

## TESTING CONDITIONAL INDEPENDENCE RESTRICTIONS

Oliver Linton<sup>1</sup> and Pedro Gozalo<sup>2</sup>

<sup>1</sup>Faculty of Economics, Cambridge University, Cambridge, UK

<sup>2</sup>Department of Community Health, Brown University, Providence, Rhode Island, USA

□ We propose a nonparametric test of the hypothesis of conditional independence between variables of interest based on a generalization of the empirical distribution function. This hypothesis is of interest both for model specification purposes, parametric and semiparametric, and for nonmodel-based testing of economic hypotheses. We allow for both discrete variables and estimated parameters. The asymptotic null distribution of the test statistic is a functional of a Gaussian process. A bootstrap procedure is proposed for calculating the critical values. Our test has power against alternatives at distance  $n^{-1/2}$  from the null; this result holding independently of dimension. Monte Carlo simulations provide evidence on size and power.

**Keywords** Conditional independence; Empirical distribution; Independence; Nonparametric; Smooth bootstrap; Test.

**JEL Classification** C12; C14; C15; C52.

### 1. INTRODUCTION

Les Godfrey has made some important contributions to hypothesis testing in econometrics. He wrote a series of papers tackling central issues in misspecification testing, culminating in his Econometric Society monograph, Godfrey (1988). This paper is about testing for conditional independence and related hypotheses.

We investigate the application of the hypothesis of conditional independence in econometrics. Let  $Y$ ,  $X$ , and  $Z$  be random variables; following Dawid (1979), we write

$$Y \perp\!\!\!\perp X \mid Z \tag{1}$$

to denote that  $Y$  is independent of  $X$  given  $Z$ . This assumption is related to the more commonly treated hypothesis that  $Y \perp\!\!\!\perp X$  ( $Y$  is

independent of  $X$ ), in that it imposes an infinite number of restrictions on the joint distribution.<sup>1</sup> These assumptions are stronger than the mean independence conditions usually employed in regression analysis: for example, that  $Y$  is mean independent of  $X$ , i.e., that  $E(Y|X) = 0$ , or that  $Y$  is mean independent of  $X$  given  $Z$ , i.e., that  $E(Y|X, Z) = E(Y|Z)$ . We now give two concrete reasons for interest in the conditional independence hypothesis.

A large literature now exists on testing parametric regression models against general alternatives, see for example Bierens and Ploberger (1997), Hong and White (1995), and Fan and Li (1996). These amount to testing a null hypothesis of mean independence of the regressors from some parametrically defined residual. Andrews (1997) extends this to testing the null hypothesis of a parametric conditional distribution against a general nonparametric alternative. There are many nonparametric tests of independence for continuous data, starting with Hoeffding (1948), including those based on empirical distribution functions such as Blum et al. (1961) and Skaug and Tjøstheim (1993) and Delgado (1996), and those based on smoothing methods like Robinson (1991) and Zheng (1997). However, there do not appear to be any fully nonparametric tests for conditional independence.<sup>2</sup>

Our first application concerns the evaluation of the impact of a social program such as a job training program. Let  $D$  denote the dummy variable such that  $D = 1$  when the person receives treatment (participates), and  $D = 0$  if not treated. Let  $Y_1$  and  $Y_0$  be the outcomes associated with the participation values  $D = 1$  and  $D = 0$ , respectively, and let  $X$  denote individual observed characteristics. A common measure of the impact of partial coverage programs, such as job training programs, is the *average treatment effects on the treated*

$$E(Y_1 - Y_0 | D = 1, X) = E(Y_1 | D = 1, X) - E(Y_0 | D = 1, X).$$

If it exceeds the appropriate measure of cost, the program should be maintained, see for example Heckman et al. (1998). The main problem in the estimation of  $E(Y_1 - Y_0 | D = 1, X)$  is that the second term,  $E(Y_0 | D = 1, X)$ , cannot be observed. Replacing it with the observable average outcomes of nonparticipants  $E(Y_0 | D = 0, X)$  leads to the presence of the self-selection bias term  $B(X) = E(Y_0 | D = 1, X) - E(Y_0 | D = 0, X)$ . One can try to characterize  $B(X)$  by using a control group (people

<sup>1</sup>See Phillips (1988) for a discussion of the difference between independence and conditional independence.

<sup>2</sup>When  $X, Y, Z$  are jointly normal with mean  $\mu$  and covariance matrix  $\Sigma = (\sigma_{ij})$ ,  $Y \perp\!\!\!\perp X$  is equivalent to  $\sigma_{YX} = 0$ , while  $Y \perp\!\!\!\perp X | Z$  is equivalent to  $\sigma^{YX} = 0$ , where the concentration matrix  $\Sigma^{-1} = (\sigma^{ij})$ . In this case, there are simple parametric tests of both independence and conditional independence. For categorical data, there are also numerous tests of independence and conditional independence, see (Agresti, 1990, p. 228).

87 that applied to participate in the program but were randomly denied  
 88 access to program) to estimate  $E(Y_0 | D = 1, X)$  and a comparison group  
 89 (eligible non-participants) to estimate  $E(Y_0 | D = 0, X)$ . The potentially  
 90 high dimension of  $X$ , however, makes direct nonparametric estimation  
 91 of  $B(X)$  problematic. Instead, the most common approach in this  
 92 literature is to use the probability of program participation given observed  
 93 characteristics,  $\Pr(D = 1 | X) = P(X)$ , also referred to as the *propensity score*,  
 94 to characterize the bias. The important role of the propensity score is often  
 95 motivated by the results of Rosenbaum and Rubin (1983). They show that  
 96 if there exists an  $X$  such that

$$97 \quad (Y_1, Y_0) \perp\!\!\!\perp D | X,$$

98  
 99 and  $0 < \Pr(D = 1 | X) < 1$  for all  $X$ , then, conditioning on  $X$  is equivalent  
 100 to conditioning on the univariate index  $P(X)$ : specifically,  $E(Y_0 | D, X) =$   
 101  $E(Y_0 | X) = E(Y_0 | P(X))$ , so that

$$102 \quad B(X) = B\{P(X)\} = 0, \quad \text{for all } X.$$

103  
 104 This index sufficiency restriction is essentially the conditional  
 105 independence restriction that treatment  $D$  is ignorable given the  
 106 observables  $X$ .

107 Our second application concerns semiparametric model specification.  
 108 Consider the semiparametric binary choice model

$$109 \quad Y = \begin{cases} 1 & \text{if } \beta^\top X \geq \varepsilon \\ 0 & \text{otherwise,} \end{cases} \quad (2)$$

110  
 111 where  $\beta$  is a vector of unknown parameters and  $\varepsilon$  is an unobservable  
 112 stochastic error term. The semiparametric literature divides into two broad  
 113 categories according to whether  $\varepsilon$  is assumed to be independent of  $X$  or  
 114 only mean (actually median) independent (see the recent review papers  
 115 by Manski, 1994, and Powell, 1994, for discussion). In the latter case,  
 116 Manski (1975) developed the maximum score procedure for estimating  $\beta$   
 117 which was subsequently shown by Kim and Pollard (1990) to converge, on  
 118 centering, at rate  $n^{1/3}$  to a non-normal limit. Horowitz (1992) suggested  
 119 a smoothed version of the maximum score procedure obtaining, under  
 120 smoothness conditions, asymptotic normality at a rate faster than  $n^{1/3}$ ,  
 121 but still less than  $n^{1/2}$ . In fact, Chamberlain (1987) showed that the  
 122 semiparametric information in this model is zero: i.e., that one cannot  
 123 estimate  $\beta$  in this model at the usual  $n^{1/2}$  rate. By contrast, in the case that  
 124  $\varepsilon$  is assumed to be independent of  $X$ , it is possible to estimate  $\beta$  with the  
 125  $n^{1/2}$  rate of convergence; the Ichimura (1993) and Klein and Spady (1993)  
 126 procedures both achieve this. If

$$127 \quad \varepsilon \perp\!\!\!\perp X, \quad (3)$$

128  
 129

130 then

$$131 \quad Y \perp\!\!\!\perp X \mid \beta^\top X, \quad (4)$$

132  
133 so that the independence condition on the unobservable random variable  
134  $\varepsilon$  implies the conditional independence of the observable quantities.  
135 The conditional independence restriction (4) is weaker than (3), which  
136 suggests that an alternative way of specifying (2) would be to assume  
137 the weaker condition (4). Note also that when  $Y$  is binary, for example,  
138 independence is equivalent to mean independence.

139 The above discussion holds much more generally in the class of  
140 transformation models considered in Han (1987),  $Y = D \circ F(\beta^\top X, \varepsilon)$ ,  
141 where  $D$  is a monotonic function, while  $F$  is monotonic in each of its  
142 arguments. This includes transformation models, binary choice, duration  
143 models and censored regression. We can also extend the discussion to  
144 include panel data models where the independence assumption is even  
145 more crucial. Suppose that  $Y_t = F(\alpha + \beta^\top X_t + \varepsilon_t)$ ,  $t = 1, \dots, \top$ , where the  
146 composite random term  $\alpha + \varepsilon_t$  is independent of  $\beta^\top X_1, \dots, \beta^\top X_\top$ . Then

$$147 \quad Y_1, \dots, Y_\top \perp\!\!\!\perp X_1, \dots, X_\top \mid \beta^\top X_1, \dots, \beta^\top X_\top.$$

148  
149 In this case, it has only been possible to consistently estimate  $\beta$  under the  
150 independence assumption, see (Powell, 1994, p. 2513).

151 Let  $\sigma(Y)$ ,  $\sigma(X)$ , and  $\sigma(Z)$  be, respectively, the sigma algebras generated  
152 by the random variables  $Y$ ,  $X$ , and  $Z$  defined on the same probability space  
153  $(\Omega, \mathcal{F}, P)$ . We say that  $\sigma(Y)$  is independent of  $\sigma(X)$  given  $\sigma(Z)$ , if for all  
154  $A \in \sigma(Y)$  and  $B \in \sigma(X)$ , we have

$$155 \quad P(A \cap B \mid \sigma(Z)) = P(A \mid \sigma(Z)) P(B \mid \sigma(Z)),$$

156 where  $P(A \mid \sigma(Z))$  is the conditional probability of the event  $A$  given  $\sigma(Z)$ .

157 When  $Y$ ,  $X$ ,  $Z$  are all continuous random variables with joint  
158 density function  $f(y, x, z)$  with respect to Lebesgue measure on  $\mathbb{R}^d$ , then  
159 conditional independence of  $Y, X$  given  $Z = z$  becomes  
160

$$161 \quad \frac{f(y, x, z)}{f(z)} = \frac{f(y, z)}{f(z)} \frac{f(x, z)}{f(z)},$$

162 for all  $y, x, z$  in the support of  $Y, X, Z$  or, equivalently,

$$163 \quad f(z)f(y, x, z) = f(y, z)f(x, z).^3 \quad (5)$$

164  
165  
166  
167  
168  
169  
170 <sup>3</sup>Similarly, when  $Y, X, Z$  are all discrete random variables with joint probability mass function  
171  $p(y, x, z)$ , then conditional independence implies that  $p(z)p(y, x, z) = p(y, z)p(x, z)$  for all  $y, x, z$ . In  
172 the sequel we will use  $f(\cdot)$  to denote both density and probability mass functions to simplify notation  
and let the continuous or discrete nature of the random variables determine its interpretation.

173 A consistent nonparametric test of conditional independence based  
 174 on the relationship between joint and marginal density functions given  
 175 by (5) would appear to require (except in the special case where all  
 176 variables are discrete) nonparametric smoothing estimation of potentially  
 177 high dimensional density functions  $f(y, x, z)$ . The main disadvantage  
 178 of this approach is the slower rate of growth of the noncentrality of  
 179 the test statistic which results in lower power against local alternatives.  
 180 In particular, a smoothing based test would not be able to detect  
 181 alternatives at distance  $n^{-1/2}$  from the null detected by parametric and  
 182 some nonparametric test statistics.

183 Our strategy to avoid this dimensionality problem is to verify the  
 184 relationship (5) over subsets (with positive Lebesgue measure) in the  
 185 support of  $Y, X, Z$  rather than at individual values, thereby replacing  
 186 densities with distribution functions. Expression (5) is extended to  
 187

$$188 \quad P(C)P(A \cap B \cap C) = P(A \cap C)P(B \cap C), \quad (6)$$

189  
 190 for all subsets  $A, B$  in the support of  $Y, X$ , and  $C$  a subset in the open  
 191 support of  $Z$ .<sup>4</sup>

192 We now come to the contribution of this paper. We first provide a  
 193 nonparametric test of (1) based on the empirical distribution functions  
 194 representation of condition (6). A key question addressed in the sequel  
 195 is how to choose the subsets  $A, B$ , and  $C$  so as to apply this principle  
 196 to mixed continuous and discrete data. We extend the usual treatment  
 197 based on quadrants to a general class of rectangles suitable for these  
 198 types of data. In our first example the central hypothesis is simple  
 199 (although it is about a function), while the second case was composite and  
 200 there were unknown nuisance parameters involved; our test is therefore  
 201 devised to allow for estimated parameters.<sup>5</sup> A key technical issue we  
 202 face in this paper is the verification of the stochastic equicontinuity  
 203 property for processes involving both discontinuous indicator functions  
 204 of rectangles and nonlinear functions of the underlying parameters. Our  
 205 test statistic is easy to compute and to analyze. Its asymptotic distribution  
 206 is a functional of a Gaussian process whose quantiles are found by the  
 207 bootstrap.

