

Testing Conditional Moment Restrictions: A Partitioning Approach

Antonio Raiola^{*}

University Carlos III of Madrid

January 8, 2024

[\[Click here for the latest version\]](#)

Abstract

This paper proposes χ^2 tests for assessing the specification of regression models or general conditional moment restrictions (CMR). The data is partitioned based on the explanatory variables into several cells, and the χ^2 tests evaluate whether the difference between the observed average of the dependent variable and its expected value under the model specification, within each cell, arises by chance. In contrast to existing omnibus procedures, χ^2 tests are asymptotically pivotal and fairly insensitive to the curse of dimensionality. The computation of statistics is straightforward and does not require bootstrapping or smoothing techniques. Importantly, the asymptotic properties of the test are invariant to sample-dependent partitions, which can be chosen to favor certain alternatives. A Monte Carlo study provides evidence of the good performance of χ^2 tests using samples of small or moderate size compared with existing omnibus alternatives, particularly in high-dimensional settings. An empirical application regarding returns to education of African American students in the US complements the finite sample study.

JEL classification: C12, C31, C52

Keywords: Goodness-of-fit tests; Regression specifications; Distribution-free tests; Regressogram; Neyman-Pearson cells.

^{*}I am grateful to Miguel Angel Delgado for his guidance and support, and to Juan Carlos Escanciano, Carlos Velasco, Juan Jose Dolado, and Nazarii Salish for the many insightful discussions. I also thank the participants at the ENTER seminar at Toulouse School of Economics, as well as the participants of the Econometrics Reading Group at the University Carlos III of Madrid for their helpful comments. I gratefully acknowledge support from the Ministerio de Ciencia e Innovación through research grants PRE2020-092409, the AEI (Spain) through research grant PID2019-110779GB-I00, and University Carlos III for the scholarship for Master's study. Email address: araiola@eco.uc3m.es.

1 Introduction

The identification and estimation of causal relationships often rely on models defined by conditional moment restrictions (CMR), typically in the form of parametric regression models. Ensuring the correct model specification is crucial for drawing valid inferences and providing reliable interpretations. This article introduces χ^2 tests to check the goodness-of-fit of CMR model specifications. Leveraging their well-established usefulness in classical goodness-of-fit contexts, these tests offer a valuable complementary approach to existing methodologies, as discussed below.

The need to ensure correct model specification lies at the core of the many model-checking proposals in the literature. These proposals, rooted in the definition of the integrated regression function (IRF), expand upon classical goodness-of-fit tests designed for cumulative density function (CDF) specifications and apply them to the domain of regression analysis. The IRF generalizes the one-to-one relationship between density and the CDF to the regression framework, enabling the extension of all the tests proposed for CDFs to the context of regression. The alternative test statistics are characterized by a functional of the standard empirical process (SEP). The SEP is defined as the difference between the empirical IRF/CDF and its restricted counterpart, based on the assumed model specification, suitably scaled to stabilize the variance.

The omnibus tests are designed to assess whether any deviations of the SEP from zero arise by chance rather than due to model misspecification. Classical test statistics for CDF specifications, which are based on a norm of the SEP, such as the well-known Kolmogorov-Smirnov, Cramér-von Mises, or Anderson-Darling tests, have been extended to the context of regression model specifications by [Bierens \(1982\)](#), [Stute \(1997\)](#), and [Andrews \(1997\)](#), among others. Several tests are based on transformations of the SEP. For instance, test statistics based on Fourier transforms of the SEP, which in fact compare characteristic functions (e.g., [Koutrouvelis and Kellermeier 1981](#)), have been adapted to the regression case by [Bierens \(1982\)](#) and [Bierens and Ploberger \(1997\)](#). Tests based on martingale transforms of the SEP ([Khmaladze 1982](#)) have been extended by [Khmaladze and Koul \(2004\)](#). Additionally, specification tests for Lebesgue densities, which rely on the convolution of the SEP and a kernel function (e.g., [Bickel and Rosenblatt 1973](#)), have

been expanded to regression models by [Hardle and Mammen \(1993\)](#). All of these tests are known as minimum-distance tests.

When the model parameters are estimated from the sample, the SEP converges in distribution to a case-dependent process, and the tests are implemented with the help of bootstrap techniques. Two exceptions exist: (i) tests based on the SEP martingale transform, which converge in distribution to a Brownian motion with tabulated critical values; and (ii) tests based on the convolution of the SEP and a kernel function, suitably centered and standardized, which exhibit a standard normal limit null distribution when the bandwidth parameter vanishes at a suitable rate as the sample size grows. However, the size properties of convolution tests depend on the bandwidth choice, and in practice, the tests are implemented using bootstrap.

It is well known that tests based on the SEP exhibit poor power properties in the direction of high-frequency alternatives, as pointed out by [Durbin and Knott \(1972\)](#) based on the study of the SEP's principal components. [Stute \(1997\)](#) generalized these results to the regression case. Tests that rely on convolutions of the SEP and a kernel, with a vanishing bandwidth, are capable of detecting high-frequency alternatives but fail to detect the standard alternatives converging to the model in the null at the parametric rate. All of these tests are affected by the curse of dimensionality, showing limited power when the number of explanatory variables is large.

In practice, χ^2 tests, such as the [Pearson \(1900\)](#) test and its subsequent modifications, are among the most popular goodness-of-fit tests for CDF model specifications. These tests, after partitioning the data into a number of cells, say L , evaluates whether the difference between the observed and the expected frequencies within each cell occur by chance. The Pearson test statistics, for instance, can be represented as a quadratic form in the vector of differences between observed and expected frequencies. Note that the χ^2 tests are not omnibus: their goal is to detect alternatives where the true probability of the variable falling into a particular cell differs from the expected value specified by the model. However, they have important advantages with respect to tests based on SEP functionals. They are distribution-free and do not require the use of bootstrapping or smoothing techniques. Additionally, these tests are invariant to sample-dependent partitions, which can then be chosen to favor certain alternatives of interest.

Surprisingly, the χ^2 tests have not been extended to check the specification of CMR. This paper aims to fill this gap by proposing to group the data according to a partition of the explanatory variables into L classes, and considering χ^2 tests given by quadratic forms in the vector of differences between observed and expected (according to the model) sample averages within each class. The Pearson test analog is a weighted sum of these squared differences, which is in fact a J test for the orthogonality conditions between regression errors and the explanatory variables in each of the L classes. The paper also consider a generalized Wald test statistic: a quadratic form in the vector of differences between observed and expected sample averages using any \sqrt{n} -consistent estimator of the parameters in the model.

Much like the classical χ^2 test statistic, these tests have a limiting null distribution which is invariant to sample-dependent partitions. I propose algorithms for the calculation of balanced cells, each designed to contain approximately the same number of observations, as well as partitions intended to prioritize specific alternatives. I consider, for instance, tests based on Neyman-Pearson (NP) cells built upon deviations of the model under the null hypothesis from a predefined alternative, see [Balakrishnan, Voinov, and Nikulin \(2013\)](#). Under the correct specification of the alternative model, the NP partition maximizes the distance between the null and alternative models, thereby enhancing the power of the tests. In cases where there are no specific alternatives to prioritize, the recommendation is to utilize NP cells that compare the model under the null hypothesis with an auxiliary, flexible specification of the regression model.

Monte Carlo simulations show the good performances of the χ^2 tests compared to omnibus tests, particularly when the covariates dimension is high. The finite-sample study is complemented with an empirical illustration in which I apply the tests to analyze the returns of attending historically black college and universities (HBCU), relative to non-HCBU, for black students in the United States.

The structure of the paper is as follows: In the next section, I introduce the tests focusing on regression specifications. This facilitates the motivation and presentation and reduces the notational burden required for the more general case. In Section 3, I discuss the asymptotic properties of the tests using sample dependent partitions. Section 4 introduces partitioning algorithms designed to enhance the power of the tests. Sec-

tion 5 analyzes Pitman's local power. In Section 6, I extend the discussion to general CMR. Section 7 presents a Monte Carlo study. In the last section, I provide a real-data application.

2 Chi-Squared Tests for CMR

Let $\{Z_i\}_{i=1}^n = \{(Y_i, X_i')'\}_{i=1}^n$ be an i.i.d sample from the \mathbb{R}^{1+d_x} -valued random vector $Z = (Y, X)'$ with distribution P , where Y is the response variable, and X is the d_x -dimensional vector of explanatory variables with support in $\mathcal{X} \subset \mathbb{R}^{d_x}$. This paper considers testing the parametric model specification of the regression function, $m(x) = \mathbb{E}[Y|X = x]$. The correct specification hypothesis is stated as,

$$H_0 : m \in \mathcal{M} \tag{1}$$

with $\mathcal{M} = \{m_\theta(\cdot) : \theta \in \Theta\}$ and $\Theta \subset \mathbb{R}^{d_\theta}$ being a family of parametric regression functions and a suitable parameter space, respectively. Thus, under H_0 , there exists a $\theta_0 \in \Theta$ such that the IRF, $M(x) = \mathbb{E}[Y\mathbb{I}_{(-\infty, x]}(X)]$, and its version imposing the null hypothesis restrictions, $M_\theta(x) = \mathbb{E}[m_\theta(X)\mathbb{I}_{(-\infty, x]}(X)]$, are equal in all the Borel sets of \mathcal{X} , i.e.

$$M\{A\} = M_{\theta_0}\{A\} \text{ for all } A \in \mathcal{B}^{d_x}, \tag{2}$$

where \mathcal{B}^{d_x} represent the Borel sigma-field of \mathbb{R}^{d_x} , $M\{A\} = \int_A M(dx)$ and $M_\theta\{A\} = \int_A M_\theta(dx)$ denote the Lebesgue-Stieltjes measure of M and M_θ over A , respectively. Recall that (2) is the definition of the regression function for the specification in H_0 (e.g., see definition 34.1 in Billingsley (2017)), which is equivalent to the following orthogonality conditions,

$$H_0 : \mathbb{E}[\varepsilon_{\theta_0}|X] = 0 \text{ a.s. for some } \theta_0 \in \Theta,$$

with $\varepsilon_\theta(z) = y - m_\theta(x)$ denoting the regression error.

Intuitively, when X takes values in a finite set and θ_0 is known (i.e., under the simple hypothesis), a test for H_0 , following the proposal by Pearson (1900), consists of comparing the mean of Y given $x \in \mathcal{X}$ with $m_{\theta_0}(x)$ for all $x \in \mathcal{X}$. In this case, such a test is indeed

omnibus (i.e. the test detect all the possible alternatives). For any type of covariates, once the data is partitioned into L cells, say, χ^2 tests assess whether the difference between the expected and observed averages in each cell arose by chance.

