NBER WORKING PAPER SERIES

A SCORE BASED APPROACH TO WILD BOOTSTRAP INFERENCE

Patrick M. Kline Andres Santos

Working Paper 16127 http://www.nber.org/papers/w16127

NATIONAL BUREAU OF ECONOMIC RESEARCH 1050 Massachusetts Avenue Cambridge, MA 02138 June 2010

We thank Justin McCrary, Graham Elliott and Michael Jansson for useful comments. The views expressed herein are those of the authors and do not necessarily reflect the views of the National Bureau of Economic Research.

 \bigcirc 2010 by Patrick M. Kline and Andres Santos. All rights reserved. Short sections of text, not to exceed two paragraphs, may be quoted without explicit permission provided that full credit, including \bigcirc notice, is given to the source.

A Score Based Approach to Wild Bootstrap Inference Patrick M. Kline and Andres Santos NBER Working Paper No. 16127 June 2010 JEL No. C01.C12

ABSTRACT

We propose a generalization of the wild bootstrap of Wu (1986) and Liu (1988) based upon perturbing the scores of M-estimators. This "score bootstrap" procedure avoids recomputing the estimator in each bootstrap iteration, making it significantly easier to implement than the conventional bootstrap, particularly in complex nonlinear models. Despite this computational advantage, for studentized statistics, the score bootstrap distribution is equivalent to that of the conventional wild bootstrap up to order $O(n^{-1})$. We establish the consistency of score bootstrapped Wald and Lagrange Multiplier type tests and tests of moment restrictions under weak regularity conditions and in the presence of potential misspecification. However, a higher order analysis reveals that certain forms of misspecification may undermine the ability of the score bootstrap to provide an Edgeworth refinement, a result which carries over to the Wild bootstrap as well. To gauge the empirical relevance of these results we conduct an extensive series of Monte Carlo experiments comparing the performance of several different bootstrap procedures in settings with clustered data. We find that variants of our proposed score based bootstrap substantially outperform analytical cluster robust methods and in some cases outperform the block bootstrap even in the presence of misspecification.

Patrick M. Kline Department of Economics UC, Berkeley 508-1 Evans Hall #3880 Berkeley, CA 94720 and NBER pkline@econ.berkeley.edu

Andres Santos Department of Economics 9500 Gilman Drive La Jolla, CA 92093-0508 a2santos@ucsd.edu

1 Introduction

Applied researchers often wish to analyze samples with few independent observations. The bootstrap of Efron (1979) has become a standard tool for conducting inference in such settings. Among the numerous variants of the original bootstrap, the so-called "wild" bootstrap of Wu (1986) and Liu (1988) has been shown to yield dramatic improvements in the ability to control the size of Wald tests of OLS regression coefficients in small samples (Mammen (1993), Horowitz (1997, 2001), Cameron, Gelbach, and Miller (2008)).

Originally proposed as an alternative to the residual bootstrap of Freedman (1981), the wild bootstrap has often been interpreted as a procedure that resamples residuals in a manner that captures any heteroscedasticity in the underlying errors. Perhaps for this reason, the applications and extensions of the wild bootstrap have largely been limited to linear models where residuals are straightforward to obtain; see for example Hardle and Mammen (1993) for nonparametric regression, You and Chen (2006) for the partially linear regression, Davidson and MacKinnon (2008) for IV regression and Cavaliere and Taylor (2008) for unit root inference.

We propose a new bootstrap procedure (the "score" bootstrap) which perturbs the fitted scores of an M-estimator conditional on a fixed Hessian. In the linear model, this procedure is numerically equivalent to the conventional wild bootstrap for unstudentized statistics and higher order equivalent for studentized ones. However, in contrast to the wild bootstrap, our approach is easily adapted to estimators without conventional residuals and avoids recomputing the estimator in each bootstrap iteration. As a result, the score bootstrap possesses an important advantage over existing bootstraps in settings where the model is computationally expensive to estimate or poorly behaved in a subset of the bootstrap draws. For example, computational problems often arise in small samples even in simple probit or logit models where, for some bootstrap draws, the estimator cannot be computed.

We provide results establishing the consistency of the score bootstrap for a broad class of test statistics under weak regularity conditions and in the presence of potential misspecification. Our framework is shown to encompass Wald and Lagrange Multiplier (LM) tests as well as tests of moment restrictions. We then examine the higher order properties of the proposed bootstrap in the specific context of a linear model. There we derive the conditions under which a score bootstrapped Wald statistic yields an Edgeworth refinement and find that the presence of a refinement is sensitive to certain forms of misspecification. This conclusion holds true for the original wild bootstrap as well and hence may be of independent interest. To assess the empirical relevance of these theoretical results, we conduct an extensive series of Monte Carlo experiments comparing the performance of several different bootstrap procedures in settings with clustered data. We find that variants of our proposed score based bootstrap substantially outperform analytical cluster robust methods and in some cases outperform the traditional block bootstrap even in the presence of misspecification. In line with our theoretical results, we find negligible differences in the performance of the score and wild bootstraps despite their large difference in computational cost.

The remainder of the paper is structured as follows: Section 2 reviews the wild bootstrap, while Section 3 introduces the score bootstrap and establishes its higher order equivalence. In Section 4 we establish the consistency of the score bootstrap under weak regularity conditions and illustrate its applicability to a variety of settings. Section 5 examines the conditions under which the score bootstrap yields an Edgeworth refinement in the linear model. Our simulation study is contained in Section 6, while Section 7 briefly concludes. All proofs are contained in the Appendix.

2 Wild Bootstrap Review

We begin by reviewing the wild bootstrap and the reasons for its consistency in the context of a linear model. A careful examination of the arguments justifying its validity provides us with the intuition necessary for developing the score bootstrap and its extension to M-estimation problems.

While there are multiple approaches to implementing the wild bootstrap, for expository purposes we focus on the original methodology developed in Liu (1988). Suppose $\{Y_i, X_i\}_{i=1}^n$ is an i.i.d. sequence of random variables, with $Y_i \in \mathbf{R}, X_i \in \mathbf{R}^m$ and satisfying the linear relationship:

$$Y_i = X_i'\beta_0 + \epsilon_i . (1)$$

Letting $\hat{\beta}$ denote the OLS estimate of β_0 and $e_i \equiv (Y_i - X'_i \hat{\beta})$ the implied residual, the wild bootstrap generates new residuals of the form $\epsilon_i^* \equiv W_i e_i$ for some randomly generated i.i.d. sequence $\{W_i\}_{i=1}^n$ that is independent of $\{Y_i, X_i\}_{i=1}^n$ and satisfies $E[W_i] = 0$ and $E[W_i^2] = 1$. Under these conditions,

$$E[\epsilon_i^*|\{Y_i, X_i\}_{i=1}^n] = 0 \qquad E[(\epsilon_i^*)^2|\{Y_i, X_i\}_{i=1}^n] = e_i^2 , \qquad (2)$$

and hence ϵ_i^* is mean independent of $\{Y_i, X_i\}_{i=1}^n$ and in addition captures the pattern of heteroscedasticity found in the original sample. This property, originally noted in Wu (1986), enables the wild bootstrap to remain consistent even in the presence of heteroscedasticity or model misspecification.¹

¹We refer to misspecification in model (1) as $E[\epsilon_i | X_i] \neq 0$ but $E[\epsilon_i X_i] = 0$.

The wild bootstrap resampling scheme consists of generating dependent variables $\{Y_i^*\}_{i=1}^n$ by

$$Y_i^* \equiv X_i'\hat{\beta} + \epsilon_i^* \tag{3}$$

and then conducting OLS on the sample $\{Y_i^*, X_i\}_{i=1}^n$ in order to obtain a bootstrap estimate $\hat{\beta}^*$. The distribution of $\sqrt{n}(\hat{\beta}^* - \hat{\beta})$ conditional on $\{Y_i, X_i\}_{i=1}^n$ (but not on $\{W_i\}_{i=1}^n$) is then used as an estimate of the unknown distribution of $\sqrt{n}(\hat{\beta} - \beta_0)$. Since the former distribution can be computed through simulation, the wild bootstrap provides a simple way to obtain critical values for inference.

We review why the wild bootstrap is consistent by drawing from arguments in Mammen (1993). First, observe that standard OLS algebra and the relationships in (1) and (3) imply that:

$$\sqrt{n}(\hat{\beta} - \beta_0) = H_n^{-1} \frac{1}{\sqrt{n}} \sum_{i=1}^n X_i \epsilon_i \qquad \sqrt{n}(\hat{\beta}^* - \hat{\beta}) = H_n^{-1} \frac{1}{\sqrt{n}} \sum_{i=1}^n X_i \epsilon_i^* , \qquad (4)$$

where $H_n \equiv n^{-1} \sum_i X_i X'_i$. Since both the true and bootstrap scores are properly centered, both expressions in (4) can be expected to converge to a normal limit. Therefore, consistency of the wild bootstrap hinges on whether this limit is the same or, equivalently, whether the asymptotic variances agree. However, as $E[W_i^2] = 1$ and $\{W_i\}_{i=1}^n$ is independent of $\{Y_i, X_i\}_{i=1}^n$, we obtain:

$$E[(\frac{1}{\sqrt{n}}\sum_{i=1}^{n}X_{i}\epsilon_{i})(\frac{1}{\sqrt{n}}\sum_{i=1}^{n}X_{i}\epsilon_{i})'] = E[X_{i}X_{i}'\epsilon_{i}^{2}]$$

$$\tag{5}$$

$$E[(\frac{1}{\sqrt{n}}\sum_{i=1}^{n}X_{i}\epsilon_{i}^{*})(\frac{1}{\sqrt{n}}\sum_{i=1}^{n}X_{i}\epsilon_{i}^{*})'|\{Y_{i},X_{i}\}_{i=1}^{n}] = \frac{1}{n}\sum_{i=1}^{n}X_{i}X_{i}'e_{i}^{2}$$
(6)

and hence the second moments indeed agree asymptotically by standard arguments. As a result, $\sqrt{n}(\hat{\beta} - \beta_0)$ and $\sqrt{n}(\hat{\beta}^* - \hat{\beta})$ converge in distribution to the same normal limit and the consistency of the wild bootstrap is immediate.

While the ability of the wild bootstrap to asymptotically match the first two moments of the scores provides the basis for establishing its validity, it does not elucidate why it often performs better than a normal approximation. Improvements occur when the bootstrap is able to additionally match higher moments of the statistics. If, for example, $E[W_i^3] = 1$, then the third moments match asymptotically and the wild bootstrap provides a refinement over the normal approximation to a studentized statistic by providing a skewness correction (Liu (1988)). Alternatively, the Rademacher distribution,² which satisfies $E[W_i] = E[W_i^3] = 0$ and $E[W_i^2] = E[W_i^4] = 1$, is able to match the first four moments for symmetric distributions and can in such cases provide an additional refinement (Liu (1988); Davidson and Flachaire (2008)).

²A Rademacher random variable puts probability one half on both one and negative one.

3 The Score Bootstrap

The wild bootstrap resampling scheme is often interpreted as a means of generating a set of bootstrap residuals mimicking the heteroscedastic nature of the true errors. An alternative interpretation is that it creates a set of bootstrap scores mimicking the heteroscedastic nature of the true scores. In this section, we develop the implications of this observation, which provides the basis for our proposed procedure.

The relationship between the wild bootstrap and the scores is transparent from the discussion of its consistency in Section 2. Since $\epsilon_i^* = e_i W_i$, we learn from (4) that the wild bootstrap may be interpreted as a perturbation of the scores $(X_i(Y_i - X'_i\beta))$ evaluated at the estimated parameter value $(\hat{\beta})$ that leaves the Hessian $(\sum_i X_i X'_i)$ unchanged.³ More precisely, a numerically equivalent way to implement the wild bootstrap would be to employ the following algorithm:

STEP 1: Obtain the OLS estimate $\hat{\beta}$ and generate the fitted scores $\{X_i(Y_i - X'_i \hat{\beta})\}_{i=1}^n$.

STEP 2: Using random weights $\{W_i\}_{i=1}^n$ independent of $\{Y_i, X_i\}_{i=1}^n$ and satisfying $E[W_i] = 0$ and $E[W_i^2] = 1$, perturb the original fitted scores to obtain a new set of scores $\{X_i(Y_i - X'_i\hat{\beta})W_i\}_{i=1}^n$.

STEP 3: Multiply the perturbed scores by the original Hessian to obtain $H_n^{-1}n^{-\frac{1}{2}}\sum_i (Y_i - X_i\hat{\beta})X_iW_i$ and use its distribution conditional on $\{Y_i, X_i\}_{i=1}^n$ as an estimate of the distribution of $\sqrt{n}(\hat{\beta} - \beta_0)$.

Unlike the residual based view of the wild bootstrap, the score interpretation is easily generalized to more complex nonlinear models. One may simply perturb the fitted scores of such a model while keeping the Hessian unchanged and, provided the perturbations satisfy $E[W_i] = 0$ and $E[W_i^2] = 1$, the first two moments of the perturbed and true scores will match asymptotically. Under the appropriate regularity conditions, this moment equivalence will suffice for establishing the consistency of the proposed bootstrap. For obvious reasons, we term this approach a "score bootstrap."

3.1 Higher Order Equivalence

In the linear model, the wild and score bootstrap statistics for $\sqrt{n}(\hat{\beta}-\beta_0)$ are numerically equivalent. However, in most instances the statistic of interest is studentized, since only in this context is a refinement over an analytical approximation available (Liu (1988), Horowitz (2001)).

Because a bootstrap estimate $\hat{\beta}^*$ is not computed under the score bootstrap, it is not practical to studentize by following the exact analogue of the full sample computation. In fact, in accord

³This is in contrast to the weighted bootstrap which perturbs the score and the Hessian (Ma and Kosorok (2005)).

with the perturbed score interpretation, it is more natural to simply employ the sample variance of the perturbed scores for studentization. For this reason, we define the bootstrap statistics:

$$T_{w,n}^* \equiv (H_n^{-1} \Sigma_n^* (\hat{\beta}^*) H_n^{-1})^{-\frac{1}{2}} \sqrt{n} (\hat{\beta}^* - \hat{\beta}) \qquad T_{s,n}^* \equiv (H_n^{-1} \Sigma_n^* (\hat{\beta}) H_n^{-1})^{-\frac{1}{2}} H_n^{-1} \frac{1}{\sqrt{n}} \sum_{i=1}^n X_i \epsilon_i^* , \quad (7)$$

where $\Sigma_n^*(\beta) \equiv n^{-1} \sum_i X_i X_i' (Y_i^* - X_i'\beta)^2$ and $T_{w,n}^*$ and $T_{s,n}^*$ are the studentized wild and score bootstrap statistics respectively. It is important to note that in the computation of $T_{s,n}^*$, the full sample estimator $\hat{\beta}$ is used in obtaining the standard errors, and hence calculation of $\hat{\beta}^*$ remains unnecessary. As a result, the score bootstrap is computationally simpler to implement than the wild bootstrap which requires use of the bootstrap estimate $\hat{\beta}^*$.

While for the statistics in (4) the wild and score bootstraps are numerically equivalent, such a relationship fails to hold for the studentized versions. An important concern then is whether the refinement of the wild bootstrap over a normal approximation (Liu (1988)) is lost due to this difference. Somewhat surprisingly, the answer is negative. The score bootstrap not only remains consistent despite not recomputing the estimator but can in addition be expected to obtain a refinement over an analytical approximation in precisely the same instances as the wild bootstrap. Such a result is the consequence of the wild and score bootstrap statistics being asymptotically equivalent up to a higher order than the refinement over the normal approximation.