208 The use of (6) over (5) has the advantage of increased power over  
 209 local alternatives for any dimension, but it introduces a difficulty. As the  
 210 conditioning set  $C$  in (6) deviates from the singleton set  $C = \{z\}$ , we will  
 211 see in the next section that the exact equivalence between conditional  
 212

213 <sup>4</sup>We need to exclude  $C$  from being the entire support since (6) would reduce to  $P(A \cap B) =$   
 214  $P(A)P(B)$  which implies  $Y \perp\!\!\!\perp X$  but it is well known that conditional independence does not  
 215 imply independence (Chow and Teicher, 1998, p. 221).

<sup>5</sup>Note that (3) itself is not directly testable because one cannot estimate  $\varepsilon$  consistently.

216 independence (1) and expression (5) is only partially preserved by (6). In  
 217 particular we will see that while (6) implies (1), the reverse is not always  
 218 true.<sup>6</sup> In other words, the set of distributions  $P$  for which (6) holds is  
 219 a subset of the set of distributions for which conditional independence  
 220 (1) holds. The practical implication of this is that a test of conditional  
 221 independence based on (6) may have a noncentrality different from zero  
 222 not only in cases of conditional dependence but also in some cases where  
 223 conditional independence holds. However, our proposed bootstrap is able  
 224 to correct this by automatically adjusting for the noncentrality present  
 225 under the null of conditional independence. Our simulations illustrate the  
 226 performance of our test in such a case where (6) does not hold exactly  
 227 under the null.

228 We remark that since this paper was first written (Linton and  
 229 Gozalo, 1995), a number of authors have proposed tests of conditional  
 230 independence, see for example Fernandes and Flores (2001), Su and  
 231 White (2008), and Song (2009). These authors all use smoothing  
 232 techniques (i.e., kernels) and obtain tests that are consistent of the full  
 233 null hypothesis against all alternatives in a large class.

234 The rest of the paper is organized as follows. In Section 2 we introduce  
 235 our test statistics, in Section 3 we make our assumptions and present the  
 236 limiting distributions (we use an i.i.d. setup suitable for cross-sectional  
 237 data); Section 4 introduces a bootstrap method for obtaining critical  
 238 values. We provide a small simulation experiment in Section 5.

239 **Notation.** We use  $\mathbf{1}(\cdot)$  for the indicator function, i.e.,  $\mathbf{1}(A) = 1$  if event  $A$   
 240 occurs and  $\mathbf{1}(A) = 0$  otherwise. Let  $\Rightarrow$  and  $\rightarrow_p$  denote weak convergence  
 241 of probability measures and convergence in probability, respectively; all  
 242 limits are taken as sample size  $n \rightarrow \infty$ .

243

## 244 2. TEST STATISTICS

245

246 The underlying population is the random vector  $U \in \mathbb{R}^q$  from which  
 247 we observe an independent and identically distributed (i.i.d.) sample  
 248  $\{U_i\}_{i=1}^n$ . Of interest are certain residual [or index] functions computed  
 249 from  $U$ , that is  $V(U; \theta) = (Y(U; \theta), X(U; \theta), Z(U; \theta)) \in \mathbb{R}^d$ , where the  
 250 parameter  $\theta \in \Theta \subset \mathbb{R}^p$  and  $d = l + m + k$ . The null hypothesis to be tested  
 251 is that  $Y(U; \theta^0)$  and  $X(U; \theta^0)$  are independent conditional on  $Z(U; \theta^0)$  for  
 252 some particular  $\theta^0$  whose value is not known.

253 We shall base our test on the equality (6) for some (separating)  
 254 class of subsets and replace the population probability measure by an  
 255 empirical measure. Most previous work has been based on quadrants,  
 256

257

258 <sup>6</sup>This was pointed out to us by Jon Wellner to whom we are grateful.

259 i.e., the empirical distribution function.<sup>7</sup> These sets apparently work  
 260 well for continuous data but, as currently applied, are unsuited for  
 261 discrete data as the following example illustrates. Suppose that  $(Y, X)$  are  
 262 binary with  $\Pr(Y = 1, X = 0) = \Pr(Y = 0, X = 1) = 1/2$ , then  $Y$  and  $X$  are  
 263 perfectly dependent with the same marginals  $\Pr(Y = 1) = \Pr(X = 1) =$   
 264  $1/2$ . Unfortunately, quadrants located at the observations will not uncover  
 265 this dependence; in fact,  $\Pr(Y \leq 1, X \leq 0) = \Pr(Y \leq 1) \Pr(X \leq 0)$ ,  $\Pr(Y \leq$   
 266  $1, X \leq 1) = \Pr(Y \leq 1) \Pr(X \leq 1)$ , etc. In view of this, we consider the  
 267 more general class of all rectangular subsets of  $\mathbb{R}^d$  of a certain (possibly  
 268 zero) width. Let  $a_\alpha, b_\alpha, \alpha = 1, \dots, d$  be given nonnegative numbers, possibly  
 269 infinite, and let

$$270 \mathfrak{B}_\alpha(v_\alpha) = [v_\alpha - a_\alpha, v_\alpha + b_\alpha]$$

271 be a rectangle in the component  $\alpha$ . Let also  $\mathfrak{B}(v) = \times_{\alpha=1}^d \mathfrak{B}_\alpha(v_\alpha)$ , and  
 272  $\mathfrak{B}(y)$ ,  $\mathfrak{B}(x)$ , and  $\mathfrak{B}(z)$  be the rectangles obtained by intersecting the  
 273 corresponding intervals.<sup>8</sup> Then let

$$274 F(v | \theta) = P(V(U; \theta) \in \mathfrak{B}(v))$$

275 be the joint rectangular distribution function of  $V$ , and denote the  
 276 corresponding probability functions of  $(Y, Z)$ ,  $(X, Z)$ , and  $Z$  by  $G(y, z | \theta)$ ,  
 277  $H(x, z | \theta)$ , and  $L(z | \theta)$ , respectively; also, let  $F(v) = F(v | \theta^0)$ ,  $G(y, z) =$   
 278  $G(y, z | \theta^0)$ ,  $H(x, z) = H(x, z | \theta^0)$ , and  $L(z) = L(z | \theta^0)$ . When  $b_\alpha = 0$  and  
 279  $a_\alpha = \infty$ , these functions correspond to the usual distribution functions.  
 280 For discrete variables, events of the form  $\{V \leq v\}$  are not a wise choice,  
 281 as discussed above, and would give zero power against some alternatives,  
 282 see Joag-Dev (1984). For these variables we shall take  $a_\alpha = b_\alpha = 0$ .<sup>9</sup> For  
 283 continuously distributed data we take  $a_\alpha > 0$  and  $b_\alpha \geq 0$ , except we also  
 284 rule out the case  $a_\alpha = b_\alpha = \infty$ . Note that the choice of rectangles can  
 285 vary with location, so that a data series with both continuous and discrete  
 286 components can be accommodated by choosing atomic rectangles at  
 287 points of discreteness but intervals elsewhere. There is, therefore, wide  
 288 latitude in choosing which rectangles to use for a given application. We  
 289 discuss this further in section 5 below.

290 <sup>7</sup>There has also been some work using multivariate half spaces, i.e., hyperplanes, see Beran  
 291 and Millar (1986).

292 <sup>8</sup>Formally speaking, the sets we examine are of the form  $A = \{V \in \mathfrak{B}(y) \times (-\infty, \infty)^{m+k}\} \in \mathbb{R}^d$ ,  
 293  $B = \{V \in \mathfrak{B}(x) \times (-\infty, \infty)^{l+k}\} \in \mathbb{R}^d$ , and  $C = \{V \in \mathfrak{B}(z) \times (-\infty, \infty)^{l+m}\} \in \mathbb{R}^d$ . Then, for example  
 294  $A \cap B \cap C = \{V \in \mathfrak{B}(v)\}$ .

295 <sup>9</sup>In our discrete example, the dependence is uncovered by this choice of events, since clearly  
 296  $\Pr(Y = 1, X = 0) \neq \Pr(Y = 1) \Pr(X = 0)$ .

302 Letting  $A(v|\theta, P) = L(z|\theta)F(v|\theta) - G(y, z|\theta)H(x, z|\theta)$ , the Eq. (6) is  
 303 equivalent to.<sup>10</sup>

$$304 \quad A(v|\theta^0, P) = 0, \quad \text{for all } v \in \mathbb{R}^d, \quad \text{some } \theta^0 \in \Theta \subset \mathbb{R}^p. \quad (7)$$

305  
 306 A number of functionals of  $A$  can be used to test conditional independence;  
 307 specifically, the Kolmogorov-Smirnov  $KS = \sup_v |A(v|\theta^0, P)|$  and the  
 308 Cramér von-Mises  $CM = \int A^2(v|\theta^0, P)d\mu(v)$  for some measure  $\mu(\cdot)$  (for  
 309 example  $\mu = F$ ). Shorack and Wellner (1986) discuss a number of  
 310 alternative test functionals in a variety of contexts. The quantities  $KS$  and  
 311  $CM$  provide a general measure of the amount of conditional independence  
 312 there is.

313 Let now  $P^{XYZ}$  be the joint distribution of  $(X, Y, Z)$ , and write  $P^{XYZ} =$   
 314  $P^{XY|Z} \cdot P^Z$ , where  $P^{XY|Z}$  and  $P^Z$  are the distribution of  $(X, Y)$  conditional  
 315 on  $Z$ , and the marginal distribution of  $Z$ , respectively. Similarly, let  
 316  $P^{X|Z}$  and  $P^{Y|Z}$  denote the distributions of  $X$  and  $Y$  conditional on  $Z$ ,  
 317 respectively. Then we define  $P_{CI}^{XYZ}$  as the joint distribution of  $(X, Y, Z)$   
 318 generated by  $P^{XYZ}$  subject to the null hypothesis  $P^{XY|Z} = P^{X|Z} \cdot P^{Y|Z}$  of  
 319 independence between  $X$  and  $Y$  conditional on  $Z$ . That is,

$$320 \quad P_{CI}^{XYZ} = P^{X|Z} \cdot P^{Y|Z} \cdot P^Z.$$

321 This is just a functional of the original probability measure  $P^{XYZ}$ . We can  
 322 then express the null hypothesis of conditional independence as

$$323 \quad \mathbf{H}_0 : P^{XYZ}(v|\theta^0) = P_{CI}^{XYZ}(v|\theta^0), \quad \text{for all } v \in \mathbb{R}^d : \lambda(\mathfrak{B}(z)) = 0,$$

$$324 \quad \text{some } \theta^0 \in \Theta \subset \mathbb{R}^p,$$

325 where  $\lambda(\cdot)$  denotes Lebesgue measure. Thus, for some parameter value  $\theta^0$ ,  
 326  $P^{XYZ}(v|\theta^0)$  and  $P_{CI}^{XYZ}(v|\theta^0)$  agree under  $\mathbf{H}_0$  over all rectangles  $\mathfrak{B}(v)$  of zero  
 327 width in the conditional variable  $Z$ . The alternative hypothesis  $\mathbf{H}_A$  is the  
 328 negation of this for each  $\theta \in \Theta \subset \mathbb{R}^p$ .

329 Consider now how (7), which is evaluated over all rectangles  $\mathfrak{B}(v)$ ,  
 330 relates to  $\mathbf{H}_0$ , which is defined over all rectangles  $\mathfrak{B}(v)$  with single-valued  
 331 conditioning sets  $\mathfrak{B}(z) = \{z\}$ . Suppose that (7) holds. If  $(X, Y, Z)$  has a  
 332 discrete distribution, (7) holding over all rectangles  $\mathfrak{B}(v)$  implies that it  
 333 also holds for single-valued sets  $\mathfrak{B}(v) = \{v\}$ . Hence (5) holds and this  
 334 implies  $\mathbf{H}_0$ . Similarly, when  $(X, Y, Z)$  has a continuous distribution, a  
 335 limiting argument for  $\mathfrak{B}(v)$  shrinking towards a single-valued set  $\{v\}$  shows  
 336 that (5) holds and this again implies  $\mathbf{H}_0$ . Therefore, (7) implies  $\mathbf{H}_0$ .

337  
 338  
 339  
 340  
 341  
 342  
 343  
 344 <sup>10</sup>This is because the class of rectangles of a given width separates probability measures. That  
 is, if two probability measures  $P_1$  and  $P_2$  agree on the class of all rectangles of given width, then  
 they agree on all Borel sets.