Consider the partition $\gamma = (\gamma_1, \dots, \gamma_L)$ of the covariates support \mathcal{X} into L cells and let $\mathbf{I}_\gamma(x) = (\mathbb{I}_{\gamma_1}(x), \dots, \mathbb{I}_{\gamma_L}(x))'$ denote the vector of indicator functions over the sets within γ . The building block of the χ^2 test statistics is the standardized vector of differences between the observed averages of Y and the corresponding averages under the specification of H_0 in each cell,

$$\hat{\Phi}_\gamma(\theta) = \sqrt{n} (\hat{\mu}_\gamma^0 - \hat{\mu}_\gamma(\theta)) = \sqrt{n} (\hat{\mu}_1^0 - \hat{\mu}_1(\theta), \dots, \hat{\mu}_L^0 - \hat{\mu}_L(\theta))', \quad (3)$$

where $\hat{\mu}_l^0 = \hat{M}\{\gamma_l\} = n^{-1} \sum_{i=1}^n Y_i \mathbb{I}_{\gamma_l}(X_i)$ represents the empirical IRF evaluated at γ_l , with $\hat{M}(x) = n^{-1} \sum_{i=1}^n Y_i \mathbb{I}_{(-\infty, x]}(X_i)$, and $\hat{\mu}_l(\theta) = \hat{M}_\theta\{\gamma_l\} = n^{-1} \sum_{i=1}^n m_\theta(X_i) \mathbb{I}_{\gamma_l}(X_i)$ represents its restricted counterpart under H_0 , with $\hat{M}_\theta(x) = n^{-1} \sum_{i=1}^n m_\theta(X_i) \mathbb{I}_{(-\infty, x]}(X_i)$.

Under a simple hypothesis, i.e. when θ_0 is known, by the central limit theorem, under the null,

$$\hat{\Sigma}_\gamma(\theta_0)^{-1/2} \hat{\Phi}_\gamma(\theta_0) \xrightarrow{d} N(0, I_L), \quad (4)$$

where $\hat{\Sigma}_\gamma(\theta) = n^{-1} \sum_{i=1}^n \varepsilon_\theta^2(Z_i) \mathbf{I}_\gamma(X_i) \mathbf{I}_\gamma(X_i)'$ estimates,

$$\Sigma_{\gamma,0} = \text{Avar} \left(\hat{\Phi}_\gamma(\theta_0) \right) = \text{diag}\{\sigma_{0,1}^2, \dots, \sigma_{0,L}^2\}, \quad (5)$$

under H_0 , with $\sigma_{0,l}^2 = \sigma_l^2(\theta_0)$, and $\sigma_l^2(\theta) = \mathbb{E}[\varepsilon_\theta^2(Z) \mathbb{I}_{\gamma_l}(X)]$. Thus, taking for granted that $\rho_l = \mathbb{E}[\mathbb{I}_{\gamma_l}(X)] > 0$ for all l , under H_0 ,

$$\hat{\chi}_{\gamma,0}^2(\theta_0) \xrightarrow{d} \chi_L^2, \quad (6)$$

where

$$\hat{\chi}_{\gamma,0}^2(\theta) = \hat{\Phi}_\gamma(\theta)' \hat{\Sigma}_\gamma(\theta_0)^{-1} \hat{\Phi}_\gamma(\theta) = n \sum_{l=1}^L \frac{(\hat{\mu}_l^0 - \hat{\mu}_l(\theta))^2}{\hat{\sigma}_{0,l}^2}, \quad (7)$$

and $\hat{\sigma}_{0,l}^2 = n^{-1} \sum_{i=1}^n \varepsilon_{\theta_0}^2(Z_i) \mathbb{I}_{\gamma_l}(X_i)$. The test statistic $\hat{\chi}_{\gamma,0}^2(\theta_0)$ extends the conventional Pearson's chi-squared test to check regression model specifications. It does so by replacing the comparison of observed frequencies with expected frequencies under the null

hypothesis with a comparison of averages.

Of course, tests based on (3) are not omnibus, but designed for detecting deviations from H_0 of the type,

$$H_1(\gamma) : \boldsymbol{\mu}_\gamma^0 \neq \boldsymbol{\mu}_\gamma(\theta) \text{ for all } \theta \in \Theta, \quad (8)$$

where $\boldsymbol{\mu}_\gamma^0 = (\mu_{\gamma_1}^0, \dots, \mu_{\gamma_L}^0)' = (M\{\gamma_1\}, \dots, M\{\gamma_L\})'$ is the vector of expected sample averages of Y in each cell and $\boldsymbol{\mu}_\gamma(\theta) = (\mu_{\gamma_1}(\theta), \dots, \mu_{\gamma_L}(\theta))' = (M_\theta\{\gamma_1\}, \dots, M_\theta\{\gamma_L\})'$ is the vector of expected averages under H_0 . The partitions, therefore, cover a fundamental role in implementing the test and provide a flexible tool to exploit out-of-sample information on the possible alternatives. Neyman-Pearson (NP) cells (Balakrishnan, Voinov, and Nikulin 2013), for instance, consist of the points in \mathcal{X} where the model under the null and a pre-specified alternative parametrization of $m(\cdot)$ meet (see Example 1 below). Under the pre-specified alternative, NP classes maximize the L2 norm of $\hat{\boldsymbol{\Phi}}_\gamma$.

Example 1 (NP Classes)

Consider testing the linear model $m_{\theta_0}(X_i) = \theta_{00} + X_i\theta_{01}$ against the alternative specification,

$$H_1 : m_{1,\theta_0^*}(X_i) = \theta_{00} + \theta_{01}X_i + \theta_{02} \sin\left(\frac{50X_i}{2\pi}\right)$$

where $\theta_0^* = (\theta_0, \theta_{02})$ (and, thus, θ_0) is a known vector. NP classes split \mathcal{X} over the points where $\theta_{02} \sin\left(\frac{50X_i}{2\pi}\right) = 0$. As a result, under H_1 , $m_{\theta_0}(\cdot)$ is strictly bigger or strictly smaller than $m_{1,\theta_0^*}(\cdot)$ within each cell, and most cell-specific errors have the same sign, implying that the average error of a single cell is larger than the average error over the union of two contiguous cells (in absolute terms). As a matter of fact, in this example, the average error of (any) two contiguous cells is close to zero. See Figure 1 for a graphical representation of the partition.

When θ_0 is unknown, the criterion in (7) suggests the following minimum distance estimator, hereafter referred to as the grouped GMM estimator:

$$\hat{\theta}_\gamma = \arg \min_{\theta \in \Theta} \chi_\gamma^2(\theta) \quad (9)$$

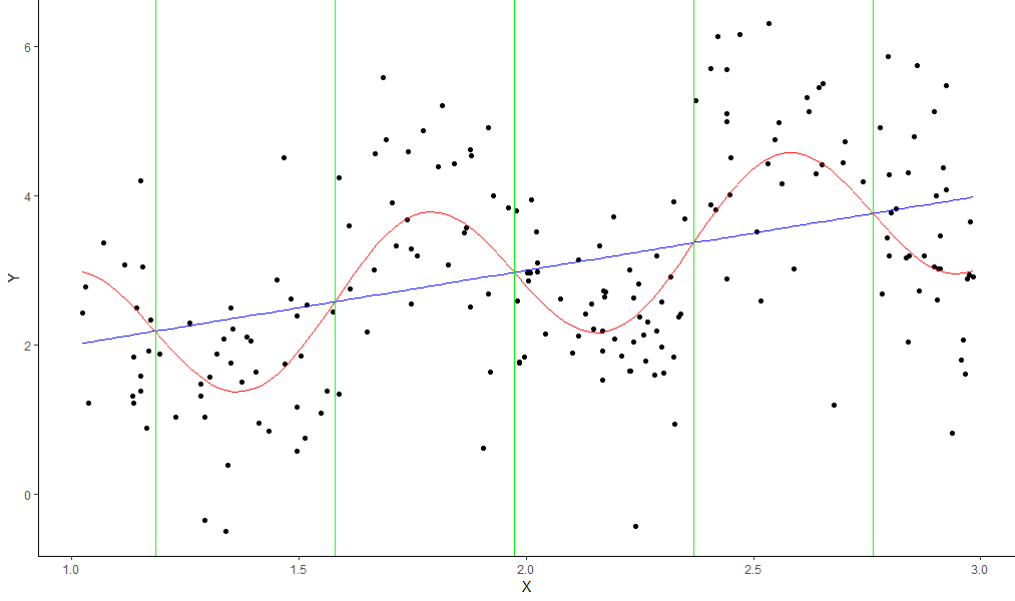


Figure 1: The graph depicts a random draw from the model under H_1 with $\theta_0^* = (1, 1, 1)$. The green lines depict the points where the model under the null (blue line) and under the alternative (red line) meet.

where

$$\hat{\chi}_\gamma^2(\theta) = \hat{\Phi}_\gamma(\theta)' \hat{\Sigma}_\gamma(\tilde{\theta})^{-1} \hat{\Phi}_\gamma(\theta) = n \sum_{l=1}^L \frac{(\hat{\mu}_l^0 - \hat{\mu}_l(\theta))^2}{\hat{\sigma}_l^2(\tilde{\theta})},$$

$\hat{\sigma}_l^2(\theta) = n^{-1} \sum_{i=1}^n \varepsilon_\theta^2(Z_i) \mathbb{I}_{\gamma_l}(X_i)$, and $\tilde{\theta}$ is some initial \sqrt{n} -consistent estimator of θ_0 . The estimator $\hat{\theta}_\gamma$ is analogous to the multinomial maximum-likelihood estimator (or minimum χ^2 estimator) in the classical case (see [Cramér 1946](#)). Under linear null hypothesis, i.e. $m_\theta(x) = x'\theta$, it is a feasible GLS estimator based on the aggregated data $\{\bar{Y}_l, \bar{X}_l\}_{l=1}^L$,

$$\hat{\theta}_\gamma = \left[\sum_{l=1}^L \frac{\bar{X}_l \bar{X}_l'}{\hat{\sigma}_l^2(\tilde{\theta})} \right]^{-1} \sum_{l=1}^L \frac{\bar{X}_l \bar{Y}_l}{\hat{\sigma}_l^2(\tilde{\theta})},$$

with $\bar{Y}_l = n^{-1} \sum_{i=1}^n Y_i \mathbb{I}_{\gamma_l}(X_i)$ and $\bar{X}_l = n^{-1} \sum_{i=1}^n X_i \mathbb{I}_{\gamma_l}(X_i)$. In the non-linear case, we can use a feasible asymptotically efficient one-step ahead Gauss-Newton estimator starting from any preliminary \sqrt{n} -consistent estimator $\tilde{\theta}$,

$$\hat{\theta}_\gamma^{(1)} = \tilde{\theta} - \left[\sum_{l=1}^L \frac{\hat{\mu}_l^*(\tilde{\theta}) \hat{\mu}_l^*(\tilde{\theta})'}{\hat{\sigma}_l^2(\tilde{\theta})} \right]^{-1} \sum_{l=1}^L \frac{\hat{\mu}_l^*(\tilde{\theta}) (\hat{\mu}_l^0 - \hat{\mu}_l(\tilde{\theta}))}{\hat{\sigma}_l^2(\tilde{\theta})}$$

with $\hat{\mu}_l^*(\theta) = n^{-1} \sum_{i=1}^n \nabla m_\theta(X_i) \mathbb{I}_{\gamma_l}(X_i)$ and $\nabla m_{\tilde{\theta}} = d/d\theta' m_\theta|_{\theta=\tilde{\theta}}$. The estimator belongs

to the class of minimum-distance estimators considered by [Cristobal, Roca, and Manteiga \(1987\)](#) and [Koul and Ni \(2004\)](#), with the main difference that the regressogram is used instead of kernels, and the weighting introduced to improve efficiency. Under suitable regularity conditions and the null H_0 ,

$$\sqrt{n}(\hat{\theta}_\gamma - \theta_0) \xrightarrow{d} N\left(0, \left[\boldsymbol{\mu}_{\gamma,0}^{*'}(\Sigma_{\gamma,0})^{-1}\boldsymbol{\mu}_{\gamma,0}^*\right]^{-1}\right),$$

where $\boldsymbol{\mu}_{\gamma,0}^* = \boldsymbol{\mu}_\gamma^*(\theta_0)$, with $\boldsymbol{\mu}_\gamma^*(\theta) = (\mu_1^*(\theta), \dots, \mu_L^*(\theta))'$ and $\mu_l^*(\theta) = \mathbb{E}[\hat{\mu}_l^*(\theta)]$ denoting the matrix of partial derivatives of $\boldsymbol{\mu}_\gamma(\theta)$.