In order to establish the higher order equivalence of $T_{n,s}^*$ and $T_{n,w}^*$, we impose the following:

Assumption 3.1. (i) $\{Y_i, X_i\}_{i=1}^n$ are *i.i.d.* $E[X_i\epsilon_i] = 0$, $E[X_iX'_i] = I$ and $E[X_iX'_i\epsilon^2_i]$ is full rank; (ii) The moments $E[||X_i||^9]$, $E[\epsilon^9_i]$ and $E[||X_i||^9\epsilon^9_i]$ are finite; (iii) $\beta_0 \in \Theta$, where $\Theta \subset \mathbf{R}^m$ is compact; (iv) $\{W_i\}_{i=1}^n$ are *i.i.d.*, independent of $\{Y_i, X_i\}_{i=1}^n$ with $E[W_i] = 0$, $E[W_i^2] = 1$ and $E[W_i^9] < \infty$.

The requirement that $E[X_iX'_i] = I$ of Assumption 3.1(i) constitutes a notationally convenient normalization for establishing the higher order expansions in Section 5. The existence of high order moments, imposed in Assumption 3.1(ii) and 3.1(iv), is necessary for the computation of the first three moments of t-statistics, but not for establishing the consistency of the procedure. The bias, variance and skewness of the t-statistics are in turn needed to explore whether a bootstrap procedure yields a refinement over an analytical approximation. In Section 4, where we establish bootstrap consistency results for a general class of M-estimators, we only require existence of the first two moments. As a special case those results imply the consistency of the score bootstrap in the context of the linear model under weaker requirements than those in Assumption 3.1.

Let P^* and E^* denote probability and expectation conditional on $\{Y_i, X_i\}_{i=1}^n$ (but not $\{W_i\}_{i=1}^n$). Under Assumption 3.1 we can then establish the higher order equivalence of $T^*_{w,n}$ and $T^*_{s,n}$. **Lemma 3.1.** Under Assumption 3.1, $T_{w,n}^* = T_{s,n}^* + O_{p^*}(n^{-1})$ almost surely.

If the conditions for an Edgeworth expansion of the bootstrap statistics $T_{w,n}^*$ and $T_{s,n}^*$ are satisfied, then Lemma 3.1 implies that they can be expected to disagree only in terms of order n^{-1} or smaller; see Chapter 2.7 in Hall (1992) for such arguments. Therefore, in settings where the wild bootstrap obtains the traditional Edgeworth refinement of order $n^{-\frac{1}{2}}$ over a normal approximation, the score bootstrap should as well. The higher order equivalence of $T_{w,n}^*$ and $T_{s,n}^*$ is at first glance unexpected since the score bootstrap appears to violate the usual plug-in approach of the standard bootstrap. However, this only introduces a smaller order error due to the residuals $\{\epsilon_i^*\}_{i=1}^n$ being mean independent of $\{X_i\}_{i=1}^n$ under the bootstrap distribution. Importantly, the higher order equivalence would fail to hold if the residuals $\{\epsilon_i^*\}$ were sampled in a manner under which they were merely uncorrelated with $\{X_i\}_{i=1}^n$ under the bootstrap distribution.

4 Inference

We turn now to establishing the validity of a score bootstrap procedure for estimating the critical values of a large class of tests. Building on our earlier discussion we consider test statistics based upon the parametric scores of M-estimators, using perturbations of those scores to estimate their sampling distribution. Since this approach does not depend upon resampling of residuals, we do not distinguish between dependent and exogenous variables and instead consider a random vector $Z_i \in \mathcal{Z} \subseteq \mathbf{R}^m$ which may contain both.

We focus on test statistics G_n that are quadratic forms of a vector valued statistic T_n :

$$G_n \equiv T'_n T_n \ . \tag{8}$$

Under the null hypothesis, the underlying statistic T_n is required to be asymptotically pivotal and allow for a linear expansion. More precisely, we require that under the null hypothesis:

$$T_n = (A_n(\theta_0)\Sigma_n(\theta_0)A_n(\theta_0)')^{-\frac{1}{2}}S_n(\theta_0) + o_p(1) \qquad S_n(\theta) \equiv A_n(\theta)\frac{1}{\sqrt{n}}\sum_{i=1}^n s(Z_i,\theta) , \qquad (9)$$

where $A_n(\theta)$ is a $r \times k$ matrix, $s(z, \theta)$ is a $k \times 1$ vector, $\Sigma_n(\theta)$ is the sample covariance matrix of $s(Z_i, \theta)$ and θ_0 is an unknown parameter vector. Under appropriate regularity conditions, T_n is therefore asymptotically normally distributed with identity covariance matrix and hence G_n is asymptotically Chi-squared distributed with degrees of freedom equal to the dimension of T_n . Though we only consider asymptotically pivotal statistics, our results readily extend to unstudentized ones as well. The bootstrap statistics employed to estimate the distributions of G_n and T_n are given by:

$$G_{n}^{*} \equiv T_{n}^{*'}T_{n}^{*} \qquad T_{n}^{*} \equiv (A_{n}(\hat{\theta})\Sigma_{n}^{*}(\hat{\theta})A_{n}(\hat{\theta})')^{-\frac{1}{2}}S_{n}^{*}(\hat{\theta}) \qquad S_{n}^{*}(\theta) \equiv A_{n}(\theta)\frac{1}{\sqrt{n}}\sum_{i=1}^{n}s(Z_{i},\theta)W_{i} \quad (10)$$

where $\Sigma_n^*(\theta)$ is the sample covariance matrix of $s(Z_i, \theta)W_i$ and $\hat{\theta}$ is a consistent estimator for θ_0 . As discussed in the previous section, implementation of the score bootstrap only requires calculation of the full sample estimator $\hat{\theta}$ and no additional optimization is needed in each bootstrap iteration.

4.1 Bootstrap Consistency

We establish the consistency of the bootstrap under the following set of assumptions:

Assumption 4.1. (i) $\hat{\theta} \xrightarrow{p} \theta_0$ with $\hat{\theta}, \theta_0 \in \Theta \subset \mathbf{R}^p$ and Θ a compact set; (ii) The limit point θ_0 satisfies $E[s(Z_i, \theta_0)s(Z_i, \theta_0)'] < \infty$ and the matrix $A(\theta_0)E[s(Z_i, \theta_0)s(Z_i, \theta_0)']A(\theta_0)'$ is invertible.

Assumption 4.2. (i) Under the null hypothesis T_n satisfies (9) and θ_0 is such that $E[s(Z_i, \theta_0)] = 0$; (ii) Under the alternative hypothesis $G_n \xrightarrow{p} \infty$.

Assumption 4.3. (i) $\{Z_i\}_{i=1}^n$ is i.i.d. (ii) $\sup_{\theta \in \Theta} ||A_n(\theta) - A(\theta)||_F = o_p(1)$ with $A(\theta)$ continuous.

Assumption 4.4. (i) $\{W_i\}_{i=1}^n$ is an i.i.d. sample, independent of $\{Z_i\}_{i=1}^n$ satisfying $E[W_i] = 0$ and $E[W_i^2] = 1$; (ii) $s(z,\theta)$ is continuously differentiable in $\theta \in \Theta$ and $\sup_{\theta \in conv(\Theta)} \|\nabla s(z,\theta)\|_F \leq F(z)$ for some function F(z) with $E[F^2(Z_i)] < \infty$.

In Assumption 4.1 we require $\hat{\theta}$ to converge in probability to some parameter vector $\theta_0 \in \Theta$ whose value may depend upon the distribution of Z_i . The compactness of the parameter space Θ is employed to verify the perturbed scores form a Donsker class. This restriction may be relaxed at the expense of a more complicated argument that exploits the consistency of $\hat{\theta}$ for a local analysis. Though in the notation we suppress such dependence, it is important to note that θ_0 may take different values under the null and alternative hypotheses. In Assumptions 4.3(ii) and 4.4(ii), $\|\cdot\|_F$ denotes the Frobenius norm. Assumptions 4.2 and 4.3, in turn enable us to establish the asymptotic behavior of G_n under the null and alternative hypotheses; see Lemma 7.3 in the Appendix. Assumption 4.4(i) imposes the only requirements on the random weights $\{W_i\}_{i=1}^n$, which are the same conditions imposed for inference on the linear model in previous wild bootstrap studies. Assumption 4.4(ii) allows us to establish that the empirical process induced by functions of the form $ws(z, \theta)$ is asymptotically tight. Differentiability is not necessary for this end, but we opt to impose it due to its ease of verification and wide applicability.⁴

⁴For non-differentiable settings, the relevant condition is that $\mathcal{F} \equiv \{ws(z,\theta) : \theta \in \Theta\}$ be a Donsker class.

Assumptions 4.1-4.4 are sufficient for establishing the consistency of the proposed score bootstrap procedure under the null hypothesis.

Theorem 4.1. Let F_n and F_n^* be the cdfs of G_n and of G_n^* conditional on $\{Z_i\}_{i=1}^n$ and suppose that Assumptions 4.1, 4.2, 4.3 and 4.4 hold. If the null hypothesis is true, it then follows that:

$$\sup_{t \in \mathbf{R}} |F_n(t) - F_n^*(t)| = o_p(1)$$

Theorem 4.1 justifies the use of quantiles from the distribution of G_n^* conditional on $\{Z_i\}_{i=1}^n$ as critical values. In order to control the size of the test at level α , we may employ:

$$\hat{c}_{1-\alpha} \equiv \inf\{t : P(G_n^* \ge t \mid \{Z_i\}_{i=1}^n) \ge 1 - \alpha\} .$$
(11)

While difficult to compute analytically, $\hat{c}_{1-\alpha}$ may easily be calculated via simulation. Employing a random number generator, B samples $\{\{W_{i1}\}_{i=1}^{n}, \ldots, \{W_{iB}\}_{i=1}^{n}\}$ may be created independently of the data and used to construct B statistics $\{G_{n1}^{*}, \ldots, G_{nB}^{*}\}$. Provided B is sufficiently large, the empirical $1 - \alpha$ quantile of $\{G_{n1}^{*}, \ldots, G_{nB}^{*}\}$ will yield an accurate approximation to $\hat{c}_{1-\alpha}$.

While Theorem 4.1 implies that the critical value $\hat{c}_{1-\alpha}$ in conjunction with the test statistic G_n delivers size control, it does not elucidate the behavior of the test under the alternative hypothesis. As in other bootstrap procedures, the test is consistent due to the bootstrap statistic G_n^* being properly centered even under the alternative. As a result, $\hat{c}_{1-\alpha}$ converges in probability to the $1-\alpha$ quantile of a Chi-squared distribution with r degrees of freedom, while G_n diverges to infinity. Therefore, under the alternative hypothesis, G_n is larger than $\hat{c}_{1-\alpha}$ with probability tending to one and the test rejects asymptotically. We summarize these findings in the following corollary:

Corollary 4.1. Under Assumptions 4.1, 4.2, 4.3 and 4.4, it follows that under the null hypothesis:

$$\lim_{n \to \infty} P(G_n \ge \hat{c}_{1-\alpha}) = 1 - \alpha ,$$

for any $0 < \alpha < 1$. Under the same assumptions, if the alternative hypothesis is instead true, then:

$$\lim_{n \to \infty} P(G_n \ge \hat{c}_{1-\alpha}) = 1 \; .$$

4.2 Parameter Tests

A principal application of the proposed bootstrap is in obtaining critical values for hypothesis tests on parametric models. We consider a general M-estimation framework in which the parameter of interest θ_M is the unique minimizer of some non-stochastic but unknown function $Q: \Theta \to \mathbf{R}$:

$$\theta_M = \arg\min_{\theta \in \Theta} Q(\theta) \ . \tag{12}$$

We examine the classic problem of conducting inference on a function of θ_M . Specifically, for some known and differentiable mapping $c: \Theta \to \mathbf{R}^l$ with $l \leq p$, the hypothesis we study is:

$$H_0: c(\theta_M) = 0 \qquad H_1: c(\theta_M) \neq 0$$
 (13)

Standard tests for this hypothesis include the Wald and Lagrange Multiplier (LM) tests. Intuitively, the Wald test examines whether the value of the function c evaluated at an unrestricted estimator $\hat{\theta}_M$ is statistically different from zero. In contrast, the LM test instead checks whether the first order condition of an estimator $\hat{\theta}_{M,R}$ computed imposing the null hypothesis is statistically different from zero. Therefore, in the nomenclature of Assumption 4.1(i), $\hat{\theta}$ equals $\hat{\theta}_M$ for the Wald test and $\hat{\theta}_{M,R}$ for the LM test. Similarly, if $\theta_{M,R}$ denotes the minimizer of Q over Θ subject to $c(\theta) = 0$, then θ_0 equals θ_M and $\theta_{M,R}$ under the Wald and LM test respectively.

We proceed to illustrate the details of the score bootstrap in this setting for both generalized method of moments (GMM) and maximum likelihood (ML) estimators. We focus on the analytical expressions $A_n(\theta)$ and $s(z, \theta)$ take in those specific settings and provide references for primitive conditions that ensure Assumptions 4.1, 4.2, 4.3 and 4.4 hold.

4.2.1 ML Estimators

For an ML estimator, the criterion function Q and its sample analogue Q_n are of the form:

$$Q_n(\theta) \equiv \frac{1}{n} \sum_{i=1}^n q(Z_i, \theta) \qquad Q(\theta) \equiv E[q(Z_i, \theta)] , \qquad (14)$$

where $q : \mathbb{Z} \times \Theta \to \mathbf{R}$ is the log-likelihood. If q is twice differentiable in θ , then we may define the Hessian $H_n(\theta) \equiv n^{-1} \sum_i \nabla^2 q(Z_i, \theta)$. For notational convenience, it is also helpful to denote the gradient of the function c evaluated at θ by $C(\theta) \equiv \nabla c(\theta)$.

Example 4.1. (Wald) The relevant Wald statistic is the studentized quadratic form of $\sqrt{nc(\theta_M)}$, which under both the null and alternative hypothesis satisfies the asymptotic expansion:

$$\sqrt{n}(c(\hat{\theta}_M) - c(\theta_M)) = -C(\theta_M)H_n^{-1}(\theta_M)\frac{1}{\sqrt{n}}\sum_{i=1}^n \nabla q(Z_i, \theta_M) + o_p(1) .$$
(15)

Therefore, the Wald statistic fits the formulation in (9) with $A_n(\theta) = -C(\theta)H_n^{-1}(\theta)$ and $s(Z_i, \theta) = \nabla q(Z_i, \theta)$. Under the alternative hypothesis, G_n diverges to infinity since $c(\theta_M) \neq 0$. Refer to Section 3.2 in Newey and McFadden (1994) for a formal justification of these arguments.

Example 4.2. (LM) In this setting, the LM statistic is the normalized quadratic form of:

$$C(\hat{\theta}_{M,R})H_n^{-1}(\hat{\theta}_{M,R})\frac{1}{\sqrt{n}}\sum_{i=1}^n \nabla q(Z_i,\hat{\theta}_{M,R}) .$$
(16)

Moreover, under conditions stated in Chapter 12.6.2 in Wooldridge (2002), we additionally have:

$$C(\hat{\theta}_{M,R})H_n^{-1}(\hat{\theta}_{M,R})\frac{1}{\sqrt{n}}\sum_{i=1}^n \nabla q(Z_i,\hat{\theta}_{M,R}) = C(\theta_{M,R})H_n^{-1}(\theta_{M,R})\frac{1}{\sqrt{n}}\sum_{i=1}^n \nabla q(Z_i,\theta_{M,R}) + o_p(1) , \quad (17)$$

under the null hypothesis. Thus, the LM statistic also fits the general formulation in (9) with $A_n(\theta) = C(\theta)H_n^{-1}(\theta)$ and score $s(z,\theta) = \nabla q(z,\theta)$. Under the alternative, $G_n \xrightarrow{p} \infty$ provided $\theta_{M,R}$ is not a local minimizer of Q, $C(\theta_{M,R})E[\nabla^2 q(Z_i,\theta_{M,R})]$ is full rank and Assumption 4.1(ii) holds.

