0,$$

where  $M = (M_l)_{l=1}^k$  is a measurable map from  $\mathcal{A} \times \mathcal{H}$  to  $\mathbb{R}^k$ , and  $\mathcal{A} \times \mathcal{H}$  are some convex subsets of  $\mathbb{R}^d \times \mathbb{R}^p$ . Here the dimension  $d$  of the target parameter  $\alpha \in \mathcal{A}$  and the number of equations  $k$  are assumed to be fixed, and the dimension  $p = p_n$  of the nuisance parameter  $\eta \in \mathcal{H}$  is allowed to be very high, potentially much larger than  $n$ . To handle the high-dimensional nuisance parameter  $\eta$ , we employ structured assumptions and selection or regularization methods appropriate for the structure to estimate  $\eta_0$ .

Given an appropriate estimator  $\hat{\eta}$ , we can construct an estimator  $\hat{\alpha}$  as an approximate solution to the estimating equation:

$$\|\hat{M}(\hat{\alpha}, \hat{\eta})\| \leq \inf_{\alpha \in \mathcal{A}} \|\hat{M}(\alpha, \hat{\eta})\| + o(n^{-1/2}),$$

where  $\hat{M} = (\hat{M}_l)_{l=1}^k$  is the empirical analog of theoretical equations  $M$ , which is a measurable map from  $\mathcal{W}^n \times \mathcal{A} \times \mathcal{H}$  to  $\mathbb{R}^k$ . We can also use  $\hat{M}(\alpha, \hat{\eta})$  to test hypotheses about  $\alpha_0$  and then invert the tests to construct confidence sets.

It is not required in the formulation above, but a typical case is when  $\hat{M}$  and  $M$  are formed as theoretical and empirical moment functions:

$$M(\alpha, \eta) := E[\psi(w_i, \alpha, \eta)], \quad \hat{M}(\alpha, \eta) := \mathbb{E}_n[\psi(w_i, \alpha, \eta)],$$

where  $\psi = (\psi_l)_{l=1}^k$  is a measurable map from  $\mathcal{W} \times \mathcal{A} \times \mathcal{H}$  to  $\mathbb{R}^k$ . Of course, there are many problems that do not fall in the moment condition framework. As illustrations of the general conditions we will provide, we show how our general conditions can be verified in the context of affine-quadratic models and use these results to give primitive conditions in the linear IV model with many instruments and many controls in Sections 4 and 5.

## 2.2. Valid Inference via Testing

A simple introduction to the inferential problem is via the testing problem in which we would like to test some hypothesis about the true parameter value  $\alpha_0$ . By inverting the test, we create a confidence set for  $\alpha_0$ . The key condition for the validity of this confidence region is adaptivity, which can be ensured by using orthogonal estimating equations and using structured assumptions on the high-dimensional nuisance parameter.<sup>2</sup>

The key condition enabling us to perform valid inference on  $\alpha_0$  is the adaptivity condition:

$$\sqrt{n}(\hat{M}(\alpha_0, \hat{\eta}) - \hat{M}(\alpha_0, \eta_0)) \rightarrow_p 0. \quad (1)$$

This condition states that using  $\sqrt{n}\hat{M}(\alpha_0, \hat{\eta})$  is as good as using  $\sqrt{n}\hat{M}(\alpha_0, \eta_0)$ , at least to the first order. This condition may hold despite using estimators  $\hat{\eta}$  that are not asymptotically linear and are nonregular. Verification of adaptivity may involve substantial work, as illustrated below. A key requirement that often arises is the orthogonality or immunization condition:

$$\partial_\eta M(\alpha_0, \eta_0) = 0. \quad (2)$$

This condition states that the equations are locally insensitive to small perturbations of the nuisance parameter around the true parameter values. In several important models, this condition is

<sup>2</sup>Readers are referred to Bickel (1982) for a definition of and introduction to adaptivity.

equivalent to the double-robustness condition (Robins & Rotnitzky 1995). Additional assumptions regarding the quality of estimation of  $\eta_0$  are also needed and are highlighted below.

The adaptivity condition immediately allows us to use the statistic  $\sqrt{n}\hat{M}(\alpha_0, \hat{\eta})$  to perform inference. Indeed, suppose we have that

$$\Omega^{-1/2}(\alpha_0)\sqrt{n}\hat{M}(\alpha_0, \eta_0) \rightsquigarrow_{P_n} \mathcal{N}(0, I_k) \quad (3)$$

for some positive-definite  $\Omega(\alpha) = \text{Var}(\sqrt{n}\hat{M}(\alpha, \eta_0))$ . This condition can be verified using central limit theorems for triangular arrays. Such theorems are available for independently and identically distributed (i.i.d.) as well as dependent and clustered data. Suppose further that there exists  $\hat{\Omega}(\alpha)$  such that

$$\hat{\Omega}^{-1/2}(\alpha_0)\Omega^{1/2}(\alpha_0) \rightarrow_{P_n} I_k. \quad (4)$$

It is then immediate that the following score statistic, evaluated at  $\alpha = \alpha_0$ , is asymptotically normal,

$$S(\alpha) := \hat{\Omega}_n^{-1/2}(\alpha)\sqrt{n}\hat{M}(\alpha, \hat{\eta}) \rightsquigarrow_{P_n} \mathcal{N}(0, I_k), \quad (5)$$

and that the quadratic form of this score statistic is asymptotically  $\chi^2$  with  $k$  degrees of freedom:

$$C(\alpha_0) = \|S(\alpha_0)\|^2 \rightsquigarrow_{P_n} \chi^2(k). \quad (6)$$

The statistic given in Equation 6 simply corresponds to a quadratic form in appropriately normalized statistics that have the desired immunization or orthogonality condition. We refer to this statistic as a generalized  $C(\alpha)$ -statistic in honor of Neyman's fundamental contributions (e.g., Neyman 1959, 1979) because, in likelihood settings, the statistic in Equation 6 reduces to Neyman's  $C(\alpha)$ -statistic and the generalized score  $S(\alpha_0)$  given in Equation 5 reduces to Neyman's orthogonalized score. We demonstrate these relationships in the special case of likelihood models in Section 3.1 and provide a generalization to GMM models in Section 3.2. Both these examples serve to illustrate the construction of appropriate statistics in different settings, but we note that the framework applies far more generally.

The following elementary result is an immediate consequence of the preceding discussion.

**Proposition 1 (valid inference after selection or regularization):** Consider a sequence  $\{P_n\}$  of sets of probability laws such that for each sequence  $\{P_n\} \in \{P_n\}$  the adaptivity condition in Equation 1, the normality condition in Equation 3, and the variance consistency condition in Equation 4 hold. Then  $\text{CR}_{1-a} = \{\alpha \in \mathcal{A} : C(\alpha) \leq c(1-a)\}$ , where  $c(1-a)$  is the  $1-a$ -quantile of a  $\chi^2(k)$ , is a uniformly valid confidence interval for  $\alpha_0$  in the sense that

$$\limsup_{n \rightarrow \infty} \sup_{P \in P_n} |P(\alpha_0 \in \text{CR}_{1-a}) - (1-a)| = 0.$$

We remark here that in order to make the uniformity claim interesting, we should insist that the sets of probability laws  $P_n$  nondecreasing in  $n$  (i.e.,  $P_{\bar{n}} \subseteq P_n$  whenever  $\bar{n} \leq n$ ).

**Proof:** For any sequence of positive constants  $\epsilon_n$  approaching 0, let  $P_n \in P_n$  be any sequence such that

$$|\mathbb{P}_n(\alpha_0 \in \text{CR}_{1-a}) - (1-a)| + \epsilon_n \geq \sup_{\mathbb{P} \in \mathbb{P}_n} |\mathbb{P}(\alpha_0 \in \text{CR}_{1-a}) - (1-a)|.$$

By the conditions in Equations 3 and 4, we have that

$$\mathbb{P}_n(\alpha_0 \in \text{CR}_{1-a}) = \mathbb{P}_n(C(\alpha_0) \leq c(1-a)) \rightarrow \mathbb{P}(\chi^2(k) \leq c(1-a)) = 1-a,$$

which implies the conclusion from the preceding display.

### 2.3. Valid Inference via Adaptive Estimation

Suppose that  $M(\alpha_0, \eta_0) = 0$  holds for  $\alpha_0 \in \mathcal{A}$ . We consider an estimator  $\hat{\alpha} \in \mathcal{A}$  that is an approximate minimizer of the map  $\alpha \mapsto \|\hat{M}(\alpha, \hat{\eta})\|$  in the sense that

$$\|\hat{M}(\hat{\alpha}, \hat{\eta})\| \leq \inf_{\alpha \in \mathcal{A}} \|\hat{M}(\alpha, \hat{\eta})\| + o(n^{-1/2}). \quad (7)$$

To analyze this estimator, we assume that the derivatives  $\Gamma_1 := \partial_{\alpha'} M(\alpha_0, \eta_0)$  and  $\partial_{\eta'} M(\alpha, \eta_0)$  exist. We assume that  $\alpha_0$  is interior relative to the parameter space  $\mathcal{A}$ ; namely, for some  $\ell_n \rightarrow \infty$  such that  $\ell_n/\sqrt{n} \rightarrow 0$ ,

$$\{\alpha \in \mathbb{R}^d : \|\alpha - \alpha_0\| \leq \ell_n/\sqrt{n}\} \subset \mathcal{A}. \quad (8)$$

We also assume that the following local-global identifiability condition holds: For some constant  $c > 0$ ,

$$2\|M(\alpha, \eta_0)\| \geq \|\Gamma_1(\alpha - \alpha_0)\| \wedge c \quad \forall \alpha \in \mathcal{A}, \quad \text{mineig}(\Gamma_1' \Gamma_1) \geq c. \quad (9)$$

Furthermore, for  $\Omega = \text{Var}(\sqrt{n}\hat{M}(\alpha_0, \eta_0))$ , we suppose that the central limit theorem,

$$\Omega^{-1/2} \sqrt{n}\hat{M}(\alpha_0, \eta_0) \rightsquigarrow_{\mathbb{P}_n} \mathcal{N}(0, I), \quad (10)$$

and the stability condition,

$$\|\Gamma_1' \Gamma_1\| + \|\Omega\| + \|\Omega^{-1}\| \lesssim 1, \quad (11)$$

hold.

Assume that for some sequence of positive numbers  $\{r_n\}$  such that  $r_n \rightarrow 0$  and  $r_n n^{1/2} \rightarrow \infty$ , the following stochastic equicontinuity and continuity conditions hold:

$$\sup_{\alpha \in \mathcal{A}} \frac{\|\hat{M}(\alpha, \hat{\eta}) - M(\alpha, \hat{\eta})\| + \|M(\alpha, \hat{\eta}) - M(\alpha, \eta_0)\|}{r_n + \|\hat{M}(\alpha, \hat{\eta})\| + \|M(\alpha, \eta_0)\|} \rightarrow_{\mathbb{P}_n} 0, \quad (12)$$

$$\sup_{\|\alpha - \alpha_0\| \leq r_n} \frac{\|\hat{M}(\alpha, \hat{\eta}) - M(\alpha, \hat{\eta}) - \hat{M}(\alpha_0, \eta_0)\|}{n^{-1/2} + \|\hat{M}(\alpha, \hat{\eta})\| + \|M(\alpha, \eta_0)\|} \rightarrow_{\mathbb{P}_n} 0. \quad (13)$$

Suppose that uniformly for all  $\alpha \neq \alpha_0$  such that  $\|\alpha - \alpha_0\| \leq r_n \rightarrow 0$ , the following conditions on the smoothness of  $M$  and the quality of the estimator  $\hat{\eta}$  hold, as  $n \rightarrow \infty$ :

$$\begin{aligned} & \left\| M(\alpha, \eta_0) - M(\alpha_0, \eta_0) - \Gamma_1[\alpha - \alpha_0] \right\| \|\alpha - \alpha_0\|^{-1} \rightarrow 0, \\ & \sqrt{n} \left\| M(\alpha, \hat{\eta}) - M(\alpha, \eta_0) - \partial_{\eta'} M(\alpha, \eta_0)[\hat{\eta} - \eta_0] \right\| \rightarrow_{\mathbb{P}_n} 0, \\ & \left\| \left\{ \partial_{\eta'} M(\alpha, \eta_0) - \partial_{\eta'} M(\alpha_0, \eta_0) \right\} [\hat{\eta} - \eta_0] \right\| \|\alpha - \alpha_0\|^{-1} \rightarrow_{\mathbb{P}_n} 0. \end{aligned} \tag{14}$$

Finally, as before, we assume that the orthogonality condition

$$\partial_{\eta'} M(\alpha_0, \eta_0) = 0 \tag{15}$$

holds.

The above conditions extend the analysis of Pakes & Pollard (1989) and Chen et al. (2003), which in turn extended Huber's (1964) classical results on Z-estimators. These conditions allow for both smooth and nonsmooth systems of estimating equations. The identifiability condition imposed above is mild and holds for broad classes of identifiable models. The equicontinuity and smoothness conditions imposed above require mild smoothness on the function  $M$  and also require that  $\hat{\eta}$  is a good-quality estimator of  $\eta_0$ . In particular, these conditions will often require that  $\hat{\eta}$  converges to  $\eta_0$  at a faster rate than  $n^{-1/4}$ , as demonstrated, for example, in the next section. However, the rate condition alone is not sufficient for adaptivity. We also need the orthogonality condition in Equation 15. In addition, it is required that  $\hat{\eta} \in \mathcal{H}_n$ , where  $\mathcal{H}_n$  is a set whose complexity does not grow too quickly with the sample size, to verify the stochastic equicontinuity condition (see, e.g., Belloni et al. 2013a,d). In Sections 4 and 5, we use the sparsity of  $\hat{\eta}$  to control this complexity. Note that the conditions in Equations 12 and 13 can be simplified by leaving only  $r_n$  and  $n^{-1/2}$  in the denominator, although this simplification would then require imposing compactness on  $\mathcal{A}$  even in linear problems.

**Proposition 2 (valid inference via adaptive estimation after selection or regularization):**

Consider a sequence  $\{\mathbb{P}_n\}$  of sets of probability laws such that for each sequence  $\{\mathbb{P}_n\} \in \{\mathbb{P}_n\}$  the conditions in Equations 7–15 hold. Then we obtain

$$\sqrt{n}(\hat{\alpha} - \alpha_0) + \left[ \Gamma_1' \Gamma_1 \right]^{-1} \Gamma_1' \sqrt{n} \hat{M}(\alpha_0, \eta_0) \rightarrow_{\mathbb{P}_n} 0.$$

In addition, for  $V_n := (\Gamma_1' \Gamma_1)^{-1} \Gamma_1' \Omega \Gamma_1 (\Gamma_1' \Gamma_1)^{-1}$ , we have that

$$\limsup_{n \rightarrow \infty} \sup_{\mathbb{P} \in \mathbb{P}_n} \sup_{R \in \mathcal{R}} \left| \mathbb{P}(V_n^{-1/2}(\hat{\alpha} - \alpha_0) \in R) - \mathbb{P}(\mathcal{N}(0, I) \in R) \right| = 0,$$

where  $\mathcal{R}$  is a collection of all convex sets. Moreover, the result continues to apply if  $V_n$  is replaced by a consistent estimator  $\hat{V}_n$  such that  $\hat{V}_n - V_n \rightarrow_{\mathbb{P}_n} 0$  under each sequence  $\{\mathbb{P}_n\}$ . Thus,  $\text{CR}_{1-a}^l = \left[ l' \hat{\alpha} \pm c(1-a/2) (l' \hat{V}_n l / n)^{1/2} \right]$ , where  $c(1-a/2)$  is the  $(1-a/2)$ -quantile of  $\mathcal{N}(0, 1)$ , is a uniformly valid confidence set for  $l' \alpha_0$ :

$$\limsup_{n \rightarrow \infty} \sup_{\mathbb{P} \in \mathbb{P}_n} \left| \mathbb{P}(l' \alpha_0 \in \text{CR}_{1-a}^l) - (1-a) \right| = 0.$$



Note that the above formulation implicitly accommodates weighting options. Suppose  $M^o$  and  $\hat{M}^o$  are the original theoretical and empirical systems of equations, and let  $\Gamma_1^o = \partial_\alpha M^o(\alpha_0, \eta_0)$  be the original Jacobian. We could consider  $k \times k$  positive-definite weight matrices  $A$  and  $\hat{A}$  such that

$$\|A^2\| + \|(A^2)^{-1}\| \lesssim 1, \quad \|\hat{A}^2 - A^2\| \rightarrow_{P_n} 0. \quad (16)$$

For example, we may wish to use the optimal weighting matrix  $A^2 = \text{Var}(\sqrt{n}\hat{M}^o(\alpha_0, \eta_0))^{-1}$ , which can be estimated by  $\hat{A}^2$  obtained using a preliminary estimator  $\hat{\alpha}^o$  resulting from solving the problem with some nonoptimal weighting matrix such as  $I$ . We can then simply redefine the system of equations and the Jacobian according to

$$M(\alpha, \eta) = AM^o(\alpha, \eta), \quad \hat{M}(\alpha, \eta) = \hat{A}\hat{M}^o(\alpha, \eta), \quad \Gamma_1 = A\Gamma_1^o. \quad (17)$$

**Proposition 3 (adaptive estimation via weighted equations):** Consider a sequence  $\{P_n\}$  of sets of probability laws such that for each sequence  $\{P_n\} \in \{P_n\}$  the conditions of Proposition 2 hold for the original pair of systems of equations  $(M^o, \hat{M}^o)$  and Equation 16 holds. Then these conditions also hold for the new pair  $(M, \hat{M})$  in Equation 17, so all the conclusions of Proposition 2 apply to the resulting approximate argmin estimator  $\hat{\alpha}$ . In particular, if we use  $A^2 = \text{Var}(\sqrt{n}\hat{M}^o(\alpha_0, \eta_0))^{-1}$  and  $\hat{A}^2 - A^2 \rightarrow_{P_n} 0$ , then the large sample variance  $V_n$  simplifies to  $V_n = (\Gamma_1' \Gamma_1)^{-1}$ .