The $\hat{\chi}^2$ test statistics is, for $L > d_\theta$,

$$\hat{\chi}_\gamma^2 = \min_{\theta \in \Theta} \hat{\chi}_\gamma^2(\theta) = \hat{\chi}_\gamma^2(\hat{\theta}_\gamma), \quad (10)$$

which is in fact a J test on the set of the L , out of the many, orthogonality conditions implied by the null,

$$\mathbb{E}[Y\mathbb{I}_l(X)] = \mathbb{E}[m_{\theta_0}(X)\mathbb{I}_l(X)] \quad \text{for all } l \in \{1, 2, \dots, L\}.$$

Thus, under H_0 , and for $L > d_\theta$,

$$\hat{\chi}_\gamma^2 \xrightarrow{d} \chi_{L-d_\theta}^2. \quad (11)$$

Is also well motivated, as suggested in classical goodness-of-fit χ^2 tests (e.g. [Nikulin 1973](#) and [Rao and Robson 1974](#)), using the Wald testing principle based on $\hat{\boldsymbol{\Phi}}_\gamma(\tilde{\theta})$, employing any \sqrt{n} -consistent estimator $\tilde{\theta}$,

$$\hat{\mathcal{W}}_\gamma(\tilde{\theta}) = \hat{\boldsymbol{\Phi}}_\gamma(\tilde{\theta}) \widehat{\text{Avar}}^- \left(\hat{\boldsymbol{\Phi}}_\gamma(\tilde{\theta}) \right) \hat{\boldsymbol{\Phi}}_\gamma(\tilde{\theta}), \quad (12)$$

where $\widehat{\text{Avar}}^- \left(\hat{\boldsymbol{\Phi}}_\gamma(\tilde{\theta}) \right)$ is a consistent estimator of some generalized inverse of $\text{Avar} \left(\hat{\boldsymbol{\Phi}}_\gamma(\tilde{\theta}) \right)$, $\text{Avar}^- \left(\hat{\boldsymbol{\Phi}}_\gamma(\tilde{\theta}) \right)$ say. Assuming $\widehat{\text{Avar}}^- \left(\hat{\boldsymbol{\Phi}}_\gamma(\tilde{\theta}) \right) \xrightarrow{p} \text{Avar}^- \left(\hat{\boldsymbol{\Phi}}_\gamma(\tilde{\theta}) \right)$, and suitable regularity conditions, under H_0 ,

$$\hat{\mathcal{W}}_\gamma(\tilde{\theta}) \xrightarrow{d} \chi_{r(\text{Avar}(\hat{\boldsymbol{\Phi}}_\gamma(\tilde{\theta})))}^2, \quad (13)$$

where for a given square matrix A , $r(A)$ denotes its rank.

Taking for granted the asymptotic linearity of the estimator (see Assumption 3 in the next section), the covariance matrix of $\hat{\Phi}_\gamma(\tilde{\theta})$ is characterized as,

$$\text{Avar}\left(\hat{\Phi}_\gamma(\tilde{\theta})\right) = \Sigma_{\gamma,0} - \boldsymbol{\mu}_{\gamma,0}^* C'_{\gamma,0} - C_{\gamma,0} \boldsymbol{\mu}_{\gamma,0}^{*'} + \boldsymbol{\mu}_{\gamma,0}^* L_0 \boldsymbol{\mu}_{\gamma,0}^{*'}, \quad (14)$$

where $C_{\gamma,0} = \mathbb{E}[\varepsilon_{\theta_0}(Z) \mathbf{I}_\gamma(X) l_{\theta_0}(Z)']$, $L_0 = \mathbb{E}[l_{\theta_0}(Z) l_{\theta_0}(Z)']$, and $l_{\theta_0}(\cdot)$ is the influence function of $\tilde{\theta}$. When $\text{Avar}\left(\hat{\Phi}_\gamma(\tilde{\theta})\right)$ is full rank (e.g., if $\varepsilon_{\theta_0}(\cdot) \mathbf{I}_\gamma(\cdot)$ and $l_{\theta_0}(\cdot)$ have linearly independent components), the Wald test can be performed on any finite splitting of the data. In this case, a valid choice of $\widehat{\text{Avar}}^-\left(\hat{\Phi}_\gamma(\tilde{\theta})\right)$ is given by the inverse of,

$$\hat{W}_\gamma(\tilde{\theta}) = \hat{\Sigma}_\gamma(\tilde{\theta}) - \hat{\boldsymbol{\mu}}_\gamma^*(\tilde{\theta}) \hat{C}_\gamma(\tilde{\theta})' - \hat{C}_\gamma(\tilde{\theta}) \hat{\boldsymbol{\mu}}_\gamma^{*'}(\tilde{\theta}) + \hat{\boldsymbol{\mu}}_\gamma^*(\tilde{\theta}) \hat{L}(\tilde{\theta}) \hat{\boldsymbol{\mu}}_\gamma^{*'}(\tilde{\theta}), \quad (15)$$

where $\hat{\boldsymbol{\mu}}_\gamma^*(\theta) = (\hat{\mu}_1^*(\theta), \dots, \hat{\mu}_L^*(\theta))'$, $\hat{C}_\gamma(\theta) = n^{-1} \sum_{i=1}^n \varepsilon_\theta(Z_i) \mathbf{I}_\gamma(X_i) l_\theta(Z_i)'$, and $\hat{L}(\theta) = n^{-1} \sum_{i=1}^n l_\theta(Z_i) l_\theta(Z_i)'$. If the covariance matrix is rank deficient and the Moore-Penrose inverse of $\hat{W}_\gamma(\tilde{\theta})$, denoted as $\hat{W}_\gamma^+(\tilde{\theta})$, has rank converging in probability to the one of $\text{Avar}\left(\hat{\Phi}_\gamma(\tilde{\theta})\right)$, then $\hat{W}_\gamma^+(\tilde{\theta}) \xrightarrow{p} \text{Avar}^+\left(\hat{\Phi}_\gamma(\tilde{\theta})\right)$ (Theorem 2 of [Andrews 1987](#)). However, this need not be the case ([Schott 2016](#), p. 222-224) and more complex methods might be required (see, e.g., [Lütkepohl and Burda 1997](#)). Consider, for instance, the nonlinear least squares (NLLS),

$$\tilde{\theta} = \arg \min_{\theta \in \Theta} \sum_{i=1}^n \varepsilon_\theta(Z_i)^2.$$

Under homoskedasticity, i.e. $\text{Var}(\varepsilon_{\theta_0}(Z)|X=x) = \sigma_0^2 = \sum_{l=1}^L \sigma_{0,l}^2$, it holds that, $l_{\theta_0}(z) = \Psi_0^{-1} \nabla m_{\theta_0}(x) \varepsilon_{\theta_0}(z)$, where $\Psi_0 = \mathbb{E}[\nabla m_{\theta_0}(X) \nabla m_{\theta_0}(X)']$. Under H_0 , when $\tilde{\theta}$ is asymptotically more efficient than $\hat{\theta}_\gamma$, i.e., when $\text{Avar}(\hat{\theta}_\gamma) - \text{Avar}(\tilde{\theta})$ is p.d., $\text{Avar}\left(\hat{\Phi}_\gamma(\tilde{\theta})\right) = \Sigma_{\gamma,0} - \sigma_0^2 \boldsymbol{\mu}_{\gamma,0}^* \Psi_0^{-1} \boldsymbol{\mu}_{\gamma,0}^{*'}$ is also p.d., and $\hat{W}(\tilde{\theta}) \xrightarrow{d} \chi_L$. A similar reasoning holds for the probit/logit maximum likelihood estimator (MLE), as discussed in Section 7. Under heteroskedasticity, however, $\text{Avar}\left(\hat{\Phi}_\gamma(\tilde{\theta})\right)$ is not necessarily an invertible matrix.

Notice that the estimation of the covariance matrix, when $\text{Avar}\left(\hat{\Phi}_\gamma(\tilde{\theta})\right)$ is full rank, can be avoided by using a random normalizing weighting matrix, as in [Kuan and Lee \(2006\)](#) (see also [Kiefer, Vogelsang, and Bunzel 2000](#) for an early reference). Similarly, the

over-identification test of [Lee, Kuan, and Hsu \(2014\)](#) provides a robust version of the $\hat{\chi}^2$ test which does not require estimating $\hat{\Sigma}_\gamma(\tilde{\theta})$. In these cases, the limit null distribution is non-standard but pivotal.

3 Data-dependent Cells

When studying the large sample behavior of the statistics, it is crucial to address the inherent influence of the data on the selection of cells ([Watson 1959](#)). [Moore and Spruill \(1975\)](#) were among the first to address this concern in the distribution model check literature providing a rigorous derivation of the null distribution of χ^2 tests with rectangular data-dependent cells. In a more general setting, [Pollard \(1979\)](#) established the result for cells of arbitrary form using uniform results for empirical processes indexed by sets, while [Andrews \(1988a,b\)](#) and [Delgado and Vainora \(2023\)](#) applied the methodology to conditional distribution testing. This section provides a similar result for the more general CMR testing framework, showing that the grouped GMM estimator and the tests keep standard limit distribution when the partition is built with data-dependent cells.

A minimal set of assumptions, consisting of smoothness conditions and restrictions on the partitioning algorithm complexity, allows to state the convergence results in a self-contained fashion.

Assumption 1 (a) $\{Z_i = (Y_i, X_i')'\}_{i=1}^n$ is a sequence of i.i.d. random vectors with $\mathbb{E}|Y_i| < \infty$; (b) $\mathbb{E}[\varepsilon_{\theta_0}^2] < C$, with $C < \infty$; (c) Θ is a compact subset of \mathbb{R}^{d_θ} and θ_0 is an interior point of Θ .