4.2.2 GMM Estimators

In the context of GMM estimation, the criterion function Q and its sample analogue Q_n are:

$$Q_n(\theta) \equiv \left[\frac{1}{n}\sum_{i=1}^n q(Z_i,\theta)'\right]\Omega_n\left[\frac{1}{n}\sum_{i=1}^n q(Z_i,\theta)\right] \qquad Q(\theta) \equiv E[q(Z_i,\theta)']\Omega E[q(Z_i,\theta)] , \qquad (18)$$

where $q: \mathbb{Z} \times \Theta \to \mathbf{R}^k$ is a known function and Ω_n , Ω are positive definite matrices such that $\Omega_n \xrightarrow{p} \Omega$. Assuming q is differentiable in θ , let $D_n(\theta) \equiv n^{-1} \sum_i \nabla q(Z_i, \theta)$ and $B_n(\theta) \equiv D_n(\theta)' \Omega_n D_n(\theta)$. As in the discussion of ML estimators, we also denote $C(\theta) \equiv \nabla c(\theta)$.

Example 4.3. (Wald) The Wald statistic for the hypothesis in (13) is given by the studentized quadratic form of $\sqrt{nc(\hat{\theta}_M)}$. In the present context we therefore obtain an expansion of the form:

$$\sqrt{n}(c(\hat{\theta}_M) - c(\theta_M)) = -C(\theta_M)B_n^{-1}(\theta_M)D_n(\theta_M)'\Omega_n \frac{1}{\sqrt{n}}\sum_{i=1}^n q(Z_i, \theta_M) + o_p(1) , \qquad (19)$$

which implies $A_n(\theta) = -C(\theta)B_n^{-1}(\theta)D_n(\theta)'\Omega_n$ and $s(z,\theta) = q(z,\theta)$ and Assumption 4.2(i) is satisfied provided $E[q(Z_i, \theta_M)] = 0.5$ Primitive conditions under which Assumptions 4.1-4.4 hold in this context can be found in Section 3.3 of Newey and McFadden (1994).

Example 4.4. (LM) In this setting, the LM test statistic is the studentized quadratic form of:

$$C(\hat{\theta}_{M,R})B_n^{-1}(\hat{\theta}_{M,R})D_n(\hat{\theta}_{M,R})'\Omega_n\frac{1}{\sqrt{n}}\sum_{i=1}^n q(Z_i,\hat{\theta}_{M,R}) , \qquad (20)$$

which, as shown in Section 9.1 of Newey and McFadden (1994), is asymptotically equivalent to:

$$C(\theta_{M,R})B_n^{-1}(\theta_{M,R})D_n(\theta_{M,R})'\Omega_n\frac{1}{\sqrt{n}}\sum_{i=1}^n q(Z_i,\theta_{M,R})$$
(21)

under the null hypothesis. Hence, $A_n(\theta) = C(\theta)B_n^{-1}(\theta)D_n(\theta)'\Omega_n$ and $s(z,\theta) = q(z,\theta)$.

⁵Notice this is trivially satisfied in a just identified system. The extension to overidentified models in which $E[q(Z_i, \theta_M)] \neq 0$ but $E[\nabla q(Z_i, \theta_M)'] \Omega E[q(Z_i, \theta_M)] = 0$ can be accomplished by letting $s(z, \theta)$ depend on n and setting $s_n(z, \theta) = D_n(\theta)' \Omega_n g(z, \theta)$. Though straightforward to establish, we do not pursue such an extension.

4.3 Moment Restrictions

An additional application of the bootstrap procedure we consider is for testing the hypothesis:

$$H_0: E[m(Z_i, \theta_M)] = 0 \qquad H_1: E[m(Z_i, \theta_M)] \neq 0 ,$$
 (22)

where $m : \mathcal{Z} \times \Theta \to \mathbf{R}^l$ is a known function and θ_M is the minimizer of some unknown nonstochastic $Q : \Theta \to \mathbf{R}$. Such restrictions arise, for example, in tests of proper model specification and hypotheses regarding average marginal effects in nonlinear models. As in Section 4.2, the specific nature of the bootstrap statistic is dependent on whether Q is as in (14) (ML) or as in (18) (GMM). For brevity, we focus on the former, though the extension to GMM can be readily derived following manipulations analogous to those in Example 4.3.

The Wald test statistic for the hypothesis in (22) is based on the studentized plug-in estimator:

$$\frac{1}{\sqrt{n}}\sum_{i=1}^{n}m(Z_i,\hat{\theta}_M) , \qquad (23)$$

where $\hat{\theta}_M$ is in this case the unconstrained minimizer of Q_n on Θ . Hence, in this setting θ_0 equals θ_M and $\hat{\theta}$ equals $\hat{\theta}_M$ in the notation of Assumption 4.1(i). Obtaining an expansion for T_n as in (9) is straightforward provided m and q are once and twice continuously differentiable in θ respectively. Defining the gradient $M_n(\theta) \equiv n^{-1} \sum_i \nabla m(Z_i, \theta)$ and Hessian $H_n(\theta) \equiv n^{-1} \sum_i \nabla^2 q(Z_i, \theta)$, standard arguments imply that under the null hypothesis:

$$\frac{1}{\sqrt{n}}\sum_{i=1}^{n}m(Z_{i},\hat{\theta}_{M}) = \frac{1}{\sqrt{n}}\sum_{i=1}^{n}m(Z_{i},\theta_{M}) - M_{n}(\theta_{M})H_{n}^{-1}(\theta_{M})\frac{1}{\sqrt{n}}\sum_{i=1}^{n}\nabla q(Z_{i},\theta_{M}) + o_{p}(1); \quad (24)$$

see Newey (1985a) for primitive conditions for (24). Thus, in this setting $s(z,\theta)$ and $A_n(\theta)$ are:

$$s(z,\theta) = \begin{pmatrix} m(z,\theta) \\ \nabla q(z,\theta) \end{pmatrix} \qquad A_n(\theta) = \begin{bmatrix} I & \vdots & -M_n(\theta)H_n^{-1}(\theta) \end{bmatrix} .$$
(25)

Moreover, if θ_M is an interior point of Θ , then $E[\nabla q(Z_i, \theta_M)] = 0$ because θ_M minimizes Q on Θ . Therefore, $G_n \xrightarrow{p} \infty$ under the alternative hypothesis due to $E[m(Z_i, \theta_M)] \neq 0$.

4.3.1 ML Specification Tests

A prominent application of hypotheses as in (22) is in model specification testing. In particular, this setting encompasses moment based specification tests ("m-tests") for maximum likelihood models, as considered in White (1982, 1994), Newey (1985b) and Tauchen (1985).⁶ Computations

⁶A bootstrap construction for the Information Matrix Equality test was also developed in Horowitz (1994).

are significantly simplified for maximum likelihood models by the generalized information matrix equality, which implies:

$$E[\nabla^2 q(Z_i, \theta_M)] = -E[\nabla q(Z_i, \theta_M) \nabla q(Z_i, \theta_M)'] \qquad E[\nabla m(Z_i, \theta_M)] = -E[m(Z_i, \theta_M) \nabla q(Z_i, \theta_M)']$$

For example, as noted in Chesher (1984) and Newey (1985b), computation of the Wald test statistic for the null hypothesis in (22) can be performed through the auxiliary regression:

$$1 = m(Z_i, \hat{\theta}_M)' \delta + \nabla q(Z_i, \hat{\theta}_M)' \gamma + \epsilon_i .$$
⁽²⁶⁾

If R^2 is the uncentered *R*-squared of the regression in (26), then under the generalized information matrix equality result in (26) the Wald test statistic is asymptotically equivalent to:⁷

$$G_n = nR^2 . (27)$$

The calculation of the score bootstrap simplifies in an analogous fashion. Under a uniform law of large numbers, we obtain that $A_n(\hat{\theta}_M)$ as defined in (24) satisfies,

$$A_{n}(\hat{\theta}_{M}) = \left[I : -\frac{1}{n} \sum_{i=1}^{n} m(Z_{i}, \hat{\theta}_{M}) \nabla q(Z_{i}, \hat{\theta}_{M})' \times \left[\frac{1}{n} \sum_{i=1}^{n} \nabla q(Z_{i}, \hat{\theta}_{M}) \nabla q(Z_{i}, \hat{\theta}_{M})'\right]^{-1}\right] + o_{p}(1) , \quad (28)$$

under the null hypothesis. As a result, the score bootstrap has a simple interpretation in terms of the multivariate regression of the moments $m(Z_i, \hat{\theta}_M)$ on the score $\nabla q(Z_i, \hat{\theta}_M)$:

$$m^{(1)}(Z_i, \hat{\theta}_M) = \nabla q(Z_i, \hat{\theta}_M)' \beta_1 + \epsilon_{1,i}$$

$$\vdots = \vdots , , \qquad (29)$$

$$m^{(l)}(Z_i, \hat{\theta}_M) = \nabla q(Z_i, \hat{\theta}_M)' \beta_l + \epsilon_{l,i}$$

where $m^{(j)}(Z_i, \hat{\theta}_M)$ is the j^{th} component of $m(Z_i, \hat{\theta}_M)$. Letting $e_{j,i} \equiv m^{(j)}(Z_i, \hat{\theta}_M) - \nabla q(Z_i, \hat{\theta}_M)' \hat{\beta}_j$ be the fitted residual of the j^{th} regression and $e_i = (e_{1,i}, \ldots, e_{l,i})'$, we obtain that:

$$S_n^*(\hat{\theta}_M) = \frac{1}{\sqrt{n}} \sum_{i=1}^n e_i W_i .$$
 (30)

Therefore, G_n^* is simply the Wald test for the null hypothesis that the mean of $e_i W_i$ equals zero.

In summary, if the generalized information matrix equality holds, then in testing (22) we may follow the simple algorithm:

STEP 1: Run the regression in (26) and compute the uncentered *R*-squared to obtain G_n as in (27). STEP 2: Regress $\{m(Z_i, \hat{\theta}_M)\}_{i=1}^n$ on the scores $\{\nabla q(Z_i, \hat{\theta}_M)\}_{i=1}^n$ to generate residual vectors $\{e_i\}_{i=1}^n$.

 $^{^7\}mathrm{See}$ also Chapter 8.2.2 in Cameron and Trivedi (2005) for a summary of these results.

STEP 3: Using random weights $\{W_i\}_{i=1}^n$ independent of $\{Y_i, X_i\}_{i=1}^n$ with $E[W_i] = 0$ and $E[W_i^2] = 1$, perturb the original residual vectors $\{e_i\}_{i=1}^n$ to obtain a new set of residual vectors $\{e_iW_i\}_{i=1}^n$.

STEP 4: Let G_n^* be the Wald test statistic for the null that $E[e_iW_i] = 0$ calculated using $\{e_iW_i\}_{i=1}^n$. To control size at level α , reject if G_n is larger than the $1 - \alpha$ quantile of G_n^* conditional on $\{Z_i\}_{i=1}^n$.

4.4 Clustered Data

Theorem 4.1 and Corollary 4.1 may be applied to clustered data provided all clusters have the same number of observations. An extension to unbalanced clusters is feasible, essentially requiring an extension of Theorem 4.1 to independent but not identically distributed observations.

Let Z_{ic} denote observation number *i* in cluster *c*, *J* be the total number of observations per cluster, *n* be the total number of clusters and $Z_c = \{Z_{1c}, \ldots, Z_{Jc}\}$. Following (9), we consider test statistics of the general form $\tilde{G}_n \equiv \tilde{T}'_n \tilde{T}_n$, where \tilde{T}_n satisfies:

$$\tilde{T}_{n} = (A_{n}(\theta_{0})\tilde{\Sigma}_{n}(\theta_{0})A_{n}(\theta_{0})')^{-\frac{1}{2}}\tilde{S}_{n}(\theta_{0}) + o_{p}(1) \qquad \tilde{S}_{n}(\theta) \equiv A_{n}(\theta)\frac{1}{\sqrt{n}}\sum_{c=1}^{n}\frac{1}{\sqrt{J}}\sum_{i=1}^{J}\tilde{s}(Z_{ic},\theta) , \quad (31)$$

where $A_n(\theta)$ is again a $r \times m$ matrix, $\tilde{s}(z, \theta)$ maps each (Z_{ic}, θ) into a $m \times 1$ vector and $\tilde{\Sigma}_n(\theta)$ is a robust covariance matrix that allows for arbitrary correlation within cluster. The Wald and LM test statistics, as well as the moment restriction tests previously discussed all extend to this setting when observations are allowed to be dependent within clusters.

The applicability of Theorem 4.1 and Corollary 4.1 to the present context is immediate once we notice that we may define $s(z, \theta)$, mapping each (Z_c, θ) into a $m \times 1$ vector, to be given by:

$$s(Z_c,\theta) = \frac{1}{\sqrt{J}} \sum_{i=1}^{J} \tilde{s}(Z_{ic},\theta) .$$
(32)

The statistics \tilde{T}_n and $\tilde{S}_n(\theta)$ are then special cases of T_n and $S_n(\theta)$ as considered in (9) but with Z_c in place of Z_i . Hence, equations (9) and (32) indicate that the relevant bootstrap statistic should perturb the data at the cluster rather than at the individual observation level. We thus define:

$$\tilde{G}_n^* \equiv \tilde{T}_n^{*'} \tilde{T}_n^* \qquad \tilde{T}_n^* \equiv (A_n(\hat{\theta}) \tilde{\Sigma}_n^*(\hat{\theta}) A_n(\hat{\theta})')^{-\frac{1}{2}} \tilde{S}_n^*(\hat{\theta}) \qquad \tilde{S}_n^*(\theta) \equiv A_n(\theta) \frac{1}{\sqrt{n}} \sum_{c=1}^n \frac{W_c}{\sqrt{J}} \sum_{i=1}^J \tilde{s}(Z_{ic}, \theta)$$

where $\tilde{\Sigma}_{n}^{*}(\theta)$ is a robust bootstrap covariance matrix for $s(Z_{ic}, \theta)W_{c}$.

Given these definitions, it is readily apparent that \tilde{G}_n^* , \tilde{T}_n^* and $\tilde{S}_n^*(\theta)$ are themselves special cases of the bootstrap statistics G_n^* , T_n^* and $S_n^*(\theta)$. The consistency of the proposed score bootstrap then follows immediately provided the clusters are i.i.d., the number of clusters tends to infinity and $s(z, \theta)$ as defined in (32) satisfies Assumption 4.1(ii), 4.2(i) and 4.4(ii).