Suppose now that  $\mathbf{H}_0$  holds. The first element of  $A(v | \theta, P)$  (omitting the dependence on  $\theta$ ) is

$$\begin{aligned} L(z)F(v) &= \int_{\mathfrak{B}(z)} f(z) dz \int_{\mathfrak{B}(z)} \int_{\mathfrak{B}(y)} \int_{\mathfrak{B}(x)} f(y, x, z) dy dx dz \\ &= \int_{\mathfrak{B}(z)} f(z) dz \int_{\mathfrak{B}(z)} \left[ \int_{\mathfrak{B}(y)} \int_{\mathfrak{B}(x)} f(y, x | z) dy dx \right] f(z) dz \\ &= \int_{\mathfrak{B}(z)} f(z) dz \int_{\mathfrak{B}(z)} \left[ \int_{\mathfrak{B}(y)} f(y | z) dy \right] \left[ \int_{\mathfrak{B}(x)} f(x | z) dx \right] f(z) dz, \end{aligned}$$

where the last equality follows by  $\mathbf{H}_0$ . But, in general,

$$\begin{aligned} &\int_{\mathfrak{B}(z)} f(z) dz \int_{\mathfrak{B}(z)} \left[ \int_{\mathfrak{B}(y)} f(y | z) dy \right] \left[ \int_{\mathfrak{B}(x)} f(x | z) dx \right] f(z) dz \\ &\neq \int_{\mathfrak{B}(z)} \left[ \int_{\mathfrak{B}(y)} f(y | z) dy \right] f(z) dz \int_{\mathfrak{B}(z)} \left[ \int_{\mathfrak{B}(x)} f(x | z) dy \right] f(z) dz \\ &= G(y, z)H(x, z), \end{aligned}$$

whenever  $\mathfrak{B}(z)$  is not a single-valued set. Therefore, in general,  $\mathbf{H}_0$  does not imply  $A(v | \theta^0, P) = 0$ .<sup>11</sup> In other words, the set of distributions  $P$  for which (7) holds is a subset of the set of distributions for which conditional independence  $\mathbf{H}_0$  holds.

Our choice of functional based on multivalued conditioning sets is therefore nonstandard in the sense that it is zero only when additional restrictions are true, and therefore can have noncentrality different from zero for some values of the null hypothesis. However, we show that this noncentrality can be estimated and that our test statistics are consistent for the null hypothesis of interest. Specifically, we will see that  $A_n(v | \theta)$  can be decomposed as

$$A_n(v | \theta) = A(v | \theta, P) + \Delta_n(v | \theta) + O_p(n^{-1}),$$

where  $\Delta_n(v | \theta)$  is a zero-mean stochastic term that determines the limiting distribution of our test. Let  $A(v | \theta, P_{CI})$  denote  $A(v | \theta, P)$  when  $P$  is replaced by the distribution it generates under  $\mathbf{H}_0$ ,  $P_{CI}$ , and let  $\mu(v | \theta, P) = A(v | \theta, P) - A(v | \theta, P_{CI})$ . Then, adding and subtracting  $A(v | \theta, P_{CI})$ , we can express  $A_n(v | \theta)$  as

$$\begin{aligned} A_n(v | \theta) &= \mu(v | \theta, P) + A(v | \theta, P_{CI}) + \Delta_n(v | \theta) + O_p(n^{-1}) \\ &= \mu(v | \theta, P) + D_n(v | \theta) + O_p(n^{-1}), \end{aligned}$$

<sup>11</sup>A case where  $\mathbf{H}_0$  would imply  $A(v | \theta^0, P) = 0$  is when  $f(z)$  is constant for all  $z$  in  $\mathfrak{B}(z)$  and  $P(\mathfrak{B}(y)|z)$  or  $P(\mathfrak{B}(x)|z)$  are constant for all  $z$  in  $\mathfrak{B}(z)$ .

say. Note that  $\mathbf{H}_0$  is true if and only if the noncentrality  $\mu(v|\theta^0, P) = 0$  for all  $v$ , while  $\mu(v|\theta, P) > 0$  under  $\mathbf{H}_A$  for each  $\theta$  and some  $v$  (which may depend on  $\theta$ ). Typical test statistics based on empirical distributions require an approximation to the case-dependent distribution of the zero-mean term  $\Delta_n(v|\theta)$ . In Section 4 below, we propose a bootstrap that approximates the distribution of the (possibly nonzero mean) term  $D_n(v|\theta)$ .

To implement the test suppose that there exists estimates  $\hat{\theta}$  of  $\theta^0$  that are root- $n$  consistent under  $\mathbf{H}_0$ , and replace  $A(v|\theta^0, P)$  by its empirical analogue  $A(v|\hat{\theta}, P_n)$  or  $A_n = L_n F_n - G_n H_n$ , suppressing dependence on  $\hat{\theta}$  and  $P_n$ , where

$$L_n(z) = n^{-1} \sum_{j=1}^n \mathbf{1}\{\widehat{Z}_j \in \mathfrak{B}(z)\}; \quad G_n(y, z) = n^{-1} \sum_{j=1}^n \mathbf{1}\{(\widehat{Y}_j, \widehat{Z}_j) \in \mathfrak{B}(y, z)\}$$

$$F_n(v) = n^{-1} \sum_{j=1}^n \mathbf{1}\{\widehat{V}_j \in \mathfrak{B}(v)\}; \quad H_n(x, z) = n^{-1} \sum_{j=1}^n \mathbf{1}\{(\widehat{X}_j, \widehat{Z}_j) \in \mathfrak{B}(x, z)\}$$

with  $\widehat{Z}_j = Z_j(U_j; \hat{\theta})$ ,  $\widehat{X}_j = X_j(U_j; \hat{\theta})$ , and  $\widehat{Y}_j = Y_j(U_j; \hat{\theta})$ . We then estimate *CM* and *KS* by

$$CM_n = n^{-1} \sum_{i=1}^n A_n^2(\widehat{V}_i); \quad KS_n = \max_{1 \leq i \leq n} |A_n(\widehat{V}_i)|. \quad (8)$$

Note that a maximum is used in  $KS_n$  instead of the usual supremum. This particular version of the Kolmogorov–Smirnov statistic has recently been suggested by Andrews (1997) in another context. It has the advantage of requiring only  $O(n^2)$  computations. Computation of both tests can be completely vectorized.<sup>12</sup> Given critical values  $\hat{c}_\alpha$ , our level- $\alpha$  test based on either test statistic  $I_n = CM_n, KS_n$  is then

$$\text{Reject if: } I_n > \hat{c}_\alpha. \quad (9)$$

<sup>12</sup>Although  $CM_n$  and  $KS_n$  are desirable from a computational point of view, they can have poor (small sample) performance for large  $d$ , because the evaluation points  $\widehat{V}_i$  are not representative enough. In practice, the following statistics may work better with large  $d$  and small  $n$ ,

$$CM_n^f = m^{-1} \sum_{i=1}^m A_n^2(t_i); \quad KS_n^f = \max_{1 \leq i \leq m} |A_n(t_i)|,$$

where  $\{t_i; i = 1, \dots, m\}$  is a fixed or random grid of points. The number of evaluation points,  $m$ , is under the control of the practitioner, but should increase with sample size, see Beran and Millar (1986) for justification of this device. In the simulations presented in Section 5, we used a random grid of points based on the observations.

431 The bootstrap critical values  $\hat{c}_\alpha$  have the property that  $\Pr[I_n > \hat{c}_\alpha | \mathbf{H}_0] \rightarrow \alpha$   
 432 and  $\Pr[I_n > \hat{c}_\alpha | \mathbf{H}_A] \rightarrow 1$ .

### 434 3. ASYMPTOTIC PROPERTIES OF THE TEST

435 We now establish the asymptotic properties of  $CM_n$  and  $KS_n$ . The main  
 436 technical difficulty here is that  $V$  is a nonlinear function of both the  
 437 data and the parameters and occurs inside an indicator which is itself a  
 438 non-smooth function. Empirical processes with estimated parameters were  
 439 first studied by Durbin (1973). There followed a number of papers that  
 440 extended his results to a variety of situations, including nonlinear and  
 441 dependent data. See the books Van der Vaart and Wellner (1996) and  
 442 Shorack and Wellner (1986) for many references. Some recent works of  
 443 special interest to econometricians include Bai (1994), Andrews (1997),  
 444 and Koul (1996).

445 First of all we introduce some notation. Define:  $\delta_1(\cdot, v | \theta) = \mathbf{1}\{V(\cdot, \theta) \in$   
 446  $\mathfrak{B}(v)\} - F(v)$ ,  $\delta_2(\cdot, v | \theta) = \mathbf{1}\{Z(\cdot, \theta) \in \mathfrak{B}(z)\} - L(z)$ ,  $\delta_3(\cdot, v | \theta) = \mathbf{1}\{(X(\cdot, \theta),$   
 447  $Z(\cdot, \theta)) \in \mathfrak{B}(x, z)\} - H(x, z)$ ,

448  $\delta_4(\cdot, v | \theta) = \mathbf{1}\{(Y(\cdot, \theta), Z(\cdot, \theta)) \in \mathfrak{B}(y, z)\} - G(y, z)$ , and

$$449 \delta_0(\cdot, v | \theta) = L(z)\delta_1(\cdot, v | \theta) + F(v)\delta_2(\cdot, v | \theta) - G(y, z)\delta_3(\cdot, v | \theta) \\ 450 - H(x, z)\delta_4(\cdot, v | \theta).$$

451 The process  $\Delta_n(v | \theta) = n^{-1} \sum_{i=1}^n \delta_0(U_i, v | \theta)$  is an approximation to  
 452  $A_n(v | \theta)$  in the sense that

$$453 A_n(v | \theta) - A(v | \theta) = \Delta_n(v | \theta) + O_p(n^{-1}), \quad (10)$$

454 where the error is uniform in both  $v$  and  $\theta$ . If  $\theta^0$  were known, then the  
 455 asymptotic distribution of the empirical process  $\Delta_n(v | \theta^0)$  determines the  
 456 limiting distribution of our test. When  $\theta^0$  is replaced by an estimate we must  
 457 also take account of its variation. For this we must calculate how  $\Delta_n(v | \theta)$   
 458 changes with movements in  $\theta$ . Letting  $\Delta_j(u | \theta) = E_{\theta^0}[\delta_j(U, V(u, \theta) | \theta)]$  for  
 459  $j = 0, \dots, 4$ , we have

$$460 \Delta_n(V(u, \theta) | \theta) = \Delta_n(V(u, \theta^0) | \theta^0) + \frac{\partial \Delta_0(u | \theta^0)}{\partial \theta^\top} (\theta - \theta^0) + O_p(|\theta - \theta^0|^2).$$

461 We make the following assumptions.

462 **Assumption A1.** Under  $\mathbf{H}_0$ ,

$$463 \sqrt{n}(\hat{\theta} - \theta^0) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \psi(U_i | \theta^0) + o_p(1),$$

464 where  $E[\psi(U_i | \theta^0)] = 0$  and  $E[\psi(U_i | \theta^0)\psi(U_i | \theta^0)^\top] < \infty$ .

474 **Assumption A2.** The function  $V(u; \theta)$  is uniformly continuous in  $u$  and twice  
 475 continuously differentiable in  $\theta$  on  $\Theta_0 = \{\theta: |\theta - \theta^0| \leq c\}$  for some  $c > 0$ , with  
 476  $E[|\frac{\partial V_\ell}{\partial \theta_k}(U_i, \theta^0)|^2] < \infty$  and for some  $\kappa \geq 2$ ,  $E[\sup_{\theta \in \Theta_0} |\frac{\partial^2 V_\ell}{\partial \theta_k \partial \theta_r}(U_i, \theta)|^\kappa] < \infty$  for  
 477  $\ell = 1, \dots, d$  and  $k, r = 1, \dots, p$ .

478

479 **Assumption A3.** The functions  $\Delta_j(\cdot | \theta)$ ,  $j = 0, \dots, 4$  are continuously  
 480 differentiable in  $\theta$  on  $\Theta_0$ , and the derivative vector  $\Gamma(\cdot | \theta) = \partial \Delta_0(\cdot | \theta) / \partial \theta^\top$   
 481 satisfies

482

$$483 \int \Gamma(U | \theta^0) \Gamma(U | \theta^0)^\top dP(U) < \infty, \quad (11)$$

484

485 where  $P$  is the distribution of  $U$ .

486

487 **Remarks.** Assumption 2 could, perhaps, be weakened to once  
 488 differentiability at the cost of a longer proof. In general these assumptions  
 489 are fairly standard and can be verified directly. For illustrative purposes,  
 490 we just consider the linear model

491

$$492 y = \beta_0^\top X + \varepsilon, \quad (12)$$

493

494 where  $\beta_0 = (1, 1)^\top$  and  $X = (X_1, X_2)^\top$ , with  $(\varepsilon, X^\top)^\top \sim N(0, I_3)$ . Consider  
 495 testing the hypothesis that  $y \perp \alpha_0^\top X | \beta_0^\top X$ , where  $\alpha_0 = (1, -1)^\top$ . For any  $\alpha$  and  
 496  $\beta$ , we have

$$497 \begin{pmatrix} y \\ \alpha^\top X \\ \beta^\top X \end{pmatrix} \sim N \left[ 0, \begin{pmatrix} 1 + \beta_0^\top \beta_0 & \beta_0^\top \alpha & \beta_0^\top \beta \\ \beta_0^\top \alpha & \alpha^\top \alpha & \alpha^\top \beta \\ \beta_0^\top \beta & \alpha^\top \beta & \beta^\top \beta \end{pmatrix} \right].$$

500

501 In order for Assumption 3 to be satisfied, this covariance matrix should be  
 502 nonsingular at  $\alpha = \alpha_0$  and  $\beta = \beta_0$ ; this certainly holds, since by construction  
 503  $\beta_0^\top \alpha_0 = 0$ . The derivatives of  $\Delta_j$  are fairly easy to compute in this case. For  
 504 example,

505

$$506 \frac{\partial}{\partial \beta} E[\mathbf{1}\{\beta^\top X \leq \beta^\top x\}] = [I - (\beta^\top \beta)^{-1} \beta \beta^\top] \frac{x}{(\beta^\top \beta)^{1/2}} \phi(\beta^\top x / (\beta^\top \beta)^{1/2}).$$

508

509 This quantity is mean zero and has finite variance with respect to the  
 510 distribution of  $x$ . Kim and Pollard (1990) carry out similar calculations for  
 511 general distributions with instead  $\beta^\top x$  replaced by 0.