Assumption 2 $m_\theta(\cdot)$ is twice continuously differentiable in a neighborhood \mathcal{N}_{θ_0} of θ_0 , with $\mathcal{N}_{\theta_0} \subset \Theta$. The gradient, $\nabla m_\theta(\cdot) = d/d\theta m_\theta(\cdot)$, is bounded by a square-integrable function $R(\cdot)$ such that $\sup_{\theta \in \Theta_0} |\nabla^{(j)} m_\theta(\cdot)| \leq R(\cdot)$ for all $j \in \{1, \dots, d_\theta\}$, where $\nabla^{(j)}$ denotes the j -th partial derivative, and $\mathbb{E}[R(X)^2] < \infty$.

Assumption 3

(a) The estimator $\tilde{\theta}$ satisfies the following asymptotic expansion under the null,

$$\sqrt{n}(\tilde{\theta} - \theta_0) = \frac{1}{\sqrt{n}} \sum_{i=1}^n l_{\theta_0}(Z_i) + o_p(1)$$

where $\mathbb{E}[l_{\theta_0}(Z)] = 0$ and $L_0 = \mathbb{E}[l_{\theta_0}(Z)l_{\theta_0}(Z)']$ is a finite and non-singular matrix.

(b) The vector-valued function $l_{\theta}(\cdot)$ is twice continuously differentiable in a neighborhood Θ_0 of θ_0 with first partial derivatives bounded by a square-integrable function $R_2(\cdot)$ such that $\sup_{\theta \in \Theta_0} |\nabla^{(j)} l_{\theta}(\cdot)| \leq R_2(\cdot)$ for all $j \in \{1, \dots, d_{\theta}\}$ and $\mathbb{E}[R_2(Z)^2] < \infty$.

Assumptions 1 and 2 are common in the model check literature (see [Stute and Zhu 2002](#), for instance). Compared to papers developing χ^2 tests based on probability model (e.g., [Tauchen 1985](#)), these require slightly higher smoothness condition of the regression function, but leave completely unrestricted the data distribution. Assumption 3(a) holds for most of the estimators used in practice, such as least squares or GMM estimators, as well as for identification-robust minimum-distance estimators (see, e.g., [Domínguez and Lobato 2004](#)). While Assumption 3(b) is a technical requirement for the consistency of the plug-in estimator $\hat{W}_{\gamma}(\tilde{\theta})$.

Let also state the necessary global identification and finite-variance conditions for the consistency and asymptotic normality of the grouped GMM estimator.

Assumption 2' (a) $\Sigma_{\gamma,0}$ is positive definite (p.d.); (b) $\mathbb{E}[Y\mathbf{I}_{\gamma}] = \mathbb{E}[m_{\theta}(X)\mathbf{I}_{\gamma}]$ if and only if $\theta = \theta_0$; (c) $\boldsymbol{\mu}_{\gamma,0}^*$ is full rank.

Following [Pollard \(1979\)](#) and [Andrews \(1988a\)](#), the data-dependent partitions are modeled as random functions over a class of properly restricted measurable sets. Specifically, let \mathbb{C} be a class of measurable sets in \mathcal{X} from which the cells of each partition are drawn, and denote as \mathbb{D} the class of partitions of \mathcal{X} comprised of L sets from \mathbb{C} (L is fixed for all n); that is,

$$\mathbb{D} = \left\{ \gamma = (\gamma_1, \dots, \gamma_L) \in \mathbb{C}^L : \bigcup_{l=1}^L \gamma_l = \mathcal{X}, \gamma_l \cap \gamma_f = \emptyset, \forall l \neq f \right\}, \quad (16)$$

where γ_l and γ_f denote sets of the partition γ . We equip \mathbb{C} with the topology generated

by the $L^2(F_x)$ semi-norm, F_x being the distribution of X under P , and give \mathbb{D} the corresponding product topology. This means that two set C_1, C_2 in \mathcal{X} are close if $F_x\{C_1 \Delta C_2\}$ is small, where Δ denotes the symmetric difference operator, $C_1 \Delta C_2 = C_1 \cup C_2 \setminus C_1 \cap C_2$.

For each sample size n the corresponding partition is determined by a measurable mapping, denoted as $\hat{\gamma} = (\hat{\gamma}_1, \dots, \hat{\gamma}_L)$, from the underlying probability space to \mathbb{D} , which converges in probability to some fixed partition of cells in \mathbb{D} . In other words, for all $\epsilon > 0$,

$$P(F_x\{\hat{\gamma}_l \Delta \gamma_l\} > \epsilon) \rightarrow 0 \text{ as } n \rightarrow \infty, \text{ for all } l = 1, 2, \dots, L.$$

Assumption 4 *Under H_0 , $\hat{\gamma} \xrightarrow{P} \gamma$ for some fixed set of cells $\gamma \in \mathbb{D}$*

Assumption 4 represents a standard requirement in the literature concerning empirical processes indexed by sets, and it is satisfied by a broad range of partitioning algorithms. For example, this assumption is met when the cells depend on $\tilde{\theta}$ in a continuous fashion, and $\tilde{\theta}$ converges in probability to a non-random vector ([Andrews 1988b](#)). However, it is less clear whether cells generated by the comparison of different estimators meet this requirement, as discussed in the next section.

It is also possible to relax Assumption 4 to cases where $\hat{\gamma}$ converges to a random element $\gamma_0 \in \mathbb{D}$. This is achieved by replacing the convergence in probability with a uniform tightness condition and ensuring asymptotic independence between $\hat{\gamma}$ and $\hat{\Phi}_{\gamma_0}(\theta_0)$. Refer to Section 6 of [Pollard \(1979\)](#) for further discussions

Crucially, deriving the limit null distribution requires bounding the complexity of the partitions employed for test construction. This is achieved by assuming that the cells are drawn from a Vapnik-Cervonenkis (VC) class.

Assumption 5 *\mathbb{C} is a VC class of sets.*

This assumption is convenient because it is independent of the data distribution P , yet it is general enough for our purposes. For example, algorithms that generate cells with a finite number of straight edges and the class of hyper ellipsoids are VC classes. Furthermore, unions, intersections, differences, and complements of VC classes are also VC classes ([Andrews 1988a](#) and [Pollard 1984](#) provide a thorough discussion, see also Section 2.6 in [Van Der Vaart 1996](#)). A less stringent condition assumes that \mathbb{C} is a

Donsker class for the underlying probability measure (Pollard 1979). While this allows for a wider range of admissible partitions, verifying Donsker assumptions may prove more challenging in practice.

The following theorems show that the asymptotic distribution of the grouped GMM estimator and of the test statistics are unaffected by data-dependent cells. All the proofs are relegated to the appendix.

Theorem 1 *Let Assumptions 1, 2, 2', 4, 5 hold. Then, under the null hypothesis H_0 ,*

$$\sqrt{n}(\hat{\theta}_{\hat{\gamma}} - \theta_0) \xrightarrow{d} N\left(0, \left[\boldsymbol{\mu}_{\gamma,0}^{*'}(\Sigma_{\gamma,0})^{-1}\boldsymbol{\mu}_{\gamma,0}^*\right]^{-1}\right),$$

Theorem 2 *Let Assumptions 1, 2, 4, 5, and the null hypothesis H_0 hold. Then,*

(a) *Under Assumption 2' and $L > d_\theta$,*

$$\hat{\chi}_{\hat{\gamma}}^2 \xrightarrow{d} \chi_{L-d_\theta}^2.$$

(b) *Under Assumption 3, and $\widehat{\text{Avar}}^-\left(\hat{\boldsymbol{\Phi}}_{\hat{\gamma}}(\tilde{\theta})\right) \xrightarrow{p} \text{Avar}^-\left(\hat{\boldsymbol{\Phi}}_{\hat{\gamma}}(\tilde{\theta})\right)$,*

$$\hat{\mathcal{W}}_{\hat{\gamma}}(\tilde{\theta}) \xrightarrow{d} \chi_{r(\text{Avar}(\hat{\boldsymbol{\Phi}}_{\hat{\gamma}}(\tilde{\theta})))}^2.$$

The upcoming section introduces a variety of partitioning procedures, discussing methods based on unsupervised clustering and model-based techniques that leverage information on the null hypothesis and pre-specified alternatives.

4 Partitioning Procedures

The partitioning algorithms can be classified into two categories: covariates-based methods and model-based methods, depending on the information used to split the sample. The first classifies the data exclusively based on the characteristics of the population. This includes simple methods such as dividing the data into cubic cells, as well as more complex methodologies like k-means clustering, hierarchical clustering, and nearest-neighbor clustering, to mention only a few (see chapter 14 of Hastie, Tibshirani, and Friedman

2009 for a review of unsupervised clustering methods). In contrast, model-based methods exploit information about the model under the null hypothesis and about pre-specified alternatives to increase the power of the tests in those directions. The objective of these procedures is to construct partitions where the distance between the vectors of average predictions under the null hypothesis model and the model under the alternative is maximized.

This section, first explore the construction of partitions based on the union of statistically equivalent blocks (SEB) (Gessaman 1970). The method provides a simple and intuitive approach for partitioning the data into balanced cells, ensuring that each cell contains an (approximately) equal number of observations. It can be directly applied to the covariate matrix or to a lower-dimensional data projection.

Next, the discussion will turn to the construction of Neyman-Pearson (NP) classes. These partitioning algorithms aim to maximize the probability of rejecting the null hypothesis by comparing the parametric fit under the null with the parametric fit under a pre-specified alternative. This method is particularly useful when the researcher's primary concern is rejecting a subset of deviations critical for the application at hand.

When the alternative is left unrestricted, the recommendation is to construct NP classes using an auxiliary flexible specification of the regression model (see Davidson, MacKinnon, et al. 2004 Ch. 15.2 for a review of tests based on artificial regressions). The procedure assumes a structured alternative model, which allows to capture deviations from the null hypothesis by exploiting the conditional mean dependence of the parametric residuals from the fitted values under the alternative. Tests built with these partitions can be seen as a formalization of the statistical practice of looking at the residual-fitted values scatter plot to capture model misspecification and heteroskedasticity (e.g., Cook 1994).

Lastly, I explore the convergence in probability of NP classes to fixed cells in \mathbb{D} . Since these partitioning methods rely on comparing estimated parameters, understanding the convergence to fixed cells is not straightforward and requires additional investigation. Covariates-based partitioning methods are already discussed in Andrews (1988a) and Andrews (1988b).

4.1 Union of Statistically Equivalent Blocks (SEB)

This algorithm offers a straightforward and intuitive approach to partition the data while ensuring that each cell contains approximately an equal number of observations. The SEB algorithm is easy to apply and aligns with statistical practice, which involves dividing the data into equiprobable cells when testing against unknown probability distributions (see, for instance, [Greenwood and Nikulin 1996](#)).

Let \mathbb{F} be an $n \times d_F$ design matrix. The procedure involves sorting observations by column values and grouping them into S blocks at each iteration, resulting in a partition with S^{d_F} cells. Since the number of cells increases exponentially with d_F , applying the procedure directly to the raw design matrix, \mathbb{X} , can be problematic for datasets with large or moderate dimensions of covariates. Even for small choices of S , the final partition may become too fine for the validity of the tests' asymptotic approximation.