## 2.4. Inference via Adaptive One-Step Estimation

We next consider a one-step estimator. To define the estimator, we start with an initial estimator  $\tilde{\alpha}$  that satisfies, for  $r_n = o(n^{-1/4})$ ,

$$P_n(\|\tilde{\alpha} - \alpha_0\| \leq r_n) \rightarrow 1. \quad (18)$$

The one-step estimator  $\hat{\alpha}$  then solves a linearized version of Equation 7:

$$\hat{\alpha} = \tilde{\alpha} - \left[ \hat{\Gamma}_1' \hat{\Gamma}_1 \right]^{-1} \hat{\Gamma}_1' \hat{M}(\tilde{\alpha}, \hat{\eta}), \quad (19)$$

where  $\hat{\Gamma}_1$  is an estimator of  $\Gamma_1$  such that

$$P_n\left(\|\hat{\Gamma}_1 - \Gamma_1\| \leq r_n\right) \rightarrow 1. \quad (20)$$

Because the one-step estimator is considerably more crude than the argmin estimator, we need to impose additional smoothness conditions. Specifically, we suppose that uniformly for all  $\alpha \neq \alpha_0$  such that  $\|\alpha - \alpha_0\| \leq r_n \rightarrow 0$ , the following strengthened conditions on stochastic equicontinuity, smoothness of  $M$ , and the quality of the estimator  $\hat{\eta}$  hold, as  $n \rightarrow \infty$ :

$$\begin{aligned} n^{1/2} \left\| \hat{M}(\alpha, \hat{\eta}) - M(\alpha, \hat{\eta}) - \hat{M}(\alpha_0, \eta_0) \right\| &\rightarrow_{P_n} 0, \\ \left\| M(\alpha, \eta_0) - M(\alpha_0, \eta_0) - \Gamma_1[\alpha - \alpha_0] \right\| \|\alpha - \alpha_0\|^{-2} &\lesssim 1, \\ \sqrt{n} \left\| M(\alpha, \hat{\eta}) - M(\alpha, \eta_0) - \partial_{\eta'} M(\alpha, \eta_0)[\hat{\eta} - \eta_0] \right\| &\rightarrow_{P_n} 0, \\ \sqrt{n} \left\| \left\{ \partial_{\eta'} M(\alpha, \eta_0) - \partial_{\eta'} M(\alpha_0, \eta_0) \right\} [\hat{\eta} - \eta_0] \right\| &\rightarrow_{P_n} 0. \end{aligned} \quad (21)$$

**Proposition 4 (valid inference via adaptive one-step estimators):** Consider a sequence  $\{\mathbf{P}_n\}$  of sets of probability laws such that for each sequence  $\{\mathbf{P}_n\} \in \{\mathbf{P}_n\}$  the conditions of Proposition 2 as well as those in Equations 18, 20, and 21 hold. Then the one-step estimator  $\tilde{\alpha}$  defined by Equation 19 is first-order equivalent to the argmin estimator  $\hat{\alpha}$ :

$$\sqrt{n}(\tilde{\alpha} - \hat{\alpha}) \rightarrow_{P_n} 0.$$

Consequently, all conclusions of Proposition 2 apply to  $\tilde{\alpha}$  in place of  $\hat{\alpha}$ .

The one-step estimator requires stronger regularity conditions than the argmin estimator. Moreover, there is finite-sample evidence (e.g., Belloni et al. 2013e) that in practical problems the argmin estimator often works much better, as the one-step estimator typically suffers from higher-order biases. This problem could be alleviated somewhat by iterating on the one-step estimator, treating the previous iteration as the crude start  $\tilde{\alpha}$  for the next iteration.

### 3. ACHIEVING ORTHOGONALITY USING NEYMAN'S ORTHOGONALIZATION

Here we describe orthogonalization ideas that go back at least to Neyman (1959) (see also Neyman 1979). Neyman's idea was to project the score that identifies the parameter of interest onto the orthocomplement of the tangent space for the nuisance parameter. This projection underlies semi-parametric efficiency theory, which is concerned particularly with the case in which  $\eta$  is infinite dimensional (see van der Vaart 1998). Here we consider finite-dimensional  $\eta$  of high dimension (for discussion of infinite-dimensional  $\eta$  in an approximately sparse setting, see Belloni et al. 2013a,d).

#### 3.1. The Classical Likelihood Case

In likelihood settings, the construction of orthogonal equations was proposed by Neyman (1959), who used them in construction of his celebrated  $C(\alpha)$ -statistic. The  $C(\alpha)$ -statistic, or the orthogonal score statistic, was first explicitly utilized for testing (and also for setting up estimation) in high-dimensional sparse models in Belloni et al. (2013d) and Belloni et al. (2013c), in the context of quantile regression, and Belloni et al. (2013e) in the context of logistic regression and other generalized linear models. More recent uses of  $C(\alpha)$ -statistics (or close variants) include those by Voorman et al. (2014), Ning & Liu (2014), and Yang et al. (2014).

Suppose that the (possibly conditional, possibly quasi-) log-likelihood function associated with observation  $w_i$  is  $\ell(w_i, \alpha, \beta)$ , where  $\alpha \in \mathcal{A} \subset \mathbb{R}^d$  is the target parameter and  $\beta \in \mathcal{B} \subset \mathbb{R}^{p_0}$  is the nuisance parameter. Under regularity conditions, the true parameter values  $\gamma_0 = (\alpha'_0, \beta_0)'$  obey

$$E[\partial_\alpha \ell(w_i, \alpha_0, \beta_0)] = 0, \quad E[\partial_\beta \ell(w_i, \alpha_0, \beta_0)] = 0. \quad (22)$$

Now consider the moment function

$$M(\alpha, \eta) = E[\psi(w_i, \alpha, \eta)], \quad \psi(w_i, \alpha, \eta) = \partial_\alpha \ell(w_i, \alpha, \beta) - \mu \partial_\beta \ell(w_i, \alpha, \beta). \quad (23)$$

Here the nuisance parameter is

$$\eta = \left( \beta', \text{vec}(\mu)' \right)' \in \mathcal{B} \times \mathcal{D} \subset \mathbb{R}^p, \quad p = p_0 + dp_0,$$

where  $\mu$  is the  $d \times p_0$  orthogonalization parameter matrix whose true value  $\mu_0$  solves the equation

$$J_{\alpha\beta} - \mu J_{\beta\beta} = 0 \quad (\text{i.e., } \mu_0 = J_{\alpha\beta} J_{\beta\beta}^{-1}), \quad (24)$$

where, for  $\gamma := (\alpha', \beta')'$  and  $\gamma_0 := (\alpha'_0, \beta'_0)'$ ,

$$J := -\partial_\gamma \text{E}[\partial_\gamma \ell(w_i, \gamma)] \Big|_{\gamma=\gamma_0} =: \begin{pmatrix} J_{\alpha\alpha} & J_{\alpha\beta} \\ J_{\beta\alpha} & J_{\beta\beta} \end{pmatrix}.$$

Note that  $\mu_0$  not only creates the necessary orthogonality but also creates the optimal score (in statistical language) or, equivalently, the optimal instrument/moment (in econometric language) for inference about  $\alpha_0$ .<sup>3</sup>

Provided  $\mu_0$  is well defined, we have by Equation 22 that

$$M(\alpha_0, \eta_0) = 0.$$

Moreover, the function  $M$  has the desired orthogonality property:

$$\partial_\eta M(\alpha_0, \eta_0) = \left[ J_{\alpha\beta} - \mu_0 J_{\beta\beta}; \quad FE[\partial_\beta \ell(w_i, \alpha_0, \beta_0)] \right] = 0, \quad (25)$$

where  $F$  is a tensor operator, such that  $Fx = \partial_\mu x / \partial \text{vec}(\mu)' \Big|_{\mu=\mu_0}$  is a  $d \times (dp_0)$  matrix for any vector  $x$  in  $\mathbb{R}^{p_0}$ . Note that the orthogonality property holds for Neyman's construction even if the likelihood is misspecified. That is,  $\ell(w_i, \gamma_0)$  may be a quasi-likelihood, and the data need not be i.i.d. and may, for example, exhibit complex dependence over  $i$ .

An alternative way to define  $\mu_0$  arises by considering that, under correct specification and sufficient regularity, the information matrix equality holds and yields

$$\begin{aligned} J &= J^0 := \text{E} \left[ \partial_\gamma \ell(w_i, \gamma) \partial_\gamma \ell(w_i, \gamma)' \right] \Big|_{\gamma=\gamma_0} \\ &= \begin{pmatrix} \text{E} \left[ \partial_\alpha \ell(w_i, \gamma) \partial_\alpha \ell(w_i, \gamma)' \right] & \text{E} \left[ \partial_\alpha \ell(w_i, \gamma) \partial_\beta \ell(w_i, \gamma)' \right] \\ \text{E} \left[ \partial_\beta \ell(w_i, \gamma) \partial_\alpha \ell(w_i, \gamma)' \right] & \text{E} \left[ \partial_\beta \ell(w_i, \gamma) \partial_\beta \ell(w_i, \gamma)' \right] \end{pmatrix} \Big|_{\gamma=\gamma_0} \\ &=: \begin{pmatrix} J_{\alpha\alpha}^0 & J_{\alpha\beta}^0 \\ J_{\beta\alpha}^0 & J_{\beta\beta}^0 \end{pmatrix}. \end{aligned}$$

Hence, define  $\mu_0^* = J_{\alpha\beta}^0 J_{\beta\beta}^{0-1}$  as the population projection coefficient of the score for the main parameter  $\partial_\alpha \ell(w_i, \gamma_0)$  on the score for the nuisance parameter  $\partial_\beta \ell(w_i, \gamma_0)$ :

$$\partial_\alpha \ell(w_i, \gamma_0) = \mu_0^* \partial_\beta \ell(w_i, \gamma_0) + \varrho, \quad \text{E} \left[ \varrho \partial_\beta \ell(w_i, \gamma_0)' \right] = 0. \quad (26)$$

We can see this construction as the nonlinear version of Frisch-Waugh's "partialling out" from the linear regression model. It is important to note that under misspecification, the information

<sup>3</sup>The connection between optimal instruments/moments and the likelihood/score has been elucidated in the fundamental work of Chamberlain (1987).

matrix equality generally does not hold, and this projection approach does not provide valid orthogonalization.

**Lemma 1 [Neyman’s orthogonalization for (quasi-)likelihood scores]:** Suppose that for each  $\gamma = (\alpha, \beta) \in \mathcal{A} \times \mathcal{B}$ , the derivative  $\partial_\gamma \ell(w_i, \gamma)$  exists, is continuous at  $\gamma$  with probability 1, and obeys the dominance condition  $E \sup_{\gamma \in \mathcal{A} \times \mathcal{B}} \|\partial_\gamma \ell(w_i, \gamma)\|^2 < \infty$ . Suppose that the condition in Equation 22 holds for some (quasi-) true value  $(\alpha_0, \beta_0)$ . Then, (a) if  $J$  exists and is finite and  $J_{\beta\beta}$  is invertible, the orthogonality condition in Equation 25 holds; (b) if the information matrix equality holds, namely  $J = J^0$ , then the orthogonality condition in Equation 25 holds for the projection parameter  $\mu_0^*$  in place of the orthogonalization parameter matrix  $\mu_0$ .

The claim follows immediately from the computations above.

With the formulations given above, Neyman’s  $C(\alpha)$ -statistic takes the form

$$C(\alpha) = \|S(\alpha)\|_2^2, \quad S(\alpha) = \hat{\Omega}^{-1/2}(\alpha, \hat{\eta}) \sqrt{n} \hat{M}(\alpha, \hat{\eta}),$$

where  $\hat{M}(\alpha, \hat{\eta}) = \mathbb{E}_n[\psi(w_i, \alpha, \hat{\eta})]$  as before,  $\Omega(\alpha, \eta_0) = \text{Var}(\sqrt{n} \hat{M}(\alpha, \eta_0))$ , and  $\hat{\Omega}(\alpha, \hat{\eta})$  and  $\hat{\eta}$  are suitable estimators based on sparsity or other structured assumptions. The estimator is then

$$\hat{\alpha} = \arg \inf_{\alpha \in \mathcal{A}} C(\alpha) = \arg \inf_{\alpha \in \mathcal{A}} \|\sqrt{n} \hat{M}(\alpha, \hat{\eta})\|,$$

provided that  $\hat{\Omega}(\alpha, \hat{\eta})$  is positive definite for each  $\alpha \in \mathcal{A}$ . If the conditions of Section 2 hold, we have

$$C(\alpha) \rightsquigarrow \chi^2(d), \quad V_n^{-1/2} \sqrt{n}(\hat{\alpha} - \alpha_0) \rightsquigarrow \mathcal{N}(0, I), \quad (27)$$

where  $V_n = \Gamma_1^{-1} \Omega(\alpha_0, \eta_0) \Gamma_1^{-1}$  and  $\Gamma_1 = J_{\alpha\alpha} - \mu_0 J_{\alpha\beta}'$ . Under correct specification and i.i.d. sampling, the variance matrix  $V_n$  further reduces to the optimal variance

$$\Gamma_1^{-1} = \left( J_{\alpha\alpha} - J_{\alpha\beta} J_{\beta\beta}^{-1} J_{\alpha\beta}' \right)^{-1}$$

of the first  $d$  components of the maximum likelihood estimator in a Gaussian shift experiment with observation  $Z \sim \mathcal{N}(b, J_0^{-1})$ . Likewise, the result in Equation 27 also holds for the one-step estimator  $\tilde{\alpha}$  of Section 2 in place of  $\hat{\alpha}$  as long as the conditions in Section 2 hold.

Provided that sparsity or its generalizations are plausible assumptions to make regarding  $\eta_0$ , the formulations above naturally lend themselves to sparse estimation. For example, Belloni et al. (2013e) used penalized and post-penalized maximum likelihood to estimate  $\beta_0$  and used the information matrix equality to estimate the orthogonalization parameter matrix  $\mu_0^*$  by employing Lasso or post-Lasso estimation of the projection equation (Equation 26). It is also possible to estimate  $\mu_0$  directly by finding approximate sparse solutions to the empirical analog of the system of equations  $J_{\alpha\beta} - \mu J_{\beta\beta} = 0$  using  $\ell_1$ -penalized estimation, as, for example, in van de Geer et al. (2014), or post- $\ell_1$ -penalized estimation.

### 3.2. Achieving Orthogonality in Generalized Method of Moments (GMM) Problems

Here we consider  $\gamma_0 = (\alpha_0', \beta_0')'$  that solve the system of equations

$$E[m(w_i, \alpha_0, \beta_0)] = 0,$$

where  $m: \mathcal{W} \times \mathcal{A} \times \mathcal{B} \mapsto \mathbb{R}^k$ ,  $\mathcal{A} \times \mathcal{B}$  is a convex subset of  $\mathbb{R}^d \times \mathbb{R}^{p_0}$ , and  $k \geq d + p_0$  is the number of moments. The orthogonal moment equation is

$$M(\alpha, \eta) = E[\psi(w_i, \alpha, \eta)], \quad \psi(w_i, \alpha, \eta) = \mu m(w_i, \alpha, \beta). \quad (28)$$

The nuisance parameter is

$$\eta = (\beta', \text{vec}(\mu)')' \in \mathcal{B} \times \mathcal{D} \subset \mathbb{R}^p, \quad p = p_0 + dk,$$

where  $\mu$  is the  $d \times k$  orthogonalization parameter matrix. The true value of  $\mu$  is

$$\mu_0 = \left( G'_\alpha \Omega_m^{-1} - G'_\alpha \Omega_m^{-1} G'_\beta (G'_\beta \Omega_m^{-1} G'_\beta)^{-1} G'_\beta \Omega_m^{-1} \right),$$

where, for  $\gamma = (\alpha', \beta')'$  and  $\gamma_0 = (\alpha'_0, \beta'_0)'$ ,

$$G_\gamma = \partial_\gamma E[m(w_i, \alpha, \beta)] \Big|_{\gamma=\gamma_0} = \left[ \partial_{\alpha'} E[m(w_i, \alpha, \beta)], \partial_{\beta'} E[m(w_i, \alpha, \beta)] \right] \Big|_{\gamma=\gamma_0} =: [G_\alpha, G_\beta],$$

and

$$\Omega_m = \text{Var}(\sqrt{n} \mathbb{E}_n[m(w_i, \alpha_0, \beta_0)]).$$

As before, we can interpret  $\mu_0$  as an operator creating orthogonality while building the optimal instrument/moment (in econometric language) or, equivalently, the optimal score function (in statistical language).<sup>4</sup> The resulting moment function has the required orthogonality property; namely, the first derivative with respect to the nuisance parameter when evaluated at the true parameter values is zero:

$$\partial_{\eta'} M(\alpha_0, \eta) \Big|_{\eta=\eta_0} = \left[ \mu_0 G_\beta, FE[m(w_i, \alpha_0, \beta_0)] \right] = 0, \quad (29)$$

where  $F$  is a tensor operator, such that  $Fx = \partial \mu x / \partial \text{vec}(\mu)' \Big|_{\mu=\mu_0}$  is a  $d \times (dk)$  matrix for any vector  $x$  in  $\mathbb{R}^k$ .

Estimation and inference on  $\alpha_0$  can be based on the empirical analog of Equation 28:

$$\hat{M}(\alpha_0, \hat{\eta}) = \mathbb{E}_n[\psi(w_i, \alpha, \hat{\eta})],$$

where  $\hat{\eta}$  is a post-selection or other regularized estimator of  $\eta_0$ . Note that the previous framework of (quasi-)likelihood is incorporated as a special case with

$$m(w_i, \alpha, \beta) = \left[ \partial_\alpha \ell(w_i, \alpha)', \partial_\beta \ell(w_i, \beta)' \right]'$$

With the formulations above, Neyman's  $C(\alpha)$ -statistic takes the form

<sup>4</sup>Readers are referred to footnote 3.

$$C(\alpha) = \|S(\alpha)\|_2^2, \quad S(\alpha) = \hat{\Omega}^{-1/2}(\alpha, \hat{\eta})\sqrt{n}\hat{M}(\alpha, \hat{\eta}),$$

where  $\hat{M}(\alpha, \hat{\eta}) = \mathbb{E}_n[\psi(w_i, \alpha, \hat{\eta})]$  as before,  $\Omega(\alpha, \eta_0) = \text{Var}(\sqrt{n}\hat{M}(\alpha, \eta_0))$ , and  $\hat{\Omega}(\alpha, \hat{\eta})$  and  $\hat{\eta}$  are suitable estimators based on structured assumptions. The estimator is then

$$\hat{\alpha} = \arg \inf_{\alpha \in \mathcal{A}} C(\alpha) = \arg \inf_{\alpha \in \mathcal{A}} \|\sqrt{n}\hat{M}(\alpha, \hat{\eta})\|,$$

provided that  $\hat{\Omega}(\alpha, \hat{\eta})$  is positive definite for each  $\alpha \in \mathcal{A}$ . If the high-level conditions of Section 2 hold, we have that

$$C(\alpha) \rightsquigarrow_{P_n} \chi^2(d), \quad V_n^{-1/2}\sqrt{n}(\hat{\alpha} - \alpha) \rightsquigarrow_{P_n} \mathcal{N}(0, I), \quad (30)$$

where  $V_n = (\Gamma_1')^{-1}\Omega(\alpha_0, \eta_0)(\Gamma_1)^{-1}$  coincides with the optimal variance for GMM; here  $\Gamma_1 = \mu_0 G_\alpha$ . Likewise, the same result in Equation 30 holds for the one-step estimator  $\tilde{\alpha}$  of Section 2 in place of  $\hat{\alpha}$  as long as the conditions in Section 2 hold. In particular, the variance  $V_n$  corresponds to the variance of the first  $d$  components of the maximum likelihood estimator in the normal shift experiment with the observation  $Z \sim \mathcal{N}(b, (G_\gamma' \Omega_m^{-1} G_\gamma)^{-1})$ .