512 The large sample properties of our test statistics are given in the  
 513 following theorem which is proved in the Appendix. Let

514

$$515 CM_n^c = n^{-1} \sum_{i=1}^n (A_n(\widehat{V}_i) - \mu(V_i))^2; \quad KS_n^c = \max_{1 \leq i \leq n} |A_n(\widehat{V}_i) - \mu(V_i)|.$$

516

517 **Theorem 1.** Suppose that Assumptions A1–A3 hold. (i) Under  $H_0$ ,

518

519

520

521

$$nCM_n^c \Rightarrow \sum_{\ell=1}^{\infty} \lambda_{\ell} \chi_{1\ell}^2, \quad (13)$$

522

523

524

525

526

527

where  $\chi_{1\ell}^2$ ,  $\ell = 1, 2, \dots$  are independent chi-squared [with one degree of freedom] random variables, while  $\{\lambda_{\ell}\}_{\ell=1}^{\infty}$  are the eigenvalues of the operator  $\tau$ , where  $\tau q(\cdot) = \int h(\cdot, y)q(y)dP(y)$  in which  $h(u_1, u_2) = \int \zeta(u_1, V(U, \theta^0) | \theta^0) \zeta(u_2, V(U, \theta^0) | \theta^0) dP(U)$  with  $\zeta(u, v | \theta) = \delta(u, v | \theta) + \Gamma(v | \theta)\psi(u | \theta)$ ; and (ii) under  $H_0$ ,

528

529

530

$$n^{1/2}KS_n^c \Rightarrow \sup_{t \in \mathbb{R}^q} |W(t)|, \quad (14)$$

531

532

533

in which  $W$  is a Gaussian process with mean zero and covariance function  $\omega(u, u') = \int \zeta(U, V(u, \theta^0) | \theta^0) \zeta(U, V(u', \theta^0) | \theta^0) dP(U)$ .

534

535

536

537

538

539

540

541

542

543

The limiting distributions are non-Gaussian. Also, there is a “correction factor” [the term  $\Gamma(v | \theta^0)\psi(u | \theta^0)$  inside  $\zeta(u, v | \theta^0)$ ] in the limiting distribution of both test statistics due to the estimation of  $\theta^0$ . When the parameters are known, this term disappears and  $\zeta = \delta$ . Similarly, when the parameters enter in a linear fashion, the correction term can be zero, see for example Pierce and Kopecky (1979). Even in this case, the null distributions of our tests are complicated functionals of a Gaussian process and depend on the underlying distribution, i.e., neither test is distribution free. This is why we use the bootstrap, see below, to construct critical values.

544

545

546

547

548

549

550

551

552

553

Consider next the power of our tests against local alternatives. Our tests should have power against all root- $n$  alternatives, just like the Bierens and Ploberger test (1997). This is true for discrete variables by virtue of the events  $\mathfrak{B}$  we have taken; this is also true regardless of the dimensionality of  $V$ .<sup>13</sup> We suppose, for simplicity, that the parameters are known. The choice of how to specify alternatives even for the (unconditional) independence test is not universally agreed on, see (Nikitin, 1995, p. 194). For simplicity we shall assume that the data are generated by a sequence of distribution functions  $\bar{F}$  shrinking towards a distribution function  $F$  that does satisfy the null hypothesis, i.e.,

554

555

556

$$\mathbf{H}_n : \bar{F}(v) = F(v) + \frac{a_V(v)}{n^{1/2}}$$

557

558

559

<sup>13</sup>Also note that the rate of convergence to the limiting distributions in Theorem 1 is  $n^{1/2}$  independently of dimensions which implies that the size distortion is of order  $n^{-1/2}$  independently of dimensions. See Csörgó and Faraway (1996).

560 for some function  $a(\cdot)$  which is not identically zero (and which makes  $\bar{F}(v)$   
 561 a probability for all  $v$  for  $n$  larger than some  $n_0$ ). This implies that

$$562 \bar{L}(z) = L(z) + \frac{a_Z(z)}{n^{1/2}}; \quad \bar{G}(y, z) = G(y, z) + \frac{a_{YZ}(y, z)}{n^{1/2}};$$

$$563 \bar{H}(x, z) = H(x, z) + \frac{a_{XZ}(x, z)}{n^{1/2}}$$

564 for functions  $a_Z(z) = a_V(\infty, \infty, z)$ ,  $a_{YZ}(y, z) = a_V(y, \infty, z)$ , and  $a_{XZ}(x, z) =$   
 565  $a_V(\infty, x, z)$ . For  $\bar{F}(v)$  to be a proper alternative hypothesis, we require that

$$566 \tau(v) = L(z)a_V(v) + F(v)a_Z(z) - G(y, z)a_{XZ}(x, z) - H(x, z)a_{YZ}(y, z)$$

567 is not identically zero. In any case, under the sequence of hypotheses  $\mathbf{H}_n$ ,  
 568 we have

$$569 n^{1/2}KS_n^c \Rightarrow \sup_{t \in \mathbb{R}^q} |W(t) + \tau(t)|.$$

570 This guarantees nontrivial power against such alternatives. A similar result  
 571 holds for the Cramér-von Mises test.

#### 572 4. BOOTSTRAP CRITICAL VALUES

573 We use the bootstrap because it performs well in many other related  
 574 situations (and can lead to some improvements, see for example Canepa  
 575 and Godfrey, 2007) and because the alternative bounding methods  
 576 suggested in Bierens and Ploberger (1997), for example, are quite  
 577 complicated to implement and unintuitive.

578 We first discuss the method for the case that the parameters  $\theta^0$   
 579 are known. The basic problem for the bootstrap is how to impose the  
 580 null hypothesis in the resampling scheme. Simple resampling from the  
 581 empirical joint distribution of  $V_i$  will not impose the null restriction. In  
 582 the independence case, one can resample from the marginal empiricals  
 583 thereby imposing independence, see for example Skaug and Tjøstheim  
 584 (1993). We essentially do the same here except that our marginals are  
 585 conditional on  $Z$ . Let  $P_n^{XYZ}$  be the joint empirical distribution, then write  
 586  $P_n^{XYZ} = P_n^{XY|Z} \cdot P_n^Z$ , where  $P_n^{XY|Z}$  and  $P_n^Z$  are the empirical distribution of  
 587  $(X, Y)$  conditional on  $Z$ , and the empirical marginal distribution of  $Z$ ,  
 588 respectively. The conditioning variable  $Z$  is an ancillary statistic, so that we  
 589 can conduct inference conditional on the sample  $\{Z_i\}_{i=1}^n$  without any loss of  
 590 information, i.e., we can work with  $P_n^{XY|Z}$ . Our proposal consists of drawing  
 591 resamples  $\{X_i^*, Y_i^*, Z_i^*\}_{i=1}^n$ , where  $Z_i^* = Z_i$ , from a conditional distribution

603  $\widehat{P}_n^{XY|Z}$  in which we impose the null hypothesis of independence between  
 604  $X$  and  $Y$  conditional on  $Z$ . That is,

$$605 \quad 606 \quad 607 \quad \widehat{P}_n^{XY|Z} = \widehat{P}_n^{X|Z} \cdot \widehat{P}_n^{Y|Z},$$

608 where  $\widehat{P}_n^{X|Z}$  and  $\widehat{P}_n^{Y|Z}$  denote the bootstrap conditional distributions of  $X$   
 609 and  $Y$ , respectively. We just explain the procedure for computing  $\widehat{P}_n^{X|Z}$ ,  
 610 since  $\widehat{P}_n^{Y|Z}$  is constructed in the same manner. Unlike with the joint  
 611 distribution  $P^{XZ}$  of  $X$  and  $Z$  where the (naive) bootstrap distribution can  
 612 be chosen to be the empirical distribution  $P_n^{XZ} = n^{-1} \sum_{i=1}^n \mathbf{1}(X = X_i)\mathbf{1}(Z =$   
 613  $Z_i)$ , the analogy does not carry over to  $\widehat{P}_n^{X|Z}$ . The reason for this is simple:  
 614 unless one has repeated observations for  $Z$  among the observed values  
 615  $\{Z_i\}_{i=1}^n$ , only one value of  $X$ , namely  $X_i$ , will be associated with each  $Z_i$ ,  
 616  $i = 1, \dots, n$ , so that

$$617 \quad 618 \quad 619 \quad 620 \quad \widehat{P}_n^{X|Z}(X_i^* | Z_i) = \frac{n^{-1} \sum_{j=1}^n \mathbf{1}(Z_j = Z_i)\mathbf{1}(X_j = X_i^*)}{n^{-1} \sum_{j=1}^n \mathbf{1}(Z_j = Z_i)},$$

621 or  $X_i^* = X_i$  with probability 1,  $i = 1, \dots, n$ . Even in the rare event that each  
 622 value of  $Z$  in our sample is associated with two or three distinct values of  
 623  $X$ , it will still be inadequate to produce a good approximation of  $P^{X|Z}$   
 624 through the empirical distribution  $P_n^{X|Z}$ .

625 One way to solve this problem is to smooth  $P_n^{X|Z}$ . We choose the  
 626 following smoothing procedure in our simulations and application below.  
 627 For any set  $A$ , including singletons, let

$$628 \quad 629 \quad 630 \quad 631 \quad 632 \quad \widehat{P}_n^{X|Z}(A | Z_i) = \frac{n^{-1} \sum_{j=1}^n K_h(\|Z_j - Z_i\|)\mathbf{1}(X_j \in A)}{n^{-1} \sum_{j=1}^n K_h(\|Z_j - Z_i\|)}, \quad (15)$$

633 where  $K_h(u) = h^{-1}K(u/h)$ , and the univariate kernel  $K$  is a symmetric,  
 634 nonnegative function that integrates to one, and is absolutely integrable.  
 635 In practice, a weighted distance is chosen to reflect the different scales of  
 636 the vector components.<sup>14</sup> We resample from (15); this involves choosing

$$637 \quad 638 \quad 639 \quad 640 \quad X_i^* = X_j \quad \text{with probability} \quad \frac{K_h(\|Z_j - Z_i\|)}{\sum_{j=1}^n K_h(\|Z_j - Z_i\|)}, \quad j = 1, \dots, n.$$

641 In practice, it will be advisable in small samples to choose the bandwidth  
 642 parameter  $h$  to be, for example, the distance from  $Z_i$  to its  $k$ th nearest  
 643

644 <sup>14</sup>See Härdle and Linton (1994) for discussion of smoothing methods and Horowitz (1995)  
 645 for background on the bootstrap.

646 neighbor ( $k$ -NN).<sup>15</sup> This guarantees that each  $X_i^*$  is drawn from at least  $k$   
 647 observations of  $X$  whose associated  $Z$  are the  $k$  closest to  $Z_i$ . We generate  
 648  $B$  bootstrap samples and with each sample compute  $CM_n^*$  and  $KS_n^*$  in  
 649 analogous fashion to  $CM_n$  and  $KS_n$ . The level  $\alpha$  critical values  $\hat{c}_\alpha$  are  
 650 computed as an approximate solution to

$$651 \Pr^*[CM_n^* > \hat{c}_\alpha] = \alpha, \quad (16)$$

652 where  $\Pr^*$  denotes probability conditional on the sample.  
 653 The consistency of this procedure,

$$654 \sup_{c \in \mathbb{R}} |\Pr^*(I_n^* \leq c) - \Pr(I_n \leq c)| = o_p(1), \quad n \rightarrow \infty,$$

655 where  $I_n$  denotes either  $CM_n$  or  $KS_n$  and  $I_n^*$  denotes either  $CM_n^*$  or  $KS_n^*$ ,  
 656 should follow from the following argument. Firstly, suppose that instead  
 657 of the fixed distribution  $P$  of  $U$ , there was a deterministic sequence  $P_n$   
 658 of probability measures which for each  $n$  satisfies the null hypothesis.  
 659 Theorem 1 can be extended to include this triangular array and, provided  
 660  $P_n \rightarrow P$ , the limiting distribution is the same as given in Theorem 1.  
 661 Secondly, one can show that  $\hat{P}_n^{X|Z}$  and  $\hat{P}_n^{Y|Z}$  are uniformly consistent,  
 662 under regularity conditions such as can be found in Härdle et al. (1988).  
 663 Thus, the bootstrap test statistic has the same asymptotic distribution.