Corollary 4.2. Under Assumptions 4.1, 4.2, 4.3 and 4.4, it follows that under the null hypothesis:

$$\lim_{n \to \infty} P(\tilde{G}_n \ge \hat{c}_{1-\alpha}) = 1 - \alpha$$

for any $0 < \alpha < 1$. Under the same assumptions, if the alternative hypothesis is instead true, then:

$$\lim_{n \to \infty} P(\tilde{G}_n \ge \hat{c}_{1-\alpha}) = 1$$

5 Higher Order Refinements

We conclude our theoretical discussion of the score bootstrap by returning to the linear model of Sections 2 and 3 and conducting an analysis of the bootstrap higher order properties in the special but important case of Wald tests. Specifically, for $\hat{\beta}$ the OLS estimate of (1) and an arbitrary $\lambda \in \mathbf{R}^m$ with $\lambda \neq 0$, we examine the studentized form of $\sqrt{n}\lambda'(\hat{\beta}-\beta_0)$ and its bootstrap counterpart:

$$T_n = \frac{\sqrt{n\lambda'}}{\hat{\sigma}} (\hat{\beta} - \beta_0) \qquad \qquad T_n^* = \frac{\sqrt{n\lambda'}}{\hat{\sigma}^*} H_n^{-1} \frac{1}{\sqrt{n}} \sum_{i=1}^n X_i \epsilon_i^* , \qquad (33)$$

where recall that in the present context $H_n = n^{-1} \sum_i X_i X'_i$ and the perturbed score is given by $X_i \epsilon_i^* = X_i (Y_i - X'_i \hat{\beta}) W_i$. The standard errors in (33) are therefore $\hat{\sigma}^2 = \lambda' H_n^{-1} \Sigma_n(\hat{\beta}) H_n^{-1} \lambda$ for the full sample estimate, and $(\hat{\sigma}^*)^2 = \lambda' H_n^{-1} \Sigma_n^*(\hat{\beta}) H_n^{-1} \lambda$ for the score bootstrap analogue.

We conduct our higher order analysis in two steps. First, we obtain stochastic expansions for T_n and T_n^* in (33) up to a remainder term of order $O_p(n^{-1})$ and $O_{p^*}(n^{-1})$ respectively. Second, we compute and compare the first three moments of such expansions and provide conditions under which they agree up to terms of order $O_p(n^{-1})$. If Edgeworth expansions are valid, then under appropriate regularity conditions the concordance of the first three moments is necessary and sufficient to ensure a first order refinement over a normal approximation; see for example Hall (1992). Because the score and wild bootstraps are equivalent to higher order as established in Lemma 3.1, the conclusions of our analysis are relevant for both bootstrap procedures.

We first define three parameters that play an important role in the higher order analysis:

$$\kappa \equiv E[(\lambda' X_i)^3 \epsilon_i^3] \qquad \gamma_0 \equiv E[(\lambda' X_i)^2 X_i \epsilon_i] \qquad \gamma_1 \equiv E[(\lambda' X_i) X_i' X_i \epsilon_i] . \tag{34}$$

The parameter κ may be interpreted as the third moment of the score, while γ_0 and γ_1 capture the possible correlation of the error ϵ_i with higher moments of the regressor X_i . If the model is properly

specified, so that ϵ_i is mean independent of X_i , then both γ_0 and γ_1 equal zero. However, as we conduct our analysis under Assumption 3.1, which only requires that the regressor and error be uncorrelated, both γ_0 and γ_1 may take nonzero values. The parameter κ may of course be nonzero under both proper specification or misspecification.

Under Assumption 3.1, we establish the following asymptotic expansion for T_n and T_n^* :

Lemma 5.1. Suppose Assumption 3.1 holds, and for $\lambda \in \mathbb{R}^m$ define the following random variables:

$$L_{n} \equiv \lambda' \{ I + \Delta_{n} \} \frac{1}{\sqrt{n\sigma}} \sum_{i=1}^{n} X_{i} \epsilon_{i} - \frac{1}{2\sigma^{3}\sqrt{n}} \sum_{i=1}^{n} (\lambda' X_{i}) \epsilon_{i} \{ (\hat{\sigma}_{R}^{2} - \sigma^{2}) - \frac{2}{n} \sum_{i=1}^{n} \gamma_{0}' X_{i} \epsilon_{i} \}$$
(35)

$$L_n^* \equiv \lambda' H_n^{-1} \frac{1}{\sqrt{n}} \sum_{i=1}^{n} X_i \epsilon_i^* \{ \frac{1}{\hat{\sigma}} - \frac{1}{2\hat{\sigma}^3} ((\hat{\sigma}^*)^2 - \hat{\sigma}^2) \}$$
(36)

where $\Delta_n \equiv I - H_n$, $\Sigma(\beta_0) = E[X_i X_i' \epsilon_i^2]$, $\hat{\sigma}_R^2 = \lambda' \Sigma_n(\beta_0) \lambda + 2\lambda' \Delta_n \Sigma(\beta_0) \lambda$ and $\sigma^2 = \lambda' \Sigma(\beta_0) \lambda$. It then follows that, $T_n = L_n + O_p(n^{-1})$ and $T_n^* = L_n^* + O_{p^*}(n^{-1})$ almost surely.

Recall that in Assumption 3.1(i) the covariance $E[X_iX'_i]$ was normalized to equal the identity matrix. Therefore $\Delta_n \equiv I - H_n$ is the estimation error in the Hessian and the first term in (35) captures the contribution to T_n of not knowing the true value of $E[X_iX'_i]$. Similarly, the contribution of having to estimate the variance is divided into two parts: (i) $2n^{-1}\sum_i \gamma'_0 X_i \epsilon_i$ which reflects use of $\hat{\beta}$ rather than β_0 in the sample variance calculations and (ii) $\hat{\sigma}_R^2 - \sigma^2$ which captures the randomness that would be present in estimating σ^2 if β_0 were known.

Under proper specification, γ_0 equals zero and as a result the contribution of employing β rather than the true parameter β_0 in estimating the variance is of smaller order than when the model is misspecified. Since under the bootstrap distribution mean independence holds, there is no term analogous to $2n^{-1}\sum_i \gamma'_0 X_i \epsilon_i$ in the bootstrap expansion (36). Finally, we also note that the expansion under the bootstrap distribution lacks a term analogous to Δ_n due to there being no randomness present in the Hessian.

Our discussion thus far suggests a discordance between the moments of L_n and L_n^* may emerge when mean independence is violated. This is indeed the case, as we establish in Lemma 5.2.

Lemma 5.2. Suppose Assumption 3.1 holds. Then $E[(L_n^2)] = E^*[(L_n^*)^2] + O_p(n^{-1})$ and in addition:

$$E[L_n] = -\frac{\kappa}{2\sigma^3\sqrt{n}} - \frac{\gamma_1}{\sigma\sqrt{n}} - \frac{2\lambda'\Sigma(\beta_0)\gamma_0}{\sigma^3\sqrt{n}} \qquad E^*[L_n] = -\frac{E[W_i^3]\hat{\kappa}}{2\hat{\sigma}^3\sqrt{n}} \\ E[(L_n)^3] = -\frac{7\kappa}{2\sigma^3\sqrt{n}} - \frac{3\gamma_1}{\sigma\sqrt{n}} - \frac{6\lambda'\Sigma(\beta_0)\gamma_0}{\sigma^3\sqrt{n}} + O(n^{-1}) \qquad E^*[(L_n^*)^3] = -\frac{7E[W_i^3]\hat{\kappa}}{2\hat{\sigma}^3\sqrt{n}} + O_p(n^{-1}) ,$$

where $\Sigma(\beta_0) \equiv E[X_i X_i' \epsilon_i^2]$, γ_0 , κ_1 and κ_2 are as defined in (34) and $\hat{\kappa} \equiv n^{-1} \sum_i (\lambda' H_n^{-1} X_i)^3 e_i^3$.

Observe first that unless the scores are symmetric ($\kappa = 0$), the score and wild bootstraps both fail to correct the first term in the bias and skewness if $E[W_i^3] \neq 1$. This property has already been noted in Liu (1988) who advocates imposing $E[W_i^3] = 1$ for precisely this reason. However, even with this restriction, two additional terms in the bias and skewness of L_n remain. These terms capture (i) the correlation between Hessian estimation error and the score, and (ii) the additional randomness of employing $\hat{\beta}$ rather than β_0 in estimating σ^2 . Both these expressions are of smaller order under mean independence, but may be present otherwise. Because the score and wild bootstraps impose mean independence in the bootstrap distribution, *both* bootstraps fail to mimic these terms. However, it is worth noting that if misspecification is local, then the wild bootstrap still provides a refinement over the normal approximation.⁸ This suggests the wild bootstrap is likely to outperform a normal approximation in datasets where misspecification is not too severe.

We conclude that under proper model specification both the score and wild bootstrap provide a refinement over the normal approximation provided the restriction $E[W_i^3] = 1$ is imposed on the perturbation weights. However, if X_i and ϵ_i are only uncorrelated instead of mean independent, then the availability of a refinement depends on whether the error is additionally uncorrelated with the third moments of the regressor (i.e. whether the Hessian is correlated with the score). To the best of our knowledge, this observation is novel to the literature on the wild bootstrap, which has primarily focused on the case of fixed regressors (Wu (1986), Liu (1988)). A notable exception is Mammen (1993), who allows for misspecification but is primarily concerned with the consistency of the bootstrap under many regressors asymptotics.

6 Simulation Evidence

To assess the small sample performance of the score bootstrap we conduct a series of Monte Carlo experiments examining the performance of bootstrap Wald and LM tests of hypotheses regarding the parameters of a linear model estimated by OLS and a nonlinear probit model estimated by maximum likelihood. We also examine the performance of a test for proper specification in the probit model. Because small sample issues often arise in settings with dependent data we work with hierarchical data generating processes (DGPs) exhibiting dependence of micro-units i within independent clusters c. We consider balanced panels with 20 observations per cluster and sampling

⁸Misspecification is local if γ_0 and γ_1 are allowed to depend on n and satisfy $\gamma_0 = O(n^{-\frac{1}{2}})$ and $\gamma_1 = O(n^{-\frac{1}{2}})$. Lemma 5.2 then implies $E[L_n] = E^*[L_n^*] + O_p(n^{-1})$ and $E[(L_n)^2] = E^*[(L_n^*)^2] = O_p(n^{-1})$ provided $E[W_i^3] = 1$.

designs ranging from 5 to 200 independent clusters.⁹

6.1 Designs

As pointed out by Chesher (1995), symmetric Monte Carlo designs are likely to yield an overly optimistic assessment of the ability of testing procedures to control size. Moreover, our theoretical analysis suggests that misspecification may impede the performance of the wild and score bootstraps. For this reason we study the performance of the different bootstrap procedures in conducting inference on a linear model under four different designs meant to reflect realistic features of microeconomic datasets. Throughout, the linear model we examine is given by:

$$Y_{ic} = X_{ic} + D_c + \eta_c + \epsilon_{ic} , \qquad (37)$$

where the regressors (X_{ic}, D_c) and cluster level error (η_c) are generated according to:

$$X_{ic} = X_c + \xi_{ic}$$
 $D_c = X_c \omega_c$ $\eta_c = (1 + D_c) v_c$. (38)

The regressor of interest is D_c , which varies only at the cluster level. Note that the cluster level random effect η_c exhibits heteroscedasticity with respect to D_c . The designs are:

Design I: (baseline) We let $(X_c, \xi_{ic}, \omega_c, \epsilon_{ic})$ be normally distributed with identity covariance matrix, and v_c independent of other variables with a *t*-distribution with six degrees of freedom.

Design II: (skewed regressor) Design I is modified to generate ω_c according to a mixture between a N(0, 1) with probability 0.9 and a N(2, 9) with probability 0.1 as in Horowitz (1997). This yields a regressor with occasional "outliers" and substantial skew and kurtosis in its marginal distribution.

Design III: (misspecification) The model estimated is still (37), but the DGP is modified to:

$$Y_{ic} = X_{ic} + D_c + .1D_c^2 + \eta_c + \epsilon_{ic} , \qquad (39)$$

and other features remain as in Design I. Hence, the quadratic term in the regressor of interest is ignored in estimation which yields an asymmetric reduced form regression error. Note that $E[D_c^3] = E[X_{ic}D_c^2] = 0$ which ensures the population regression coefficient on D_c is still one.

Design IV: (skew and misspecification) Design III is modified so that ω_c is as in Design II.

Our baseline design for the linear model exhibits fat tails in the random cluster effect but no skew in the score. Design II introduces skewness into the experiment by modifying the regressor

⁹In unreported results we found our results to be insensitive to variation in the number of observations per cluster.

(and hence the reduced form error) to contain outliers. Finally, in Designs III and IV, we explore the effects of misspecification, which may be important in lieu of Lemma 5.2.

To study the performance of the score bootstrap in a nonlinear model we consider probit estimation of the following DGP:

$$Y_{ic} = 1\{X_{ic} + D_c + \eta_c + \epsilon_{ic} \ge 0\} \qquad X_{ic} = X_c + \xi_{ic} \qquad D_c = X_c \omega_c .$$
(40)

This is essentially a latent variable representation of the model in (37) without heteroscedasticity in the group error η_c . We consider the following two designs for our probit analysis:

Design V: (baseline probit) In (40), we let $(X_{ic}, \xi_{ic}, \omega_c) \sim N(0, I_3)$ and $(\eta_c, \epsilon_{ic}) \sim N(0, I_2/2)$.¹⁰

Design VI: (skew probit) We modify Design V by generating ω_c according to a mixture distribution as in Design III, so that the regressor of interest is heavily skewed.

Finally, we illustrate the methods of Section 4.3 by testing the following moment restrictions implied by the probit model:

$$E[e_{ic}D_c^2] = E[e_{ic}D_c^3] = E[e_{ic}X_{ic}^2] = E[e_{ic}X_{ic}^3] = E[e_{ic}X_{ic}D_c] = 0$$
(41)

where $e_{ic} = [Y_{ic} - p_{ic}]\phi(X_{ic} + D_c)/[p_{ic}(1 - p_{ic})]$ is a generalized residual and $p_{ic} = \Phi(X_{ic} + D_c)$ is the conditional probability that Y_{ic} equals one given D_c and X_{ic} .¹¹ A test of these five moment conditions examines the probit model for unmodeled nonlinearities in the response function.

6.2 Results

Table 1 provides empirical false rejection rates from 10,000 Monte Carlo repetitions of Wald and LM tests of the null that the population least squares coefficient on D_c in (37) is one. All tests have a nominal size of 5% and are studentized using a recentered variance matrix estimator.¹² We consider implementations of the score bootstrap using both Rademacher weights and the skew correcting weights suggested by Mammen (1993).¹³ For comparison with the various score bootstraps we also compute the empirical rejection rates of Wald and LM tests based upon analytical clustered standard errors, the conventional wild bootstrap, and the pairs-based block bootstrap. All bootstrap tests

 $^{^{10}}$ Though the DGP contains a cluster level random effect, the marginal model for the outcome given covariates is

a standard probit ensuring that conventional maximum likelihood estimation is consistent.

¹¹Note that the ML probit scores are of the form $e_{ic}[1, X_{ic}, D_c]$.

¹²We also make a finite sample degrees of freedom correction of n/(n-1) to all variance estimators.

¹³Rademacher weights impose $E[W_i^4] = 1$ while Mammen's weights impose $E[W_i^3] = 1$.