The above is a generic outline of the properties that are expected for inference using orthogonalized GMM equations under structured assumptions. The problem of inference in GMM under sparsity is a very delicate matter owing to the complex form of the orthogonalization parameters. One approach to the problem is developed in Chernozhukov et al. (2014).

## 4. ACHIEVING ADAPTIVITY IN AFFINE-QUADRATIC MODELS VIA APPROXIMATE SPARSITY

Here we take orthogonality as given and explain how we can use approximate sparsity to achieve the adaptivity property in Equation 1.

### 4.1. The Affine-Quadratic Model

We analyze the case in which  $\hat{M}$  and  $M$  are affine in  $\alpha$  and affine quadratic in  $\eta$ . Specifically, we suppose that for all  $\alpha$ ,

$$\hat{M}(\alpha, \eta) = \hat{\Gamma}_1(\eta)\alpha + \hat{\Gamma}_2(\eta), \quad M(\alpha, \eta) = \Gamma_1(\eta)\alpha + \Gamma_2(\eta),$$

where the orthogonality condition holds,

$$\partial_\eta M(\alpha_0, \eta_0) = 0,$$

and  $\eta \mapsto \hat{\Gamma}_j(\eta)$  and  $\eta \mapsto \Gamma_j(\eta)$  are affine quadratic in  $\eta$  for  $j = 1$  and  $j = 2$ . That is, we will have that all second-order derivatives of  $\hat{\Gamma}_j(\eta)$  and  $\Gamma_j(\eta)$  for  $j = 1$  and  $j = 2$  are constant over the convex parameter space  $\mathcal{H}$  for  $\eta$ .

This setting is both useful, including most widely used linear models as a special case, and pedagogical, permitting simple illustration of the key issues that arise in treating the general problem. The derivations given below easily generalize to more complicated models, but we defer the details to the interested reader.

The estimator in this case is

$$\hat{\alpha} = \arg \min_{\alpha \in \mathbb{R}^d} \left\| \hat{M}(\alpha, \hat{\eta}) \right\|^2 = - \left[ \hat{\Gamma}_1(\hat{\eta})' \hat{\Gamma}_1(\hat{\eta}) \right]^{-1} \hat{\Gamma}_1(\hat{\eta})' \hat{\Gamma}_2(\hat{\eta}), \quad (31)$$

provided the inverse is well defined. It follows that

$$\sqrt{n}(\hat{\alpha} - \alpha_0) = - \left[ \hat{\Gamma}_1(\hat{\eta})' \hat{\Gamma}_1(\hat{\eta}) \right]^{-1} \hat{\Gamma}_1(\hat{\eta})' \sqrt{n} \hat{M}(\alpha_0, \hat{\eta}). \quad (32)$$

This estimator is adaptive if, for  $\Gamma_1 := \Gamma_1(\eta_0)$ ,

$$\sqrt{n}(\hat{\alpha} - \alpha_0) + \left[ \Gamma_1' \Gamma_1 \right]^{-1} \Gamma_1' \sqrt{n} \hat{M}(\alpha_0, \eta_0) \rightarrow_{P_n} 0,$$

which occurs under the conditions in Equations 10 and 11 if

$$\sqrt{n}(\hat{M}(\alpha_0, \hat{\eta}) - \hat{M}(\alpha_0, \eta_0)) \rightarrow_{P_n} 0, \quad \hat{\Gamma}_1(\hat{\eta}) - \hat{\Gamma}_1(\eta_0) \rightarrow_{P_n} 0. \quad (33)$$

Therefore, the problem of the adaptivity of the estimator is directly connected to the problem of the adaptivity of testing hypotheses about  $\alpha_0$ .

**Lemma 2 (adaptive testing and estimation in affine-quadratic models):** Consider a sequence  $\{\mathbf{P}_n\}$  of sets of probability laws such that for each sequence  $\{\mathbf{P}_n\} \in \{\mathbf{P}_n\}$ , conditions stated in the first paragraph of Section 4.1, the condition in Equation 33, the asymptotic normality condition in Equation 10, the stability condition in Equation 11, and the condition in Equation 4 hold. Then all the conditions of Propositions 1 and 2 hold. Moreover, the conclusions of Proposition 1 hold, and the conclusions of Proposition 2 hold for the estimator  $\hat{\alpha}$  in Equation 31.

## 4.2. Adaptivity for Testing via Approximate Sparsity

Assuming the orthogonality condition holds, we follow Belloni et al. (2012) in using approximate sparsity to achieve the adaptivity property in Equation 1 for the testing problem in the affine-quadratic models.

We can expand each element  $\hat{M}_j$  of  $\hat{M} = (\hat{M}_j)_{j=1}^k$  as follows:

$$\sqrt{n}(\hat{M}_j(\alpha_0, \hat{\eta}) - \hat{M}_j(\alpha_0, \eta_0)) = T_{1,j} + T_{2,j} + T_{3,j}, \quad (34)$$

where

$$\begin{aligned} T_{1,j} &:= \sqrt{n} \partial_{\eta} M_j(\alpha_0, \eta_0)' (\hat{\eta} - \eta_0), \\ T_{2,j} &:= \sqrt{n} (\partial_{\eta} \hat{M}_j(\alpha_0, \eta_0) - \partial_{\eta} M_j(\alpha_0, \eta_0))' (\hat{\eta} - \eta_0), \\ T_{3,j} &:= \sqrt{n} 2^{-1} (\hat{\eta} - \eta_0)' \partial_{\eta} \partial_{\eta}' \hat{M}_j(\alpha_0) (\hat{\eta} - \eta_0). \end{aligned} \quad (35)$$

The term  $T_{1,j}$  vanishes precisely because of orthogonality; that is,

$$T_{1,j} = 0.$$

However, terms  $T_{2,j}$  and  $T_{3,j}$  need not vanish. To show that they are asymptotically negligible, we need to impose further structure on the problem.

**4.2.1. Structure 1 (exact sparsity).** We first consider the case of using an exact sparsity structure in which  $\|\eta_0\|_0 \leq s$  and  $s = s_n \geq 1$  can depend on  $n$ . We then use estimators  $\hat{\eta}$  that exploit the sparsity structure.

Suppose that the following bounds hold with probability  $1 - o(1)$  under  $P_n$ :

$$\begin{aligned} \|\hat{\eta}\|_0 &\lesssim s, \quad \|\eta_0\|_0 \leq s, \\ \|\hat{\eta} - \eta_0\|_2 &\lesssim \sqrt{(s/n)\log(pn)}, \quad \|\hat{\eta} - \eta_0\|_1 \lesssim \sqrt{(s^2/n)\log(pn)}. \end{aligned} \quad (36)$$

These conditions are typical performance bounds that are well known to hold for many sparsity-based estimators, such as Lasso, post-Lasso, and their extensions (see, e.g., Belloni & Chernozhukov 2011).

We suppose further that the moderate deviation bound

$$\bar{T}_{2,j} = \left\| \sqrt{n} \left( \partial_{\eta'} \hat{M}_j(\alpha_0, \eta_0) - \partial_{\eta'} M_j(\alpha_0, \eta_0) \right) \right\|_{\infty} \lesssim_{P_n} \sqrt{\log(pn)} \quad (37)$$

holds and that the sparse norm of the second-derivatives matrix is bounded:

$$\bar{T}_{3,j} = \left\| \partial_{\eta} \partial_{\eta'} \hat{M}_j(\alpha_0) \right\|_{\text{sp}(\ell_n s)} \lesssim_{P_n} 1, \quad (38)$$

where  $\ell_n \rightarrow \infty$  but  $\ell_n = o(\log n)$ .

Following Belloni et al. (2012), we can verify the condition in Equation 37 using the moderate deviation theory for self-normalized sums (e.g., Jing et al. 2003), which allows us to avoid making highly restrictive sub-Gaussian or Gaussian tail assumptions. Likewise, following Belloni et al. (2012), we can verify the second condition using laws of large numbers for large matrices acting on sparse vectors, as in Rudelson & Vershynin (2008) and Rudelson & Zhou (2011) (see Lemma 7). Indeed, the condition in Equation 38 holds if

$$\left\| \partial_{\eta} \partial_{\eta'} \hat{M}_j(\alpha_0) - \partial_{\eta} \partial_{\eta'} M_j(\alpha_0) \right\|_{\text{sp}(\ell_n s)} \rightarrow_{P_n} 0, \quad \left\| \partial_{\eta} \partial_{\eta'} M_j(\alpha_0) \right\|_{\text{sp}(\ell_n s)} \lesssim 1.$$

The above analysis immediately implies the following elementary result.

**Lemma 3 (elementary adaptivity for testing via sparsity):** Let  $\{P_n\}$  be a sequence of probability laws. Assume that (a)  $\eta \mapsto \hat{M}(\alpha_0, \eta)$  and  $\eta \mapsto M(\alpha_0, \eta)$  are affine quadratic in  $\eta$ , and the orthogonality condition holds; (b) the conditions on sparsity and the quality of estimation in Equation 36 hold, and the sparsity index obeys

$$s^2 \log(pn)^2 / n \rightarrow 0; \quad (39)$$

(c) the moderate deviation bound in Equation 37 holds; and (d) the sparse norm of the second-derivatives matrix is bounded as in Equation 38. Then the adaptivity condition in Equation 1 holds for the sequence  $\{P_n\}$ .

We note that Equation 39 requires that the true value of the nuisance parameter sufficiently sparse. We can relax this condition in some special cases to the requirement  $s \log(pn)^c / n \rightarrow 0$ , for some constant  $c$ , by using sample-splitting techniques (see Belloni et al. 2012). However, this requirement seems unavoidable in general.



**Proof:** We note above that  $T_{1,j} = 0$  by orthogonality. Under Equations 36 and 37, if  $s^2 \log(pn)^2/n \rightarrow 0$ , then  $T_{2,j}$  vanishes in probability, as by Hölder's inequality

$$T_{2,j} \leq \bar{T}_{2,j} \|\hat{\eta} - \eta_0\|_1 \lesssim_{P_n} \sqrt{s^2 \log(pn)^2/n} \rightarrow_{P_n} 0.$$

Also, if  $s^2 \log(pn)^2/n \rightarrow 0$ , then  $T_{3,j}$  vanishes in probability, as by Hölder's inequality and for sufficiently large  $n$

$$T_{3,j} \leq \bar{T}_{3,j} \|\hat{\eta} - \eta_0\|^2 \lesssim_{P_n} \sqrt{ns} \log(pn)/n \rightarrow_{P_n} 0.$$

The conclusion follows from Equation 34.

**4.2.2. Structure 2 (approximate sparsity).** Following Belloni et al. (2012), we next consider an approximate sparsity structure. Approximate sparsity imposes that, given a constant  $c > 0$ , we can decompose  $\eta_0$  into a sparse component  $\eta_0^m$  and a small nonsparse component  $\eta_0^r$ :

$$\begin{aligned} \eta_0 &= \eta_0^m + \eta_0^r, \text{ support}(\eta_0^m) \cap \text{support}(\eta_0^r) = \emptyset, \\ \|\eta_0^m\|_0 &\leq s, \quad \|\eta_0^r\|_2 \leq c\sqrt{s/n}, \quad \|\eta_0^r\|_1 \leq c\sqrt{s^2/n}. \end{aligned} \quad (40)$$

This condition allows for much more realistic and richer models than can be accommodated under exact sparsity. For example,  $\eta_0$  needs not have any zero components at all under approximate sparsity. In Section 5, we provide an example in which Equation 40 arises from a more primitive condition that the absolute values  $\{|\eta_{0j}|, j = 1, \dots, p\}$ , sorted in decreasing order, decay at a polynomial speed with respect to  $j$ .

Suppose that we have an estimator  $\hat{\eta}$  such that with probability  $1 - o(1)$  under  $P_n$  the following bounds hold:

$$\|\hat{\eta}\|_0 \lesssim s, \quad \|\hat{\eta} - \eta_0^m\|_2 \lesssim \sqrt{(s/n)\log(pn)}, \quad \|\hat{\eta} - \eta_0^m\|_1 \lesssim \sqrt{(s^2/n)\log(pn)}. \quad (41)$$

This condition is again a standard performance bound expected to hold for sparsity-based estimators under approximate sparsity conditions (see Belloni et al. 2012). Note that by the approximate sparsity condition, we also have that, with probability  $1 - o(1)$  under  $P_n$ ,

$$\|\hat{\eta} - \eta_0\|_2 \lesssim \sqrt{(s/n)\log(pn)}, \quad \|\hat{\eta} - \eta_0\|_1 \lesssim \sqrt{(s^2/n)\log(pn)}. \quad (42)$$

We can employ the same moderate deviation and bounded sparse norm conditions as in the previous subsection. In addition, we require the pointwise norm of the second-derivatives matrix to be bounded. Specifically, for any deterministic vector  $a \neq 0$ , we require

$$\left\| \partial_\eta \partial_\eta \hat{M}_j(\alpha_0) \right\|_{pw(a)} \lesssim_{P_n} 1. \quad (43)$$

This condition can be easily verified using ordinary laws of large numbers.

**Lemma 4 (elementary adaptivity for testing via approximate sparsity):** Let  $\{P_n\}$  be a sequence of probability laws. Assume that (a)  $\eta \mapsto \hat{M}(\alpha_0, \eta)$  and  $\eta \mapsto M(\alpha_0, \eta)$  are affine quadratic in  $\eta$ , and the orthogonality condition holds; (b) the conditions on

approximate sparsity in Equation 40 and the quality of estimation in Equation 41 hold, and the sparsity index obeys

$$s^2 \log(pn)^2 / n \rightarrow 0;$$

(c) the moderate deviation bound in Equation 37 holds; (d) the sparse norm of the second-derivatives matrix is bounded as in Equation 38; and (e) the pointwise norm of the second-derivatives matrix is bounded as in Equation 43. Then the adaptivity condition in Equation 1 holds:

$$\sqrt{n}(\hat{M}(\alpha_0, \hat{\eta}) - \hat{M}(\alpha_0, \eta_0)) \rightarrow_{P_n} 0.$$

### 4.3. Adaptivity for Estimation via Approximate Sparsity

We work with the approximate sparsity setup and the affine-quadratic model introduced in the previous subsections. In addition to the previous assumptions, we impose the following conditions on the components  $\partial_\eta \Gamma_{1,ml}$  of  $\partial_\eta \Gamma_1$ , where  $m = 1, \dots, k$  and  $l = 1, \dots, d$ . First, we need the following deviation and boundedness condition: For each  $m$  and  $l$ , we need that

$$\left\| \partial_\eta \hat{\Gamma}_{1,ml}(\eta_0) - \partial_\eta \Gamma_{1,ml}(\eta_0) \right\|_\infty \lesssim_{P_n} 1, \quad \left\| \partial_\eta \Gamma_{1,ml}(\eta_0) \right\|_\infty \lesssim 1. \quad (44)$$

Second, we require the sparse and pointwise norms of the following second-derivatives matrices to be stochastically bounded: For each  $m$  and  $l$ , we need that

$$\left\| \partial_\eta \partial_{\eta'} \hat{\Gamma}_{1,ml} \right\|_{\text{sp}(\ell, ns)} + \left\| \partial_\eta \partial_{\eta'} \hat{\Gamma}_{1,ml} \right\|_{\text{pw}(a)} \lesssim_{P_n} 1, \quad (45)$$

where  $a \neq 0$  is any deterministic vector. Both these conditions are mild. They can be verified using self-normalized moderate deviation theorems and using laws of large numbers for matrices, as discussed in the previous subsections.

**Lemma 5 (elementary adaptivity for estimation via approximate sparsity):** Consider a sequence  $\{P_n\}$  for which the conditions of Lemma 4 hold. In addition, assume that the deviation bound in Equation 44 holds and that the sparse norm and pointwise norms of the second-derivatives matrices are stochastically bounded as in Equation 45. Then the adaptivity condition in Equation 33 holds for the testing and estimation problem in the affine-quadratic model.

## 5. ANALYSIS OF THE INSTRUMENTAL VARIABLES MODEL WITH VERY MANY CONTROL AND INSTRUMENTAL VARIABLES

Consider the linear IV model with response variable

$$y_i = d_i' \alpha_0 + x_i' \beta_0 + \varepsilon_i, \quad E[\varepsilon_i] = 0, \quad \varepsilon_i \perp (z_i, x_i), \quad (46)$$

where here and below we write  $w \perp v$  to denote  $\text{Cov}(w, v) = 0$ ,  $y_i$  is the response variable, and  $d_i = (d_{ik})_{k=1}^{p^d}$  is a  $p^d$ -vector of endogenous variables, such that

$$\begin{aligned}
d_{i1} &= x_i' \gamma_{01} + z_i' \delta_{01} + u_{i1}, & E[u_{i1}] &= 0, & u_{i1} &\perp (z_i, x_i), \\
\vdots & & \vdots & & \vdots & \\
d_{ip^d} &= x_i' \gamma_{0p^d} + z_i' \delta_{0p^d} + u_{ip^d}, & E[u_{ip^d}] &= 0, & u_{ip^d} &\perp (z_i, x_i).
\end{aligned} \tag{47}$$

Here  $x_i = (x_{ij})_{j=1}^{p^x}$  is a  $p^x$ -vector of exogenous control variables, including a constant, and  $z_i = (z_i)_{i=1}^{p^z}$  is a  $p^z$ -vector of IV. We will have  $n$  i.i.d. draws of  $w_i = (y_i, d_i', x_i', z_i)'$  obeying this system of equations. We also assume that  $\text{Var}(w_i)$  is finite throughout so that the model is well defined.