664 Suppose now that  $\theta^0$  is replaced by the estimated value  $\hat{\theta}$ . In some  
 665 special cases the correction term due to parameter estimation is zero,  
 666 i.e.,  $E(\Gamma) = 0$ . This occurs in the linear index model when  $X$  is mean  
 667 zero. In this case, one can use the above algorithm without re-estimating  
 668  $\theta$  each time. In general, however,  $E(\Gamma) \neq 0$ . We suppose that there is a  
 669 well defined inversion mapping  $r(Y, X, Z) = U$  [which is certainly the  
 670 case in linear index models]. In this case, we recommend the following  
 671 procedure:

- 672 1. With the original data and  $\hat{\theta}$  estimate  $\hat{P}_n^{X|Z}$  and  $\hat{P}_n^{Y|Z}$  but also  $\hat{P}_n^Z$  (for the  
 673 latter just take the unsmoothed empirical distribution function). Now  
 674 draw a random sample  $\{Y_i^*, X_i^*, Z_i^*\}_{i=1}^n$  from the joint distribution  $\hat{P}_n^{X|Z} \cdot$   
 675  $\hat{P}_n^{Y|Z} \cdot \hat{P}_n^Z$ .
- 676 2. Compute  $U_i^* = r(Y_i^*, X_i^*, Z_i^*)$  and re-estimate  $\hat{\theta}^*$  using the bootstrap  
 677 sample  $\{U_i^*\}_{i=1}^n$ .
- 678 3. Compute  $\hat{Y}_i^* = Y(U_i^*, \hat{\theta}^*)$ ,  $\hat{X}_i^* = X(U_i^*, \hat{\theta}^*)$ , and  $\hat{Z}_i^* = Z(U_i^*, \hat{\theta}^*)$ ,  $i =$   
 679  $1, \dots, n$ .
- 680 4. Compute  $CM_n^*$  and  $KS_n^*$  using the bootstrap sample  $\{\hat{Y}_i^*, \hat{X}_i^*, \hat{Z}_i^*\}_{i=1}^n$ .

681 <sup>15</sup>For such  $h$ , and  $K$  the uniform distribution, we get a  $k$ -nearest neighbor smooth distribution  
 682 with  $X_i^* = X_j$  with probability  $1/k$  for all  $X_j$ ,  $j = 1, \dots, n$ , such that  $Z_j \in \mathcal{N}_k(Z_i)$ , where  $\mathcal{N}_k(Z_i)$   
 683 denotes a  $k$ -neighborhood of  $Z_i$ .



689 Repeat the above  $B$  times and compute  $\hat{c}_\alpha$  as in (16).

690 In the cases where there is no an inversion mapping  $r$ , one could  
 691 instead use the approach described in Horowitz (1995), where the  
 692 recentred test statistic  $CM_n^* - CM_n$  is used and the resample does not need  
 693 to impose the null hypothesis.

694

695

696

## 5. SIMULATIONS

697

698

699

700

We evaluate the performance of our test in a binary choice model to test the conditional independence restriction (4). Specifically, we take the following designs:

701

$$D1. \quad Y = \mathbf{1}(\beta_1 X_1 + \beta_2 X_2 > \varepsilon),$$

702

$$D2. \quad Y = \mathbf{1}(\beta_1 X_1 + \beta_2 X_2 + 10n^{-1/2}(X_1^2 + X_2^2) > \varepsilon),$$

703

$$D3. \quad Y = \mathbf{1}(\beta_1 X_1 + \beta_2 X_2 + (X_1^2 + X_2^2) > \varepsilon),$$

704

705

$$D4. \quad Y = \mathbf{1}(\beta_1 X_1 + \beta_2 X_2 > \varepsilon(X_1^2 + X_2^2)^{1/2}),$$

706

707

708

709

710

711

712

713

714

715

716

where in all cases  $\beta_1 = \beta_2 = 1$ ,  $X = (X_1, X_2)$  is bivariate standard normal, and  $\varepsilon$  is standard normal independent of  $X$ .<sup>16</sup> The first design satisfies the null hypothesis of conditional independence of  $Y$  from  $X$  given the index  $Z = \beta_1 X_1 + \beta_2 X_2$ . The second is an order  $n^{-1/2}$  local alternative to this hypothesis. The third and fourth designs represent global alternatives arising from location and scale shifts, respectively. Note that conditioning on the index  $Z = \beta_1 X_1 + \beta_2 X_2$  implies that  $Y$  is independent of  $(X_1, X_2)$  given  $\beta_1 X_1 + \beta_2 X_2$  reduces to testing  $H_0: Y$  is independent of  $X_1$  given  $\beta_1 X_1 + \beta_2 X_2$ .

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

In order to implement the tests of index sufficiency we first need to estimate the index. There are two estimators to consider; the probit estimator  $\hat{\beta}_p$  and the Klein–Spady semiparametrically efficient estimate  $\hat{\beta}_{KS}$ . The implementation of the Klein–Spady estimator is perhaps problematic in that one has to select bandwidth and trimming parameters for which there is very little theoretical guidance as yet. Furthermore, this procedure is also quite time consuming for our purposes, since we have to re-compute this quantity for each bootstrap sample. We could have, therefore, based our simulations on  $\hat{\beta}_p$ . This, however, leaves one open to the criticism that an estimate too efficient under the null is being used relative to  $\hat{\beta}_{KS}$ , the central semiparametric estimator here. To address all these issues, we used the following first order approximation to  $\hat{\beta}_{KS}$ . Consider the

<sup>16</sup>Additional simulations with  $(Y, X, Z)$  trivariate normal are reported in Linton and Gozalo (1995).

732 local alternative in which  $\Pr(Y = 1|X) = \Phi\{\beta^{0\top}X + n^{-1/2}g(X)\}$  for some  
 733 function  $g(\cdot)$ . Let

$$734 \tilde{\beta}_{KS} = \beta^0 - \left\{ n^{-1} \sum_{i=1}^n \frac{\phi_i^2 (X_i - \bar{X}_i) (X_i - \bar{X}_i)^\top}{\Phi_i(1 - \Phi_i)} \right\}^{-1}$$

$$735 \times \left\{ n^{-1} \sum_{i=1}^n \frac{Y_i - \Phi_i - \phi_i \bar{g}_i / n^{1/2}}{\Phi_i(1 - \Phi_i)} \phi_i (X_i - \bar{X}_i) \right\},$$

736 with  $\Phi_i = \Phi(\beta^{0\top} X_i)$ ,  $\phi_i = \phi(\beta^{0\top} X_i)$ ,  $\bar{X}_i = E(X_i | \beta^{0\top} X_i)$ ,  $g_i = g(X_i)$ , and  $\bar{g}_i =$   
 737  $E(g_i | \beta^{0\top} X_i)$ .<sup>17</sup> Then, under both  $D1$  and  $D2$ ,  $n^{1/2}(\tilde{\beta}_{KS} - \hat{\beta}_{KS}) = o_p(1)$ . In  
 738 fact this result holds under the global alternatives  $D3$  and  $D4$ , except that  
 739 in those misspecified cases  $\tilde{\beta}_{KS}$  nor  $\hat{\beta}_{KS}$  will converge to  $\beta^0$ .

740 In the construction of the test we used zero width rectangles for  
 741 the discrete variable  $Y$ ,  $\mathfrak{B}(x) = (-\infty, x]$  for the continuous variable  $X$ ,  
 742 and  $\mathfrak{B}(z) = [z - \hat{\sigma}_{50}/4, z + \hat{\sigma}_{50}/4]$  for the index  $Z$ , where  $\hat{\sigma}_{50}$  denotes the  
 743 estimated standard deviation of  $Z$  from one fix sample of size  $n = 50$ .  
 744 Depending on the location of the evaluation point  $z$ , the interval  $\mathfrak{B}(z)$   
 745 will contain different number of observations. The interval width of  $\mathfrak{B}(z)$   
 746 was kept fix and independent of  $n$  throughout the simulations. We did  
 747 not experiment greatly with the choice of rectangles and presumably one  
 748 could achieve superior performance via tuning of this choice.

749 For each estimated index value  $\tilde{Z}_i = \tilde{\beta}_1 X_{1,i} + \tilde{\beta}_2 X_{2,i}$ , the dependence  
 750 score  $A_n = L_n F_n - G_n H_n$  was evaluated at the  $m_i$  sample points  $V_j =$   
 751  $(Y_j, X_j, \tilde{Z}_i)$  for observations  $j$  whose  $\tilde{Z}_j$  is in the interval  $\mathfrak{B}(\tilde{Z}_i)$ . This results in  
 752  $m = \sum_{i=1}^n m_i$  evaluation points. For the sample sizes considered of  $n = 50$ ,  
 753  $n = 100$ , and  $n = 500$ , the average value of  $m_i$  was approximately 7.5, 14.7,  
 754 and 70.2, respectively, resulting on  $m = 375$ , 1470, and 35100 evaluation  
 755 points, respectively.<sup>18</sup>

756 We conducted 500 replications of the Cramér-von Mises and  
 757 Kolmogorov–Smirnov type tests under the null, and 100 under each  
 758 alternative design. We used 100 bootstrap samples in each replication to  
 759 calculate the critical values at significance levels  $\alpha = 1\%$ ,  $5\%$ ,  $10\%$ , and  
 760  $20\%$ . To compute the bootstrap test, we used (15) to obtain the bootstrap  
 761 observations  $X_i^*$  with  $Z_i$  and  $Z_j$  replaced by the estimated indexes  $\tilde{Z}_i$  and  
 762  $\tilde{Z}_j$ , Epanechnikov kernel  $K(\cdot)$ , and bandwidth  $h$  equal to the distance  
 763 from  $\tilde{Z}_i$  to its  $k$ th nearest neighbor. The values of  $k$  chosen were  $k = 5$ ,  
 764

771 <sup>17</sup>Note that  $E(X_j | \sum_{i=1}^p X_i) = p^{-1} \sum_{i=1}^p X_i$ , for  $j = 1, \dots, p$ , and  $E(\sum_{i=1}^p X_i^2 | \sum_{i=1}^p X_i) = (p -$   
 772  $1) + p^{-1}(\sum_{i=1}^p X_i)^2$ .

773 <sup>18</sup>Given the large value of  $m$  using this procedure for  $n = 500$ , we decided to evaluated the  
 774 test at only 10% of the points in each interval  $\mathfrak{B}(\tilde{Z}_i)$ . This cut  $m$  to a more manageable 3510  
 points on average.

775 **TABLE 1** Size and power with estimated conditioning variable  $Z$ .

776

777

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

		Percentage rejections of null hypothesis ( $Z = \bar{\beta}_{KS}^T X$ )					
Design	$\alpha$ (%)	Cramér-von mises			Kolmogorov–Smirnov		
		$n = 50$	$n = 100$	$n = 500$	$n = 50$	$n = 100$	$n = 500$
D1	20	15.4	20.6	20.4	24.0	25.6	20.0
	10	5.8	8.6	10.2	12.4	12.4	10.6
	5	2.8	3.4	5.4	5.0	6.2	5.8
D2	1	1.2	1.0	1.2	1.2	1.4	2.8
	20	19.0	37.0	86.0	32.0	38.0	60.0
	10	11.0	15.0	60.0	20.0	21.0	36.0
D3	5	2.0	3.0	38.0	9.0	9.0	19.0
	1	1.0	0.0	12.0	3.0	1.0	6.0
	20	15.0	37.0	97.0	24.0	38.0	90.0
D4	10	6.0	15.0	92.0	15.0	21.0	66.0
	5	3.0	3.0	74.0	6.0	9.0	44.0
	1	1.0	0.0	45.0	1.0	1.0	23.0
D4	20	11.0	22.0	38.0	14.0	19.0	32.0
	10	6.0	11.0	24.0	7.0	9.0	18.0
	5	2.0	4.0	18.0	1.0	3.0	13.0
	1	0.0	1.0	4.0	0.0	1.0	2.0

793

794

795 10, and 10 for  $n = 50, 100, \text{ and } 500$ , respectively. Similarly for  $Y_i^*$ . The

796 bootstrap sample  $(Y_i^*, X_i^*)$ ,  $i = 1, \dots, n$ , was then used to obtain new

797 parameter estimates  $\beta^*$  with which we form a new set of index values  $\tilde{Z}_i^*$ .

798

799

800 **TABLE 2** Size and power with known conditioning variable  $Z$

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

		Percentage rejections of null hypothesis ( $Z = \bar{\beta}_{KS}^T X$ )					
Design	$\alpha$ (%)	Cramér-von mises			Kolmogorov–Smirnov		
		$n = 50$	$n = 100$	$n = 500$	$n = 50$	$n = 100$	$n = 500$
D1	20	18.4	19.6	20.2	23.2	21.6	20.4
	10	7.8	11.0	10.2	11.6	11.0	10.2
	5	3.8	5.4	5.2	5.8	6.2	5.4
D2	1	0.4	1.0	0.8	1.4	1.4	1.8
	20	20.0	51.0	94.0	26.0	40.0	61.0
	10	8.0	22.0	80.0	13.0	20.0	41.0
D3	5	3.0	9.0	52.0	4.0	11.0	23.0
	1	1.0	2.0	18.0	2.0	3.0	8.0
	20	20.0	51.0	100.0	28.0	40.0	92.0
D4	10	9.0	22.0	100.0	18.0	20.0	76.0
	5	3.0	9.0	97.0	8.0	11.0	63.0
	1	0.0	2.0	82.0	2.0	3.0	31.0
D4	20	14.0	33.0	59.0	20.0	30.0	42.0
	10	8.0	21.0	37.0	8.0	11.0	26.0
	5	3.0	8.0	28.0	2.0	4.0	17.0
	1	0.0	2.0	9.0	1.0	2.0	4.0

818 Finally,  $(Y_i^*, X_i^*, \tilde{Z}_i^*)$ ,  $i = 1, \dots, n$ , is used to compute the bootstrap tests  
 819  $CM_n^*$  and  $KS_n^*$ .

820 Our results using the Klein–Spady index estimate are given in Table 1.

821 The Cramér–von Mises test appears to have good size, even in relatively  
 822 small samples, while the Kolmogorov–Smirnov test requires a larger sample  
 823 size to achieve values close to the nominal values. Both tests have power  
 824 against the root- $n$  local alternative design  $D2$ , and have power against the  
 825 global alternatives  $D3$  and  $D4$  of shifts in the location and the scale of the  
 826 distribution of the error term (particularly against  $D3$ ). The Cramér–von  
 827 Mises test has higher power against all alternatives except for  $n = 50$ .