	Normal Regressor					Mixture Regressor							
Wald Tests	n = 5	n = 10	n = 20	n = 50	n = 200	-	n = 5	n = 10	n = 20	n = 50	n = 200		
Analytical	0.442	0.328	0.240	0.153	0.083	-	0.467	0.398	0.317	0.240	0.140		
Pairs	0.020	0.088	0.079	0.054	0.040		0.030	0.134	0.110	0.080	0.052		
Wild Rademacher	0.243	0.185	0.128	0.078	0.052		0.273	0.250	0.193	0.127	0.075		
Wild Mammen	0.252	0.187	0.146	0.105	0.060		0.282	0.240	0.187	0.138	0.094		
Score Rademacher	0.263	0.194	0.142	0.091	0.048		0.270	0.223	0.188	0.142	0.091		
Score Mammen	0.288	0.220	0.162	0.104	0.053		0.292	0.245	0.206	0.156	0.096		
		Normal Regressor						Mixture Regressor					
LM Tests	n = 5	n = 10	n = 20	n = 50	n = 200	_	n = 5	n = 10	n = 20	n = 50	n = 200		
Analytical	0.001	0.023	0.030	0.037	0.038	-	0.001	0.021	0.024	0.026	0.030		
Pairs	0.051	0.061	0.051	0.043	0.039		0.054	0.057	0.047	0.045	0.035		
Wild Rademacher	0.103	0.065	0.039	0.039	0.046		0.112	0.067	0.036	0.031	0.038		
Wild Mammen	0.165	0.103	0.068	0.057	0.051		0.177	0.102	0.065	0.049	0.046		
Score Rademacher	0.105	0.077	0.062	0.053	0.048		0.121	0.088	0.060	0.052	0.052		
Score Mammen	0.084	0.034	0.026	0.026	0.033		0.097	0.036	0.024	0.017	0.026		

Table 1: EMPIRICAL REJECTION RATES, OLS (PROPERLY SPECIFIED)

were computed using 200 bootstrap repetitions. Stata code for our Monte Carlo experiments is available online.

The standard clustered Wald test severely over-rejects in samples with few clusters, with performance further degrading when the regressors are generated according to a mixture distribution. A conventional pairs bootstrap of the Wald test yields dramatic improvements in size control though its performance degrades somewhat when the regressor of interest exhibits outliers. Wild bootstrapping the Wald test yields improvements over analytical methods but under performs relative to pairs regardless of whether Mammen or Rademacher weights are used. As suggested by our theoretical results, the score bootstrap yields results roughly in line with those of the corresponding Wild bootstrap.

In contrast to the Wald tests, the clustered LM tests appear to perform well across a range of sample sizes and regardless of the distribution of the regressors. While the analytical LM test yields mild underrejection with few clusters, its Wild bootstrapped analogue actually yields slight over-rejection.¹⁴ The score bootstrapped LM tests perform as well as or better than the wild

¹⁴We note that the wild bootstrapped LM test is similar to the Wild bootstrap procedure of Cameron, Gelbach,

	Normal Regressor					Mixture Regressor						
Wald Tests	n = 5	n = 10	n = 20	n = 50	n = 200	-	n = 5	n = 10	n = 20	n = 50	n = 200	
Analytical	0.448	0.333	0.248	0.162	0.086	-	0.473	0.411	0.331	0.262	0.157	
Pairs	0.022	0.091	0.084	0.059	0.042		0.033	0.135	0.110	0.073	0.042	
Wild Rademacher	0.249	0.192	0.135	0.078	0.051		0.278	0.257	0.198	0.135	0.068	
Wild Mammen	0.254	0.195	0.150	0.105	0.060		0.286	0.247	0.191	0.146	0.090	
Score Rademacher	0.253	0.184	0.135	0.087	0.045		0.259	0.214	0.185	0.144	0.101	
Score Mammen	0.277	0.210	0.154	0.099	0.047		0.281	0.234	0.200	0.156	0.102	
	Normal Regressor						Mixture Regressor					
LM Tests	n = 5	n = 10	n = 20	n = 50	n = 200	-	n = 5	n = 10	n = 20	n = 50	n = 200	
Analytical	0.001	0.022	0.032	0.037	0.040	-	0.001	0.021	0.026	0.024	0.026	
Pairs	0.051	0.062	0.053	0.045	0.040		0.055	0.062	0.049	0.039	0.031	
Wild Rademacher	0.106	0.062	0.042	0.041	0.044		0.116	0.062	0.033	0.032	0.037	
Wild Mammen	0.167	0.103	0.072	0.059	0.048		0.181	0.099	0.063	0.052	0.042	
Score Rademacher	0.110	0.082	0.061	0.057	0.051		0.118	0.090	0.061	0.055	0.051	
Score Mammen	0.094	0.033	0.026	0.027	0.034		0.094	0.038	0.025	0.021	0.024	

Table 2: EMPIRICAL REJECTION RATES, OLS (MISSPECIFIED)

bootstrapped LM tests under both regressor designs. They also perform comparably to the pairs bootstrapped Wald tests. However the pairs bootstrapped LM tests yield the best performance of the group, with coverage rates closest to nominal levels across a range of sample sizes.

Table 2 examines the performance of Wald and LM tests when the model is misspecified. Again the performance of the analytical clustered Wald test appears to be very poor in small samples or when the regressor of interest exhibits outliers. Correcting the critical values of the Wald test with the pairs bootstrap yields much improved though still sometimes unsatisfactory performance. As before, the Wild bootstrap improves on the performance of analytical Wald tests but still overrejects substantially. Score bootstrapping the Wald statistic yields results mimicking those of the Wild bootstrap.

With misspecification and an asymmetric regressor, the clustered LM test underrejects substantially in small samples. Wild or score bootstrapping the LM test leads to slight overrejection in small samples. The score bootstrap with Rademacher weights seems to perform particularly well. and Miller (2008) who impose the null $\hat{\beta} = \beta_0$ when generating the bootstrap distribution of outcomes as in (3). In results not shown we found the results from this procedure (which is akin to comparing the bootstrap critical values of an LM statistic to a full sample Wald) to be quite similar to those found in wild or score bootstrap LM tests.

	Normal Regressor						Mixture Regressor					
Wald Tests	n = 5	n = 10	n = 20	n = 50	n = 200	-	n = 5	n = 10	n = 20	n = 50	n = 200	
Analytical	0.319	0.192	0.132	0.094	0.063	-	0.320	0.210	0.137	0.092	0.063	
Pairs	0.032	0.065	0.073	0.064	0.053		0.037	0.062	0.075	0.064	0.052	
Score Rademacher	0.272	0.153	0.097	0.064	0.037		0.283	0.174	0.102	0.062	0.039	
Score Mammen	0.303	0.179	0.110	0.067	0.038		0.316	0.204	0.115	0.063	0.038	
	Normal Regressor						Mixture Regressor					
LM Tests	n = 5	n = 10	n = 20	n = 50	n = 200	-	n = 5	n = 10	n = 20	n = 50	n = 200	
Analytical	0.004	0.070	0.079	0.080	0.060	-	0.004	0.080	0.089	0.079	0.062	
Pairs ¹⁵	n.a.	0.140	0.125	0.096	0.065		n.a.	0.144	0.138	0.091	0.059	
Score Rademacher	0.101	0.124	0.099	0.085	0.060		0.111	0.142	0.108	0.085	0.064	
Score Mammen	0.079	0.054	0.077	0.081	0.063		0.089	0.063	0.085	0.082	0.061	

Table 3: Empirical Rejection Rates, Probit (Properly Specified)

Table 3 examines the performance of Wald and LM tests in the probit model. Here both Wald and LM tests tend to overreject when asymptotic critical values are used. Use of the pairs bootstrap corrects for this overrejection though in small samples we were sometimes unable to compute the bootstrap distribution.¹⁶ Score bootstrapping the Wald test yields improvements over analytical clustered standard errors but substantial overrejection remains in small samples. Score bootstrapping the LM test with Mammen weights, on the other hand, yields size control roughly on par with the pairs bootstrap.

Table 4: EMPIRICAL REJECTION RATES, m-TEST (PROBIT)

	Normal Regressor					_	Mixture Regressor					
Wald Tests	n = 5	n = 10	n = 20	n = 50	n = 200	-	n = 5	n = 10	n = 20	n = 50	n = 200	
Analytical	0.766	0.593	0.341	0.190	0.130		0.776	0.584	0.348	0.208	0.131	
Score Rademacher	0.067	0.135	0.174	0.159	0.138		0.068	0.129	0.180	0.177	0.142	
Score Mammen	0.049	0.034	0.092	0.134	0.141		0.046	0.032	0.093	0.147	0.143	

Finally, Table 4 examines the performance of tests for proper specification of the probit model via the restrictions in (41). Because the information matrix equality holds under both DGPs we use the outer product version of the test described in 4.3.1 generalized to allow for clustering

 $^{^{15}}$ We were unable to compute the LM statistic in the majority of pairs draws with 5 clusters.

¹⁶We discarded bootstrap draws for which we were unable to compute a maximum likelihood estimate.

and a recentered variance matrix estimator. We see that the analytical m-test procedure overrejects substantially in small samples and continues to exhibit poor control over size even with 200 clusters. Surprisingly the score bootstrapped versions of the test work well in small samples, but appear to degrade slightly as the number of clusters increase. With 200 clusters, the analytical and bootstrap approaches appear to work equally well.

7 Conclusion

Score bootstrap tests provide a computational advantage over conventional wild and pairs bootstraps and may easily be applied to estimators that lack conventional residuals. Both our theoretical and Monte Carlo results suggest the wild bootstrap possesses no inferential advantage over the score bootstrap despite the potentially significant increase in computational cost. And though our theoretical analysis indicates the performance of the score bootstrap may be sensitive to misspecification, we found little effect of introducing mild specification errors into our Monte Carlo designs.

Like Moreira, Porter, and Suarez (2009) we find that bootstrapping Lagrange Multiplier type tests yields improved small sample size control in a number of difficult testing environments of substantial applied interest. Economists have typically shied away from bootstrap LM tests perhaps due to the difficulty of constructing confidence intervals by test inversion. The score bootstrap methods developed here substantially reduce the cost of such an exercise and may enable researchers to conduct inference in a wider range of small sample environments than previously contemplated.

An interesting extension would be to consider a modification of the score bootstrap which additionally perturbs the Hessian with random weights. This may provide a refinement in models where the Hessian and score are correlated either due to misspecification or nonlinearity. However, in nonlinear models such a perturbation would need to account for the influence of parameter uncertainty in the Hessian which, if dealt with via an additional linearization, would require computation of third derivatives. This would likely prove to be overly burdensome in all but a few special cases.

APPENDIX

PROOF OF LEMMA 3.1: First notice that by Markov's inequality, $E[W_i^2] = 1$ and the i.i.d. assumption

$$P^*(\|\frac{1}{\sqrt{n}}\sum_{i=1}^n X_i\epsilon_i^*\| > C) \le \frac{1}{nC^2}E^*[(\sum_{i=1}^n X_i\epsilon_i^*)'(\sum_{i=1}^n X_i\epsilon_i^*)] = \frac{1}{nC^2}\sum_{i=1}^n X_i'X_ie_i^2.$$
(42)

Since $n^{-1} \sum_i X'_i X_i e_i^2 \xrightarrow{a.s.} E[X'_i X_i \epsilon_i^2] < \infty$ and $H_n \xrightarrow{a.s.} I$, we obtain from (42) that almost surely:

$$\|\sqrt{n}(\hat{\beta}^* - \hat{\beta})\| \le \|H_n^{-1}\|_F \times \|\frac{1}{\sqrt{n}} \sum_{i=1}^n X_i \epsilon_i^*\| = O_{p^*}(1) , \qquad (43)$$

where $\|\cdot\|_F$ denotes the Frobenius norm. Next observe that for $\|\cdot\|_o$ the operator norm, we have:

$$\begin{aligned} \| (H_n^{-1} \Sigma_n^*(\hat{\beta}^*) H_n^{-1})^{-1} - (H_n^{-1} \Sigma_n^*(\hat{\beta}) H_n^{-1})^{-1} \|_o \\ & \leq \| (H_n^{-1} \Sigma_n^*(\hat{\beta}) H_n^{-1})^{-1} \|_o \times \| H_n^{-1} (\Sigma_n^*(\hat{\beta}) - \Sigma_n^*(\hat{\beta}^*)) H_n^{-1} \|_o \times \| (H_n^{-1} \Sigma_n^*(\hat{\beta}^*) H_n^{-1})^{-1} \|_o . \end{aligned}$$
(44)

Let $X_i^{(k)}$ denote the k^{th} element of the vector X_i . Arguing as in (42), it is straightforward to show that $n^{-\frac{1}{2}} \sum_i X_i^{(k)} X_i^{(l)} X_i^{(s)} \epsilon_i^* = O_{p^*}(1)$ almost surely for any indices k, l, s. Therefore, since $\|\cdot\|_o \leq \|\cdot\|_F$ we conclude from (43) and direct calculation that we must have:

$$\begin{split} \|\Sigma_{n}^{*}(\hat{\beta}) - \Sigma_{n}^{*}(\hat{\beta}^{*})\|_{o} &\leq \|\frac{1}{n}\sum_{i=1}^{n}X_{i}X_{i}'\{(Y_{i}^{*} - X_{i}'\hat{\beta})^{2} - (Y_{i}^{*} - X_{i}'\hat{\beta}^{*})^{2}\}\|_{F} \\ &= \|\frac{1}{n}\sum_{i=1}^{n}X_{i}X_{i}'\{2\epsilon_{i}^{*}(X_{i}'(\hat{\beta} - \hat{\beta}^{*})) + (X_{i}'(\hat{\beta} - \hat{\beta}^{*}))^{2}\}\|_{F} = O_{p^{*}}(n^{-1}) \quad a.s. \quad (45) \end{split}$$

Moreover, since $E[(\epsilon_i^*)^k] = E[W_i^k]e_i^k$, we can also obtain from the i.i.d. assumption that:

$$E^*[\|\frac{1}{n}\sum_{i=1}^n X_i X_i'\{(\epsilon_i^*)^2 - e_i^2\}\|_F^2] = \sum_{l=1}^m \sum_{s=1}^m E^*[(\frac{1}{n}\sum_{i=1}^n X_i^{(l)}X_i^{(s)}\{(\epsilon_i^*)^2 - e_i^2\})^2]$$
$$= \frac{1}{n}\sum_{l=1}^m \sum_{s=1}^m \frac{1}{n}\sum_{i=1}^n (X_i^{(l)}X_i^{(s)})^2\{(E[W_i^4] - 1)e_i^4\} = o_{a.s.}(1) . \quad (46)$$