To avoid creating excessively fine partitions, one can divide the data by applying the SEB algorithm to a lower-dimensional projection of the data or to some known function of the covariates. In the Monte Carlo simulations of section 7, for example, the data is partitioned using the SEB algorithm on the first q principal components of \mathbb{X} (where $q < d_x$) and on the fitted values of the parametric model under the null hypothesis. I refer to these methods as **PSEB** and **FSEB**, respectively. It's worth noting that when the design matrix is one-dimensional (i.e., $d_F = 1$), the algorithm splits the sorted vector into S blocks.

Algorithm 1 (SEB Algorithm)

- (1) *Sort the observations based on the first column of \mathbb{F} and split them into S blocks, $S > 1$, containing $\lfloor n/S \rfloor$ observations each.*
- (2) *For each block, sort the observations based on the next column of \mathbb{F} and split them into S blocks containing $\lfloor n/S^2 \rfloor$ observations each.*
- (3) *Iterate until the last column of \mathbb{F} is reached.*

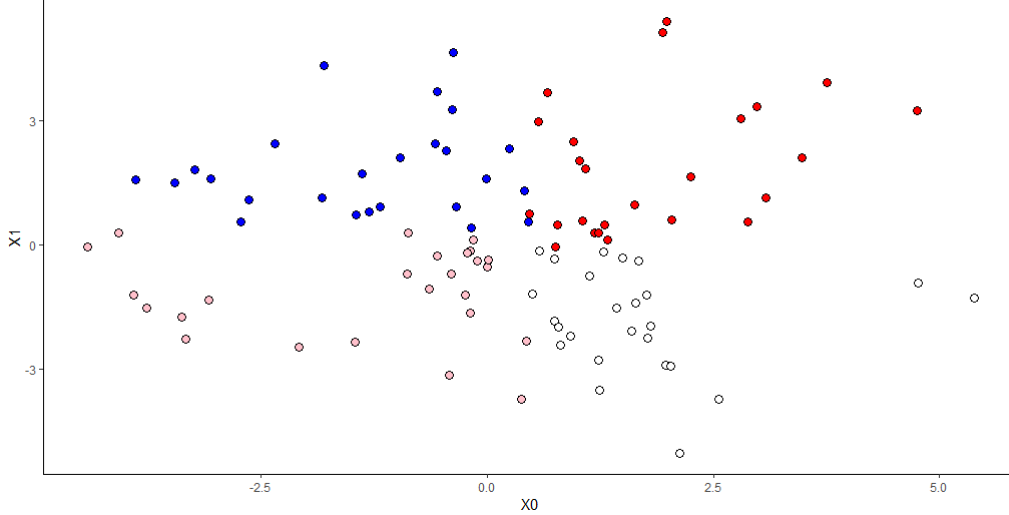


Figure 2: SEB partitioning on raw data with $d_F = 2$ and $S = 2$. The sample is extracted from a bivariate standard normal.

4.2 Neyman-Pearson (NP) Classes

Frequently, the researcher's primary concern is rejecting a subset of deviations critical for the application at hand. For instance, when estimating Mincer's earning regression, it is common to compare it with alternative specifications of the log-income profile (e.g. [Polachek et al. 2008](#)). Similarly, in the [Cox \(1972\)](#) model, a parsimonious parametric specification of the baseline hazard is often compared to more flexible models (e.g. [Seetharaman and Chintagunta 2003](#)).

In this case, a procedure that splits the data where the deviations from the null toward the alternative hypothesis are the largest is preferable. To this end, consider the construction of Neyman-Pearson (NP) classes built over the intersection points between the null and the alternative specification. These classes, in the probability distribution model checking literature, increase (or even maximize) Pearson's measure of discrepancy between the null and the alternative hypothesis (see Remark 3.3 in [Balakrishnan, Voinov, and Nikulin 2013](#)).

Let H_1 be a given alternative parametric specification of the regression function,

$$H_1 : m \in \mathcal{M}_1$$

with $\mathcal{M}_1 = \{m_{1,\theta^*}(\cdot) : \theta^* \in \Theta^*\}$ being a family of parametric regression functions and $\Theta^* \subset \mathbb{R}^{d_{\theta^*}}$ a suitable parameter space. Thus, under H_1 , $m_{1,\theta_0^*}(X) = m(X)$ a.s. for some

$\theta_0^* \in \Theta^*$. Neyman-Pearson classes partition the data splitting over the points where the models under H_0 and under H_1 meet, i.e. where $m_{1,\tilde{\theta}^*}(\cdot) = m_{\tilde{\theta}}(\cdot)$ a.s., with $\tilde{\theta}^* = \theta_0^* + o_p(1)$ under H_1 and $\tilde{\theta}$ the corresponding consistent estimator of θ_0 under H_0 . By doing so, the difference between the average prediction under the null and the alternative model in each class is the largest possible, making it easier to detect deviations from the null toward H_1 .

If X is a univariate variable, finding the intersection set, $\{x \in \mathcal{X} : m_{1,\tilde{\theta}^*}(x) - m_{\tilde{\theta}}(x) = 0\}$, and dividing the data accordingly is a straightforward process. However, when X is multivariate, the intersection set includes surfaces, and splitting \mathcal{X} across different dimensions becomes significantly more complex. In line with the original proposal of [Balakrishnan, Voinov, and Nikulin \(2013\)](#), rather than dividing the data at these intersections, \mathcal{X} is partitioned into two classes based on whether the difference between the models is greater or less than zero.

$$\begin{aligned}\hat{\gamma}_1 &= \{x \in \mathcal{X} : m_{1,\tilde{\theta}^*}(x) - m_{\tilde{\theta}}(x) > 0\}, \\ \hat{\gamma}_2 &= \{x \in \mathcal{X} : m_{1,\tilde{\theta}^*}(x) - m_{\tilde{\theta}}(x) \leq 0\},\end{aligned}$$

Notice that, if the alternative model is correctly specified, i.e. $\mathbb{E}[Y|X] = m_{1,\theta_0^*}(X)$ a.s. for some $\theta_0^* \in \Theta^*$, then $\hat{\gamma}_1$ is the set of points in \mathcal{X} where the null model underestimates the regression function, and $\hat{\gamma}_2$ is the set of points in \mathcal{X} where the null model overestimates the regression function. In this case, χ^2 tests with NP classes detect deviations of the type,

$$H_1(\gamma) : \mathbb{E}[(m_{1,\theta^*}(X) - m_\theta(X)) \mathbb{I}_{\gamma_l}(X)] \neq 0 \text{ for some } l \in \{1, 2\} \text{ and for all } \theta \in \Theta.$$

where $\gamma_1 = \{x \in \mathcal{X} : m_{1,\tilde{\theta}^*}(x) - m_\theta(x) > 0\}$ and $\gamma_2 = \{x \in \mathcal{X} : m_{1,\tilde{\theta}^*}(x) - m_\theta(x) \leq 0\}$. That is, if the alternative model is correct, χ^2 tests with Neyman-Pearson classes are capable of detecting any deviations from it, at least asymptotically.

It is convenient to further divide the data into additional cells to ensure $L > d_\theta$ and perform the J test. This is done by employing the SEB procedure, using a fixed value of $S = L/2$, on the respective vector of fitted values $\{m_{\tilde{\theta}}(X_i)\}_{i \in \gamma_l}$ within each cell. The advantages of using SEB on the fitted values are duals: on the one hand, we

reduce the splitting to a uni-dimensional problem; on the other hand, we exploit the dependence between the residuals and the fitted values under the alternative to generate large aggregate squared residuals.

Algorithm 2 (Parametric Neyman-Pearson (PNP) Algorithm)

- (0) Fix an even number of cells, L , and a minimum number of observations in each cell, n_{min} .
- (1) Compute the fitted values under the null $\{m_{\tilde{\theta}}(X_i)\}_{i=1}^n$.
- (2) Compute the fitted values under the alternative $\{m_{1,\tilde{\theta}^*}(X_i)\}_{i=1}^n$.
- (3) Split the data into two cells, $\hat{\gamma}_1$ and $\hat{\gamma}_2$, where $\hat{\gamma}_1 = \{x \in \mathcal{X} : m_{1,\tilde{\theta}^*} - m_{\tilde{\theta}}(x) > 0\}$ and $\hat{\gamma}_2 = \{x \in \mathcal{X} : m_{1,\tilde{\theta}^*} - m_{\tilde{\theta}}(x) \leq 0\}$.
- (4) If the number of observations in either cell is less than n_{min} , merge them.
- (5) If $L > 2$, split each cell $\hat{\gamma}_l$, $l = 1, 2$, into $L/2$ cells by applying the SEB algorithm on the fitted values $\{m_{\tilde{\theta}}(X_i)\}_{i \in \hat{\gamma}_l}$ of each cell.

This is not the only solution to determine additional splits. The key point for the NP procedure is that, in each cell, the differences between the two models have the same sign. Once this is granted, the the two-cell NP partition can be split with any other procedure.

Example 2 (PNP Classes)

Consider the same scenario as in Example 1, but now under a composite hypothesis, where the model parameters are unknown both under the null hypothesis and under the alternative. After obtaining the OLS estimates for both the null and alternative models, the Neyman-Pearson critical regions are defined based on the points $x \in \mathcal{X}$ where the equation $(\tilde{\theta}_{00} - \tilde{\theta}_{00}^*) + (\tilde{\theta}_{01} - \tilde{\theta}_{01}^*)x - \tilde{\theta}_{02}^* \sin(50x/2\pi) = 0$ holds. The resulting partition is depicted in Figure 3.

4.3 Flexible Neyman-Pearson (FNP) Classes

When considering the unrestricted alternative, Neyman-Pearson classes correspond to subsets of points where the model under the null hypothesis consistently overestimates and

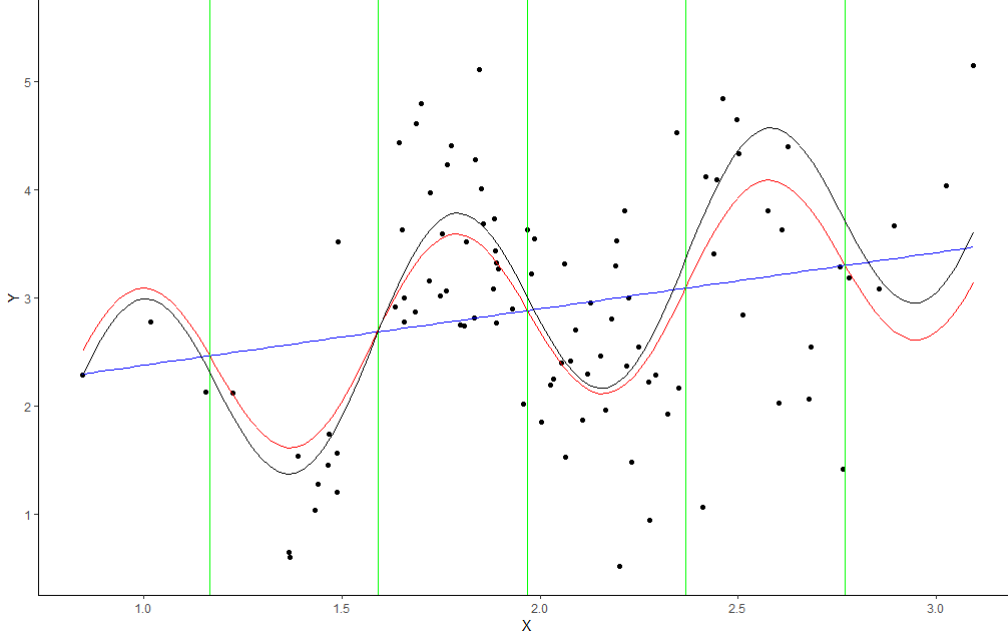


Figure 3: Neyman-Pearson classes with parametric estimate of the regression function. The data is generated under H_1 (black line) with $\theta^* = (1, 1, 1)$. The green line depicts the points where the estimates under the null (blue line) and the estimates under the alternative (red line) meet.

underestimates the (unknown) true regression model. In other words, these classes help us distinguish situations where $m(x)$ is bigger or smaller than $m_{\theta_0}(x)$, which characterize the alternative scenario.