828 To evaluate the size/power loss due to having to estimate the index, we  
 829 computed the two tests with  $Z = \beta_1 X_1 + \beta_2 X_2$  assumed known. The results  
 830 are given in Table 2. There is not much difference in size performance  
 831 between the two tests. The power has increased for all designs and sample  
 832 sizes, as expected, but particularly for the scale-shift design  $D4$ .

833

## 834 6. CONCLUSIONS

835

836 We have proposed a test of conditional independence or rather a  
 837 stronger condition than conditional independence. The test has the  
 838 advantage that it does not involve an explicit bandwidth parameter and  
 839 has good statistical properties. One key question is how to choose the  
 840 rectangles used in the construction of our test. The default rectangles are  
 841 the half spaces used in the usual construction of c.d.f.’s and these may work  
 842 well for most applications, but may be improved in certain cases. We leave  
 843 this for future work.

844

## 845 A. APPENDIX

846

### 847 A.1. Preliminaries

848 Let  $\Theta_n(c) = \{\theta: \sqrt{n}|\theta - \theta^0| \leq c\}$ . Since  $\sqrt{n}(\hat{\theta} - \theta^0) = O_p(1)$ , for all  
 849  $\epsilon > 0$  there exists a  $c_\epsilon$  such that  $\Pr[\hat{\theta} \in \Theta_n(c_\epsilon)] \geq 1 - \epsilon$ . Let  $\mathcal{A}_\theta$  be  
 850 any event depending on  $\theta$ , and let  $\mathcal{B}$  be the event that  $\hat{\theta} \in \Theta_n(c_\epsilon)$ .  
 851 Then  $\Pr[\mathcal{A}_{\hat{\theta}}] = \Pr[\mathcal{A}_{\hat{\theta}} \cap \mathcal{B}] + \Pr[\mathcal{A}_{\hat{\theta}} \cap \mathcal{B}^c] \leq \Pr[\mathcal{A}_{\hat{\theta}} \cap \mathcal{B}] + \Pr[\mathcal{B}^c] = \Pr[\mathcal{A}_{\hat{\theta}} \cap$   
 852  $\mathcal{B}] + o(1)$ . Finally, we can replace  $\theta$  by the fixed value in  $\Theta_n(c)$  that makes  
 853  $\Pr[\mathcal{A}_\theta \cap \Theta_n(c)]$  the largest.

854 We begin by providing some background results concerning the  
 855 processes:

856

$$857 v_n(\theta, v) = \frac{1}{\sqrt{n}} \sum_{i=1}^n [\mathbf{1}\{V(U_i, \theta) \in \mathfrak{B}(v)\} - E\{\mathbf{1}\{V(U_i, \theta) \in \mathfrak{B}(v)\}\}],$$

858

$$859 \theta \in \Theta_n, v \in \mathbb{R}^d$$

860

$$v'_n(\theta, u) = \frac{1}{\sqrt{n}} \sum_{i=1}^n [\mathbf{1}\{V(U_i, \theta) \in \mathfrak{B}(V(u, \theta))\} - E\{\mathbf{1}\{V(U_i, \theta) \in \mathfrak{B}(V(u, \theta))\}\}],$$

$$\theta \in \Theta_n, u \in \mathbb{R}^q.$$

The same results hold for the corresponding empirical processes involving subvectors of  $V(U_i, \theta)$ , but for convenience we just state results for  $v_n(\theta, v)$  and  $v'_n(\theta, v)$ . Note that for given  $\theta$ , these are standard empirical processes with the slight difference driven by the shape of the sets  $\mathfrak{B}(v)$ . Define the pseudo-metric

$$\rho((\theta, v), (\theta', v')) = E \left( [\mathbf{1}\{V(U_i, \theta) \in \mathfrak{B}(v)\} - \mathbf{1}\{V(U_i, \theta') \in \mathfrak{B}(v')\}]^2 \right),$$

on  $\Theta_0 \times \mathbb{R}^d$ . Likewise, define the pseudo-metric  $\rho'((\theta, u), (\theta', u'))$  on  $\Theta_0 \times \mathbb{R}^q$ . Under these metrics, the parameter spaces  $\Theta_0 \times \mathbb{R}^d$  and  $\Theta_0 \times \mathbb{R}^q$  are totally bounded. In the sequel we shall just use the generic notation  $\rho(\cdot, \cdot)$  for a metric.

We are only interested in the behavior of these processes as  $\theta$  varies in the small set  $\Theta_n$ . By writing  $\theta = \theta^0 + \gamma n^{-1/2}$ , we shall make a reparameterization to  $v_n(\gamma, v)$  and  $v'_n(\gamma, u)$ , where  $\gamma \in \Gamma(c) \subset \mathbb{R}^b$ . We establish the following results:

$$\sup_{\gamma \in \Gamma, v \in \mathbb{R}^d} |v_n(\gamma, v) - v_n(0, v)| = o_p(1) \quad (17)$$

$$\sup_{\gamma \in \Gamma, u \in \mathbb{R}^d} |v'_n(\gamma, u) - v'_n(0, u)| = o_p(1). \quad (18)$$

To prove (17) and (18) it is sufficient to show a pointwise law of large numbers, e.g.,  $v_n(\gamma, v) - v_n(0, v) = o_p(1)$  for any  $\gamma \in \Gamma, v \in \mathbb{R}^d$ , and stochastic equicontinuity of the processes  $v_n$  and  $v'_n$ . The pointwise result is immediate because the random variables are sums of bounded i.i.d. random variables with zero mean; the limits of  $v_n(\gamma, v)$  and  $v_n(0, v)$  are the same by smoothness of the expected value in  $\gamma$ . To complete the proof of (17) and (18) we shall use the following lemma, proved below, which states that  $v_n$  and  $v'_n$  are stochastically equicontinuous in  $\theta$  and  $v$  and  $\theta$  and  $u$  respectively. First, recall definition (2.3) of Andrews (1994).

**Definition.** A process  $v_n(\cdot)$  is stochastically equicontinuous if for all  $\epsilon > 0$  and  $\eta > 0$ , there exists  $\delta > 0$  such that

$$\limsup_{n \rightarrow \infty} \Pr \left[ \sup_{\rho(t_1, t_2) < \delta} |v_n(t_1) - v_n(t_2)| > \eta \right] < \epsilon.$$

904 **Lemma SE.** *Under the above assumptions, the processes  $v_n(\gamma, v)$  and  $v'_n(\gamma, u)$*   
 905 *are stochastically equicontinuous.*

906

907 Before we give the proof of Lemma SE, we state and prove a lemma  
 908 which is useful in the sequel.

909

910 **Lemma C.** *If a stochastic process  $v_n(t)$  is stochastically equicontinuous in  $t$ , and*  
 911 *if  $t = g(s)$ , where  $g$  is a uniformly continuous function on its domain, then the*  
 912 *induced process  $v_n^*(s) = v_n(g(s))$  is stochastically equicontinuous in  $s$ .*

913

914 *Proof.* Let  $\varepsilon, \eta > 0$  be given. Take  $\delta > 0$  for which

915

$$916 \limsup_{n \rightarrow \infty} \Pr \left[ \sup_{\rho(t_1, t_2) < \delta} |v_n(t_1) - v_n(t_2)| > \eta \right] < \varepsilon.$$

917

918 By the definition of continuity: there exists  $\zeta > 0$  such that  $\rho(s_1, s_2) < \zeta \Rightarrow$   
 919  $\rho(g(s_1), g(s_2)) < \delta$ . But this implies that

920

$$921 \limsup_{n \rightarrow \infty} \Pr \left[ \sup_{\rho(s_1, s_2) < \zeta} |v_n^*(s_1) - v_n^*(s_2)| > \eta \right] < \varepsilon,$$

922

923 which satisfies the definition of stochastic equicontinuity for the process  
 924  $v_n^*(\cdot)$ . □

925

926 We now give the proof of Lemma SE.

927

928 *Proof of Lemma SE.* We first prove the stochastic equicontinuity of  
 929  $v_n(\gamma, v)$ . Make a Taylor series expansion of  $V(U_i, \theta)$  about  $V(U_i, \theta^0)$

930

$$931 V_\ell(U_i, \theta^0 + \gamma n^{-1/2}) = V_\ell(U_i, \theta^0) + \frac{1}{\sqrt{n}} \sum_{k=1}^p \frac{\partial V_\ell}{\partial \theta_k}(U_i, \theta^0) \gamma_k$$

$$932 + \frac{1}{n} \sum_{k=1}^p \sum_{r=1}^p \frac{\partial^2 V_\ell}{\partial \theta_k \partial \theta_r}(U_i; \bar{\theta}_\ell) \gamma_k \gamma_r \quad (19)$$

933

934 for some intermediate points  $\bar{\theta}_\ell$ . Define the processes:

935

$$936 v_{n1}(\gamma, v) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \left[ \mathbf{1} \left\{ \Upsilon \left( U_i; \theta^0 + \frac{\gamma}{\sqrt{n}} \right) \in \mathfrak{B}(v) \right\} \right.$$

$$937 \left. - E \left\{ \mathbf{1} \left( \Upsilon \left( U_i; \theta^0 + \frac{\gamma}{\sqrt{n}} \right) \in \mathfrak{B}(v) \right) \right\} \right]$$

938

939

940

941

942

943

944

945

946

$$v_{n2}(\gamma, v) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \left[ \mathbf{1} \left\{ V \left( U_i; \theta^0 + \frac{\gamma}{\sqrt{n}} \right) \in \mathfrak{B}(v) \right\} - \mathbf{1} \left\{ \tau \left( U_i; \theta^0 + \frac{\gamma}{\sqrt{n}} \right) \in \mathfrak{B}(v) \right\} \right],$$

where  $\tau = (\tau_1, \dots, \tau_d)^\top$  with  $\tau_\ell(U_i; \theta) = V_\ell(U_i, \theta^0) + \sum_{k=1}^p \frac{\partial V_\ell}{\partial \theta_k}(U_i, \theta^0)(\theta_k - \theta_k^0)$ ,  $\ell = 1, \dots, d$ , being the linear approximation to  $V_\ell(U_i; \theta)$ . Define also the deterministic centering term

$$m_{n3}(\gamma, v) = \sqrt{n} E \left[ \mathbf{1} \left\{ V \left( U_i, \theta^0 + \frac{\gamma}{\sqrt{n}} \right) \in \mathfrak{B}(v) \right\} - E \left[ \mathbf{1} \left\{ \tau \left( U_i, \theta^0 + \frac{\gamma}{\sqrt{n}} \right) \in \mathfrak{B}(v) \right\} \right] \right].$$

Then, since  $v_n(\gamma, v) = v_{n1}(\gamma, v) + v_{n2}(\gamma, v) + m_{n3}(\gamma, v)$ , for stochastic equicontinuity of  $v_n(\gamma, v)$  it suffices to establish as follows:

- (a) The process  $v_{n1}(\gamma, v)$  is stochastically equicontinuous in  $\gamma, v$ ;
- (b) The process  $v_{n2}(\gamma, v)$  is stochastically equicontinuous in  $\gamma, v$ ;
- (c) The process  $m_{n3}(\gamma, v)$  is equicontinuous in  $\gamma, v$ .

**Proof of (a).** Our argument is very similar to that contained in Sherman (1993) because we basically have a linear index structure in this part. We show that the following class  $\mathcal{F}$  is Euclidean for the envelope 1,  $\mathcal{F} = \{f(\cdot, \tau), \tau \in \Gamma \times \mathbb{R}^d\}$ , where for each  $U$  and  $\tau$ ,

$$f(U, \tau) = \prod_{\ell=1}^d \mathbf{1} \left\{ V_\ell(U, \theta^0) + \sum_{k=1}^p \frac{\partial V_\ell}{\partial \theta_k}(U, \theta^0) \tau_k \leq v_\ell + b_\ell \right\} \times \mathbf{1} \left\{ V_\ell(U, \theta^0) + \sum_{k=1}^p \frac{\partial V_\ell}{\partial \theta_k}(U, \theta^0) \tau_k \geq v_\ell - a_\ell \right\}.$$

For each  $U$ , define

$$g(U, v, r, \boldsymbol{\varkappa}_1, \boldsymbol{\varkappa}_2, \boldsymbol{\varkappa}_3, \boldsymbol{\varkappa}_4) = \boldsymbol{\varkappa}_1 r + \sum_{\ell=1}^d \boldsymbol{\varkappa}_{2\ell} v_\ell + \sum_{\ell=1}^d \boldsymbol{\varkappa}_{3\ell} V_\ell(U, \theta^0) + \sum_{\ell=1}^d \sum_{k=1}^p \boldsymbol{\varkappa}_{2\ell k} \frac{\partial V_\ell}{\partial \theta_k}(U, \theta^0)$$

990 and

$$991 \quad \mathcal{G} = \left\{ g(\cdot, \cdot, \cdot, \cdot, \boldsymbol{\varkappa}_1, \boldsymbol{\varkappa}_2, \boldsymbol{\varkappa}_3, \boldsymbol{\varkappa}_4) : \boldsymbol{\varkappa}_1 \in \mathbb{R}, \boldsymbol{\varkappa}_2 \in \mathbb{R}^d, \boldsymbol{\varkappa}_3 \in \mathbb{R}^d, \boldsymbol{\varkappa}_4 \in \mathbb{R}^{dp} \right\}.$$