The parameter value  $\alpha_0$  is our target. We allow  $p^x = p_n^x \gg n$  and  $p^z = p_n^z \gg n$ , but we maintain that  $p^d$  is fixed in our analysis. This model includes the case of many instruments and a small number of controls considered by Belloni et al. (2012) as a special case, and the analysis readily accommodates the case of many controls and no instruments (i.e., the linear regression model) considered by Belloni et al. (2013b, 2014a) and Zhang & Zhang (2014). For the latter, we simply set  $p_n^z = 0$  and impose the additional condition  $\varepsilon_i \perp u_i$  for  $u_i = (u_{ij})_{j=1}^{p^d}$ , which together with  $\varepsilon_i \perp x_i$  implies that  $\varepsilon_i \perp d_i$ . We also note that the condition  $\varepsilon_i \perp x_i, z_i$  is weaker than the condition  $E[\varepsilon_i | x_i, z_i] = 0$ , which allows for some misspecification of the model.

We may have that  $z_i$  and  $x_i$  are correlated so that  $z_i$  are valid instruments only after controlling for  $x_i$ ; specifically, we let  $z_i = \Pi x_i + \zeta_i$ , for  $\Pi$  a  $p_n^z \times p_n^x$  matrix and  $\zeta_i$  a  $p_n^z$ -vector of unobservables with  $x_i \perp \zeta_i$ . Substituting this expression for  $z_i$  as a function of  $x_i$  into Equation 46 gives a system for  $y_i$  and  $d_i$  that depends only on  $x_i$ :

$$\begin{aligned}
y_i &= x_i' \theta_0 + \rho_i^y, & E[\rho_i^y] &= 0, & \rho_i^y &\perp x_i, \\
d_{i1} &= x_i' \vartheta_{01} + \rho_{i1}^d, & E[\rho_{i1}^d] &= 0, & \rho_{i1}^d &\perp x_i, \\
\vdots & & \vdots & & \vdots & \\
d_{ip^d} &= x_i' \vartheta_{0p^d} + \rho_{ip^d}^d, & E[\rho_{ip^d}^d] &= 0, & \rho_{ip^d}^d &\perp x_i.
\end{aligned} \tag{48}$$

Because the dimension  $p = p_n$  of

$$\eta_0 = \left( \theta_0', (\vartheta_{0k}', \gamma_{0k}', \delta_{0k}')_{k=1}^{p^d} \right)'$$

may be larger than  $n$ , informative estimation and inference about  $\alpha_0$  are impossible without imposing restrictions on  $\eta_0$ .

To state our assumptions, we fix a collection of positive constants  $(a, A, c, C)$ , where  $a > 1$ , and a sequence of constants  $\delta_n \searrow 0$  and  $\ell_n \nearrow \infty$ . These constants will not vary with  $P$ ; rather, we will work with collections of  $P$  defined by these constants.

**Condition AS 1:** We assume that  $\eta_0$  is approximately sparse, namely that the decreasing rearrangement  $(|\eta_0|_j^*)_{j=1}^p$  of absolute values of coefficients  $(|\eta_{0j}|)_{j=1}^p$  obeys

$$|\eta_0|_j^* \leq A j^{-a}, \quad a > 1, \quad j = 1, \dots, p. \tag{49}$$

Given this assumption, we can decompose  $\eta_0$  into a sparse component  $\eta_0^m$  and small nonsparse component  $\eta_0^r$ :

$$\begin{aligned} \eta_0 &= \eta_0^m + \eta_0^r, \quad \text{support}(\eta_0^m) \cap \text{support}(\eta_0^r) = \emptyset, \\ \|\eta_0^m\|_0 &\leq s, \quad \|\eta_0^r\|_2 \leq c\sqrt{s/n}, \quad \|\eta_0^r\|_1 \leq c\sqrt{s^2/n}, \\ s &= cn^{\frac{1}{2\alpha}}, \end{aligned} \tag{50}$$

where the constant  $c$  depends only on  $(\mathbf{a}, \mathbf{A})$ .

**Condition AS 2:** We assume that

$$s^2 \log(pn)^2 / n \leq o(1). \tag{51}$$

We perform inference on  $\alpha_0$  using the empirical analog of theoretical equations:

$$\mathbf{M}(\alpha_0, \eta_0) = 0, \quad \mathbf{M}(\alpha, \eta) := \mathbb{E}[\psi(w_i, \alpha, \eta)], \tag{52}$$

where  $\psi = (\psi_k)_{k=1}^p$  is defined by

$$\psi_k(w_i, \alpha, \eta) := \left( y_i - x_i' \theta - \sum_{\bar{k}=1}^{p-d} (d_{i\bar{k}} - x_i' \vartheta_{\bar{k}}) \alpha_{\bar{k}} \right) (x_i' \gamma_k + z_i' \delta_k - x_i' \vartheta_k).$$

We can verify that the following orthogonality condition holds:

$$\partial_{\eta'} \mathbf{M}(\alpha_0, \eta) \Big|_{\eta=\eta_0} = 0. \tag{53}$$

This means that missing the true value  $\eta_0$  by a small amount does not invalidate the moment condition. Therefore, the moment condition will be relatively insensitive to nonregular estimation of  $\eta_0$ .

We denote the empirical analog of Equation 52 as

$$\hat{\mathbf{M}}(\hat{\alpha}, \hat{\eta}) = 0, \quad \hat{\mathbf{M}}(\alpha, \eta) := \mathbb{E}_n[\psi_i(\alpha, \eta)]. \tag{54}$$

Inference based on this condition can be shown to be immunized against small selection mistakes by virtue of orthogonality.

The above formulation is a special case of the linear-affine model. Indeed, here we have

$$\begin{aligned} \mathbf{M}(\alpha, \eta) &= \Gamma_1(\eta)\alpha + \Gamma_2(\eta), \quad \hat{\mathbf{M}}(\alpha, \eta) = \hat{\Gamma}_1(\eta)\alpha + \hat{\Gamma}_2(\eta), \\ \Gamma_1(\eta) &= \mathbb{E}[\psi^a(w_i, \eta)], \quad \hat{\Gamma}_1(\eta) = \mathbb{E}_n[\psi^a(w_i, \eta)], \\ \Gamma_2(\eta) &= \mathbb{E}[\psi^b(w_i, \eta)], \quad \hat{\Gamma}_2(\eta) = \mathbb{E}_n[\psi^b(w_i, \eta)], \end{aligned}$$

where

$$\begin{aligned} \psi_{k,\bar{k}}^a(w_i, \eta) &= - \left( d_{i\bar{k}} - x_i' \vartheta_{\bar{k}} \right) (x_i' \gamma_k + z_i' \delta_k - x_i' \vartheta_k), \\ \psi_k^b(w_i, \eta) &= (y_i - x_i' \theta) (x_i' \gamma_k + z_i' \delta_k - x_i' \vartheta_k). \end{aligned}$$

Consequently we can use the results of the previous section. To do so, we need to provide a suitable estimator for  $\eta_0$ . Here we use the Lasso and post-Lasso estimators, as defined in Belloni et al. (2012), to deal with nonnormal errors and heteroscedasticity.

**Algorithm 1 (estimation of  $\eta_0$ ):** (a) For each  $k$ , do Lasso or post-Lasso regression of  $d_{ik}$  on  $x_i, z_i$  to obtain  $\hat{\gamma}_k$  and  $\hat{\delta}_k$ . (b) Do Lasso or post-Lasso regression of  $y_i$  on  $x_i$  to get  $\hat{\theta}$ . (c) Do Lasso or post-Lasso regression of  $\hat{d}_{ik} = x_i' \hat{\gamma}_k + z_i' \hat{\delta}_k$  on  $x_i$  to get  $\hat{\vartheta}_k$ . The estimator of  $\eta_0$  is given by  $\hat{\eta} = \left( \hat{\theta}', \left( \hat{\vartheta}'_k, \hat{\gamma}'_{0k}, \hat{\delta}'_k \right)_{k=1}^{p^d} \right)'$ .

We then use

$$\hat{\Omega}(\alpha, \hat{\eta}) = \mathbb{E}_n \left[ \psi(w_i, \alpha, \hat{\eta}) \psi(w_i, \alpha, \hat{\eta})' \right]$$

to estimate the variance matrix  $\Omega(\alpha, \eta_0) = \mathbb{E}_n \left[ \psi(w_i, \alpha, \eta_0) \psi(w_i, \alpha, \eta_0)' \right]$ . We formulate the orthogonal score statistic and the  $C(\alpha)$ -statistic,

$$S(\alpha) := \hat{\Omega}_n^{-1/2}(\alpha, \hat{\eta}) \sqrt{n} \hat{M}(\alpha, \hat{\eta}), \quad C(\alpha) = \|S(\alpha)\|^2, \quad (55)$$

as well as our estimator  $\hat{\alpha}$ :

$$\hat{\alpha} = \arg \min_{\alpha \in \mathcal{A}} \left\| \sqrt{n} \hat{M}(\alpha, \hat{\eta}) \right\|^2.$$

Note also that  $\hat{\alpha} = \arg \min_{\alpha \in \mathcal{A}} C(\alpha)$  under mild conditions, as we work with exactly identified systems of equations. We also need to specify a variance estimator  $\hat{V}_n$  for the large sample variance  $V_n$  of  $\hat{\alpha}$ . We set  $\hat{V}_n = \left( \hat{\Gamma}_1(\hat{\eta})' \right)^{-1} \hat{\Omega}(\hat{\alpha}, \hat{\eta}) \left( \hat{\Gamma}_1(\hat{\eta}) \right)^{-1}$ .

To estimate the nuisance parameter, we impose the following condition. Let  $f_i := (f_{ij})_{j=1}^{p_l} := (x_i', z_i)'$ ;  $b_i := (b_{il})_{l=1}^{p_b} := (y_i, d_i', \bar{d}_i)'$ , where  $\bar{d}_i = (\bar{d}_{ik})_{k=1}^{p^d}$  and  $\bar{d}_{ik} := x_i' \gamma_{0k} + z_i' \delta_{0k}$ ;  $v_i = (v_{il})_{l=1}^{p_b} := (\varepsilon_i, \rho_i^y, \rho_i^d, \varrho_i)'$ , with  $\varrho_i = (\varrho_{ik})_{k=1}^{p^d}$  and  $\varrho_{ik} := d_{ik} - \bar{d}_{ik}$ . Let  $\tilde{b}_i := b_i - \mathbb{E}[b_i]$ .

**Condition RF:** (a) The eigenvalues of  $\mathbb{E}[f_i f_i']$  are bounded from above by  $C$  and from below by  $c$ . For all  $j$  and  $l$ , (b)  $\mathbb{E}[b_{il}^2] + \mathbb{E}[|f_{ij}^2 \tilde{b}_{il}^2|] + 1/\mathbb{E}[f_{ij}^2 v_{il}^2] \leq C$  and  $\mathbb{E}[|f_{ij}^2 v_{il}^2|] \leq \mathbb{E}[|f_{ij}^2 \tilde{b}_{il}^2|]$ , (c)  $\mathbb{E}[|f_{ij}^3 v_{il}^3|]^2 \log^3(pn)/n \leq \delta_n$ , and (d)  $s \log(pn)/n \leq \delta_n$ . With probability no less than  $1 - \delta_n$ , we have that (e)  $\max_{i \leq n, j} f_{ij}^2 [s^2 \log(pn)] / n \leq \delta_n$ ,  $\max_{l, j} |(\mathbb{E}_n - \mathbb{E})[f_{ij}^2 v_{il}^2]| + |(\mathbb{E}_n - \mathbb{E})[f_{ij}^2 \tilde{b}_{il}^2]| \leq \delta_n$ , and (f)  $\left\| \mathbb{E}_n[f_i f_i'] - \mathbb{E}[f_i f_i'] \right\|_{\text{sp}(\ell_n s)} \leq \delta_n$ .

The conditions are motivated by those given in Belloni et al. (2012). The current conditions are made slightly stronger to account for the fact that we use zero covariance conditions in formulating the moments. Some conditions could be easily relaxed at a cost of more complicated exposition.

To estimate the variance matrix and establish asymptotic normality, we also need the following condition. Let  $q > 4$  be a fixed constant.

**Condition SM:** For each  $l$  and  $k$ , (a)  $\mathbb{E}[|b_{il}|^q] + \mathbb{E}[|v_{il}|^q] \leq C$ , (b)  $c \leq \mathbb{E}[\varepsilon_i^2 | x_i, z_i] \leq C$ ,  $c < \mathbb{E}[\varrho_{ik}^2 | x_i, z_i] \leq C$  almost surely, and (c)  $\sup_{\alpha \in \mathcal{A}} \|\alpha\|_2 \leq C$ .

Under the conditions set forth above, we have the following result on the validity of post-selection and post-regularization inference using the  $C(\alpha)$ -statistic and estimators derived from it.

**Proposition 5 [valid inference in large linear models using  $C(\alpha)$ -statistics]:** Let  $\mathbf{P}_n$  be the collection of all  $\mathbf{P}$  such that Conditions AS 1 and 2, RF, and SM hold for the given  $n$ . Then uniformly in  $\mathbf{P} \in \mathbf{P}_n$ , we find that  $S(\alpha_0) \rightsquigarrow \mathcal{N}(0, I)$  and  $C(\alpha_0) \rightsquigarrow \chi^2(p^d)$ . As

a consequence, the confidence set  $CR_{1-a} = \{\alpha \in \mathcal{A} : C(\alpha) \leq c(1-a)\}$ , where  $c(1-a)$  is the  $1-a$ -quantile of a  $\chi^2(p^d)$ , is uniformly valid for  $\alpha_0$ , in the sense that

$$\lim_{n \rightarrow \infty} \sup_{P \in \mathcal{P}_n} |\mathbb{P}(\alpha_0 \in CR_{1-a}) - (1-a)| = 0.$$

Furthermore, for  $V_n = (\Gamma_1)^{-1} \Omega(\alpha_0, \eta_0) (\Gamma_1)^{-1}$ , we have that

$$\lim_{n \rightarrow \infty} \sup_{P \in \mathcal{P}_n} \sup_{R \in \mathcal{R}} |\mathbb{P}(V_n^{-1/2}(\hat{\alpha} - \alpha_0) \in R) - \mathbb{P}(\mathcal{N}(0, I) \in R)| = 0,$$

where  $\mathcal{R}$  is the collection of all convex sets. Moreover, the result continues to apply if  $V_n$  is replaced by  $\hat{V}_n$ . Thus,  $CR_{1-a}^l = [l\hat{\alpha} \pm c(1-a/2)(l\hat{V}_n l/n)^{1/2}]$ , where  $c(1-a/2)$  is the  $(1-a/2)$ -quantile of an  $\mathcal{N}(0, 1)$ , provides a uniformly valid confidence set for  $l\alpha_0$ :

$$\lim_{n \rightarrow \infty} \sup_{P \in \mathcal{P}_n} |\mathbb{P}(l\alpha_0 \in CR_{1-a}^l) - (1-a)| = 0.$$

The proof of Proposition 5 is given in the **Supplemental Appendix** (follow the Supplemental Material link from the Annual Reviews home page at <http://www.annualreviews.org>).

### 5.1. Simulation Illustration

In this section, we provide results from a small Monte Carlo simulation to illustrate the performance of the estimator resulting from the application of Algorithm 1 in a small sample setting. As comparison, we report results from two commonly used unprincipled alternatives for which uniformly valid inference over the class of approximately sparse models does not hold. Simulation parameters were chosen so that approximate sparsity holds but exact sparsity is violated in such a way that we expect the unprincipled procedures to perform poorly.

For our simulation, we generate data as  $n$  i.i.d. draws from the model:

$$\begin{matrix} y_i = \alpha d_i + x_i' \beta + 2\varepsilon_i \\ d_i = x_i' \gamma + z_i' \delta + u_i \\ z_i = \Pi x_i + 0.125 \zeta_i \end{matrix} \quad \left| \quad \begin{matrix} \varepsilon_i \\ u_i \\ \zeta_i \\ x_i \end{matrix} \right. \sim \mathcal{N} \left( 0, \begin{pmatrix} 1 & 0.6 & 0 & 0 \\ 0.6 & 1 & 0 & 0 \\ 0 & 0 & I_{p_n^z} & 0 \\ 0 & 0 & 0 & \Sigma \end{pmatrix} \right),$$

where  $\Sigma$  is a  $p_n^x \times p_n^x$  matrix with  $\Sigma_{kj} = (0.5)^{|j-k|}$  and  $I_{p_n^z}$  is a  $p_n^z \times p_n^z$  identity matrix. We set the number of potential control variables ( $p_n^x$ ) to 200, the number of instruments ( $p_n^z$ ) to 150, and the number of observations ( $n$ ) to 200. For model coefficients, we set  $\alpha = 0$ ,  $\beta = \gamma$  as  $p_n^x$ -vectors with entries  $\beta_j = \gamma_j = 1/(9\nu)$ ,  $\nu = 4/9 + \sum_{j=5}^{p_n^x} 1/j^2$  for  $j \leq 4$  and  $\beta_j = \gamma_j = 1/(j^2\nu)$  for  $j > 4$ ,  $\delta$  as a  $p_n^z$ -vector with entries  $\delta_j = 3/j^2$ , and  $\Pi = [I_{p_n^z}, 0_{p_n^z \times (p_n^x - p_n^z)}]$ . We report results based on 1,000 simulation replications.

We provide results for four different estimators: an infeasible oracle estimator that knows the nuisance parameter  $\eta$ , two naive estimators, and the proposed double-selection estimator. The results for the proposed double-selection procedure are obtained following Algorithm 1 using post-Lasso at every step. To obtain the oracle results, we run standard IV regression of  $y_i - E[y_i|x_i]$  on  $d_i - E[d_i|x_i]$  using the single instrument  $\zeta_i' \delta$ . The expected values are obtained from the model above, and  $\zeta_i' \delta$  provides the information in the instruments that is unrelated to the controls.