Additionally, differences between the parametric model and the regression function become evident when we examine how the parametric residual relates to various function of the covariates vector, such as the model fitted values. This connection is harnessed with the adoption of a structured alternative "index" model of the form $\mathbb{E}[\varepsilon_{\theta_0}(Z)|X] = f(m_{\theta_0}(X))$ a.s., for some $\theta_0 \in \Theta$, where $f(\cdot)$ represents an unknown smooth function (see, e.g., [Han, Ma, Ren, and Wang 2023](#) for a similar approach). Consequently, under the null hypothesis, the regression error is independent in mean from the assumed regression function, i.e. $\mathbb{E}[\varepsilon_{\theta_0}(Z)|m_{\theta_0}(X)] = 0$. In contrast, under the alternative hypothesis, a systematic dependence in the means of the errors emerges.

This motivates the adoption of an auxiliary regression of $\varepsilon_{\theta_0}(Z)$ on a fixed-order poly-

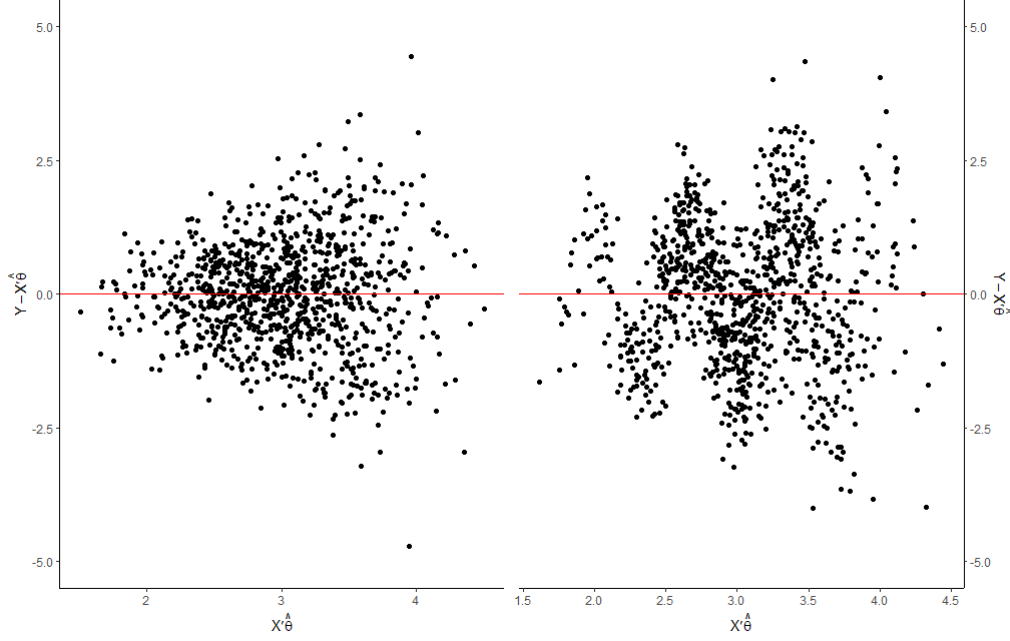


Figure 4: Residual-Fitted values scatter plot under the null (left panel) and under the alternative (right panel). Data generated using the model in (27) under homoskedasticity, $a = 0$, with $c = 0$ (left panel) and $c = 50$ (right panel).

nomial expansion of the fitted values,

$$\varepsilon_{\tilde{\theta}}(Z_i) = \beta_0 + \sum_{j=1}^q \beta_j m_{\tilde{\theta}}(X_i)^j + \epsilon_i,$$

and use it to split where the predicted residuals are greater or smaller than zero. The procedure consists of using the auxiliary model

$$\mathcal{M}_{BP} = \left\{ \beta_0 + \beta_1 m_{\tilde{\theta}}(x) + \beta_2 m_{\tilde{\theta}}(x)^2 + \cdots + \beta_q m_{\tilde{\theta}}(x)^q : \beta \in \mathbb{B} \subset \mathbb{R}^q \right\},$$

to split the sample. Denote as $\tilde{\beta} = (\tilde{\beta}_1, \dots, \tilde{\beta}_q)$ the OLS estimates of β . Under the null hypothesis, $\tilde{\theta} = \theta_0 + o_p(1)$, and $\tilde{\beta}_0 = \tilde{\beta}_1 = \cdots = \tilde{\beta}_q = o_p(1)$, implying that the fitted values do not have predictive power over the residuals. Under the structured alternative, some of the estimated coefficients converge to a non-zero limit, and the fitted values exhibit predictive power on the residuals. The dependence, then, is used to split cases based on the predicted residuals to generate cells with residuals of the same sign and, thus, larger aggregate residuals.

Algorithm 3 (Flexible Neyman-Pearson (FNP) Algorithm)

- (0) Fix an even number of cells, L , and a minimum number of observations in each cell, n_{min} .
- (1) Compute the fitted values under the null $\{m_{\hat{\theta}}(X_i)\}_{i=1}^n$ and obtain the residuals, $\{\varepsilon_{\hat{\theta}}(Z_i)\}_{i=1}^n$.
- (2) Perform an auxiliary regression of $\varepsilon_{\hat{\theta}}(Z_i)$ on a fixed-order polynomial expansion of the fitted values and obtain the vector of predicted residuals, $\{\hat{\varepsilon}_{\hat{\theta}}(Z_i)\}_{i=1}^n$.
- (3) Split the data into two cells, $\hat{\gamma}_1$ and $\hat{\gamma}_2$, where $\hat{\gamma}_1 = \{x \in \mathcal{X} : \hat{\varepsilon}_{\hat{\theta}}(x) > 0\}$ and $\hat{\gamma}_2 = \{x \in \mathcal{X} : \hat{\varepsilon}_{\hat{\theta}}(x) \leq 0\}$.
- (4) If the number of observations in either cell is less than n_{min} , merge them.
- (5) If $L > 2$, split each cell $\hat{\gamma}_l$, $l = 1, 2$, into $L/2$ cells by applying the SEB algorithm on the fitted values $\{m_{\hat{\theta}}(X_i)\}_{i \in \hat{\gamma}_l}$ of each cell.

4.4 Convergence of NP Classes

I examine the conditions under which NP (PNP and FNP) partitions converge to fixed cells in \mathbb{D} . The validity of Assumption 4 for NP partitions may not be immediately apparent, considering that the splitting points in this case are determined through the comparison of different estimators of the regression function. However, as discussed below, a significant class of NP partitions converges to fixed cells in \mathbb{D} .

First, consider the following setting, where the null hypothesis is given by $H_0 : m(X) = \theta'_0 X$ a.s., for some $\theta_0 \in \mathbb{R}^{d_x}$, whereas the alternative of interest is $H_1 : m(X) = \theta_0'^* X + \theta_1^* f(X)$ a.s., for some $(\theta_0^*, \theta_1^*) \in \mathbb{R}^{d_x+1}$, with $f : \mathbb{R}^{d_x} \rightarrow \mathbb{R}$ being a known function of the regressors. Parametric NP partitions are constructed around the points where the two models fit are equal; that is,

$$x \in \mathcal{X} : x'(\tilde{\theta}_0 - \tilde{\theta}_0^*) = f(x)\tilde{\theta}_1^*,$$

where $\tilde{\theta}_0$ is the OLS estimator of Y on X , and $(\tilde{\theta}_0^*, \tilde{\theta}_1^*)$ is the OLS estimator of Y on $\tilde{X} = (X, f(X))$. This is equivalent to stating that the set of solutions, x_0 , satisfies

$$x'_0(\tilde{\theta}_0 - \tilde{\theta}_0^*) = f(x_0)\tilde{\theta}_1^*.$$

Working out the OLS algebra, we obtain the following representations under H_0 ,

$$\begin{aligned}\sqrt{n}(\tilde{\theta}_0 - \tilde{\theta}_0^*) &= A_n \frac{1}{\sqrt{n}} \sum_{i=1}^n \tilde{X}_i e_i, \\ \sqrt{n}\tilde{\theta}_1^* &= B_n \frac{1}{\sqrt{n}} \sum_{i=1}^n \tilde{X}_i e_i,\end{aligned}$$

where $A_n = d_n^{-1} \begin{bmatrix} A_{1,n} & A_{2,n} \end{bmatrix}$, $B_n = d_n^{-1} \begin{bmatrix} 1 & A'_{1,n} \end{bmatrix}$, and,

$$A_{1,n} = \left(\frac{1}{n} \sum_{i=1}^n X_i X_i' \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n X_i f(X_i) \right),$$

$$\begin{aligned}A_{2,n} &= - \left(\frac{1}{n} \sum_{i=1}^n X_i X_i' \right) \left(\frac{1}{n} \sum_{i=1}^n f(X_i) X_i' \right) \left(\frac{1}{n} \sum_{i=1}^n X_i f(X_i) \right) \left(\frac{1}{n} \sum_{i=1}^n X_i X_i' \right)^{-1}, \\ d_n &= \left(\frac{1}{n} \sum_{i=1}^n f(X_i)^2 \right) - \left(\frac{1}{n} \sum_{i=1}^n f(X_i) X_i' \right) \left(\frac{1}{n} \sum_{i=1}^n f(X_i) X_i' \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n X_i f(X_i) \right).\end{aligned}$$

It follows that the set of splitting points is given by the solutions of the following equation,

$$(x'_0 A_n - f(x_0) B_n) \frac{1}{\sqrt{n}} \sum_{i=1}^n \tilde{X}_i e_i = 0,$$

which holds for any $x_0 \in \mathcal{X}$ such that $x'_0 A_n = f(x_0) B_n$. Thus, x_0 converge in probability to the fixed solutions of $x'_0 A = f(x_0) B$, where A and B are the probability limits of A_n and B_n , respectively.

This example illustrates that even when the partition depends on the comparison of estimated parameters, the splitting points can converge to fixed points in \mathcal{X} . The following proposition generalizes this result by incorporating specific conditions on the behavior of the parameter estimators under the null and alternative hypotheses to ensure

the convergence of NP classes.