Therefore, since $n^{-1} \sum_{i=1}^{n} X_i X'_i e_i^2 \xrightarrow{a.s.} E[X_i X'_i \epsilon_i^2]$ and $H_n^{-1} \xrightarrow{a.s.} I$, results (45) and (46) establish:

$$\|H_n^{-1}\Sigma_n^*(\hat{\beta})H_n^{-1} - E[X_iX_i'\epsilon_i^2]\|_F = o_{p^*}(1) \qquad \|H_n^{-1}\Sigma_n^*(\hat{\beta}^*)H_n^{-1} - E[X_iX_i'\epsilon_i^2]\|_F = o_{p^*}(1)$$
(47)

almost surely. Next, for any normal matrix A, let $\xi(A)$ denote its smallest eigenvalue. Since $|\xi(A) - \xi(B)| \le ||A - B||_F$ by Corollary III.2.6 in Bhatia (1997), it then follows from (47) that:

$$\xi(H_n^{-1}\Sigma_n^*(\hat{\beta})H_n^{-1}) = \xi(E[X_iX_i'\epsilon_i^2]) + o_{p^*}(1) \qquad \qquad \xi(H_n^{-1}\Sigma_n^*(\hat{\beta}^*)H_n^{-1}) = \xi(E[X_iX_i'\epsilon_i^2]) + o_{p^*}(1)$$
(48)

almost surely. However, since for any normal matrix A, we have $||A^{-1}||_o = \xi(A)$, result (48) and Assumption 3.1(i) imply $||H_n^{-1}\Sigma_n^*(\hat{\beta}^*)H_n^{-1}||_o = O_{p^*}(1)$ and $||H_n^{-1}\Sigma_n^*(\hat{\beta})H_n^{-1}||_o = O_{p^*}(1)$ almost surely. The claim of the Lemma then follows by combining results (43), (44) and (45). PROOF OF THEOREM 4.1: Let $\Sigma(\theta) = E[s(Z_i, \theta)s(Z_i, \theta)']$. As argued in (66), the matrix $s(z, \theta)s(z, \theta)'$ has an integrable envelope. Hence, since $s(z, \theta)s(z, \theta)'$ is continuous in θ for all z by Assumption 4.4(ii), the dominated convergence theorem implies $\Sigma(\theta)$ is continuous in θ . Therefore, by Lemma 7.1 and Assumption 4.1(i), we obtain $\Sigma_n^*(\hat{\theta}) = \Sigma(\theta_0) + o_p(1)$. In addition, $A_n(\hat{\theta}) = A(\theta_0) + o_p(1)$ by Assumption 4.3(ii) and hence Assumption 4.1(ii) and $\sup_{\theta \in \Theta} ||n^{-\frac{1}{2}} \sum_i s(Z_i, \theta)W_i|| = O_p(1)$ by Lemma 7.2 imply:

$$(A_{n}(\hat{\theta})\Sigma_{n}^{*}(\hat{\theta})A_{n}(\hat{\theta})')^{-\frac{1}{2}}A_{n}(\hat{\theta})\frac{1}{\sqrt{n}}\sum_{i=1}^{n}s(Z_{i},\hat{\theta})W_{i}$$

= $(A(\theta_{0})\Sigma(\theta_{0})A(\theta_{0})')^{-\frac{1}{2}}A(\theta_{0})\frac{1}{\sqrt{n}}\sum_{i=1}^{n}s(Z_{i},\hat{\theta})W_{i} + o_{p}(1)$
= $(A(\theta_{0})\Sigma(\theta_{0})A(\theta_{0})')^{-\frac{1}{2}}A(\theta_{0})\frac{1}{\sqrt{n}}\sum_{i=1}^{n}s(Z_{i},\theta_{0})W_{i} + o_{p}(1)$, (49)

where the second equality follows by Assumption 4.1(i) and Lemma 7.2. Let BL_c be the set of Lipschitz real valued functions whose Lipschitz constant and level are less than c. For two random variables Y, V:

$$||Y - V||_{BL_1} \equiv \sup_{f \in BL_1} |E[f(Y)] - E[f(V)]| , \qquad (50)$$

metrizes weak convergence, see for example Theorem 1.12.4 in van der Vaart and Wellner (1996). Define:

$$\bar{T}_{n}^{*} \equiv (A(\theta_{0})\Sigma(\theta_{0})A(\theta_{0})')^{-\frac{1}{2}}A(\theta_{0})\frac{1}{\sqrt{n}}\sum_{i=1}^{n}s(Z_{i},\theta_{0})W_{i} .$$
(51)

Using that all $f \in BL_1$ are bounded in level and Lipschitz constant by one, we obtain for any $\eta > 0$:

$$\sup_{f \in BL_1} |E[f(\bar{T}_n^*)|\{Z_i\}_{i=1}^n] - E[f(T_n^*)|\{Z_i\}_{i=1}^n]| \le \eta P(|\bar{T}_n^* - T_n^*| \le \eta |\{Z_i\}_{i=1}^n) + 2P(|\bar{T}_n^* - T_n^*| > \eta |\{Z_i\}_{i=1}^n) .$$

$$(52)$$

However, by the law of iterated expectations and (49), we have that $P(|\bar{T}_n^* - T_n^*| > \eta | \{Z_i\}_{i=1}^n)$ converges to zero in mean, and hence in probability. As a result, since η is arbitrary, result (52) in fact implies:

$$\sup_{f \in BL_1} |E[f(\bar{T}_n^*)|\{Z_i\}_{i=1}^n] - E[f(T_n^*)|\{Z_i\}_{i=1}^n]| = o_p(1) .$$
(53)

Let $T^*_{\infty} \sim N(0, I)$. Since $\|\cdot\|_{BL_1}$ metrizes weak convergence, Assumptions 4.3(i) and 4.4(i) together with Lemma 2.9.5 in van der Vaart and Wellner (1996) in turn let us conclude that:

$$\sup_{f \in BL_1} |E[f(\bar{T}_n^*)|\{Z_i\}_{i=1}^n] - E[f(T_\infty^*)]| = o_p(1) .$$
(54)

For any M > 0, define the map $g_M : \mathbf{R}^r \to \mathbf{R}$ to be given by $g_M(a) = \min\{a'a, M\}$ and notice that for any $a, b \in \mathbf{R}^r$ we have $|g_M(a) - g_M(b)| \le 2\sqrt{M} ||a - b||$ and $g_M(a) \le M$ so that for $M \ge 4$ we have $g_M \in BL_M$. As a result, for any $f \in BL_1$, $f \circ g_M \in BL_M$ and $M^{-1}f \circ g_M \in BL_1$, which implies:

$$\sup_{f \in BL_1} |E[f(g_M(T_n^*))|\{Z_i\}_{i=1}^n] - E[f(g_M(T_\infty^*))]| \le M \sup_{f \in BL_1} |E[f(T_n^*)|\{Z_i\}_{i=1}^n] - E[f(T_\infty^*)]| = o_p(1) , \quad (55)$$

where the final result follows by (53) and (54). Since $G_n^* = T_n^{*'}T_n^*$ and every $f \in BL_1$ is bounded by one,

$$\sup_{f \in BL_1} |E[f(G_n^*) - f(g_M(T_n^*))| \{Z_i\}_{i=1}^n]| \le 2P(T_n^{*'}T_n^* > M|\{Z_i\}_{i=1}^n) .$$
(56)

By (49) and the continuous mapping theorem, $T_n^{*'}T_n^* \xrightarrow{L} \mathcal{X}_r^2$ unconditionally and hence is asymptotically tight. For an arbitrary $\eta > 0$ it then follows by Markov's inequality that for M sufficiently large:

$$\limsup_{n \to \infty} P(2P(T_n^{*'}T_n^* > M | \{Z_i\}_{i=1}^n) > \eta) \le \limsup_{n \to \infty} \frac{2}{\eta} P(T_n^{*'}T_n^* > M) < \eta .$$
(57)

Similarly, let $G_{\infty}^* \sim \mathcal{X}_r^2$ and notice that by selecting M appropriately large we may also obtain:

$$\sup_{f \in BL_1} |E[f(G_{\infty}^*) - f(g_M(T_{\infty}^*))]| \le 2P(T_{\infty}^{*'}T_{\infty}^* > M) < \eta .$$
(58)

Since η is arbitrary, results (55), (56), (57) and (58) in turn allow us to conclude that:

$$\sup_{f \in BL_1} |E[f(G_n^*)| \{Z_i\}_{i=1}^n] - E[f(G_\infty^*)]| = o_p(1) , \qquad (59)$$

which establishes the weak convergence of the distribution of G_n^* conditional on $\{Z_i\}_{i=1}^n$ to that of G_∞^* in probability. Letting F be the cdf of G_∞^* , we obtain by the Portmanteau theorem, G_∞^* having a continuous distribution, result (59) and Lemma 7.3 that for any $c \in \mathbf{R}$, $F_n^*(c) = F(c) + o_p(1)$ and $F_n(c) = F(c) + o(1)$. To establish the Theorem observe that the convergence is in fact uniform in $c \in \mathbf{R}$ by Lemma 2.11 in van der Vaart (1999).

PROOF OF COROLLARY 4.1: Let F denote the cdf of a \mathcal{X}_r^2 random variable and $c_{1-\alpha}$ be its $1-\alpha$ quantile. As argued following (59), $\sup_c |F_n^*(c) - F(c)| = o_p(1)$, and hence by Lemma 7.4 it follows that $\hat{c}_{1-\alpha} = c_{1-\alpha} + o_p(1)$ provided $0 < \alpha < 1$. The first claim of the Corollary then follows by Lemma 7.3 and the continuous mapping theorem.

For the second claim of the Corollary, observe that the bootstrap statistic $S_n^*(\hat{\theta})$ remains properly centered. In fact, (59) was established without appealing to Assumption 4.2(i). Therefore, $\hat{c}_{1-\alpha} = c_{1-\alpha} + o_p(1)$ under the alternative hypothesis as well. However, under the alternative hypothesis $G_n \xrightarrow{p} \infty$ by Lemma 7.3 and therefore the second claim of the Corollary follows.

PROOF OF COROLLARY 4.2: Given the definitions, this is a special case of Corollary 4.1. ■

PROOF OF LEMMA 5.1: We first establish the expansion for the full sample statistic. Note that since $||I - H_n|| = o_p(1)$, we obtain that with probability tending to one $||I - H_n|| < 1$ and hence we expand:

$$H_n^{-1} \frac{1}{\sqrt{n}} \sum_{i=1}^n X_i \epsilon_i = \{I + \Delta_n + \sum_{k \ge 2} \Delta_n^k\} \frac{1}{\sqrt{n}} \sum_{i=1}^n X_i \epsilon_i$$
(60)

For a symmetric matrix S, let $\xi(S)$ denote its largest eigenvalue. Since $E[X_i X'_i] = I$ and $E[||X_i||^2] < \infty$, $||\Delta_n||_F = O_p(n^{-\frac{1}{2}})$. Therefore, $\xi(\Delta_n^2) \leq \operatorname{trace}\{\Delta_n^2\} = O_p(n^{-1})$ and $\xi(S^k) = \xi^k(S)$ imply:

$$\|\{\sum_{k\geq 2}\Delta_n^k\}\frac{1}{\sqrt{n}}\sum_{i=1}^n X_i\epsilon_i\| \le \sum_{k\geq 2}\xi^k(\Delta_n) \times \|\frac{1}{\sqrt{n}}\sum_{i=1}^n X_i\epsilon_i\| = \frac{\xi^2(\Delta_n)}{1-\xi(\Delta_n)} \times O_p(1) = O_p(n^{-1}) .$$
(61)

By standard arguments, Assumption 3.1(i)-(iii) imply $\Sigma_n(\hat{\beta}) = \Sigma_n(\beta_0) + O_p(n^{-\frac{1}{2}}) = \Sigma(\beta_0) + O_p(n^{-\frac{1}{2}})$. As the calculations in (60), (61) show $H_n^{-1} = I + \Delta_n + O_p(n^{-1})$, and $\|\Delta_n\|_F = O_p(n^{-\frac{1}{2}})$ we obtain:

$$\hat{\sigma}^2 = \lambda' H_n^{-1} \Sigma_n(\hat{\beta}) H_n^{-1} \lambda = \lambda' \{ I + \Delta_n \} \Sigma_n(\hat{\beta}) \{ I + \Delta_n \} \lambda + O_p(n^{-1}) = \lambda' \Sigma_n(\hat{\beta}) \lambda + 2\lambda' \Delta_n \Sigma(\beta_0) \lambda + O_p(n^{-1}) .$$
(62)

In turn, since $H_n^{-1} = I + O_p(n^{-\frac{1}{2}})$ implies $(\hat{\beta} - \beta_0) = n^{-1} \sum_{i=1}^n X_i \epsilon_i + O_p(n^{-1})$, and $n^{-1} \sum_i \epsilon_i (\lambda' X_i)^2 X_i = E[\epsilon_i (\lambda' X_i)^2 X_i] + O_p(n^{-\frac{1}{2}})$ by the central limit theorem, expanding the square we can then obtain that:

$$\lambda' \Sigma_n(\hat{\beta}) \lambda = \frac{1}{n} \sum_{i=1}^n (\lambda' X_i)^2 (Y_i - X_i' \hat{\beta})^2 = \frac{1}{n} \sum_{i=1}^n (\lambda' X_i)^2 \{ \epsilon_i^2 + 2\epsilon_i X_i' (\beta_0 - \hat{\beta}) + (X_i' (\beta_0 - \hat{\beta}))^2 \}$$
$$= \lambda' \Sigma_n(\beta_0) \lambda - \frac{2}{n} \sum_{i=1}^n E[\epsilon_i (\lambda' X_i)^2 X_i'] X_i \epsilon_i + O_p(n^{-1}) .$$
(63)

The first claim of the Lemma then follows by a Taylor expansion and results (60), (61), (62) and (63):

$$\frac{\sqrt{n\lambda'}}{\hat{\sigma}}(\hat{\beta}-\beta_0) = \frac{\lambda'}{\sqrt{n}}H_n^{-1}\sum_{i=1}^n X_i\epsilon_i\{\frac{1}{\sigma} - \frac{1}{2\sigma^3}(\hat{\sigma}^2 - \sigma^2)\} + O_p(n^{-1}) = L_n + O_p(n^{-1}) .$$
(64)

Next, for notational simplicity let $R_{in} = \lambda' H_n^{-1} X_i (Y_i - X'_i \hat{\beta})$ and apply Markov's inequality to obtain:

$$P^*(|(\hat{\sigma}^*)^2 - \hat{\sigma}^2| > \frac{C}{\sqrt{n}}) = P^*(|\frac{1}{n}\sum_{i=1}^n R_{in}^2(W_i^2 - 1)| > \frac{C}{\sqrt{n}})$$

$$\leq \frac{n}{C^2}E^*[(\frac{1}{n}\sum_{i=1}^n R_{in}^2(W_i^2 - 1))^2] = \frac{1}{Cn}\sum_{i=1}^n R_{in}^4E[(W_i^2 - 1)^2] . \quad (65)$$

However, under our moment assumptions, $n^{-1} \sum_i R_{in}^4 E[(W_i^2 - 1)^2] \xrightarrow{a.s.} E[(\lambda' X_i)^4 \epsilon_i^4] E[(W_i^2 - 1)^2] < \infty$, and therefore from (65) it follows that $(\hat{\sigma}^*)^2 = \hat{\sigma}^2 + O_{p^*}(n^{-\frac{1}{2}})$ almost surely. The second claim of the Lemma then follows from a second order Taylor expansion of $(\hat{\sigma}^*)^{-1}$.

PROOF OF LEMMA 5.2: Follows immediately from Lemmas 7.5, 7.6 and 7.7. ■

Lemma 7.1. Let $\{W_i\}_{i=1}^n$ be an *i.i.d.* sample independent of $\{Z_i\}_{i=1}^n$ satisfying $E[W_i^2] = 1$. If Assumptions 4.1, 4.3(*i*) and 4.4(*ii*) hold, then the class $\mathcal{F} = \{s(z, \theta)s(z, \theta)'w^2 : \theta \in \Theta\}$ is Glivenko-Cantelli.

PROOF: By Assumption 4.4(ii), $s(z,\theta)w$ is continuous in $\theta \in \Theta$, and hence so is $s(z,\theta)s(z,\theta)'w^2$. Let $s^{(l)}(z,\theta)$ be the l^{th} component of the vector $s(z,\theta)$. By the mean value theorem and Assumption 4.4(ii):

$$|s^{(l)}(z,\theta)| \le |s^{(l)}(z,\theta) - s^{(l)}(z,\theta_0)| + |s^{(l)}(z,\theta_0)| \le F(z) \|\theta - \theta_0\| + |s^{(l)}(z,\theta_0)| .$$
(66)

Hence, for $D = \operatorname{diam}(\Theta)$ we obtain $|s^{(l)}(z,\theta)s^{(k)}(z,\theta)w^2| \le w^2(F(z)D + |s^{(l)}(z,\theta_0)|)(F(z)D + |s^{(k)}(z,\theta_0)|)$, which is integrable for all $1 \le i \le j \le k$ due to Assumption 4.4(ii) and 4.1(ii). We conclude that \mathcal{F} has an integrable envelope, and the Lemma follows by Example 19.8 in van der Vaart (1999).