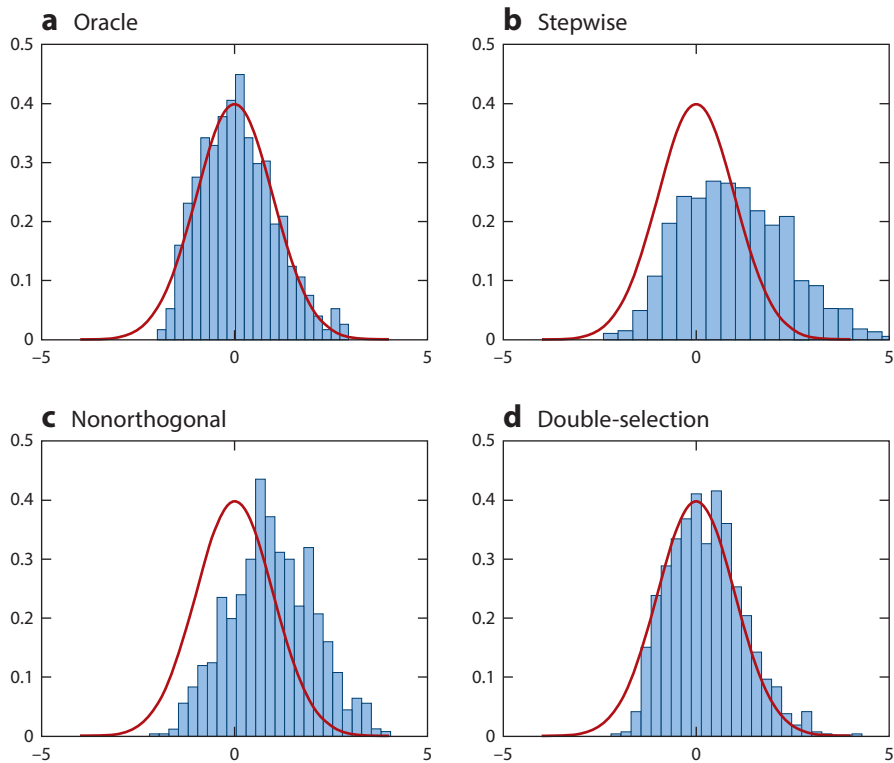
The two naive alternatives offer unprincipled, although potentially intuitive, alternatives. The first naive estimator follows Algorithm 1 but replaces Lasso/post-Lasso with stepwise regression with a  $p$  value for entry of 0.05 and a  $p$  value for removal of 0.10 (stepwise). The second naive estimator (nonorthogonal) corresponds to the use of a moment condition that does not satisfy the orthogonality condition described previously but will produce valid inference when perfect model selection in the regression of  $d$  on  $x$  and  $z$  is possible or when perfect model selection in the regression of  $y$  on  $x$  is possible and an instrument is selected in the regression of  $d$  on  $x$  and  $z$ .<sup>5</sup>

All of the Lasso and post-Lasso estimates are obtained using the data-dependent penalty level from Belloni & Chernozhukov (2013). This penalty level depends on a standard deviation that is estimated by adapting the iterative algorithm described in Belloni et al. (2012, appendix A) using post-Lasso at each iteration. For inference in all cases, we use standard  $t$ -tests based on conventional homoscedastic IV standard errors obtained from the final IV step performed in each strategy.

We display the simulation results in **Figure 1**, and we report the median bias, median absolute deviation, and size of 5% level tests for each procedure in **Table 1**. For each estimator, we plot the simulation estimate of the sampling distribution of the estimator centered around the true parameter and scaled by the estimated standard error. With this standardization, usual asymptotic approximations would suggest that these curves should line up with an  $\mathcal{N}(0, 1)$  density function, which is displayed as the red solid line in the figure. We can see that the oracle estimator and the double-selection estimator are centered correctly and line up reasonably well with the  $\mathcal{N}(0, 1)$  density function, although both estimators exhibit some mild skewness. It is interesting that the sampling distributions of the oracle and double-selection estimators are very similar, as predicted by the theory. In contrast, both the naive estimators are centered far from zero, and it is clear that the asymptotic approximation provides a very poor guide to the finite-sample distribution of these estimators in the design considered.

The poor inferential performance of the two naive estimators is driven by different phenomena. The unprincipled use of stepwise regression fails to control spurious inclusion of irrelevant variables, which leads to the inclusion of many essentially irrelevant variables, resulting in many-instrument-type problems (e.g., Chao et al. 2012). In addition, the spuriously included variables are those most highly correlated to the noise within the sample, which adds an additional type of endogeneity bias. The failure of the nonorthogonal method is driven by the fact that perfect model selection is not possible within the present design: Here we have model selection mistakes in which control variables that are correlated to the instruments but only moderately correlated to the outcome and endogenous variable are missed. Such exclusions result in standard omitted variables bias in the estimator for the parameter of interest and substantial size distortions. The additional step in the double-selection procedure can be viewed as a way to guard against such mistakes. Overall, the results illustrate the uniformity claims made in the preceding section. The feasible double-selection procedure following from Algorithm 1 performs similarly to the semiparametrically efficient infeasible oracle method. We obtain good inferential properties, with the asymptotic approximation providing a fairly good guide to the behavior of the estimator despite

<sup>5</sup>Specifically, for the second naive alternative (nonorthogonal), we first do Lasso regression of  $d$  on  $x$  and  $z$  to obtain Lasso estimates of the coefficients  $\gamma$  and  $\delta$ . Denote these estimates as  $\hat{\gamma}_L$  and  $\hat{\delta}_L$ , and denote the indices of the coefficients estimated to be nonzero as  $\hat{I}_x^d = \{j: \hat{\gamma}_{Lj} \neq 0\}$  and  $\hat{I}_z^d = \{j: \hat{\delta}_{Lj} \neq 0\}$ . We then run Lasso regression of  $y$  on  $x$  to learn the identities of controls that predict the outcome. We denote the Lasso estimates as  $\hat{\theta}_L$  and keep track of the indices of the coefficients estimated to be nonzero as  $\hat{I}_x^y = \{j: \hat{\theta}_{Lj} \neq 0\}$ . We then take the union of the controls selected in either step  $\hat{I}_x = \hat{I}_x^y \cup \hat{I}_x^d$ . The estimator of  $\alpha$  is then obtained as the usual 2SLS estimator of  $y_i$  on  $d_i$  using all selected elements from  $x_i$ ,  $x_{ij}$  such that  $j \in \hat{I}_x$ , as controls and the selected elements from  $z_i$ ,  $z_{ij}$  such that  $j \in \hat{I}_z^d$ , as instruments.



**Figure 1**

Histograms of the estimator from each method centered around the true parameters and scaled by the estimated standard error from the simulation experiment: (a) oracle, (b) stepwise, (c) nonorthogonal, and (d) double selection. The red curve is the probability density function of a standard normal, which will correspond to the sampling distribution of the estimator under the asymptotic approximation.

working in a setting in which perfect model selection is impossible. Although simply illustrative of the theory, the results are reassuring and in line with extensive simulations in the linear model with many controls provided in Belloni et al. (2014a), in the IV model with many instruments and a small number of controls provided in Belloni et al. (2012), and in linear panel data models provided in Belloni et al. (2014b).

### 5.2. Empirical Illustration: Logit Demand Estimation

As further illustration of the approach, we provide a brief empirical example in which we estimate the coefficients in a simple logit model of demand for automobiles using market share data. Our example is based on the data and most basic strategy from Berry et al. (1995). Specifically, we estimate the parameters from the model

$$\begin{aligned} \log(s_{it}) - \log(s_{0t}) &= \alpha_0 p_{it} + x'_{it} \beta_0 + \varepsilon_{it}, \\ p_{it} &= z'_{it} \delta_0 + x'_{it} \gamma_0 + u_{it}, \end{aligned}$$

where  $s_{it}$  is the market share of product  $i$  in market  $t$  with product zero denoting the outside option,  $p_{it}$  is the price and is treated as endogenous,  $x_{it}$  are observed included product characteristics, and



**Table 1 Summary of simulation results for the estimation of  $\alpha$**

Method	Median bias	Median absolute deviation	Size
Oracle	0.015	0.247	0.043
Stepwise	0.282	0.368	0.261
Nonorthogonal	0.084	0.112	0.189
Double selection	0.069	0.243	0.053

This table summarizes the simulation results from a linear instrumental variables model with many instruments and controls. Estimators include an infeasible oracle as a benchmark, two naive alternatives (stepwise and nonorthogonal) described in the text, and our proposed feasible valid procedure (double selection). Size is for 5% level tests.

$z_{it}$  are instruments. One could also adapt the proposed variable selection procedures to extensions of this model such as the nested logit model or models allowing for random coefficients (see, e.g., Gillen et al. 2014 for an example with a random coefficient).

In our example, we use the same set of product characteristics ( $x$  variables) as used in obtaining the basic results in Berry et al. (1995). Specifically, we use five variables in  $x_{it}$ : a constant, an air conditioning dummy, horsepower divided by weight, miles per dollar, and vehicle size. We refer to these five variables as the baseline set of controls.

We also adopt the argument from Berry et al. (1995) to form our potential instruments. Berry et al. (1995) argued that characteristics of other products will satisfy an exclusion restriction,  $E[\varepsilon_{it}|x_{j\tau}] = 0$  for any  $\tau$  and  $j \neq i$ , and thus that any function of characteristics of other products may be used as an instrument for price. This condition leaves a very high-dimensional set of potential instruments, as any combination of functions of  $\{x_{j\tau}\}_{j \neq i, \tau \geq 1}$  may be used to instrument for  $p_{it}$ . To reduce the dimensionality, Berry et al. (1995) used intuition and an exchangeability argument to motivate the consideration of a small number of these potential instruments formed by taking sums of product characteristics formed by summing over products excluding product  $i$ . Specifically, we form baseline instruments by taking

$$z_{k,it} = \left( \sum_{r \neq i, r \in \mathcal{I}_f} x_{k,rt}, \sum_{r \neq i, r \notin \mathcal{I}_f} x_{k,rt} \right),$$

where  $x_{k,it}$  is the  $k$ -th element of vector  $x_{it}$ , and  $\mathcal{I}_f$  denotes the set of products produced by firm  $f$ . This choice yields a vector  $z_{it}$  consisting of 10 instruments. We refer to this set of instruments as the baseline instruments.

Although the choice of the baseline instruments and controls is motivated by good intuition and economic theory, we note that theory does not clearly state which product characteristics or instruments should be used in the model. Theory also fails to indicate the functional form with which any such variables should enter the model. The high-dimensional methods outlined in this article offer one strategy to help address these concerns that complements the economic intuition motivating the baseline controls and instruments. As an illustration, we consider an expanded set of controls and instruments. We augment the set of potential controls with all first-order interactions of the baseline variables, quadratics, and cubics in all continuous baseline variables, and a time trend that yields a total of 24  $x$  variables. We refer to these as the augmented controls. We then take sums of these characteristics as potential instruments following the original strategy that yields 48 potential instruments.

**Table 2** Estimates of price coefficient

	Price coefficient	Standard error	Number inelastic
<b>Estimates without selection</b>			
Baseline OLS	-0.089	0.004	1,502
Baseline 2SLS	-0.142	0.012	670
Augmented OLS	-0.099	0.005	1,405
Augmented 2SLS	-0.127	0.014	874
<b>2SLS estimates with double selection</b>			
Baseline 2SLS selection	-0.185	0.014	139
Augmented 2SLS selection	-0.221	0.015	12

This table reports estimates of the coefficient on price along with the estimated standard error obtained using different sets of controls and instruments. The rows Baseline OLS and Baseline 2SLS, respectively, provide ordinary least-squares (OLS) and two-stage least-squares (2SLS) results using the baseline set of variables (5 controls and 10 instruments) described in the text. The rows Augmented OLS and Augmented 2SLS are defined similarly but use the augmented set of variables described in the text (24 controls and 48 instruments). The rows Baseline 2SLS with Selection and Augmented 2SLS with Selection apply the double-selection approach developed in this article to select a set of controls and instruments and perform valid post-selection inference about the estimated price coefficient in which selection occurs considering only the baseline variables. For each procedure, we also report the point estimate of the number of products for which demand is estimated to be inelastic in the column Number inelastic.

We present estimation results in **Table 2**. We report results obtained by applying the method outlined in Algorithm 1 using just the baseline set of five product characteristics and 10 instruments in the row labeled “Baseline 2SLS selection” and results obtained by applying the method to the augmented set of 24 controls and 48 instruments in the row labeled “Augmented 2SLS selection.” In each case, we apply the method outlined in Algorithm 1 using post-Lasso in each step and forcing the intercept to be included in all models. We employ the heteroscedasticity robust version of post-Lasso of Belloni et al. (2012) following the implementation algorithm provided in their appendix A. For comparison, we also report ordinary least-squares (OLS) and two-stage least-squares (2SLS) estimates using only the baseline variables, and we report OLS and 2SLS estimates using the augmented variable set. All standard errors are conventional heteroscedasticity robust standard errors.

Considering first estimates of the price coefficient, we see that the estimated price coefficient increases in magnitude as we move from OLS to 2SLS and then to the selection-based results. After selection using only the original variables, we estimate the price coefficient to be  $-0.185$  with an estimated standard error of 0.014 compared to an OLS estimate of  $-0.089$  with an estimated standard error of 0.004 and a 2SLS estimate of  $-0.142$  with an estimated standard error of 0.012. In this case, all five controls are selected in the log-share on controls regression, all five controls but only four instruments are selected in the price on controls and instruments regression, and four of the controls are selected for the price on controls relationship. The difference between the baseline results is thus largely driven by the difference in instrument sets. The change in the estimated coefficient is consistent with the wisdom from the many instrument literature that the inclusion of irrelevant instruments biases 2SLS toward OLS.

With the larger set of variables, our post-model selection estimator of the price coefficient is  $-0.221$  with an estimated standard error of  $0.015$  compared to the OLS estimate of  $-0.099$  with an estimated standard error of  $0.005$  and 2SLS estimate of  $-0.127$  with an estimated standard error of  $0.014$ . Here, we see some evidence that the original set of controls may have been overly parsimonious as we select some terms that were not included in the baseline variable set. We also see closer agreement between the OLS estimate and 2SLS estimate without selection, which is likely driven by the larger number of instruments considered and the usual bias toward OLS seen in 2SLS with many weak or irrelevant instruments. In the log-share on controls regression, we have eight control variables selected, and we have seven controls and only four instruments selected in the price on controls and instrument regression. We also have 13 variables selected for the price on controls relationship. The selection of these additional variables suggests that there is important nonlinearity missed by the baseline set of variables.

The most interesting feature of the results is that estimates of own-price elasticities become more plausible as we move from the baseline results to the results based on variable selection with a large number of controls. Recall that facing inelastic demand is inconsistent with profit-maximizing price choice within the present context, so theory would predict that demand should be elastic for all products. However, the baseline point estimates imply inelastic demand for 670 products. When we use the larger set of instruments without selection, the number of products for which we estimate inelastic demand increases to 874, with the increase generated by the 2SLS coefficient estimate moving back toward the OLS estimate. The use of the variable selection results provides results closer to the theoretical prediction. The point estimates based on selection from only the baseline variables imply inelastic demand for 139 products, and we estimate inelastic demand for only 12 products using the results based on selection from the larger set of variables. Thus, the new methods provide the most reasonable estimates of own-price elasticities.

We conclude by noting that the simple specification above suffers from the usual drawbacks of the logit demand model. However, the example illustrates how the application of the methods outlined may be used in the estimation of structural parameters in economics and adds to the plausibility of the resulting estimates. In this example, we see that we obtain more sensible estimates of key parameters with at most a modest cost in increased estimation uncertainty after applying the methods in this article while considering a flexible set of variables.

## 6. OVERVIEW OF RELATED LITERATURE

Inference following model selection or regularization more generally has been an active area of research in econometrics and statistics for the past several years. In this section, we provide a brief overview of this literature highlighting some key developments. This review is necessarily selective because of the large number of papers available and the rapid pace at which new papers are appearing. We choose to focus on papers that deal specifically with high-dimensional nuisance parameter settings and note that the ideas in these papers apply in low-dimensional settings as well.

Early work on inference in high-dimensional settings focused on inference based on the so-called perfect recovery property (see, e.g., Fan & Li 2001 for an early paper, Fan & Lv 2010 for a more recent review, and Bühlmann & van de Geer 2011 for a textbook treatment). A consequence of this property is that model selection does not impact the asymptotic distribution of the parameters estimated in the selected model. This feature allows one to do inference using standard approximate distributions for the parameters of the selected model ignoring that model selection was done. Although convenient and fruitful in many applications (e.g., signal processing), such results effectively rely on strong conditions that imply that one will be able to perfectly select the correct model. For example, such results in linear models require the so-called

beta-min condition (Bühlmann & van de Geer 2011) that all but a small number of coefficients are exactly zero and the remaining nonzero coefficients are bounded away from zero, effectively ruling out variables that have small, nonzero coefficients. Such conditions seem implausible in many applications, especially in econometrics, and relying on such conditions produces asymptotic approximations that may provide very poor approximations to finite-sample distributions of estimators as they are not uniformly valid over sequences of models that include even minor deviations from conditions implying perfect model selection. The concern about the lack of uniform validity of inference based on oracle properties was raised in a series of papers (e.g., Leeb & Pötscher 2008a,b), and the more recent work on post-model selection inference has been focused on offering procedures that provide uniformly valid inference over interesting (large) classes of models that include cases in which perfect model selection will not be possible.

To our knowledge, the first work to formally and expressly address the problem of obtaining uniformly valid inference following model selection is by Belloni et al. (2010) who considered inference about parameters on a low-dimensional set of endogenous variables following selection of instruments from among a high-dimensional set of potential instruments in a homoscedastic, Gaussian IV model. The approach does not rely on implausible beta-min conditions that imply perfect model selection but instead relies on the fact that the moment condition underlying IV estimation satisfies the orthogonality condition in Equation 2 and the use of high-quality variable selection methods. Belloni et al. (2012) further developed these ideas in the context of providing uniformly valid inference about the parameters on endogenous variables in the IV context with many instruments to allow non-Gaussian heteroscedastic disturbances. These principles have also been applied by Belloni et al. (2013b), who developed approaches for regression and IV models with Gaussian errors; Belloni et al. (2014a), who covered the estimation of the parametric components of the partially linear model, and the estimation of average treatment effects, and provided a formal statement of the orthogonality condition in Equation 2; Farrell (2014), who covered average treatment effects with discrete, multivalued treatments; Kozbur (2014), who covered additive nonparametric models; and Belloni et al. (2014b), who extended the IV and partially linear model results to allow for fixed effects panel data and clustered dependence structures. The most recent, general approach has been provided by Belloni et al. (2013a), who analyzed inference about parameters defined by a continuum of orthogonalized estimating equations with infinite-dimensional nuisance parameters and developed positive results on inference. The framework in Belloni et al. (2013a) is general enough to cover the aforementioned papers and many other parametric and semiparametric models considered in economics.