Proposition 1 *Let $\mathcal{M} = \{m_\theta(\cdot) : \theta \in \Theta\}$ and $\mathcal{M}_1 = \{m_{1,\theta^*}(\cdot) : \theta^* \in \Theta^*\}$ represent the models under H_0 and H_1 , respectively. Further, denote with $\tilde{\theta}$ and $\tilde{\theta}^*$ the respective \sqrt{n} -consistent estimators of θ_0 and θ_0^* under H_0 and H_1 , where $m_{\theta_0} = m$ a.s. under H_0 , and $m_{1,\theta_0^*} = m$ a.s. under H_1 . Then, if under H_0 :*

(a) *There exists a θ_1^* , such that $\tilde{\theta}^* \xrightarrow{p} \theta_1^*$ under H_0 .*

(b) *$m_{\theta_0} = m_{1,\theta_1^*}$ a.s.*

(c) *$\sqrt{n}(\tilde{\theta} - \theta_0) = Cn^{-1/2} \sum_{i=1}^n h(X_i, \varepsilon_{\theta_0}) + o_p(1)$, $\sqrt{n}(\tilde{\theta}^* - \theta_1^*) = Dn^{-1/2} \sum_{i=1}^n h(X_i, \varepsilon_{\theta_0}) + o_p(1)$, where C and D are $c \times d_h$ and $d \times d_h$ constant matrices, and $h(\cdot, \cdot)$ is some \mathbb{R}^{d_h} -valued function of regressors and errors such that $n^{-1/2} \sum_{i=1}^n h(X_i, \varepsilon_{\theta_0}) = O_p(1)$.*

Then, the NP partition splitting points converge in probability to fixed points in \mathcal{X} .

The most crucial conditions for the convergence to fixed cells are Conditions (a) and (b), which require that the model under the alternative encompasses the model under the null hypothesis. These conditions are typically satisfied in the context of both PNP partitions and FNP partitions. In the case of FNP partitions this is particularly easy to check, as the alternative regression model corresponds to a polynomial expansion of the fitted values, thus nesting the model under the null.

In the appendix, I provide Monte Carlo evidence demonstrating NP partitions converging to fixed cells in \mathbb{D} .

5 Local Power

To analyze the local power of the χ^2 tests under the alternative hypothesis, let's consider the following sequence of local alternatives,

$$H_{1,n} : m(x) = m_{\theta_0}(x) + \frac{1}{\sqrt{n}}h(x) \text{ a.s.}, \quad (17)$$

where $h(X)$ is a random variable representing departures from the null hypothesis with $\mathbb{E}[h(X)^2] < \infty$ and $0 \leq P(h(X) = 0) < 1$.

Denote as $\delta_\gamma = (\delta_1, \dots, \delta_L)$, with $\delta_l = \mathbb{E}[h(X)\mathbb{I}_{\gamma_l}(X)]$ for $l = 1, \dots, L$. Under $H_{1,n}$ and θ_0 known, the test statistics converge to a noncentral chi-squared distribution with L degrees of freedom,

$$\begin{aligned}\hat{\chi}_{\hat{\gamma},0}^2 &\xrightarrow{d} \chi_L^2(\lambda), \\ \hat{\mathcal{W}}_{\hat{\gamma}}(\theta_0) &\xrightarrow{d} \chi_L^2(\lambda),\end{aligned}\tag{18}$$

where,

$$\lambda = \delta'_\gamma (\Sigma_{\gamma,0})^{-1} \delta_\gamma \tag{19}$$

is the non-centrality parameter.

Equations (18) and (19) highlight an important trade-off when choosing L . On one hand, the inequality,

$$\frac{\delta_l^2}{\sigma_{0,l}^2} + \frac{\delta_f^2}{\sigma_{0,f}^2} \geq \frac{(\delta_l + \delta_f)^2}{\sigma_{0,l}^2 + \sigma_{0,f}^2},$$

shows that λ is non-decreasing for nested partitions. That is, as the partition becomes finer with more cells, the test becomes more capable of detecting smaller deviations from the null hypothesis, resulting in increased power. On the other hand, the (global) power of χ^2 tests under the alternative hypothesis $H_{1,n}$ decreases as the number of cells increases. This decline in power is due to the higher variability of the limit distribution in (18) associated with a larger number of cells. As the number of cells grows, the distribution becomes more spread out, leading to a higher chance of observing test statistics falling in less extreme regions and, thus, not rejecting. Similar trade-offs have already been noted in the classical goodness-of-fit testing literature (see, e.g., [Kallenberg, Oosterhoff, and Schriever 1985](#)).

When θ_0 is unknown, under $H_{1,n}$, the estimators of θ_0 typically follow the asymptotic expansion in Assumption 3' below (see, e.g., [Newey 1985](#)).

Assumption 3' *Under $H_{1,n}$, the estimator $\tilde{\theta}$ satisfies,*

$$\sqrt{n}(\tilde{\theta} - \theta_0) = \frac{1}{\sqrt{n}} \sum_{i=1}^n l_{\theta_0}(Z_i) + \delta + o_p(1)$$

where the function $l_\theta(\cdot)$ is as defined in Assumption 3 and $\delta \in \mathbb{R}^{d_\theta}$ is a non-zero vector.

It's easy to show, for instance, that Assumption 3' is satisfied by the grouped GMM estimator, with

$$\delta = - \left[\boldsymbol{\mu}_{\gamma,0}^* (\Sigma_{\gamma,0})^{-1} \boldsymbol{\mu}_{\gamma,0}^{*'} \right]^{-1} \boldsymbol{\mu}_{\gamma,0}^* (\Sigma_{\gamma,0})^{-1} \delta_{\gamma}.$$

It follows that under $H_{1,n}$ and the same assumption of Theorem 2,

$$\begin{aligned} \hat{\chi}_{\hat{\gamma}}^2 &\xrightarrow{d} \chi_{L-d_{\theta}}^2(\lambda_1), \\ \hat{\mathcal{W}}_{\hat{\gamma}}(\tilde{\theta}) &\xrightarrow{d} \chi_{r(\text{Avar}(\hat{\Phi}_{\gamma}(\tilde{\theta})))}^2(\lambda_2), \end{aligned} \quad (20)$$

where

$$\lambda_1 = \delta_{\gamma}' \left[(\Sigma_{\gamma,0})^{-1} - (\Sigma_{\gamma,0})^{-1} \boldsymbol{\mu}_{\gamma,0}^* (\boldsymbol{\mu}_{\gamma,0}^{*'} (\Sigma_{\gamma,0})^{-1} \boldsymbol{\mu}_{\gamma,0}^*)^{-1} \boldsymbol{\mu}_{\gamma,0}^{*'} (\Sigma_{\gamma,0})^{-1} \right] \delta_{\gamma} \quad (21)$$

$$\lambda_2 = (\delta_{\gamma} - \boldsymbol{\mu}_{\gamma,0}^* \delta)' \text{Avar}^{-}(\hat{\Phi}_{\gamma}(\tilde{\theta})) (\delta_{\gamma} - \boldsymbol{\mu}_{\gamma,0}^* \delta) \quad (22)$$

are the respective non-centrality parameters. Comparing the power of these two tests is not a straightforward task, given the differences in both degrees of freedom and non-centrality parameters in their respective distributions. In fact, scenarios can be constructed where the $\hat{\chi}^2$ test exhibits greater power than the Wald test, and vice versa, reasoning as Moore and Spruill (1975) and Moore (1977), for the classical case.

Instead of the non-centrality parameters, consider now the squared L2 norm of the drift,

$$\|\delta_{\gamma}\|_2^2 = \delta_{\gamma}' \delta_{\gamma} = \sum_{l=1}^L \delta_l^2 = \sum_{l=1}^L \mathbb{E} [h(X) \mathbb{I}_{\gamma_l}(X)]^2.$$

Then, for any pair of δ_l and δ_f such that $\text{sgn}(\delta_l) = \text{sgn}(\delta_f)$, the inequality $(\delta_l + \delta_f)^2 \geq \delta_l^2 + \delta_f^2$ shows that the optimal partitioning for $\|\delta_{\gamma}\|_2^2$ consists of two cells: one containing the points where $h(\cdot) = \sqrt{n}(m(\cdot) - m_{\theta_0}(\cdot))$ takes only positive values, and the other where it takes only negative values. Essentially, the optimal partitioning of $\|\delta_{\gamma}\|_2^2$ is achieved with two Neyman-Pearson classes.

Proposition 2 *The euclidean norm of the drifts, $\|\delta_{\gamma}\|_2$, is maximized by two Neyman-Pearson classes, $\gamma^* = \{\gamma_i^*\}_{i=1}^2$,*

$$\gamma_1^* = \{x \in \mathcal{X} : h(x) \geq 0\} \quad \gamma_2^* = \{x \in \mathcal{X} : h(x) < 0\}$$

This suggests that, in some sense, Neyman-Pearson classes correspond to optimization criteria for the non-centrality parameter. Of course, a rigorous argument to justify efficiency should account for the dependence of the non-centrality parameter denominators on the cell boundaries.

6 General CMR

The analysis is extended to general moment restrictions with the introduction of a response variables vector, Y , taking values in $\mathcal{Y} \subset \mathbb{R}^{d_y}$, $d_y \geq 1$, and a generalized residual vector (Wooldridge 1990), $\boldsymbol{\varepsilon}_\theta : (\mathcal{Y}, \mathcal{X}) \rightarrow \mathbb{R}^{d_\varepsilon}$ with $\boldsymbol{\varepsilon}_\theta(\cdot) = (\varepsilon_{1,\theta}(\cdot), \dots, \varepsilon_{d_\varepsilon,\theta}(\cdot))'$, defining parametric relationships between Y and X . The null hypothesis is defined as before, i.e. $H_0 : \mathbb{E}[\boldsymbol{\varepsilon}_{\theta_0}|X] = 0$ a.s. for some $\theta_0 \in \Theta$.

The generality of this framework allows testing for a wide range of econometric models such as simultaneous equation model identified by instrumental variables (Newey 1990) or nonlinear in parameters and endogenous variables models, e.g. Box-Cox transform.

When the dimension of the generalized residual is bigger than one, it might be optimal to consider a partition for each component of $\boldsymbol{\varepsilon}_\theta(\cdot)$. In particular, for each $j \in \{1, \dots, d_\varepsilon\}$, let \mathbb{D}_j be a class of partitions of \mathcal{X} comprised of L_j sets from \mathbb{C} (L_j is fixed for all n); that is,

$$\mathbb{D}_j = \left\{ \boldsymbol{\gamma}_j = (\gamma_{j,1}, \dots, \gamma_{j,L_j})' \in \mathbb{C}^{L_j} : \bigcup_{l=1}^{L_j} \gamma_{j,l} = \mathcal{X}, \gamma_{j,l} \cap \gamma_{j,f} = \emptyset, \forall l \neq f \right\}. \quad (23)$$

The partition corresponding to the j -th component of the generalized residual is an element of \mathbb{D}_j , $\boldsymbol{\gamma}_j \in \mathbb{D}_j$, or a random element $\hat{\boldsymbol{\gamma}}_j \in \mathbb{D}_j$, with probability limit $\boldsymbol{\gamma}_j \in \mathbb{D}_j$.