Lemma 7.2. Under Assumptions 4.1(i), 4.3(i) and 4.4(i)-(ii), $\mathcal{F} \equiv \{ws(z,\theta) : \theta \in \Theta\}$ is Donsker.

PROOF: Let $\|\cdot\|_o$ and $\|\cdot\|_F$ denote the operator and Frobenious norms. Using $\|\cdot\|_o \leq \|\cdot\|_F$, Assumption 4.4(ii) and the mean value theorem, we obtain that for some $\bar{\theta}$ a convex combination of θ_1 and θ_2 :

$$\|ws(z,\theta_1) - ws(z,\theta_2)\| = \|w\| \times \|\nabla s(z,\bar{\theta})(\theta_1 - \theta_2)\| \le \|w\| \times \|\nabla s(z,\bar{\theta})\|_o \times \|\theta_1 - \theta_2\| \le \|w\| \times F(z) \times \|\theta_1 - \theta_2\| .$$
(67)

Hence, the class \mathcal{F} is Lipschitz in $\theta \in \Theta$, and by Theorem 2.7.11 in van der Vaart and Wellner (1996):

$$N_{[]}(2\epsilon \|\tilde{F}\|_{L^2}, \mathcal{F}, \|\cdot\|_{L^2}) \le N(\epsilon, \Theta, \|\cdot\|) , \qquad (68)$$

where $\tilde{F}(w,z) \equiv |w|F(z)$. Let $D \equiv \operatorname{diam}(\Theta)$ and $M^2 \equiv E[\tilde{F}^2(W_i, Z_i)]$ and notice that Assumptions 4.4(i)-(ii) imply $M < \infty$. Since by (67), the diameter of \mathcal{F} under $\|\cdot\|_{L^2}$ is less than or equal to MD,

$$\int_{0}^{\infty} \sqrt{\log N_{[]}(\epsilon, \mathcal{F}, \|\cdot\|_{L^{2}})} d\epsilon \leq \int_{0}^{MD} \sqrt{\log N_{[]}(\epsilon, \mathcal{F}, \|\cdot\|_{L^{2}})} d\epsilon = 2M \int_{0}^{\frac{D}{2}} \sqrt{\log N_{[]}(2Mu, \mathcal{F}, \|\cdot\|_{L^{2}})} du$$
$$\leq 2M \int_{0}^{\frac{D}{2}} \sqrt{\log N(u, \Theta, \|\cdot\|)} du \leq 2M \int_{0}^{\frac{D}{2}} \sqrt{p \log(D/u)} du < \infty$$
(69)

where in the first equality we made a change of variables $u = \epsilon/2M$, the second inequality follows from (68) and the third by $N(u, \Theta, \|\cdot\|) \leq (\operatorname{diam}(\Theta)/u)^p$. The claim of the Lemma then follows from (69), $E[\tilde{F}^2(Z_i, W_i)] < \infty$ and Theorem 2.5.6 in van der Vaart and Wellner (1996).

Lemma 7.3. Suppose Assumptions 4.1, 4.2, 4.3 and 4.4(ii) hold. If the null hypothesis is true, it then follows that $G_n \xrightarrow{L} \mathcal{X}_r^2$. On the other hand, if the alternative hypothesis is true, then $G_n \xrightarrow{p} \infty$.

PROOF: We first study the limiting behavior of G_n under the null hypothesis. For this purpose, notice that Assumption 4.3(ii) implies that $A_n(\theta_0) = A(\theta_0) + o_p(1)$. Therefore, by Assumptions 4.2(i), 4.1(ii), 4.3(i), the central limit theorem and the continuous mapping theorem, we conclude that:

$$T_n = A_n(\theta_0) \frac{1}{\sqrt{n}} \sum_{i=1}^n s(Z_i, \theta_0) + o_p(1) \xrightarrow{L} N(0, A(\theta_0) \Sigma(\theta_0) A(\theta_0)') , \qquad (70)$$

where $\Sigma(\theta) \equiv E[s(Z_i, \theta)s(Z_i, \theta)']$. Lemma 7.1 applied to $W_i = 1$ with probability one in turn implies that $\Sigma_n(\theta_0) = \Sigma(\theta_0) + o_p(1)$. By Assumption 4.3(ii) and the continuous mapping theorem we obtain:

$$A_n(\theta_0)\Sigma_n(\theta_0)A_n(\theta_0)' = A(\theta_0)\Sigma(\theta_0)A(\theta_0)' + o_p(1) .$$
(71)

It follows that $A_n(\theta_0)\Sigma_n(\theta_0)A_n(\theta_0)'$ is then invertible with probability tending to one by Assumption 4.1(ii) and the first claim of the Lemma follows immediately from (70) and the continuous mapping theorem. The second claim of the Lemma was assumed in Assumption 4.2(ii).

Lemma 7.4. Let $F_n : \mathbf{R} \to [0,1]$, $F : \mathbf{R} \to [0,1]$ be monotonic, $\sup_{c \in \mathbf{R}} |F_n(c) - F(c)| = o_p(1)$ and define:

$$c_{\alpha} \equiv \inf\{c: F(c) \ge \alpha\}$$
 $c_{n,\alpha} \equiv \inf\{c: F_n(c) \ge \alpha\}$.

If F is strictly increasing at c_{α} , it then follows that $c_{n,\alpha} = c_{\alpha} + o_p(1)$.

PROOF: Fix $\epsilon > 0$. Since by hypothesis F is strictly increasing at c_{α} it follows by definition of c_{α} :

$$F(c_{\alpha} - \epsilon) < \alpha < F(c_{\alpha} + \epsilon) .$$
(72)

Moreover, since $F_n(c_{\alpha} + \epsilon) > \alpha$ implies that $c_{n,\alpha} \le c_{\alpha} + \epsilon$ and $F_n(c_{\alpha} - \epsilon) < \alpha$ implies that $c_{n,\alpha} > c_{\alpha} - \epsilon$,

$$\lim_{n \to \infty} P(|c_{\alpha} - c_{n,\alpha}| \le \epsilon) \ge \lim_{n \to \infty} P(F_n(c_{\alpha} - \epsilon) < \alpha < F_n(c_{\alpha} + \epsilon)) = 1$$
(73)

where the final equality follows from (72) and $\sup_c |F_n(c) - F(c)| = o_p(1)$ by hypothesis.

Lemma 7.5. Let $\Sigma(\beta_0) \equiv E[X_i X_i' \epsilon_i^2], \gamma_0 \equiv E[(\lambda' X_i)^2 X_i \epsilon_i]$ and Assumption 3.1 hold. It then follows that:

$$E[L_n] = -\frac{E[(\lambda' X_i)^3 \epsilon_i^3]}{2\sigma^3 \sqrt{n}} - \frac{E[(\lambda' X_i) X_i' X_i \epsilon_i]}{\sigma \sqrt{n}} + \frac{2\lambda' \Sigma(\beta_0) \gamma_0}{\sigma^3 \sqrt{n}}$$
$$E^*[L_n^*] = -\frac{E[W_i^3]}{2\hat{\sigma}^3 \sqrt{n}} \times \frac{1}{n} \sum_{i=1}^n (\lambda' H_n^{-1} X_i)^3 e_i^3$$

PROOF: We first derive an expression for $E[L_n]$. Note that $E[X_iX'_i] = I$ and $E[X_i\epsilon_i] = 0$ imply:

$$E[\lambda'\Delta_n \frac{1}{\sigma\sqrt{n}} \sum_{i=1}^n X_i \epsilon_i] = \lambda' E[\frac{1}{n} \sum_{i=1}^n (I - X_i X_i') \frac{1}{\sigma\sqrt{n}} \sum_{i=1}^n X_i \epsilon_i] = -\frac{1}{\sigma\sqrt{n}} E[(\lambda' X_i) X_i' X_i \epsilon_i]$$
(74)

due to the i.i.d. assumption. Similarly, exploiting the i.i.d. assumption and $E[(\lambda' X_i)\epsilon_i] = E[\Delta_n] = 0$:

$$E[\frac{1}{2\sigma^3\sqrt{n}}\sum_{i=1}^n (\lambda'X_i)\epsilon_i(\hat{\sigma}_R^2 - \sigma^2)] = E[\frac{1}{2\sigma^3\sqrt{n}}\sum_{i=1}^n (\lambda'X_i)\epsilon_i\{\lambda'(\Sigma_n(\beta_0) - \Sigma(\beta_0))\lambda + 2\lambda'\Delta_n\Sigma(\beta_0)\lambda\}]$$
$$= \frac{1}{2\sigma^3\sqrt{n}}\{E[(\lambda'X_i)^3\epsilon_i^3] - 2E[\epsilon_i(\lambda'X_i)^2X_i']\Sigma(\beta_0)\lambda\}.$$
(75)

The expression for $E[L_n]$ can then be obtained from (74), (75) and by analogous arguments concluding:

$$E\left[\frac{1}{2\sigma^3\sqrt{n}}\sum_{i=1}^n (\lambda'X_i)\epsilon_i \times \frac{2}{n}\sum_{i=1}^n \gamma'_0 X_i\epsilon_i\right] = \frac{\lambda'\Sigma(\beta_0)\gamma_0}{\sigma^3\sqrt{n}} .$$
(76)

In order to compute $E^*[L_n^*]$, observe that W_i independent of (Y_i, X_i) and $E[W_i^2] = 1$ implies that:

$$E^*[L_n^*] = -\frac{1}{2\hat{\sigma}^3}E^*[\frac{\lambda'H_n^{-1}}{\sqrt{n}}\sum_{i=1}^n X_i\epsilon_i^*\frac{1}{n}\sum_{i=1}^n \lambda'H_n^{-1}X_iX_i'H_n^{-1}\lambda e_i^2(W_i^2-1)] = -\frac{E[W_i^3]}{2\hat{\sigma}^3\sqrt{n}} \times \frac{1}{n}\sum_{i=1}^n (\lambda'H_n^{-1}X_i)^3e_i^3$$
(77)

which establishes the second claim of the Lemma. \blacksquare

Lemma 7.6. Under Assumption 3.1, the second moments of L_n and L_n^* satisfy,

$$E[L_n^2] = 1 + O(n^{-1})$$
 $E^*[(L_n^*)^2] = 1 + O_p(n^{-1})$

PROOF: To calculate $E[L_n^2]$, first note that $E[X_iX'_i] = I$, $E[X_i\epsilon_i] = 0$ and direct calculations yield:

$$E[(\lambda'\Delta_n \frac{1}{\sqrt{n\sigma}} \sum_{i=1}^n X_i \epsilon_i)^2] = E[(\frac{\lambda'}{n} \sum_{i=1}^n (I - X_i X_i') \frac{1}{\sqrt{n\sigma}} \sum_{i=1}^n X_i \epsilon_i)^2] = \frac{1}{\sigma^2 n^2} E[(\lambda' (I - X_i X_i') (\sum_{k=1}^n X_k \epsilon_k))^2] + \frac{(n-1)}{\sigma^2 n^2} E[\{\lambda' (I - X_i X_i) \sum_{k=1}^n X_k \epsilon_k\} \{\lambda' (I - X_j X_j') \sum_{k=1}^n X_k \epsilon_k\}] = O(n^{-1}).$$
(78)

Similarly, exploiting the i.i.d. assumption together with $E[X_i\epsilon_i] = 0$ and $E[I - X_iX'_i] = 0$ we obtain:

$$E[(\frac{1}{\sqrt{n\sigma}}\sum_{i=1}^{n}\lambda'X_{i}\epsilon_{i})(\lambda'\Delta_{n}\frac{1}{\sqrt{n\sigma}}\sum_{i=1}^{n}X_{i}\epsilon_{i})] = \frac{1}{n^{2}\sigma^{2}}E[(\sum_{i=1}^{n}\lambda'X_{i}\epsilon_{i})(\lambda'\sum_{i=1}^{n}(I-X_{i}X'_{i}))(\sum_{i=1}^{n}X_{i}\epsilon_{i})]$$
$$= \frac{1}{n\sigma^{2}}E[(\lambda'X_{i}\epsilon_{i})(\lambda'X_{i}\epsilon_{i}-\lambda'X_{i}X'_{i}X_{i}\epsilon_{i})] = O(n^{-1}). \quad (79)$$

Exploiting identical arguments to (78) on the squares of the remaining terms of L_n and the Cauchy-Schwarz inequality and arguments identical to those in (79) to address cross terms arising from expanding the square, it is then straightforward to establish that:

$$E[L_n^2] = E[(\frac{1}{\sigma\sqrt{n}}\sum_{i=1}^n \lambda' X_i \epsilon_i)^2] + O(n^{-1}) = \frac{\lambda' E[X_i X_i' \epsilon_i^2]\lambda}{\sigma^2} + O(n^{-1}) = 1 + O(n^{-1}) .$$
(80)

For notational simplicity, let $R_{in} = \lambda' H_n^{-1} X_i$. To compute $E^*[(L_n^*)^2]$, first note that the i.i.d. assumption together with $E^*[(\epsilon_i^*)^4] = e_i^4 E[W_i^4]$, $E^*[(\epsilon_i^*)^2] = e_i^2$ and $E^*[\epsilon_i^*] = 0$ imply that:

$$\frac{1}{\hat{\sigma}^4 n^2} E^* \left[\left(\sum_{i=1}^n R_{in} \epsilon_i^* \right)^2 \left(\sum_{i=1}^n R_{in}^2 \left\{ (\epsilon_i^*)^2 - e_i^2 \right\} \right) \right] = \frac{1}{\hat{\sigma}^4 n^2} \sum_{i=1}^n R_{in}^4 e_i^4 \left(E[W_i^4] - 1 \right) = O_p(n^{-1}) .$$
(81)

Next, also note that by direct calculations, $\{W_i\}_{i=1}^n$ being i.i.d. and $E[(\epsilon_i^*)^3] = e_i^3 E[W_i^3]$ we may establish:

$$\frac{1}{4\hat{\sigma}^{6}n^{3}}E^{*}[(\sum_{i=1}^{n}R_{in}\epsilon_{i}^{*})^{2}(\sum_{i=1}^{n}R_{in}^{2}\{(\epsilon_{i}^{*})^{2}-e_{i}^{2}\})^{2}] \\
= \frac{1}{4\hat{\sigma}^{6}n^{3}}\{\sum_{i=1}^{n}E^{*}[R_{in}^{2}(\epsilon_{i}^{*})^{2}(\sum_{k=1}^{n}R_{kn}^{2}\{(\epsilon_{k}^{*})^{2}-e_{k}^{2}\})^{2}] + \sum_{i=1}^{n}\sum_{j\neq i}E^{*}[(R_{in}\epsilon_{i}^{*})(R_{jn}\epsilon_{j}^{*})(\sum_{k=1}^{n}R_{kn}^{2}\{(\epsilon_{k}^{*})^{2}-e_{k}^{2}\})^{2}]\} \\
= \frac{1}{4\hat{\sigma}^{6}n^{3}}\{\sum_{i=1}^{n}\sum_{k=1}^{n}R_{in}^{2}R_{kn}^{4}E^{*}[(\epsilon_{i}^{*})^{2}\{(\epsilon_{k}^{*})^{2}-e_{k}^{2}\}^{2}] + 2\sum_{i=1}^{n}\sum_{j\neq i}R_{in}^{3}e_{i}^{3}R_{jn}^{3}e_{j}^{3}(E[W_{i}^{3}])^{2}\}.$$
(82)