As noted above, providing uniformly valid inference following model selection is closely related to the use of Neyman's  $C(\alpha)$ -statistic. Valid confidence regions can be obtained by inverting tests based on these statistics, and minimizers of  $C(\alpha)$ -statistics may be used as point estimators. The use of  $C(\alpha)$  statistics for testing and estimation in high-dimensional approximately sparse models was first explored in the context of high-dimensional quantile regression in Belloni et al. (2013c,d) and in the context of high-dimensional logistic regression and other high-dimensional generalized linear models in Belloni et al. (2013e). More recent uses of  $C(\alpha)$ -statistics (or close variants, under different names) include those by Voorman et al. (2014), Ning & Liu (2014), and Yang et al. (2014).

There have also been parallel developments based upon *ex post* debiasing of estimators. This approach is mathematically equivalent to doing classical one-step corrections in the general framework of Section 2. Indeed, although at first glance this debiasing approach may appear distinct from that taken in the papers listed above in this section, it is the same as approximately solving—by doing one Gauss-Newton step—orthogonal estimating equations satisfying Equation 2. The general results of Section 2 suggest that these approaches, the exact solving and one-step

solving, are generally first-order asymptotically equivalent, although higher-order differences may persist. To the best of our knowledge, the one-step correction approach was first employed in high-dimensional sparse models by Zhang & Zhang (2014), who covered the homoscedastic linear model, as well as in several of their follow-up works. This approach has been further used by van de Geer et al. (2014), who covered homoscedastic linear models and some generalized linear models, and Javanmard & Montanari (2014), who offered a related, although somewhat different, approach. Belloni et al. (2013d,e) also offered results on one-step corrections as part of their analysis of estimation and inference based on the orthogonal estimating equations. We would not expect the use of orthogonal estimating equations or the use of one-step corrections to dominate each other in all cases, although computational evidence from Belloni et al. (2013e) suggests that the use of exact solutions to orthogonal estimating equations may be preferable to approximate solutions obtained from one-step corrections in the contexts they considered.

Another branch of the recent literature takes a complementary, but logically distinct, approach that aims at doing valid inference for the parameters of a pseudo-true model that results from the use of a model selection procedure (see Berk et al. 2013). Specifically, this approach conditions on a model selected by a data-dependent rule and then attempts to do inference—conditional on the selection event—for the parameters of the selected model, which may deviate from the true model that generated the data. Related developments within this approach appear in G'Sell et al. (2013), Lee et al. (2013), Lee & Taylor (2014), Lockhart et al. (2014), Loftus & Taylor (2014), Taylor et al. (2014), and Fithian et al. (2014). It seems intellectually very interesting to combine the developments of the present article (and other preceding papers cited above) with developments in this literature.

The previously mentioned work focuses on doing inference for low-dimensional parameters in the presence of high-dimensional nuisance parameters. There have also been developments on performing inference for high-dimensional parameters. Belloni & Chernozhukov (2011) proposed inverting a Lasso performance bound in order to construct a simultaneous, Scheffé-style confidence band on all parameters. An interesting feature of this approach is that it uses weaker design conditions than many other approaches but requires the data analyst to supply explicit bounds on restricted eigenvalues. Gautier & Tsybakov (2011) and Chernozhukov et al. (2013) employed similar ideas while also working with various generalizations of restricted eigenvalues. van de Geer & Nickl (2013) constructed confidence ellipsoids for the entire parameter vector using sample splitting ideas. Somewhat related to this literature are the results of Belloni et al. (2013d), who used the orthogonal estimating equations framework with infinite-dimensional nuisance parameters and constructed a simultaneous confidence rectangle for many target parameters in which the number of target parameters could be much larger than the sample size. They relied on the high-dimensional central limit theorems and bootstrap results established by Chernozhukov et al. (2013).

Most of the aforementioned results rely on (approximate) sparsity and related sparsity-based estimators. Some examples of the use of alternative regularization schemes are available in the many instrument literature in econometrics. For example, Chamberlain & Imbens (2004) used a shrinkage estimator resulting from the use of a Gaussian random coefficients structure over first-stage coefficients, and Okui (2011) employed ridge regression for estimating the first-stage regression in a framework in which the instruments may be ordered in terms of relevance. Carrasco (2012) employed a different strategy based on directly regularizing the inverse that appears in the definition of the 2SLS estimator, allowing for a number of moment conditions that are larger than the sample size (see also Carrasco & Tchuente 2015). The theoretical development in Carrasco (2012) relies on restrictions on the covariance structure of the instruments rather than on the coefficients of the instruments. Hansen & Kozbur (2014) considered a combination of ridge

regularization and the jackknife to provide a procedure that is valid, allowing for the number of instruments to be greater than the sample size under weak restrictions on the covariance structure of the instruments and the first-stage coefficients. In all cases, the orthogonality condition holds, allowing root- $n$ -consistent and asymptotically normal estimation of the main parameter  $\alpha$ .

Many other interesting procedures beyond those mentioned in this review have been developed for estimating high-dimensional models (see, e.g., Hastie et al. 2009 for a textbook review). Developing new techniques for estimation in high-dimensional settings is also still an active area of research, so the list of methods available to researchers continues to expand. The use of these procedures and the impact of their use on inference about low-dimensional target parameters of interest are interesting research directions to explore. It seems likely that many of these procedures will provide sufficiently high-quality estimates that they may be used for estimating the high-dimensional nuisance parameter  $\eta$  in the present setting.

## APPENDIX A: THE LASSO AND POST-LASSO ESTIMATORS IN THE LINEAR MODEL

Suppose we have data  $\{y_i, x_i\}$  for individuals  $i = 1, \dots, n$ , where  $x_i$  is a  $p$ -vector of predictor variables, and  $y_i$  is an outcome of interest. Suppose that we are interested in a linear prediction model for  $y_i$ ,  $y_i = x_i' \eta + \varepsilon_i$ , and define the usual least-squares criterion function:

$$\hat{Q}(\eta) := \frac{1}{n} \sum_{i=1}^n (y_i - x_i' \eta)^2.$$

The Lasso estimator is defined as a solution of the following optimization program:

$$\hat{\eta}_L \in \arg \min_{\eta \in \mathbb{R}^p} \hat{Q}(\eta) + \frac{\lambda}{n} \sum_{j=1}^p |\psi_j \eta_j|, \tag{56}$$

where  $\lambda$  is the penalty level, and  $\{\psi_j\}_{j=1}^p$  are covariate specific penalty loadings. The covariate specific penalty loadings are used to accommodate data that may be non-Gaussian, heteroscedastic, and/or dependent and also help ensure basic equivariance of coefficient estimates to rescaling of the covariates.

The post-Lasso estimator is defined as the ordinary least-squares regression applied to the model  $\hat{I}$  selected by Lasso:<sup>6</sup>

$$\hat{I} = \text{support}(\hat{\eta}_L) = \left\{ j \in \{1, \dots, p\} : |\hat{\eta}_{Lj}| > 0 \right\}.$$

The post-Lasso estimator  $\hat{\eta}_{PL}$  is then

$$\hat{\eta}_{PL} \in \arg \min \left\{ \hat{Q}(\eta) : \eta \in \mathbb{R}^p \text{ such that } \eta_j = 0 \text{ for all } j \notin \hat{I} \right\}. \tag{57}$$

In other words, this estimator is OLS using only the regressors whose coefficients were estimated to be nonzero by Lasso.

Lasso and post-Lasso are motivated by the desire to predict the target function well without overfitting. The Lasso estimator is a computationally attractive alternative to some other classic

<sup>6</sup>We note that we can also allow the set  $\hat{I}$  to contain additional variables not selected by Lasso, but we do not consider that here.

approaches, such as model selection based on information criteria, because it minimizes a convex function. Moreover, under suitable conditions, the Lasso estimator achieves near-optimal rates in estimating the regression function  $x'\eta$ . However, Lasso does suffer from the drawback that the regularization by the  $\ell_1$ -norm employed in Equation 56 naturally shrinks all estimated coefficients toward zero, causing a potentially significant shrinkage bias. The post-Lasso estimator is meant to remove some of this shrinkage bias and achieves the same rate of convergence as Lasso under sensible conditions.

Practical implementation of the Lasso requires setting the penalty parameter and loadings. Verifying good properties of the Lasso typically relies on having these parameters set so that the penalty dominates the score in the sense that

$$\frac{\psi_j \lambda}{n} \geq \max_{i \leq p} 2c \left| \frac{1}{n} \sum_{i=1}^n x_{j,i} \varepsilon_i \right| \text{ or, equivalently, } \frac{\lambda}{\sqrt{n}} \geq \max_{i \leq p} 2c \left| \frac{\frac{1}{\sqrt{n}} \sum_{i=1}^n x_{j,i} \varepsilon_i}{\psi_j} \right|$$

for some  $c > 1$  with high probability. Heuristically, we would have the term inside the absolute values behaving approximately like a standard normal random variable if we set  $\psi_j = \text{Var} \left[ \frac{1}{\sqrt{n}} \sum_{i=1}^n x_{j,i} \varepsilon_i \right]$ . We could then get the desired domination by setting  $\lambda / (2c\sqrt{n})$  large enough to dominate the maximum of  $p$  standard normal random variables with high probability, for example, by setting  $\lambda = 2c\sqrt{n} \Phi^{-1}(1 - 0.1/[2p \log(n)])$ , where  $\Phi^{-1}(\cdot)$  denotes the inverse of the standard normal cumulative distribution function. Verifying that this heuristic argument holds with large  $p$  and data that may not be i.i.d. Gaussian requires careful and delicate arguments, as by, for example, Belloni et al. (2012), who covered heteroscedastic non-Gaussian data, or Belloni et al. (2014b), who covered panel data with within-individual dependence. The choice of the penalty parameter  $\lambda$  can also be refined, as done by Belloni et al. (2011). Finally, feasible implementation requires that  $\psi_j$  be estimated, which can be done through the iterative procedures suggested by Belloni et al. (2012) or Belloni et al. (2014b).

## APPENDIX B: PROOFS

### B.1. Proof of Proposition 2

Consider any sequence  $\{P_n\}$  in  $\{P_n\}$ .

**Step 1 ( $r_n$  rate).** Here we show that  $\|\hat{\alpha} - \alpha_0\| \leq r_n$  wp  $\rightarrow 1$ . We have by the identifiability condition, in particular the assumption  $\text{mineig}(\Gamma_1' \Gamma_1) \geq c$ , that

$$P_n \left( \|\hat{\alpha} - \alpha_0\| > r_n \right) \leq P_n \left( \|\mathbf{M}(\hat{\alpha}, \eta_0)\| \geq \iota(r_n) \right), \quad \iota(r_n) := 2^{-1} \left( \{\sqrt{c} r_n\} \wedge c \right).$$

Hence, it suffices to show that wp  $\rightarrow 1$ ,  $\|\mathbf{M}(\hat{\alpha}, \eta_0)\| < \iota(r_n)$ . By the triangle inequality, we obtain

$$\begin{aligned} I_1 &= \|\mathbf{M}(\hat{\alpha}, \eta_0) - \mathbf{M}(\hat{\alpha}, \hat{\eta})\|, \\ \|\mathbf{M}(\hat{\alpha}, \eta_0)\| &\leq I_1 + I_2 + I_3, \quad I_2 = \|\mathbf{M}(\hat{\alpha}, \hat{\eta}) - \hat{\mathbf{M}}(\hat{\alpha}, \hat{\eta})\|, \\ I_3 &= \|\hat{\mathbf{M}}(\hat{\alpha}, \hat{\eta})\|. \end{aligned}$$

By the assumption in Equation 12, wp  $\rightarrow 1$ , we have

$$I_1 + I_2 \leq o(1) \{r_n + I_3 + \|\mathbf{M}(\hat{\alpha}, \boldsymbol{\eta}_0)\|\}.$$

Hence, we obtain

$$\|\mathbf{M}(\hat{\alpha}, \boldsymbol{\eta}_0)\| \left(1 - o(1)\right) \leq o(1)(r_n + I_3) + I_3.$$

By construction of the estimator, we have

$$I_3 \leq o\left(n^{-1/2}\right) + \inf_{\alpha \in \mathcal{A}} \|\hat{\mathbf{M}}(\alpha, \hat{\boldsymbol{\eta}})\| \lesssim_{P_n} n^{-1/2},$$

which follows because

$$\inf_{\alpha \in \mathcal{A}} \|\hat{\mathbf{M}}(\alpha, \hat{\boldsymbol{\eta}})\| \leq \|\hat{\mathbf{M}}(\bar{\alpha}, \hat{\boldsymbol{\eta}})\| \lesssim_{P_n} n^{-1/2}, \quad (58)$$

where  $\bar{\alpha}$  is the one-step estimator defined in Step 3, as shown in Equation 59. Hence,  $\text{wp} \rightarrow 1$

$$\|\mathbf{M}(\hat{\alpha}, \boldsymbol{\eta}_0)\| \leq o(r_n) < \iota(r_n),$$

where to obtain the last inequality we have used the assumption  $\text{mineig}(\Gamma_1' \Gamma_1) \geq c$ .

**Step 2 ( $n^{-1/2}$  rate).** Here we show that  $\|\hat{\alpha} - \alpha_0\| \lesssim_{P_n} n^{-1/2}$ . By the condition in Equation 14 and the triangle inequality,  $\text{wp} \rightarrow 1$ , we find that

$$\|\mathbf{M}(\hat{\alpha}, \boldsymbol{\eta}_0)\| \geq \|\Gamma_1(\hat{\alpha} - \alpha_0)\| - o(1)\|\hat{\alpha} - \alpha_0\| \geq (\sqrt{c} - o(1))\|\hat{\alpha} - \alpha_0\| \geq \sqrt{c}/2\|\hat{\alpha} - \alpha_0\|.$$

Therefore, it suffices to show that  $\|\mathbf{M}(\hat{\alpha}, \boldsymbol{\eta}_0)\| \lesssim_{P_n} n^{-1/2}$ . We have that

$$\begin{aligned} II_1 &= \|\mathbf{M}(\hat{\alpha}, \boldsymbol{\eta}_0) - \mathbf{M}(\hat{\alpha}, \hat{\boldsymbol{\eta}})\|, \\ \|\mathbf{M}(\hat{\alpha}, \boldsymbol{\eta}_0)\| &\leq II_1 + II_2 + II_3, \quad II_2 = \|\mathbf{M}(\hat{\alpha}, \hat{\boldsymbol{\eta}}) - \hat{\mathbf{M}}(\hat{\alpha}, \hat{\boldsymbol{\eta}}) - \hat{\mathbf{M}}(\alpha_0, \boldsymbol{\eta}_0)\|, \\ II_3 &= \|\hat{\mathbf{M}}(\hat{\alpha}, \hat{\boldsymbol{\eta}})\| + \|\hat{\mathbf{M}}(\alpha_0, \boldsymbol{\eta}_0)\|. \end{aligned}$$

Then, by the orthogonality  $\partial_{\boldsymbol{\eta}'} \mathbf{M}(\alpha_0, \boldsymbol{\eta}_0) = 0$  and the condition in Equation 14,  $\text{wp} \rightarrow 1$ , we find that

$$\begin{aligned} II_1 &\leq \left\| \mathbf{M}(\hat{\alpha}, \hat{\boldsymbol{\eta}}) - \mathbf{M}(\hat{\alpha}, \boldsymbol{\eta}_0) - \partial_{\boldsymbol{\eta}'} \mathbf{M}(\hat{\alpha}, \boldsymbol{\eta}_0)[\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}_0] \right\| + \left\| \partial_{\boldsymbol{\eta}'} \mathbf{M}(\hat{\alpha}, \boldsymbol{\eta}_0)[\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}_0] \right\| \\ &\leq o(1)n^{-1/2} + o(1)\|\hat{\alpha} - \alpha_0\| \\ &\leq o(1)n^{-1/2} + o(1)(2/\sqrt{c})\|\mathbf{M}(\hat{\alpha}, \boldsymbol{\eta}_0)\|. \end{aligned}$$

Then, by the condition in Equation 13 and by  $I_3 \lesssim_{P_n} n^{-1/2}$ , we obtain

$$\begin{aligned} II_2 &\leq o(1) \left\{ n^{-1/2} + \|\hat{\mathbf{M}}(\hat{\alpha}, \hat{\boldsymbol{\eta}})\| + \|\mathbf{M}(\hat{\alpha}, \boldsymbol{\eta}_0)\| \right\} \\ &\lesssim_{P_n} o(1) \left\{ n^{-1/2} + n^{-1/2} + \|\mathbf{M}(\hat{\alpha}, \boldsymbol{\eta}_0)\| \right\}. \end{aligned}$$



Because  $III_3 \lesssim_{P_n} n^{-1/2}$  by Equation 58 and  $\|\hat{M}(\alpha_0, \eta_0)\| \lesssim_{P_n} n^{-1/2}$ , it follows that  $\text{wp} \rightarrow 1$ ,  $(1 - o(1))\|\hat{M}(\hat{\alpha}, \hat{\eta})\| \lesssim_{P_n} n^{-1/2}$ .