992 The vector space of real-valued functions  $\mathcal{G}$  is of dimension  $dp + 2d + 1$ .  
 993 For each  $\tau$ , we have  $\text{subgraph}[f(\cdot, \tau)] = \{(U, r) : 0 < r < f(U, \tau)\}$ . This can  
 994 be written as the set of all  $(U, r)$  for which the following product is equal  
 995 to one  
 996

$$997 \quad \prod_{\ell=1}^d \mathbf{1} \left\{ V_{\ell}(U, \theta^0) + \sum_{k=1}^p \frac{\partial V_{\ell}}{\partial \theta_k}(U, \theta^0) \gamma_k \leq v_{\ell} + b_{\ell} \right\}$$

$$998 \quad \times \prod_{\ell=1}^d \mathbf{1} \left\{ V_{\ell}(U, \theta^0) + \sum_{k=1}^p \frac{\partial V_{\ell}}{\partial \theta_k}(U, \theta^0) \gamma_k \geq v_{\ell} - a_{\ell} \right\} \mathbf{1} \{r \geq 1\}^c \mathbf{1} \{r > 0\}^c$$

1000 which can be written as the set of all  $(U, r)$  for which the following quantity  
 1001  $\prod_{\ell=1}^{2d} \mathbf{1} \{g_{\ell} \geq 0\} \mathbf{1} \{g_{d+1} \geq 1\}^c \mathbf{1} \{g_{d+2} > 0\}^c$  is equal to one for some choice of  
 1002  $g_1, \dots, g_{2d+2} \in \mathcal{G}$ . Thus the subgraph of  $f(\cdot, \tau)$  is the intersection of  $2d + 2$   
 1003 sets each of which belongs to a polynomial class (by Lemma 2.4 in Pakes  
 1004 and Pollard, 1989, the class of sets of the form  $\{g \geq a\}$  or  $\{g > a\}$  with  
 1005  $g \in \mathcal{G}$  and  $a \in \mathbb{R}$  is a VC class). Therefore,  $\{\text{subgraph}(f), f \in \mathcal{F}\}$  forms  
 1006 a VC class of sets. Finally, one can apply Lemma 2.12 in Pakes and Pollard  
 1007 (1989).  
 1008  
 1009  
 1010  
 1011  
 1012  
 1013  
 1014  
 1015

1016 **Proof of (b).** By virtue of assumption A2, for any  $\delta > 0$ , we can find an  
 1017  $\varepsilon > 0$  such that  $\max_{\ell=1, \dots, d; k, r=1, \dots, p} E[\sup_{\theta \in \Theta_n} |\frac{\partial^2 V_{\ell}}{\partial \theta_k \partial \theta_r}(U_i, \theta)|^2] / \varepsilon^2 \leq \delta / dp^2$ . Then,  
 1018 by the Bonferroni and Chebychev inequalities,  
 1019

$$1020 \quad \Pr \left[ \frac{1}{\sqrt{n}} \max_{\substack{1 \leq i \leq n \\ \ell=1, \dots, d \\ k, r=1, \dots, p}} \sup_{\theta \in \Theta_n} \left| \frac{\partial^2 V_{\ell}}{\partial \theta_k \partial \theta_r}(U_i, \theta) \right| > \varepsilon \right]$$

$$1021 \quad \leq n \sum_{\substack{\ell=1, \dots, d \\ k, r=1, \dots, p}} \Pr \left[ \frac{1}{\sqrt{n}} \sup_{\theta \in \Theta_n} \left| \frac{\partial^2 V_{\ell}}{\partial \theta_k \partial \theta_r}(U_i, \theta) \right| > \varepsilon \right]$$

$$1022 \quad \leq \frac{\sum_{\substack{\ell=1, \dots, d \\ k, r=1, \dots, p}} E \left[ \sup_{\theta \in \Theta_n} \left| \frac{\partial^2 V_{\ell}}{\partial \theta_k \partial \theta_r}(U_i, \theta) \right|^2 \right]}{\varepsilon^2} \leq \delta.$$

1032



1033 Therefore, with probability tending to one

1034

1035

1036

1037

$$\sup_{\gamma \in \Gamma(c)} \max_{1 \leq i \leq n} \left| \frac{1}{n} \sum_{k=1}^p \sum_{r=1}^p \frac{\partial^2 V_\ell}{\partial \theta_k \partial \theta_r} (U_i; \bar{\theta}_\ell) \gamma_k \gamma_r \right| \leq \frac{\bar{\pi}}{\sqrt{n}}$$

1038

1039

for some  $\bar{\pi} < \infty$ . Therefore, it suffices to show that the process

1040

1041

1042

1043

1044

1045

1046

$$v_{n3}(\gamma, \pi, v) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \left[ \mathbf{1} \left( \left( U_i; \theta^0 + \frac{\gamma}{\sqrt{n}} \right) + \frac{\pi}{\sqrt{n}} \in \mathfrak{B}(v) \right) - E \left\{ \mathbf{1} \left( \left( U_i; \theta^0 + \frac{\gamma}{\sqrt{n}} \right) + \frac{\pi}{\sqrt{n}} \in \mathfrak{B}(v) \right) \right\} \right],$$

1047

1048

where  $\pi \in \Pi$  a compact set, is stochastically equicontinuous in  $\gamma, \pi, v$ , and the deterministic centering term

1049

1050

1051

1052

1053

1054

1055

$$m_{n4}(\gamma, \pi, v) = \sqrt{n} \left[ E \left\{ \mathbf{1} \left( \left( U_i; \theta^0 + \frac{\gamma}{\sqrt{n}} \right) + \frac{\pi}{\sqrt{n}} \in \mathfrak{B}(v) \right) \right\} - E \left\{ \mathbf{1} \left( \left( U_i; \theta^0 + \frac{\gamma}{\sqrt{n}} \right) \in \mathfrak{B}(v) \right) \right\} \right]$$

1056

1057

1058

1059

1060

is also equicontinuous. The process  $v_{n3}(\gamma, \pi, v)$  is stochastically equicontinuous by an obvious modification of the argument given in part (a) because the parameter  $\pi$  enters in a linear fashion. The centering term is handled by Taylor expansion:

1061

1062

1063

1064

1065

1066

1067

1068

1069

$$\begin{aligned} & |m_{n4}(\gamma_1, \pi_1, v_1) - m_{n4}(\gamma_2, \pi_2, v_2)| \\ & \leq \sup_{\gamma, \pi, v} \left| \frac{\partial m_{n4}}{\partial \gamma}(\gamma, \pi, v) \right| |\gamma_1 - \gamma_2| \\ & \quad + \sup_{\gamma, \pi, v} \left| \frac{\partial m_{n4}}{\partial \pi}(\gamma, \pi, v) \right| |\pi_1 - \pi_2| + \sup_{\gamma, \pi, v} \left| \frac{\partial m_{n4}}{\partial v}(\gamma, \pi, v) \right| |v_1 - v_2| \\ & \rightarrow 0 \end{aligned}$$

1070

1071

1072

1073

as  $|\gamma_1 - \gamma_2| + |\pi_1 - \pi_2| + |v_1 - v_2| \rightarrow 0$ , since  $\sup_{\gamma, \pi, v} |\partial m_{n4}(\gamma, \pi, v) / \partial \phi| < \infty$  with  $\phi = (\gamma, \pi, v)$  by assumption A.

1074

1075

**Proof of (c).** It is straightforward to see that  $|m_{n3}(\gamma, v)| \leq c(\gamma, v)/n^{1/2}$  for some bounded continuous function  $c$ , so that (c) follows.

We now argue that the stochastic equicontinuity of  $v'_n(\gamma, u)$  is a consequence of the result for  $v_n(\gamma, v)$  combined with the uniform continuity of the function  $(u, \vartheta) \mapsto v = V(u, \vartheta)$ . Define the new process

$$v''_n(\theta, \vartheta, u) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \left[ \mathbf{1}\{V(U_i, \theta) \in \mathfrak{B}(V(u, \vartheta))\} - E \{ \mathbf{1}\{V(U_i, \theta) \in \mathfrak{B}(V(u, \vartheta))\} \} \right],$$

where  $\theta, \vartheta \in \Theta_n, u \in \mathbb{R}^q$ . Lemma C implies that  $v''_n(\theta, \vartheta, u)$  is stochastically equicontinuous. Therefore, so is  $v'_n(\theta, u)$ .

## A.2 Proof of Theorem 1

(i) Write  $CM_n = n^{-1} \sum_{i=1}^n \{A_n(\widehat{V}_i | \widehat{\theta})\}^2 = \int \{A_n(V(U, \widehat{\theta}) | \widehat{\theta})\}^2 dP_n(U)$ , where  $P_n(\cdot) = n^{-1} \sum \mathbf{1}(U_i \leq \cdot)$  is the empirical measure of  $\{U_i\}_{i=1}^n$ , and let  $CM_n^* = \int A_n^2(V(U, \widehat{\theta}) | \widehat{\theta}) dP(U)$ ,  $CM_n^{**} = \int \Delta_n^2(V(U, \widehat{\theta}) | \widehat{\theta}) dP(U)$ , and  $CM_n^{***} = n^{-2} \sum_{i=1}^n \sum_{j=1}^n h(U_i, U_j)$ , where the function  $h$  was defined in the theorem. By the triangle inequality,

$$n |CM_n - CM_n^{***}| \leq n |CM_n - CM_n^*| + n |CM_n^{**} - CM_n^*| + n |CM_n^{**} - CM_n^{***}|.$$

We establish the following results:

$$n(CM_n - CM_n^*) \xrightarrow{p} 0 \tag{20}$$

$$n(CM_n^* - CM_n^{**}) \xrightarrow{p} 0 \tag{21}$$

$$n(CM_n^{**} - CM_n^{***}) \xrightarrow{p} 0. \tag{22}$$

Then, the result (i) follows because

$$nCM_n^{***} \Rightarrow \sum_{\ell=1}^{\infty} \lambda_{\ell} \chi_{1, \ell},$$

by standard U-statistic theory, see Skaug and Tjøstheim (1993).

*Proof of (21).* We have for any  $\theta$ ,

$$\begin{aligned} n(CM_n^* - CM_n^{**}) &= \int \{A_n^2(V(U, \theta) | \theta) - \Delta_n^2(V(U, \theta) | \theta)\} dP(U) \\ &\leq 2 \sup_{v \in \mathbb{R}^d} |\Delta_n(v | \theta)| \sup_{v \in \mathbb{R}^d} |A_n(v | \theta) - \Delta_n(v | \theta)| \\ &\quad + \sup_{v \in \mathbb{R}^d} |A_n(v | \theta) - \Delta_n(v | \theta)|^2, \end{aligned}$$

1119 since the probability measure  $P \leq 1$ . We show that

$$1120 \sqrt{n} \sup_{\theta \in \Theta_n} \sup_{v \in \mathbb{R}^d} |A_n(v | \theta) - \Delta_n(v | \theta)| \rightarrow_p 0 \quad (23)$$

$$1122 \sqrt{n} \sup_{\theta \in \Theta_n} \sup_{v \in \mathbb{R}^d} |\Delta_n(v | \theta)| = O_p(1). \quad (24)$$

1125 Since  $\Pr(\hat{\theta} \in \Theta_n^c)$  can be made arbitrarily small, (21) will then follow.  
1126 Firstly,

$$1128 A_n(v | \theta) - \Delta_n(v | \theta) = \{F_n(v | \theta) - F(v)\} \{L_n(z | \theta) - L(z)\} \\ 1129 - \{G_n(y, z | \theta) - G(y, z)\} \{H_n(x, z | \theta) - H(x, z)\},$$

1131 so that (23) will follow if  $n^{1/4} \sup |G_n(y, z | \theta) - E[G_n(y, z | \theta)]| \rightarrow_p 0$ ;  
1132  $n^{1/4} \sup |E[G_n(y, z | \theta)] - G(y, z)| \rightarrow_p 0$ ;  $n^{1/4} \sup |H_n(x, z | \theta) - E[H_n(x, z | \theta)]|$   
1133  $\rightarrow_p 0$ ;  $n^{1/4} \sup |E[H_n(x, z | \theta)] - H(x, z)| \rightarrow_p 0$ ;  $n^{1/4} \sup |F_n(v | \theta) - E[F_n$   
1134  $(v | \theta)]| \rightarrow_p 0$ ;  $n^{1/4} \sup |E[F_n(v | \theta)] - F(v)| \rightarrow_p 0$ ;  $n^{1/4} \sup |L_n(z | \theta) - E[L_n$   
1135  $(z | \theta)]| \rightarrow_p 0$ ; and  $n^{1/4} \sup |E[L_n(z | \theta)] - L(z)| \rightarrow_p 0$ , where the suprema  
1136 are taken over  $\theta \in \Theta_n$  and  $v \in \mathbb{R}^d$ . We just show

$$1138 n^{1/4} \sup_{\theta \in \Theta_n} \sup_{v \in \mathbb{R}^d} |F_n(v | \theta) - E[F_n(v | \theta)]| \rightarrow_p 0 \quad (25)$$

$$1139 n^{1/4} \sup_{\theta \in \Theta_n} \sup_{v \in \mathbb{R}^d} |E[F_n(v | \theta)] - F(v)| \rightarrow_p 0, \quad (26)$$

1142 since the argument is the same for the other terms.