Therefore, expanding the square, noting that $n^{-1} \sum_{i=1}^{n} R_{in}^2 e_i^2 = \hat{\sigma}^2$ and exploiting (81) and (82):

$$E^*[(L_n^*)^2] = \frac{1}{n\hat{\sigma}^2} E^*[(\sum_{i=1}^n R_{in}\epsilon_i^*)^2] + O_p(n^{-1}) = 1 + O_p(n^{-1}) , \qquad (83)$$

which establishes the second and final claim of the Lemma. \blacksquare

Lemma 7.7. Let $\Sigma(\beta_0) \equiv E[X_i X_i' \epsilon_i^2], \ \gamma_0 \equiv E[(\lambda' X_i)^2 X_i \epsilon_i]$ and Assumption 3.1 hold. It then follows that:

$$E[(L_n)^3] = -\frac{7}{2\sigma^3\sqrt{n}}E[(\lambda'X_i)^3\epsilon_i^3] - \frac{3}{\sigma\sqrt{n}}E[(\lambda'X_i)(X_i'X_i)\epsilon_i] - \frac{6\lambda'\Sigma(\beta_0)\gamma_0}{\sigma^3\sqrt{n}} + O(n^{-1})$$
$$E^*[(L_n^*)^3] = -\frac{7E[W_i^3]}{2\hat{\sigma}^3\sqrt{n}} \times \frac{1}{n}\sum_{i=1}^n (\lambda'H_n^{-1}X_i)^3e_i^3 + O_p(n^{-1})$$

PROOF: The calculations are cumbersome and for brevity we provide only the essential steps. Define:

$$\Gamma_n \equiv \lambda' \Delta_n \frac{1}{\sigma \sqrt{n}} \sum_{i=1}^n X_i \epsilon_i - \frac{1}{2\sigma^3 \sqrt{n}} \sum_{i=1}^n (\lambda' X_i) \epsilon_i \{ (\hat{\sigma}_R^2 - \sigma^2) - \frac{2}{n} \sum_{i=1}^n \gamma_0' X_i \epsilon_i \} .$$
(84)

Notice that $L_n = \sigma^{-1} n^{-\frac{1}{2}} \lambda' \sum_i X_i \epsilon_i + \Gamma_n$. Under Assumption 3.1(ii), it is possible to establish $E[\Gamma_n^3] = O(n^{-\frac{3}{2}})$ and $E[(n^{-\frac{1}{2}} \sum_i \lambda' X_i \epsilon_i)^3] = O(n^{-\frac{1}{2}})$. Therefore, by direct calculation and Holder's inequality:

$$E[(L_n)^3] = E[(\frac{1}{\sigma\sqrt{n}}\sum_{i=1}^n (\lambda'X_i)\epsilon_i)^3] + 3E[(\frac{1}{\sigma\sqrt{n}}\sum_{i=1}^n (\lambda'X_i)\epsilon_i)^2\Gamma_n] + 3E[(\frac{1}{\sigma\sqrt{n}}\sum_{i=1}^n (\lambda'X_i)\epsilon_i)\Gamma_n^2] + E[\Gamma_n^3]$$

$$\leq E[(\frac{1}{\sigma\sqrt{n}}\sum_{i=1}^n (\lambda'X_i)\epsilon_i)^3] + 3E[(\frac{1}{\sigma\sqrt{n}}\sum_{i=1}^n (\lambda'X_i)\epsilon_i)^2\Gamma_n] + O(n^{-1}).$$
(85)

Hence, we can establish the first claim of the Lemma by analyzing the remaining terms in (85). Note that

$$E[\left(\frac{1}{\sigma\sqrt{n}}\sum_{i=1}^{n}(\lambda'X_{i})\epsilon_{i}\right)^{3}] = \frac{1}{\sigma^{3}\sqrt{n}}E[(\lambda'X_{i})^{3}\epsilon_{i}^{3}], \qquad (86)$$

by the i.i.d. assumption and $E[X_i \epsilon_i] = 0$. Similarly, by direct calculation we can also obtain the expression:

$$E[\left(\frac{1}{\sigma\sqrt{n}}\sum_{i=1}^{n}(\lambda'X_{i})\epsilon_{i}\right)^{2}\frac{\lambda'\Delta_{n}}{\sqrt{n}\sigma}\sum_{i=1}^{n}X_{i}\epsilon_{i}]$$

$$=\frac{1}{\sigma^{3}n^{\frac{5}{2}}}E[\{\sum_{i=1}^{n}(\lambda'X_{i})^{2}\epsilon_{i}^{2}+\sum_{i=1}^{n}(\lambda'X_{i})\epsilon_{i}\sum_{j\neq i}(\lambda'X_{j})\epsilon_{j}\}\sum_{k=1}^{n}\lambda'(I-X_{k}X_{k}')\sum_{l=1}^{n}X_{l}\epsilon_{l}]$$

$$=-\frac{\lambda'\Sigma(\beta_{0})\lambda}{\sigma^{3}\sqrt{n}}E[(\lambda'X_{i})(X_{i}'X_{i})\epsilon_{i}]-\frac{2}{\sigma^{3}\sqrt{n}}E[(\lambda'X_{i})(\gamma_{0}'X_{i})\epsilon_{i}^{2}]+O(n^{-\frac{3}{2}}).$$
(87)

By analogous arguments we can compute the remaining terms in $E[(\sigma^{-1}n^{-\frac{1}{2}}\sum_{i}\lambda' X_i\epsilon_i)^2\Gamma_n]$ and obtain:

$$\frac{1}{2\sigma^5}E[(\frac{1}{\sqrt{n}}\sum_{i=1}^n (\lambda'X_i)\epsilon_i)^3\lambda'\{\Sigma_n(\beta_0) - \Sigma(\beta_0)\}\lambda] = \frac{3\lambda'\Sigma(\beta_0)\lambda}{2\sigma^5\sqrt{n}}E[(\lambda'X_i)^3\epsilon_i^3] + O(n^{-\frac{3}{2}})$$
(88)

$$\frac{1}{\sigma^5} E[(\frac{1}{\sqrt{n}} \sum_{i=1}^n (\lambda' X_i) \epsilon_i)^3 \{\lambda' \Delta_n \Sigma(\beta_0) \lambda\}] = -\frac{3\lambda' \Sigma(\beta_0) \lambda}{\sigma^5 \sqrt{n}} \gamma_0' \Sigma(\beta_0) \lambda + O(n^{-\frac{3}{2}})$$
(89)

$$\frac{1}{\sigma^5} E[(\frac{1}{\sqrt{n}} \sum_{i=1}^n (\lambda' X_i) \epsilon_i)^3 \{ \frac{1}{n} \sum_{i=1}^n \gamma_0' X_i \epsilon_i \}] = \frac{3\lambda' \Sigma(\beta_0) \lambda}{\sigma^5 \sqrt{n}} \lambda' \Sigma(\beta_0) \gamma_0 + O(n^{-\frac{3}{2}}) .$$
(90)

The first claim of the Lemma then follows by combining the results from (85)-(90).

Letting $R_{in} = \lambda' H_n^{-1} X_i$ and employing Assumption 3.1(ii), it can then be shown that:

$$E^*[(\frac{1}{\sqrt{n}}\sum_{i=1}^n R_{in}\epsilon_i^*)^3(\frac{1}{2\hat{\sigma}^3}\{(\hat{\sigma}^*)^2 - \hat{\sigma}^2\})^2] = O_p(n^{-\frac{3}{2}}) \qquad E^*[(\frac{1}{\sqrt{n}}\sum_{i=1}^n R_{in}\epsilon_i^*)^3(\frac{1}{2\hat{\sigma}^3}\{(\hat{\sigma}^*)^2 - \hat{\sigma}^2\})^3] = O_p(n^{-\frac{3}{2}})$$

Therefore, expanding the cube and exploiting that $W_i \perp (Y_i, X_i)$ and $E[(\epsilon_i^*)^k] = E[W_i^k]e_i^k$, it follows that:

$$E^*[(L_n^*)^3] = E^*[(\frac{1}{\sqrt{n}}\sum_{i=1}^n R_{in}\epsilon_i^*)^3\{\frac{1}{\hat{\sigma}^3} - \frac{3((\hat{\sigma}^*)^2 - \hat{\sigma}^2)}{2\hat{\sigma}^5} + \frac{3((\hat{\sigma}^*)^2 - \hat{\sigma}^2)^2}{4\hat{\sigma}^7} - \frac{((\hat{\sigma}^*)^2 - \hat{\sigma}^2)^3}{8\hat{\sigma}^9}\}]$$
$$= \frac{E[W_i^3]}{\hat{\sigma}^3\sqrt{n}} \times \frac{1}{n}\sum_{i=1}^n R_{in}^3 e_i^3 - \frac{3}{2\hat{\sigma}^5}E^*[(\frac{1}{\sqrt{n}}\sum_{i=1}^n R_{in}\epsilon_i^*)^3\{(\hat{\sigma}^*)^2 - \hat{\sigma}^2\}] + O_p(n^{-\frac{3}{2}}).$$
(91)

Moreover, also note that by analogous arguments and direct calculations we further obtain:

$$E^{*}[(\frac{1}{\sqrt{n}}\sum_{i=1}^{n}R_{in}\epsilon_{i}^{*})^{3}\{\frac{3}{2\hat{\sigma}^{5}n}\sum_{i=1}^{n}R_{in}^{2}\{(\epsilon_{i}^{*})^{2}-e_{i}^{2}\}\}]$$

$$=\frac{3}{2\hat{\sigma}^{5}n^{\frac{3}{2}}}\times\frac{1}{n}\sum_{i=1}^{n}R_{in}^{5}E^{*}[(\epsilon_{i}^{*})^{3}\{(\epsilon_{i}^{*})^{2}-e_{i}^{2}\}]+\frac{9}{2\hat{\sigma}^{5}n^{\frac{5}{2}}}E^{*}[\{\sum_{i=1}^{n}R_{in}(\epsilon_{i}^{*})\sum_{j\neq i}R_{jn}^{2}(\epsilon_{j}^{*})^{2}\}\sum_{k=1}^{n}R_{kn}^{2}\{(\epsilon_{k}^{*})^{2}-e_{k}^{2}\}]$$

$$=\frac{9}{2\hat{\sigma}^{5}\sqrt{n}}\times\frac{1}{n}\sum_{i=1}^{n}R_{in}^{2}e_{i}^{2}\times\frac{E[W_{i}^{3}]}{n}\sum_{i=1}^{n}R_{in}^{3}e_{i}^{3}+O_{p}(n^{-\frac{3}{2}}).$$
(92)

The second claim of the Lemma is then established by (91) and (92). \blacksquare

References

- BHATIA, R. (1997): Matrix Analysis. Springer, New York.
- CAMERON, A. C., J. B. GELBACH, AND D. L. MILLER (2008): "Bootstrap-Based Improvements for Inference with Clustered Errors," *Review of Economics and Statistics*, 90, 414–427.
- CAMERON, A. C., AND P. K. TRIVEDI (2005): *Microeconometrics Methods and Applications*. Cambridge University Press, New York.
- CAVALIERE, G., AND A. M. R. TAYLOR (2008): "Bootstrap Unit Root Tests for Time Series with Nonstationary Volatility," *Econometric Theory*, 24, 43–71.
- CHESHER, A. (1984): "Testing for Neglected Heterogeneity," Econometrica, 52, 865–872.
- (1995): "A Mirror Image Invariance for M-estimators," *Econometrica*, 63, 207–211.
- DAVIDSON, R., AND E. FLACHAIRE (2008): "The Wild Bootstrap, Tamed at Last," *Journal of Economet*rics, 146, 162–169.
- DAVIDSON, R., AND J. G. MACKINNON (2008): "Wild Bootstrap Tests for IV Regression," Working paper, Queen's Economics Department.
- EFRON, B. (1979): "Bootstrap Methods: Another Look at the Jacknife," *The Annals of Statistics*, 7(1), 1–26.
- FREEDMAN, D. A. (1981): "Bootstrapping Regression Models," The Annals of Statistics, 9(6), 1218–1228.
- HALL, P. (1992): The Bootstrap and Edgeworth Expansion. Springer-Verlag, New York.
- HARDLE, W., AND E. MAMMEN (1993): "Comparing Nonparametric Versus Parametric Regression Fits," The Annals of Statistics, 21(4), 1926–1947.
- HOROWITZ, J. L. (1994): "Bootstrap-based Critical Values for the Information Matrix Test," Journal of Econometrics, 61, 395–411.
 - (1997): "Bootstrap Methods in Econometrics: Theory and Numerical Performance," in Advances in Economics and Econometrics: Theory and Applications, Seventh World Congress, ed. by D. M. Kreps, and K. F. Wallis, vol. 3. Cambridge University Press.
- (2001): "The Bootstrap," in *Handbook of Econometrics*, ed. by J. J. Heckman, and E. Leamer, vol. 5, chap. 52. Elsevier.

- LIU, R. Y. (1988): "Bootstrap Procedures under some Non-I.I.D. Models," The Annals of Statistics, 16(4), 1696–1708.
- MA, S., AND M. R. KOSOROK (2005): "Robust Semiparametric M-estimation and the Weighted Bootstrap," Journal of Multivariate Analysis, 96, 190–217.
- MAMMEN, E. (1993): "Bootstrap and Wild Bootstrap for High Dimensional Linear Models," *The Annals of Statistics*, 21(1), 255–285.
- MOREIRA, M. J., J. R. PORTER, AND G. A. SUAREZ (2009): "Bootstrap Validity for the Score Test when Instruments are Weak," *Journal of Econometrics*, 149, 52–64.
- NEWEY, W. K. (1985a): "Generalized Method of Moments Specification Testing," Journal of Econometrics, 29, 229–256.
- (1985b): "Maximum Likelihood Specification Testing and Conditional Moment Tests," Econometrica, 53(5), 1047–1070.
- NEWEY, W. K., AND D. L. MCFADDEN (1994): "Large Sample Estimation and Hypothesis Testing," in Handbook of Econometrics, ed. by R. F. Engle, and D. L. McFadden, vol. IV, pp. 2113–2245. Elsevier Science B.V.
- TAUCHEN, G. (1985): "Diagnostic Testing and Evaluation of Maximum Likelihood Models," Journal of Econometrics, 30, 415–443.
- VAN DER VAART, A. (1999): Asymptotic Statistics. Cambridge University Press, New York.
- VAN DER VAART, A. W., AND J. A. WELLNER (1996): Weak Convergence and Empirical Processes: with Applications to Statistics. Springer, New York.
- WHITE, H. (1982): "Maximum Likelihood Estimation of Misspecified Models," *Econometrica*, 50, 1–25.
- (1994): Estimation, Inference and Specification Analysis. Cambridge University Press, New York.
- WOOLDRIDGE, J. M. (2002): Econometric Analysis of Cross Section and Panel Data. The MIT Press, Cambridge.
- WU, C. F. J. (1986): "Jacknife, Bootstrap, and other Resampling Methods in Regression Analysis," Annals of Statistics, 14(4), 1261–1295.
- YOU, J., AND G. CHEN (2006): "Wild Bootstrap Estimation in Partially Linear Models with Heteroscedasticity," *Statistics and Probability Letters*, 76, 340–348.