**Step 3 (linearization).** Define the linearization map  $\alpha \mapsto \hat{L}(\alpha)$  by  $\hat{L}(\alpha) := \hat{M}(\alpha_0, \eta_0) + \Gamma_1(\alpha - \alpha_0)$ . Then we obtain

$$\begin{aligned} III_1 &= \|\hat{M}(\hat{\alpha}, \hat{\eta}) - \hat{M}(\hat{\alpha}, \eta_0)\|, \\ \|\hat{M}(\hat{\alpha}, \hat{\eta}) - \hat{L}(\hat{\alpha})\| &\leq III_1 + III_2 + III_3, \quad III_2 = \|\hat{M}(\hat{\alpha}, \eta_0) - \Gamma_1(\hat{\alpha} - \alpha_0)\|, \\ III_3 &= \|\hat{M}(\hat{\alpha}, \hat{\eta}) - \hat{M}(\hat{\alpha}, \eta_0) - \hat{M}(\alpha_0, \eta_0)\|. \end{aligned}$$

Then, using the assumptions in Equations 13 and 14, conclude that

$$\begin{aligned} III_1 &\leq \left\| \hat{M}(\hat{\alpha}, \hat{\eta}) - \hat{M}(\hat{\alpha}, \eta_0) - \partial_{\eta'} \hat{M}(\hat{\alpha}, \eta_0)[\hat{\eta} - \eta_0] \right\| + \left\| \partial_{\eta'} \hat{M}(\hat{\alpha}, \eta_0)[\hat{\eta} - \eta_0] \right\| \\ &\leq o(1)n^{-1/2} + o(1)\|\hat{\alpha} - \alpha_0\|, \\ III_2 &\leq o(1)\|\hat{\alpha} - \alpha_0\|, \\ III_3 &\leq o(1)\left(n^{-1/2} + \|\hat{M}(\hat{\alpha}, \hat{\eta})\| + \|\hat{M}(\hat{\alpha}, \eta_0)\|\right) \\ &\leq o(1)\left(n^{-1/2} + n^{-1/2} + III_2 + \|\Gamma_1(\hat{\alpha} - \alpha_0)\|\right). \end{aligned}$$

Conclude that  $\text{wp} \rightarrow 1$ , as  $\|\Gamma_1' \Gamma_1\| \lesssim 1$  by the assumption in Equation 11,

$$\|\hat{M}(\hat{\alpha}, \hat{\eta}) - \hat{L}(\hat{\alpha})\| \lesssim_{P_n} o(1)\left(n^{-1/2} + \|\hat{\alpha} - \alpha_0\|\right) = o\left(n^{-1/2}\right).$$

Also consider the minimizer of the map  $\alpha \mapsto \|\hat{L}(\alpha)\|$ , namely,

$$\bar{\alpha} = \alpha_0 - \left(\Gamma_1' \Gamma_1\right)^{-1} \Gamma_1' \hat{M}(\alpha_0, \eta_0),$$

which obeys  $\|\sqrt{n}(\bar{\alpha} - \alpha_0)\| \lesssim_{P_n} n^{-1/2}$  under the conditions of the proposition. We can repeat the argument above to conclude that  $\text{wp} \rightarrow 1$ ,  $\|\hat{M}(\bar{\alpha}, \hat{\eta}) - \hat{L}(\bar{\alpha})\| \lesssim_{P_n} o(n^{-1/2})$ . This implies, as  $\|\hat{L}(\bar{\alpha})\| \lesssim_{P_n} n^{-1/2}$ , that

$$\|\hat{M}(\bar{\alpha}, \hat{\eta})\| \lesssim_{P_n} n^{-1/2}. \quad (59)$$

This also implies that  $\|\hat{L}(\hat{\alpha})\| = \|\hat{L}(\bar{\alpha})\| + o_{P_n}(n^{-1/2})$ , as  $\|\hat{L}(\bar{\alpha})\| \leq \|\hat{L}(\hat{\alpha})\|$  and

$$\|\hat{L}(\hat{\alpha})\| - o_{P_n}(n^{-1/2}) \leq \|\hat{M}(\hat{\alpha}, \hat{\eta})\| \leq \|\hat{M}(\bar{\alpha}, \hat{\eta})\| + o(n^{-1/2}) = \|\hat{L}(\bar{\alpha})\| + o_{P_n}(n^{-1/2}).$$

The former assertion implies that  $\|\hat{L}(\hat{\alpha})\|^2 = \|\hat{L}(\bar{\alpha})\|^2 + o_{P_n}(n^{-1})$ , so that

$$\|\hat{L}(\hat{\alpha})\|^2 - \|\hat{L}(\bar{\alpha})\|^2 = \left\| \Gamma_1(\hat{\alpha} - \bar{\alpha}) \right\|^2 = o_{P_n}(n^{-1}),$$

from which we can conclude that  $\sqrt{n}\|\hat{\alpha} - \bar{\alpha}\| \rightarrow_{P_n} 0$ .

**Step 4 (conclusion).** Given the conclusion of the previous step, the remaining claims are standard and follow from the continuous mapping theorem and Lemma 8.

### B.2. Proof of Proposition 3

We have  $\text{wp} \rightarrow 1$  that, for some constants  $0 < u < l < 1$ ,  $l\|x\| \leq \|Ax\| \leq u\|x\|$  and  $l\|x\| \leq \|\hat{A}x\| \leq u\|x\|$ . Hence, we obtain

$$\begin{aligned} & \sup_{\alpha \in \mathcal{A}} \frac{\|\hat{A}\hat{M}^o(\alpha, \hat{\eta}) - AM^o(\alpha, \hat{\eta})\| + \|AM^o(\alpha, \hat{\eta}) - AM^o(\alpha, \eta_0)\|}{r_n + \|\hat{A}\hat{M}^o(\alpha, \hat{\eta})\| + \|AM^o(\alpha, \eta_0)\|} \\ & \leq \sup_{\alpha \in \mathcal{A}} \frac{u}{l} \frac{\|\hat{M}^o(\alpha, \hat{\eta}) - M^o(\alpha, \hat{\eta})\| + \|M^o(\alpha, \hat{\eta}) - M^o(\alpha, \eta_0)\|}{(r_n/l) + \|\hat{M}^o(\alpha, \hat{\eta})\| + \|M^o(\alpha, \eta_0)\|} \\ & + \sup_{\alpha \in \mathcal{A}} \frac{\|\hat{A} - A\| \|\hat{M}^o(\alpha, \hat{\eta})\|}{r_n + l \|\hat{M}^o(\alpha, \hat{\eta})\|} \lesssim_{P_n} o(1) + \|\hat{A} - A\|/l \rightarrow_{P_n} 0. \end{aligned}$$

The proof that the rest of the conditions hold is analogous and is therefore omitted.

### B.3. Proof of Proposition 4

**Step 1.** We define the feasible and infeasible one-step estimators

$$\begin{aligned} \tilde{\alpha} &= \tilde{\alpha} - \hat{F}\hat{M}(\tilde{\alpha}, \hat{\eta}), \quad \hat{F} = \left(\hat{\Gamma}'_1 \hat{\Gamma}_1\right)^{-1} \hat{\Gamma}'_1, \\ \bar{\alpha} &= \alpha_0 - F\hat{M}(\alpha_0, \eta_0), \quad F = \left(\Gamma'_1 \Gamma_1\right)^{-1} \Gamma'_1. \end{aligned}$$

We deduce by Equations 11 and 20 that

$$\|\hat{F}\| \lesssim_{P_n} 1, \quad \|\hat{F}\Gamma_1 - I\| \lesssim_{P_n} r_n, \quad \|\hat{F} - F\| \lesssim_{P_n} r_n.$$

**Step 2.** By Step 1 and by the condition in Equation 21, we have that

$$\begin{aligned} D &= \|\hat{F}\hat{M}(\tilde{\alpha}, \hat{\eta}) - \hat{F}\hat{M}(\alpha_0, \eta_0) - \hat{F}\Gamma_1(\tilde{\alpha} - \alpha_0)\| \\ &\leq \|\hat{F}\| \|\hat{M}(\tilde{\alpha}, \hat{\eta}) - \hat{M}(\alpha_0, \eta_0) - \Gamma_1(\tilde{\alpha} - \alpha_0)\| \\ &\lesssim_{P_n} \|\hat{M}(\tilde{\alpha}, \hat{\eta}) - M(\tilde{\alpha}, \hat{\eta}) - \hat{M}(\alpha_0, \eta_0)\| + D_1 \lesssim_{P_n} o(n^{-1/2}) + D_1, \end{aligned}$$

where  $D_1 := \|M(\tilde{\alpha}, \hat{\eta}) - \Gamma_1(\tilde{\alpha} - \alpha_0)\|$ .

Moreover, we have  $D_1 \leq IV_1 + IV_2 + IV_3$ , where  $\text{wp} \rightarrow 1$  by the condition in Equation 21 and  $r_n^2 = o(n^{-1/2})$

$$\begin{aligned}
IV_1 &:= \left\| \mathbf{M}(\tilde{\alpha}, \eta_0) - \Gamma_1(\tilde{\alpha} - \alpha_0) \right\| \lesssim \left\| \tilde{\alpha} - \alpha_0 \right\|^2 \lesssim r_n^2 = o\left(n^{-1/2}\right), \\
IV_2 &:= \left\| \mathbf{M}(\tilde{\alpha}, \hat{\eta}) - \mathbf{M}(\tilde{\alpha}, \eta_0) - \partial_{\eta'} \mathbf{M}(\tilde{\alpha}, \eta_0)[\hat{\eta} - \eta_0] \right\| \lesssim o\left(n^{-1/2}\right), \\
IV_3 &:= \left\| \partial_{\eta'} \mathbf{M}(\tilde{\alpha}, \eta_0)[\hat{\eta} - \eta_0] \right\| \lesssim o\left(n^{-1/2}\right).
\end{aligned}$$

Conclude that  $n^{1/2}D \rightarrow_{P_n} 0$ .

**Step 3.** We have by the triangle inequality and Steps 1 and 2 that

$$\begin{aligned}
\sqrt{n} \left\| \hat{\alpha} - \bar{\alpha} \right\| &\leq \sqrt{n} \left\| (I - \hat{F}\Gamma_1)(\tilde{\alpha} - \alpha_0) \right\| + \sqrt{n} \left\| (\hat{F} - F)\hat{\mathbf{M}}(\alpha_0, \eta_0) \right\| + \sqrt{n}D \\
&\leq \sqrt{n} \left\| (I - \hat{F}\Gamma_1) \right\| \left\| \tilde{\alpha} - \alpha_0 \right\| + \left\| \hat{F} - F \right\| \left\| \sqrt{n}\hat{\mathbf{M}}(\alpha_0, \eta_0) \right\| + \sqrt{n}D \\
&\lesssim_{P_n} \sqrt{nr_n^2} + o(1) = o(1).
\end{aligned}$$

Thus, we have  $\sqrt{n} \left\| \hat{\alpha} - \bar{\alpha} \right\| \rightarrow_{P_n} 0$ , and  $\sqrt{n} \left\| \hat{\alpha} - \hat{\alpha} \right\| \rightarrow_{P_n} 0$  follows from the triangle inequality and the fact that  $\sqrt{n} \left\| \hat{\alpha} - \bar{\alpha} \right\| \rightarrow_{P_n} 0$ .

### B.4. Proof of Lemma 2

The conditions of Proposition 1 are clearly satisfied, and thus the conclusions of Proposition 1 immediately follow. We also have that, for  $\hat{\Gamma}_1 = \hat{\Gamma}_1(\hat{\eta})$ ,

$$\begin{aligned}
\sqrt{n}(\hat{\alpha} - \alpha_0) &= -\hat{F}\sqrt{n}\hat{\mathbf{M}}(\alpha_0, \hat{\eta}), \quad \hat{F} = \left(\hat{\Gamma}'_1\hat{\Gamma}_1\right)^{-1}\hat{\Gamma}_1, \\
\sqrt{n}(\bar{\alpha} - \alpha_0) &:= -F\sqrt{n}\hat{\mathbf{M}}(\alpha_0, \eta_0), \quad F = \left(\Gamma'_1\Gamma_1\right)^{-1}\Gamma_1.
\end{aligned}$$

We deduce by Equations 11 and 33 that  $\left\| \hat{F} \right\| \lesssim_{P_n} 1$  and  $\left\| \hat{F} - F \right\| \rightarrow_{P_n} 0$ . Hence, we have by the triangle and Hölder inequalities and the condition in Equation 33 that

$$\sqrt{n} \left\| \hat{\alpha} - \bar{\alpha} \right\| \leq \left\| \hat{F} \right\| \left\| \sqrt{n} \left\| \hat{\mathbf{M}}(\alpha_0, \hat{\eta}) - \hat{\mathbf{M}}(\alpha_0, \eta_0) \right\| \right\| + \left\| \hat{F} - F \right\| \left\| \sqrt{n} \left\| \hat{\mathbf{M}}(\alpha_0, \eta_0) \right\| \right\| \rightarrow_{P_n} 0.$$

The conclusions regarding the uniform validity of inference using  $\hat{\alpha}$  of the form stated in the conclusions of Proposition 2 follow from the conclusions regarding the uniform validity of inference using  $\bar{\alpha}$ , which follow from the continuous mapping theorem, Lemma 8, and the assumed stability conditions in Equation 11. This establishes the second claim of the lemma. Verification of the conditions of Proposition 2 is omitted.

### B.5. Proof of Lemmas 3 and 4

The proof of Lemma 3 is given in the main text. As in the proof of Lemma 3, we can expand:

$$\sqrt{n}(\hat{\mathbf{M}}_j(\alpha_0, \hat{\eta}) - \hat{\mathbf{M}}_j(\alpha_0, \eta_0)) = T_{1,j} + T_{2,j} + T_{3,j}, \tag{60}$$

where the terms  $(T_{l,j})_{l=1}^3$  are as defined in the main text. We can further bound  $T_{3,j}$  as follows:

$$\begin{aligned}
T_{3,j} &\leq T_{3,j}^m + T_{4,j}, \\
T_{3,j}^m &:= \sqrt{n} \left| (\hat{\eta} - \eta_0^m)' \partial_{\eta} \partial_{\eta'} \hat{\mathbf{M}}_j(\alpha_0) (\hat{\eta} - \eta_0^m) \right|, \\
T_{4,j} &:= \sqrt{n} \left| \eta_0^r' \partial_{\eta} \partial_{\eta'} \hat{\mathbf{M}}_j(\alpha_0) \eta_0^r \right|.
\end{aligned} \tag{61}$$

Then  $T_{1,j} = 0$  by orthogonality,  $T_{2,j} \rightarrow_{P_n} 0$  as in the proof of Lemma 3. Given that  $s^2 \log(pn)^2/n \rightarrow 0$ ,  $T_{3,j}^m$  vanishes in probability because, by Hölder's inequality and for sufficiently large  $n$ ,

$$T_{3,j}^m \leq \sqrt{n} \bar{T}_{3,j} \left\| \hat{\eta} - \eta_0^m \right\|^2 \lesssim_{P_n} \sqrt{ns} \log(pn)/n \rightarrow_{P_n} 0.$$

Also, if  $s^2 \log(pn)^2/n \rightarrow 0$ ,  $T_{4,j}$  vanishes in probability because

$$T_{4,j} \leq \sqrt{n} \left\| \partial_{\eta} \partial_{\eta'} \hat{M}_j(\alpha_0) \right\|_{\text{pw}(\eta_0^r)} \left\| \eta_0^r \right\|^2 \lesssim_{P_n} \sqrt{ns} \log(pn)/n \rightarrow_{P_n} 0,$$

where the inequalities follow by Hölder's inequality and Equation 43. The conclusion follows from Equation 60.

### B.6. Proof of Lemma 5

For  $m = 1, \dots, k$  and  $l = 1, \dots, d$ , we can bound each element  $\hat{\Gamma}_{1,ml}(\eta)$  of matrix  $\hat{\Gamma}_1(\eta)$  as follows:

$$\begin{aligned} T_{1,ml} &:= \left| \partial_{\eta} \Gamma_{1,ml}(\eta_0)' (\hat{\eta} - \eta_0) \right|, \\ T_{2,ml} &:= \left| \left( \partial_{\eta} \hat{\Gamma}_{1,ml}(\eta_0) - \partial_{\eta} \Gamma_{1,ml}(\eta_0) \right)' (\hat{\eta} - \eta_0) \right|, \\ T_{3,ml} &:= \left| (\hat{\eta} - \eta_0^m)' \partial_{\eta} \partial_{\eta'} \hat{\Gamma}_{1,ml}(\hat{\eta} - \eta_0^m) \right|, \\ T_{4,ml} &:= \left| \eta_0^r' \partial_{\eta} \partial_{\eta'} \hat{\Gamma}_{1,ml} \eta_0^r \right|. \end{aligned}$$

$$\left| \hat{\Gamma}_{1,ml}(\hat{\eta}) - \hat{\Gamma}_{1,ml}(\eta_0) \right| \leq \sum_{k=1}^4 T_{k,ml},$$

Under the conditions in Equations 44 and 45, we have that  $\text{wp} \rightarrow 1$ ,

$$\begin{aligned} T_{1,ml} &\leq \left\| \partial_{\eta} \Gamma_{1,ml}(\eta_0) \right\|_{\infty} \left\| \hat{\eta} - \eta_0 \right\|_1 \lesssim_{P_n} \sqrt{s^2 \log(pn)/n} \rightarrow 0, \\ T_{2,ml} &\leq \left\| \partial_{\eta} \hat{\Gamma}_{1,ml}(\eta_0) - \partial_{\eta} \Gamma_{1,ml}(\eta_0) \right\|_{\infty} \left\| \hat{\eta} - \eta_0 \right\|_1 \lesssim_{P_n} \sqrt{s^2 \log(pn)/n} \rightarrow 0, \\ T_{3,ml} &\leq \left\| \partial_{\eta} \partial_{\eta'} \hat{\Gamma}_{1,ml} \right\|_{\text{sp}(\ell_n s)} \left\| \hat{\eta} - \eta_0^m \right\|^2 \lesssim_{P_n} s \log(pn)/n \rightarrow 0, \\ T_{4,ml} &\leq \left\| \partial_{\eta} \partial_{\eta'} \hat{\Gamma}_{1,ml} \right\|_{\text{pw}(\eta_0^r)} \left\| \eta_0^r \right\|^2 \lesssim_{P_n} s \log(pn)/n \rightarrow 0. \end{aligned}$$

The claim follows from the assumed growth conditions, as  $d$  and  $k$  are bounded.

### APPENDIX C: KEY TOOLS

Let  $\Phi$  and  $\Phi^{-1}$  denote the distribution and quantile function of  $\mathcal{N}(0, 1)$ . Note that, in particular,  $\Phi^{-1}(1-a) \leq \sqrt{2 \log(a)}$  for all  $a \in (0, 1/2)$ .

**Lemma 6 (moderate deviation inequality for the maximum of a vector):** Suppose that  $\mathcal{S}_j := \sum_{i=1}^n U_{ij} / \sqrt{\sum_{i=1}^n U_{ij}^2}$ , where  $U_{ij}$  are independent random variables across  $i$  with mean zero and finite third-order moments. Then, we obtain

$$\mathbb{P}\left(\max_{1 \leq j \leq p} |S_j| > \Phi^{-1}(1 - \gamma/2p)\right) \leq \gamma \left(1 + \frac{A}{\ell_n^3}\right),$$

where  $A$  is an absolute constant, provided for  $\ell_n > 0$ ,

$$0 \leq \Phi^{-1}(1 - \gamma/(2p)) \leq \frac{n^{1/6}}{\ell_n} \min_{1 \leq j \leq p} M_j^2 - 1, \quad M_j := \frac{\left(\frac{1}{n} \sum_{i=1}^n \mathbb{E}[U_{ij}^2]\right)^{1/2}}{\left(\frac{1}{n} \sum_{i=1}^n \mathbb{E}[|U_{ij}|^3]\right)^{1/3}}.$$

This result is essentially due to Jing et al. (2003). The proof of this result, given by Belloni et al. (2012), follows from a simple combination of union bounds with their result.