1144 The proof of (25) is a consequence of the facts that:  $n^{1/4} |F_n(v | \theta) -$   
1145  $E[F_n(v | \theta)]| \rightarrow_p 0$  for any  $\theta \in \Theta_n$  and  $v \in \mathbb{R}^d$  by a simple law of large  
1146 numbers, along with the stochastic equicontinuity result Lemma SE  
1147 established above. By the mean value theorem, we have for some  
1148 intermediate vector  $\bar{\theta}$ ,

$$1149 n^{1/4} |E[F_n(v | \theta)] - F(v)| = n^{1/4} \left| \frac{\partial F}{\partial \theta^T}(v | \bar{\theta})(\theta - \theta^0) \right| \\ 1150 \leq n^{-1/4} \sup_{\theta \in \Theta_n} \left| \frac{\partial F}{\partial \theta^T}(v | \theta) \right| \sup_{\theta \in \Theta_n} |n^{1/2}(\theta - \theta^0)| \\ 1151 \rightarrow 0, \\ 1152$$

1156 by assumption (A3).

1157 The result (24) is also a consequence of Lemma SE (see the un-named  
1158 proposition given in Andrews, 1994, p. 2251).

1160 **Proof of (22).** This follows from the stochastic equicontinuity of the  
1161 empirical process  $v_n(\gamma, v)$  and Taylor expansion of the mean, see Andrews

1162 (1994). Write  $0 = n^{1/2}\Delta_0(u, \theta^0) = n^{1/2}\Delta_0(u, \hat{\theta}) + [\partial\Delta_0(u, \theta^*)/\partial\theta]n^{1/2}(\hat{\theta} - \theta^0)$   
 1163 by the mean value theorem, where  $\theta^*$  are intermediate between  $\hat{\theta}$  and  $\theta^0$ .  
 1164 By the uniform continuity of  $\partial\Delta_0(u, \theta)/\partial\theta$  near  $\theta^0$ , we can replace  $\theta^*$  by  $\theta^0$ .  
 1165 Specifically, writing  $\Delta_0(u, \theta) = \Delta_n(u, \theta) - \{\Delta_n(u, \theta) - \Delta_0(u, \theta)\}$ , we obtain  
 1166

$$1167 \sqrt{n}\Delta_n(u, \theta) = \sqrt{n}\Delta_n(u, \theta^0) + \left\{ \frac{\partial\Delta_0(u, \theta^0)}{\partial\theta} + o(1) \right\} \sqrt{n}(\theta - \theta^0)$$

$$1168$$

$$1169 + \sqrt{n}\{\Delta_n(u, \theta) - \Delta_0(u, \theta)\} - \sqrt{n}\{\Delta_n(u, \theta^0) - \Delta_0(u, \theta^0)\}$$

$$1170$$

1171 for any  $\theta \in \Theta_n$ . We now invoke (18) and the triangle inequality to argue  
 1172 that the second line is  $o_p(1)$  uniformly in  $u \in \mathbb{R}^q$  and  $\theta \in \Theta_n$ . Finally,  
 1173 substituting in the expansion for  $n^{1/2}(\theta - \theta^0)$  we are done. Therefore,  
 1174

$$1175 CM_n^{**} = \int \left\{ \frac{1}{n} \sum_{i=1}^n \zeta(U_i, V(U, \theta^0) | \theta^0) \right\}^2 dP(U) + o_p(n^{-1}).$$

$$1176$$

$$1177$$

$$1178$$

1179 The result follows by interchanging summation and integration.

1180 Proof of (20). This follows from the fact that  $\sqrt{n} \sup_{u \in \mathbb{R}^q} |P_n(u) -$   
 1181  $P(u)| = O_p(1)$  and (21) and (22).  
 1182

1183 *Proof of (ii).* We have already shown that  $A_n(v | \theta)$  can be approximated  
 1184 by  $\Delta_n(v | \theta)$  with error of order smaller than  $n^{-1/2}$ . Also use the argument  
 1185 given in (22).  $\square$   
 1186

## 1187 ACKNOWLEDGEMENTS

1188 We both thank seminar participants for their helpful comments. The  
 1189 first author would like to thank Tilburg University for its hospitality  
 1190 and Pierre Chaussé for research assistance. Financial support from the  
 1191 National Science Foundation and the North Atlantic Treaty Organization  
 1192 is gratefully acknowledged.  
 1193

## 1194 REFERENCES

- 1197 Agresti, A. (1990). *Categorical Data Analysis*. New York: John Wiley.
- 1198 Andrews, D. W. K. (1994). Empirical process methods in econometrics. In: *Handbook of Econometrics*.  
 1199 Vol. IV. Engle, R. F., McFadden, D. L., eds. Elsevier Science B.V. Q1
- 1200 Andrews, D. W. K. (1995). Nonparametric kernel estimation for semiparametric models. *Econometric*  
 1201 *Theory* 11:560–596. Q2
- 1202 Andrews, D. W. K. (1997). A conditional Kolmogorov test. *Econometrica* 65:1097–1128.
- 1203 Bai, J. (1994). Weak convergence of the sequential residual empirical process in ARMA models.  
 1204 *The Annals of Statistics* 22:2051–2061.
- 1203 Beran, R., Millar, P. W. (1986). Confidence sets for a multivariate distribution. *The Annals of Statistics*  
 1204 14:431–443.

- 1205 Bierens, H. J., Ploberger, W. (1997). Asymptotic theory of integrated conditional moment tests.  
1206 *Econometrica* 65:1129–1151.
- 1207 Blum, J. R., Kiefer, J., Rosenblatt, M. (1961). Distribution free tests of independence based on the  
sample distribution function. *Annals of Mathematical Statistics* 32:485–498.
- 1208 Canepa, A., Godfrey, L. (2007). Improvement of the quasi-likelihood ratio test in ARMA models:  
1209 some results for bootstrap methods. *Journal of Times Series Analysis* 28:434–453.
- 1210 Chamberlain, G. (1987). Asymptotic efficiency in estimation with conditional moment restrictions.  
*Journal of Econometrics* 34:305–334.
- 1211 Chow, Y. S., Teicher, H. (1988). *Probability Theory*. 2nd ed. Berlin: Springer Texts in Statistics.
- 1212 Csörgő, S., Faraway, J. J. (1996). The exact and asymptotic distribution of Cramér von Mises  
1213 statistics. *Journal of the Royal Statistical Society Series B* 58:221–234.
- 1214 Dawid, A. P. (1979). Conditional independence in statistical theory. *Journal of the Royal Statistical  
1215 Society, Series B*. 41:1–31.
- 1216 Delgado, M. (1996). Testing serial independence using the sample distribution function. *Journal of  
1217 Time Series Analysis* 17:271–287.
- 1218 Durbin, J. (1973). Weak convergence of the sample distribution function when parameters are  
1219 estimated. *The Annals of Statistics* 1:279–290.
- 1218 Fan, Y., Li, Q. (1996). Consistent model specification tests: Omitted variables and semiparametric  
functional forms. *Econometrica* 64:865–890.
- 1219 Fernandes, M., Flores, R. G. (2001). Tests for Conditional Independence, Markovian Dynamics, and  
1220 Noncausality. <http://www.vwl.uni-mannheim.de/brownbag/flores.pdf> Q3
- 1221 Hall, P., Horowitz, J. (1996). Bootstrap critical values for tests based on Generalized Methods of  
1222 Moments estimation. *Econometrica*. Q2, Q4
- 1222 Godfrey, L. (1988). *Misspecification Tests in Econometrics: The Lagrange Multiplier Principle and Other  
1223 Approaches*. ix + pp. 252. Econometric Society Monograph No. 16. Cambridge: Cambridge  
1224 University Press, 1988, reprinted 1989.
- 1225 Linton, O., Gozalo, P. (1997). *Conditional Independence Restrictions: Testing and Estimation*. Discussion  
1226 Paper 1140, Cowles Foundation for Research in Economics, Yale Univ.
- 1227 Granger, C. W. J., Thomson, P. J. (1987). Predictive consequences of using conditioning or causal  
1228 variables. *Econometric Theory* 3:150–152. Q2
- 1228 Härdle, W., Janssen, P., Serfling, R. (1988). Strong uniform consistency rates for estimators of  
1229 conditional functionals. *Annals of Statistics* 16:1428–1449.
- 1229 Härdle, W., Linton, O. B. (1994). Applied nonparametric methods. *The Handbook of Econometrics*.  
1230 Vol. IV. McFadden, D. F., Engle III, R. F., eds. North Holland. Q5
- 1231 Han, A. K. (1987). A non-parametric analysis of transformations. *Journal of Econometrics* 35:191–209.
- 1232 Heckman, J. J., Ichimura, H., Smith, J., Todd, P. (1998). *Characterizing Selection Bias Using Experimental  
1233 Data*. *Econometrica*. Vol. 66, No. 5 (Sep., 1998), pp. 1017–1098.
- 1233 Hiemstra, C., Jones, J. D. (1994). Testing for linear and nonlinear Granger Causality in the Stock  
1234 Price-Volume Relation. *The Journal of Finance* 44:1639–1664. Q2
- 1234 Hoeffding, W. (1948). A non-parametric test of independence. *The Annals of Mathematical Statistics*  
1235 58:546–557.
- 1236 Hong, Y., White, H. (1995). Consistent specification testing via nonparametric series regression.  
*Econometrica* 63:1133–1159.
- 1237 Horowitz, J. (1992). A smoothed maximum score estimator for the binary response model.  
*Econometrica* 60:505–531.
- 1238 Horowitz, J. (1995). Bootstrap methods in econometrics: Theory and numerical performance. In:  
1239 Kreps, D., Wallis, K. W., eds. *Advances in Economics and Econometrics: 7th World Congress*,  
1240 Cambridge: Cambridge University Press, forthcoming. Q6
- 1241 Ichimura, H. (1993). Semiparametric least squares (SLS) and weighted SLS estimation of single-  
1242 index models. *Journal of Econometrics* 58:71–120.
- 1243 Joag-Dev, K. (1984). Measures of dependence. In: Krishnaiah, P. R., Sen, P. K., eds. *Handbook of  
1244 Statistics*, Vol. 4. Q5
- 1244 Kim, J., Pollard, D. (1990). Cube root asymptotics. *The Annals of Statistics* 18:191–219.
- 1245 Klein, R. W., Spady, R. H. (1993). An efficient semiparametric estimator for discrete choice models.  
*Econometrica* 61:387–421.
- 1246 Koul, H. L. (1996). Asymptotics of some estimators and sequential residual empiricals in nonlinear  
1247 time series models. *The Annals of Statistics* 24:380–404.

- 1248 Linton, O. B., Gozalo, P. L. (1995). A nonparametric test of conditional independence. *Cowles*  
 1249 *Foundation Discussion Paper* no 1106.
- 1250 Manski, C. F. (1975). The maximum score estimation of the stochastic utility model of choice.  
 1251 *Journal of Econometrics* 3:205–228.
- 1252 Manski, C. (1994). Analog estimation of econometric models. In: Engle III, R. F., McFadden, D. F.,  
 1253 eds. *The Handbook of Econometrics, vol. IV*. North Holland. **Q5**
- 1254 Nikitin, Y. (1995). *Asymptotic Efficiency of Nonparametric Tests*. Cambridge: Cambridge University Press.
- 1255 Pakes, A., Pollard, D. (1989). Simulation and the asymptotics of optimization estimators. *Econometrica*  
 1256 57:1027–1057.
- 1257 Phillips, P. C. B. (1988). Conditional and unconditional statistical independence. *Journal of*  
 1258 *Econometrics* 38:341–348.
- 1259 Pierce, D. A., Kopecky, K. J. (1979). Testing goodness of fit for the distribution of errors in  
 1260 regression model. *Biometrika* 66:1–5.
- 1261 Powell, J. L. (1994). Estimation in semiparametric models. In: Engle III, R. F., McFadden, D. F.,  
 1262 eds. *The Handbook of Econometrics*. Vol. IV. North Holland. **Q5**
- 1263 Robinson, P. M. (1991). Consistent nonparametric entropy-based testing. *Review of Economic Studies*  
 1264 58:437–453.
- 1265 Rosenbaum, P. R., Rubin, D. B. (1983). The central role of the propensity score in observational  
 1266 studies for causal effects. *Biometrika* 70:41–55.
- 1267 Sherman, R. P. (1993). The limiting distribution of the maximum rank correlation estimator.  
 1268 *Econometrica* 61:123–138.
- 1269 Shorack, G. R., Wellner, J. A. (1986). *Empirical Processes with Applications to Statistics*. New York: John  
 1270 Wiley.
- 1271 Skaug, H. J., Tjøstheim, D. (1993). A nonparametric test of serial independence based on the  
 1272 empirical distribution function. *Biometrika* 80:591–602.
- 1273 Song, K. (2009). Testing conditional independence via rosenblatt transforms. *The Annals of Statistics*  
 1274 37(6B):4011–4045.
- 1275 Su, L., White, H. (2008). A nonparametric Hellinger metric test for conditional independence.  
 1276 *Econometric Theory* 24:829–864.
- 1277 Van der Vaart, A. W., Wellner, J. A. (1996). *Weak Convergence and Empirical Processes*. Springer-Verlag:  
 1278 Berlin.
- 1279 Zheng, J. Z. (1997). A consistent specification test of independence. *Journal of Nonparametric Statistics*  
 1280 7(4). **Q7**
- 1281
- 1282
- 1283
- 1284
- 1285
- 1286
- 1287
- 1288
- 1289
- 1290