**Lemma 7 (laws of large numbers for large matrices in sparse norms):** Let  $s_n, p_n, k_n$ , and  $\ell_n$  be sequences of positive constants such that  $\ell_n \rightarrow \infty$  but  $\ell_n/\log n \rightarrow 0$  and  $c_1$  and  $c_2$  be fixed positive constants. Let  $(x_i)_{i=1}^n$  be i.i.d. vectors such that  $\|\mathbb{E}[x_i x_i']\|_{\text{sp}(s_n, \log n)} \leq c_1$ , and either one of the following holds: (a)  $x_i$  is a sub-Gaussian random vector with  $\sup_{\|u\| \leq 1} \|x_i' u\|_{\psi_2, \mathbb{P}} \leq c_2$ , where  $\|\cdot\|_{\psi_2, \mathbb{P}}$  denotes the  $\psi_2$ -Orlicz norm of a random variable, and  $s_n(\log n)(\log(p_n \vee n))/n \rightarrow 0$ , or (b)  $\|x_i\|_\infty \leq k_n$  almost surely and  $k_n^2 s_n (\log^4 n) \log(p_n \vee n)/n \rightarrow 0$ . Then there is  $o(1)$  term such that with probability  $1 - o(1)$ ,  $\|\mathbb{E}_n[x_i x_i'] - \mathbb{E}[x_i x_i']\|_{\text{sp}(s_n, \ell_n)} \leq o(1)$ ,  $\|\mathbb{E}_n[x_i x_i']\|_{\text{sp}(s_n, \ell_n)} \leq c_1 + o(1)$ .

Under (a), the result follows from theorem 3.2 of Rudelson & Zhou (2011), and under (b), the result follows from Rudelson & Vershynin (2008), as shown in the supplemental material of Belloni & Chernozhukov (2013).

**Lemma 8 (useful implications of the central limit theorem in  $\mathbb{R}^m$ ):** Consider a sequence of random vectors  $Z_n$  in  $\mathbb{R}^m$  such that  $Z_n \rightsquigarrow Z = \mathcal{N}(0, I_m)$ . The elements of the sequence and the limit variable need not be defined on the same probability space. Then we obtain

$$\lim_{n \rightarrow \infty} \sup_{R \in \mathcal{R}} |\mathbb{P}(Z_n \in R) - \mathbb{P}(Z \in R)| = 0,$$

where  $\mathcal{R}$  is the collection of all convex sets in  $\mathbb{R}^m$ .

**Proof:** Let  $R$  denote a generic convex set in  $\mathbb{R}^m$ . Let  $R^\epsilon = \{z \in \mathbb{R}^m : d(z, R) \leq \epsilon\}$  and  $R^{-\epsilon} = \{z \in \mathbb{R}^m : B(z, \epsilon) \subset R\}$ , where  $d$  is the Euclidean distance and  $B(z, \epsilon) = \{y \in \mathbb{R}^m : d(y, z) \leq \epsilon\}$ . The set  $R^\epsilon$  may be empty. By theorem 11.3.3 in Dudley (2002), we find that  $\epsilon_n := \rho(Z_n, Z) \rightarrow 0$ , where  $\rho$  is the Prohorov metric. The definition of the metric implies that  $\mathbb{P}(Z_n \in R) \leq \mathbb{P}(Z \in R^{\epsilon_n}) + \epsilon_n$ . By the reverse isoperimetric inequality (Chen & Fang 2011, proposition 2.5), we obtain  $|\mathbb{P}(Z \in R^{\epsilon_n}) - \mathbb{P}(Z \in R)| \leq m^{1/2} \epsilon_n$ . Hence,  $\mathbb{P}(Z_n \in R) \leq \mathbb{P}(Z \in R) + \epsilon_n(1 + m^{1/2})$ . Furthermore, for any convex set  $R$ , we find that  $(R^{-\epsilon_n})^{\epsilon_n} \subset R$  (interpreting the expansion of an empty set as an empty set). Hence, for any convex  $R$ , we have  $\mathbb{P}(Z \in R^{-\epsilon_n}) \leq \mathbb{P}(Z_n \in R) + \epsilon_n$  by definition of Prohorov's metric. By the reverse isoperimetric inequality, we obtain  $|\mathbb{P}(Z \in R^{-\epsilon_n}) - \mathbb{P}(Z \in R)| \leq m^{1/2} \epsilon_n$ . Conclude that  $\mathbb{P}(Z_n \in R) \geq \mathbb{P}(Z \in R) - \epsilon_n(1 + m^{1/2})$ .

## DISCLOSURE STATEMENT

The authors are not aware of any affiliations, memberships, funding, or financial holdings that might be perceived as affecting the objectivity of this review.

## ACKNOWLEDGMENTS

We thank Denis Chetverikov, Mert Demirer, Anna Mikusheva, seminar participants and the discussant Susan Athey at the AEA Session on Machine Learning in Economics and Econometrics, participants of the CEME conference on Non-Standard Problems in Econometrics, participants of the Berlin Statistics Seminar, and the students from MIT's 14.387 Applied Econometrics Class for useful comments.

## LITERATURE CITED

- Belloni A, Chen D, Chernozhukov V, Hansen C. 2012. Sparse models and methods for optimal instruments with an application to eminent domain. *Econometrica* 80:2369–429
- Belloni A, Chernozhukov V. 2011. High-dimensional sparse econometric models: an introduction. In *Inverse Problems and High-Dimensional Estimation: Stats in the Chateau Summer School, August 31–September 2, 2009*, ed. P Alquier, E Gautier, G Stoltz, pp. 121–56. New York: Springer
- Belloni A, Chernozhukov V. 2013. Least squares after model selection in high-dimensional sparse models. *Bernoulli* 19:521–47
- Belloni A, Chernozhukov V, Fernández-Val I, Hansen C. 2013a. Program evaluation with high-dimensional data. arXiv:1311.2645 [math.ST]
- Belloni A, Chernozhukov V, Hansen C. 2010. LASSO methods for Gaussian instrumental variables models. arXiv:1012.1297 [stat.ME]
- Belloni A, Chernozhukov V, Hansen C. 2013b. Inference for high-dimensional sparse econometric models. In *Advances in Economics and Econometrics: 10th World Congress*, Vol. 3: *Econometrics*, ed. D Acemoglu, M Arellano, E Dekel, pp. 245–95. Cambridge, UK: Cambridge Univ. Press
- Belloni A, Chernozhukov V, Hansen C. 2014a. Inference on treatment effects after selection amongst high-dimensional controls. *Rev. Econ. Stud.* 81:608–50
- Belloni A, Chernozhukov V, Hansen C, Kozbur D. 2014b. Inference in high dimensional panel models with an application to gun control. arXiv:1411.6507 [stat.ME]
- Belloni A, Chernozhukov V, Kato K. 2013c. Robust inference in approximately sparse quantile regression models (with an application to malnutrition). arXiv:1312.7186 [math.ST]
- Belloni A, Chernozhukov V, Kato K. 2013d. Uniform post selection inference for LAD regression models and other Z-estimation problems. arXiv:1304.0282 [math.ST]
- Belloni A, Chernozhukov V, Wang L. 2011. Square-root-LASSO: pivotal recovery of sparse signals via conic programming. *Biometrika* 98:791–806
- Belloni A, Chernozhukov V, Wei Y. 2013e. Honest confidence regions for logistic regression with a large number of controls. arXiv:1304.3969 [stat.ME]
- Berk R, Brown L, Buja A, Zhang K, Zhao L. 2013. Valid post-selection inference. *Ann. Stat.* 41:802–37
- Berry S, Levinsohn J, Pakes A. 1995. Automobile prices in market equilibrium. *Econometrica* 63:841–90
- Bickel PJ. 1982. On adaptive estimation. *Ann. Statist.* 10:647–71
- Bickel PJ, Ritov Y, Tsybakov AB. 2009. Simultaneous analysis of Lasso and Dantzig selector. *Ann. Stat.* 37:1705–32
- Bühlmann P, van de Geer S. 2011. *Statistics for High-Dimensional Data: Methods, Theory and Applications*. New York: Springer
- Candès E, Tao T. 2007. The Dantzig selector: statistical estimation when  $p$  is much larger than  $n$ . *Ann. Stat.* 35:2313–51
- Carrasco M. 2012. A regularization approach to the many instruments problem. *J. Econom.* 170:383–98

- Carrasco M, Tchuente G. 2015. Regularized LIML with many instruments. *J. Econom.* 186:427–42
- Chamberlain G. 1987. Asymptotic efficiency in estimation with conditional moment restrictions. *J. Econom.* 34:305–34
- Chamberlain G, Imbens G. 2004. Random effects estimators with many instrumental variables. *Econometrica* 72:295–306
- Chao JC, Swanson NR, Hausman JA, Newey WK, Woutersen T. 2012. Asymptotic distribution of JIVE in a heteroskedastic IV regression with many instruments. *Econom. Theory* 28:42–86
- Chen LHY, Fang X. 2011. Multivariate normal approximation by Stein’s method: the concentration inequality approach. arXiv:1111.4073 [math.PR]
- Chen X, Linton O, Keilegom IV. 2003. Estimation of semiparametric models when the criterion function is not smooth. *Econometrica* 71:1591–608
- Chernozhukov V, Chetverikov D, Kato K. 2013. Gaussian approximations and multiplier bootstrap for maxima of sums of high-dimensional random vectors. *Ann. Stat.* 41:2786–819
- Chernozhukov V, Liu H, Lu J, Ning Y. 2014. *Statistical inference in high-dimensional sparse models using generalized method of moments*. Unpublished manuscript, Mass. Inst. Technol., Cambridge, MA, Princeton Univ., Princeton, NJ
- Dudley RM. 2002. *Real Analysis and Probability*. Cambridge, UK: Cambridge Univ. Press
- Fan J, Li R. 2001. Variable selection via nonconcave penalized likelihood and its oracle properties. *J. Am. Stat. Assoc.* 96:1348–60
- Fan J, Lv J. 2010. A selective overview of variable selection in high dimensional feature space. *Stat. Sin.* 20:101–48
- Farrell MH. 2014. Robust inference on average treatment effects with possibly more covariates than observations. arXiv:1309.4686 [math.ST]
- Fithian W, Sun D, Taylor J. 2014. Optimal inference after model selection. arXiv:1410.2597v1 [math.ST]
- Frank IE, Friedman JH. 1993. A statistical view of some chemometrics regression tools. *Technometrics* 35:109–35
- Gautier E, Tsybakov AB. 2011. High-dimensional instrumental variables regression and confidence sets. arXiv:1105.2454v4 [math.ST]
- Gillen BJ, Shum M, Moon HR. 2014. Demand estimation with high-dimensional product characteristics. *Adv. Econom.* 34:301–23
- G’Sell MG, Taylor J, Tibshirani R. 2013. Adaptive testing for the graphical lasso. arXiv:1307.4765 [math.ST]
- Hansen C, Kozbur D. 2014. Instrumental variables estimation with many weak instruments using regularized JIVE. *J. Econom.* 182:290–308
- Hastie T, Tibshirani R, Friedman J. 2009. *Elements of Statistical Learning: Data Mining, Inference, and Prediction*. New York: Springer
- Huber PJ. 1964. The behavior of maximum likelihood estimates under nonstandard conditions. *Proc. 5th Berkeley Symp.*, ed. J Neyman, pp. 221–23. Berkeley: Univ. Calif. Press
- Javanmard A, Montanari A. 2014. Confidence intervals and hypothesis testing for high-dimensional regression. arXiv:1306.3171v2 [stat.ME]
- Jing B-Y, Shao Q-M, Wang Q. 2003. Self-normalized Cramer-type large deviations for independent random variables. *Ann. Probab.* 31:2167–215
- Kozbur D. 2014. *Inference in nonparametric models with a high-dimensional component*. Work. Pap., ETH Zürich
- Lee JD, Sun DL, Sun Y, Taylor JE. 2013. Exact post-selection inference, with application to the lasso. arXiv:1311.6238 [math.ST]
- Lee JD, Taylor JE. 2014. Exact post model selection inference for marginal screening. arXiv:1402.5596 [stat.ME]
- Leeb H, Pötscher BM. 2008a. Recent developments in model selection and related areas. *Econom. Theory* 24:319–22
- Leeb H, Pötscher BM. 2008b. Sparse estimators and the oracle property, or the return of Hodges’ estimator. *J. Econom.* 142:201–11
- Lockhart R, Taylor JE, Tibshirani RJ, Tibshirani R. 2014. A significance test for the lasso. *Ann. Stat.* 42:413–68

- Loftus JR, Taylor JE. 2014. A significance test for forward stepwise model selection. arXiv:1405.3920 [stat.ME]
- Meinshausen N, Yu B. 2009. Lasso-type recovery of sparse representations for high-dimensional data. *Ann. Stat.* 37:2246–70
- Neyman J. 1959. Optimal asymptotic tests of composite statistical hypotheses. In *Probability and Statistics: The Harald Cramer Volume*, ed. U Grenander, pp. 213–34. New York: Wiley
- Neyman J. 1979.  $C(\alpha)$  tests and their use. *Sankhya* 41:1–21
- Ning Y, Liu H. 2014. SPARC: optimal estimation and asymptotic inference under semiparametric sparsity. arXiv:1412.2295 [stat.ML]
- Okui R. 2011. Instrumental variable estimation in the presence of many moment conditions. *J. Econom.* 165:70–86
- Pakes A, Pollard D. 1989. Simulation and asymptotics of optimization estimators. *Econometrica* 57:1027–57
- Robins JM, Rotnitzky A. 1995. Semiparametric efficiency in multivariate regression models with missing data. *J. Am. Stat. Assoc.* 90:122–29
- Rudelson M, Vershynin R. 2008. On sparse reconstruction from Fourier and Gaussian measurements. *Commun. Pure Appl. Math.* 61:1025–45
- Rudelson M, Zhou S. 2011. Reconstruction from anisotropic random measurements. arXiv:1106.1151 [math.ST]
- Taylor J, Lockhart R, Tibshirani R, Tibshirani R. 2014. Exact post-selection inference for forward stepwise and least angle regression. arXiv:1401.3889 [stat.ME]
- Tibshirani R. 1996. Regression shrinkage and selection via the lasso. *J. R. Stat. Soc. B* 58:267–88
- van de Geer S, Bühlmann P, Ritov Y, Dezeure R. 2014. On asymptotically optimal confidence regions and tests for high-dimensional models. *Ann. Stat.* 42:1166–202
- van de Geer S, Nickl R. 2013. Confidence sets in sparse regression. *Ann. Stat.* 41:2852–76
- van der Vaart AW. 1998. *Asymptotic Statistics*. Cambridge, UK: Cambridge Univ. Press
- Voorman A, Shojaie A, Witten D. 2014. Inference in high dimensions with the penalized score test. arXiv:1401.2678 [stat.ME]
- Yang Z, Ning Y, Liu H. 2014. On semiparametric exponential family graphical models. arXiv:1412.8697 [stat.ML]
- Zhang C-H, Zhang SS. 2014. Confidence intervals for low dimensional parameters in high dimensional linear models. *J. R. Stat. Soc. B* 76:217–42





# Contents

Knowledge-Based Hierarchies: Using Organizations to Understand the Economy <i>Luis Garicano and Esteban Rossi-Hansberg</i> . . . . .	1
Beyond Ricardo: Assignment Models in International Trade <i>Arnaud Costinot and Jonathan Vogel</i> . . . . .	31
The Roots of Gender Inequality in Developing Countries <i>Seema Jayachandran</i> . . . . .	63
Reconciling Micro and Macro Labor Supply Elasticities: A Structural Perspective <i>Michael Keane and Richard Rogerson</i> . . . . .	89
International Trade, Multinational Activity, and Corporate Finance <i>C. Fritz Foley and Kalina Manova</i> . . . . .	119
Policy Implications of Dynamic Public Finance <i>Mikhail Golosov and Aleh Tsyvinski</i> . . . . .	147
Media and Politics <i>David Strömberg</i> . . . . .	173
Forecasting in Nonstationary Environments: What Works and What Doesn't in Reduced-Form and Structural Models <i>Raffaella Giacomini and Barbara Rossi</i> . . . . .	207
Political Decentralization <i>Dilip Mookherjee</i> . . . . .	231
Household Debt: Facts, Puzzles, Theories, and Policies <i>Jonathan Zinman</i> . . . . .	251
Making Progress on Foreign Aid <i>Nancy Qian</i> . . . . .	277

Credit, Financial Stability, and the Macroeconomy <i>Alan M. Taylor</i> . . . . .	309
Job Creation, Job Destruction, and Productivity Growth: The Role of Young Businesses <i>John Haltiwanger</i> . . . . .	341
The Evolution of Social Norms <i>H. Peyton Young</i> . . . . .	359
Crime and Economic Incentives <i>Mirko Draca and Stephen Machin</i> . . . . .	389
Entrepreneurship and Financial Frictions: A Macroeconomic Perspective <i>Francisco J. Buera, Joseph P. Kaboski, and Yongseok Shin</i> . . . . .	409
The US Electricity Industry After 20 Years of Restructuring <i>Severin Borenstein and James Bushnell</i> . . . . .	437
Methods of Identification in Social Networks <i>Bryan S. Graham</i> . . . . .	465
Affirmative Action in Undergraduate Education <i>Peter Arcidiacono, Michael Lovenheim, and Maria Zhu</i> . . . . .	487
Is College a Worthwhile Investment? <i>Lisa Barrow and Ofer Malamud</i> . . . . .	519
The Schumpeterian Growth Paradigm <i>Philippe Aghion, Ufuk Akcigit, and Peter Howitt</i> . . . . .	557
Climate and Conflict <i>Marshall Burke, Solomon M. Hsiang, and Edward Miguel</i> . . . . .	577
The Gains from Market Integration <i>Dave Donaldson</i> . . . . .	619
Valid Post-Selection and Post-Regularization Inference: An Elementary, General Approach <i>Victor Chernozhukov, Christian Hansen, and Martin Spindler</i> . . . . .	649

## Indexes

Cumulative Index of Contributing Authors, Volumes 3–7	689
Cumulative Index of Article Titles, Volumes 3–7	692

## Errata

An online log of corrections to *Annual Review of Economics* articles may be found at <http://www.annualreviews.org/errata/economics>