Denote as $\mathcal{E}_\theta(\cdot)$ the $\bar{L} \times \bar{L}$ block diagonal matrix of generalized residuals with main diagonal elements given by $\{\varepsilon_{j,\theta}(\cdot)I_{L_j}\}_{j=1}^{d_\varepsilon}$, where $\bar{L} = \sum_j L_j$. If $L_1 = L_2 = \dots = L_{d_\varepsilon} = L$, then $\mathcal{E}_\theta(\cdot) = \text{diag}[\boldsymbol{\varepsilon}_\theta(\cdot)] \otimes I_L$, where $\text{diag}[\boldsymbol{\varepsilon}_\theta(\cdot)] = \text{diag}\{\varepsilon_{1,\theta}(\cdot), \dots, \varepsilon_{d_\varepsilon,\theta}(\cdot)\}$ is the $d_\varepsilon \times d_\varepsilon$ diagonal matrix with the components of $\boldsymbol{\varepsilon}_\theta(\cdot)$ on the main diagonal and \otimes denotes the Kronecker product.

The χ^2 test statistics can be expressed as quadratic forms of

$$\hat{\Phi}_\gamma(\theta) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathcal{E}_\theta(Z_i) \mathbf{I}_\gamma(X_i), \quad (24)$$

where $\mathbf{I}_\gamma(\cdot) = (\mathbf{I}'_{\gamma_1}, \dots, \mathbf{I}'_{\gamma_{d_\varepsilon}})'$ is the vector of indicator functions over all the partitions. While the covariance matrices of $\hat{\Phi}_\gamma(\theta_0)$ and $\hat{\Phi}_\gamma(\tilde{\theta})$ under the null are given by,

$$\Sigma_{\gamma,0} = \mathbb{E}[\mathcal{E}_{\theta_0}(Z) \mathbf{I}_\gamma(X) \mathbf{I}_\gamma(X)' \mathcal{E}_{\theta_0}(Z)'], \quad (25)$$

and,

$$\text{Avar}\left(\hat{\Phi}_\gamma(\tilde{\theta})\right) = \Sigma_{\gamma,0} - \boldsymbol{\mu}_{\gamma,0}^* C'_{\gamma,0} - C_{\gamma,0} \boldsymbol{\mu}_{\gamma,0}^{*'} + \boldsymbol{\mu}_{\gamma,0}^* L_0 \boldsymbol{\mu}_{\gamma,0}^{*'}, \quad (26)$$

where $\boldsymbol{\mu}_{\gamma,0}^* = \mathbb{E}[\nabla \mathcal{E}_{\theta_0}(Z) \mathbf{I}_\gamma(X)]$ is the Jacobian matrix, $C'_{\gamma,0} = \mathbb{E}[\mathcal{E}_{\theta_0}(Z) \mathbf{I}_\gamma(X) l_{\theta_0}(Z)']$, and L_0 is defined as before.

In general, $\Sigma_{\gamma,0}$ is not diagonal, unless $\boldsymbol{\varepsilon}_\theta(\cdot)$ consists of orthogonal components or $d_\varepsilon = 1$. However, when $\gamma_1 = \gamma_2 = \dots = \gamma_{d_\varepsilon}$, $\Sigma_{\gamma,0}$ is block diagonal with L blocks of size $d_\varepsilon \times d_\varepsilon$ corresponding to the covariance matrix of $\boldsymbol{\varepsilon}_\theta(\cdot)$ in each cell of the partition.

7 Monte Carlo Study

This section presents a Monte Carlo study to evaluate the finite sample performance of the χ^2 tests. This study is structured into two distinct parts. The first part involves testing the null hypothesis of linearity within a regression model, while the second part focuses on assessing the proper specification of a probit model for a binary response variable. To establish a benchmark for comparison with existing tests for H_0 , the χ^2 tests in the first part are compared with the tests proposed by [Stute \(1997\)](#) and [Stute and Zhu \(2002\)](#) for regression specifications based on marked residuals processes. In the second part, the χ^2 tests are compared with the specification test for propensity score specifications proposed by [Sant'Anna and Song \(2019\)](#) (hereafter, SS).

Linear Model

The data generating process (DGP) is given by,

$$Y_i = \sum_{j=1}^{d_x} X_{j,i} + b \sin \left(\frac{c \sum_{j=1}^{d_x} X_{j,i}}{2\pi} \right) + \sigma(X_i) \epsilon_i \quad (27)$$

where $X_i = (X_{1,i}, \dots, X_{d_x,i})$ is a vector of mutually independent covariates distributed uniformly over $[0,1]$, $X_{i,j} \sim U[0,1]$, the error distributes normally and independently from X , $\epsilon|X \sim N(0,1)$,

$$\sigma^2(X) = \frac{g(X)}{\mathbb{E}[g(X, a)]},$$

and $g(X, a) = e^{aX_1}$, with $\mathbb{E}[\sigma^2(X)] = 1$. The model under the null corresponds to $b = 0$, while for the alternative models $b = 0.5$. The parameter a controls the heteroskedasticity severity, with $a = 0$ corresponding to homoskedasticity, while c governs the extent of departures from linearity. In the first part of the simulations, I provide evidence for the size of the test for $a \in \{0, 3\}$, I then fix $a = 3$ and proceeds to examine the tests power for $c \in \{10, 50\}$. Notably, with $c = 50$, the deviations manifest at higher frequencies making them more challenging to discern from sampling error. The generated samples have size $n \in \{100, 200, 500, 1000\}$ with a dimension of the covariate vector $d_x \in \{5, 10\}$. The number of Monte Carlo repetition, R , is set at $R = 5000$ for $n = 100$, $R = 2500$ for $n = 200$, $R = 1250$ for $n = 500$, $R = 625$ for $n = 1000$. The rejection rates are reported at the 5% nominal level, results at 1% and 10% are similar.

The $\hat{\chi}^2$ test is built as described in Section 2 using the GMM estimator,

$$\hat{\theta}_\gamma = \left[\sum_{l=1}^L \frac{\bar{X}_l \bar{X}_l'}{\hat{\sigma}_l^2(\tilde{\theta})} \right]^{-1} \sum_{l=1}^L \frac{\bar{X}_l \bar{Y}_l}{\hat{\sigma}_l^2(\tilde{\theta})}.$$

where $\tilde{\theta}$ is the OLS estimator. Given the data-generating process, one can easily show that the influence function of $\tilde{\theta}$, $l_{\theta_0}(z) = \mathbb{E}[XX']^{-1} x \varepsilon_{\theta_0}(z)$, and $\varepsilon_{\theta_0}(z) \mathbf{I}_\gamma(x)$ have linearly independent components (however, the independence does not hold when X includes a constant term). Thus, $\text{Avar} \left(\hat{\Phi}_\gamma(\tilde{\theta}) \right)$, is full rank and the Wald test can be implemented

using the plug-in estimator $\hat{W}_\gamma(\tilde{\theta})$,

$$\begin{aligned}\hat{W}_\gamma(\tilde{\theta}) &= \mathbb{E}_n [\varepsilon_{\tilde{\theta}}(Z_i)^2 \mathbf{I}_\gamma(X_i) \mathbf{I}_\gamma(X_i)'] - \mathbb{E}_n [\mathbf{I}_\gamma(X_i) X_i'] \mathbb{E}_n [X_i X_i']^{-1} \mathbb{E}_n [\varepsilon_{\tilde{\theta}}(Z_i)^2 X_i \mathbf{I}_\gamma(X_i)'] - \\ &\quad - \mathbb{E}_n [\varepsilon_{\tilde{\theta}}(Z_i)^2 \mathbf{I}_\gamma(X_i) X_i'] \mathbb{E}_n [X_i X_i']^{-1} \mathbb{E}_n [X_i \mathbf{I}_\gamma(X_i)'] + \\ &\quad + \mathbb{E}_n [\mathbf{I}_\gamma(X_i) X_i'] \mathbb{E}_n [X_i X_i']^{-1} \mathbb{E}_n [\varepsilon_{\tilde{\theta}}(Z_i)^2 X_i X_i'] \mathbb{E}_n [X_i X_i']^{-1} \mathbb{E}_n [X_i \mathbf{I}_\gamma(X_i)'].\end{aligned}$$

I assess the test's performance under various partitioning methods, including SEB on the first principal component (**PSEB**), SEB on the vector of fitted values (**FSEB**), parametric NP (**PNP**), and Flexible NP splitting (**FNP**). When $d_x = 5$, the partitioning is done using $L \in \{4, 6, 8\}$ cells, whereas when $d_x = 10$, partitions with $L \in \{8, 10, 12\}$ cells are examined. In the initial two-cell partition of the **FNP** and **PNP** algorithms, the minimum number of observations within each cell is set to $n_{min} = n/5$. The parametric **PNP** cells are constructed using the correctly specified alternative model (with known parameter c). Tests constructed in this manner utilize additional information about the alternative specifications of the regression function, including them allows for a comprehensive comparison with the other feasible partitioning methods. The **FNP** partition is built using the predicted residuals of a linear regression of the parametric residuals on a polynomial expansion of the fitted values of order $q = 3$.

The performance of the χ^2 tests are compared with two minimum-distance tests for regression specifications based on marked residuals processes indexed by real vectors and real numbers; cf. [Stute \(1997\)](#) and [Stute and Zhu \(2002\)](#), respectively.

$$\begin{aligned}R_{1,n}(x_1) &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \left(Y_i - X_i' \tilde{\theta} \right) \mathbb{I}\{X_i \leq x_1\}, \quad x_1 \in \mathbb{R}^{d_x} \\ R_{2,n}(x_2) &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \left(Y_i - X_i' \tilde{\theta} \right) \mathbb{I}\{X_i' \tilde{\theta} \leq x_2\}, \quad x_2 \in \mathbb{R}.\end{aligned}$$

Incidentally, the marked residual process indexed by partitions generated with the **FSEB** method,

$$\hat{\Phi}_{\hat{\gamma}_l}(\tilde{\theta}) = n^{-1/2} \sum_{i=1}^n \varepsilon_{\tilde{\theta}}(Z_i) \mathbb{I}\{X_i' \tilde{\theta} \in C_l\},$$

where $\hat{\gamma}_l = \{x \in \mathcal{X} : x' \tilde{\theta} \in C_l\}$ and $\bigcup_{l=1}^L C_l = \mathbb{R}$, is a finite cells version of the process

used in the [Stute and Zhu \(2002\)](#) test. The test statistics consist of a functional of $R_{1,n}(\cdot)$ and $R_{2,n}(\cdot)$. In these simulations, we only consider the Kolmogorov-Smirnov functional,

$$KS1 = \sup_{x_1 \in \mathbb{R}^{d_x}} |R_{1,n}(x_1)|,$$

$$KS2 = \sup_{x_2 \in \mathbb{R}} |R_{2,n}(x_2)|.$$

Since these tests have non-pivotal limiting distribution, the critical values are estimated with Wild bootstrap (see [Stute, Manteiga, and Quindimil 1998](#)) using $B = 500$ bootstrap repetitions for each replication.

The simulation results for the linear model presented in [Tables 1](#) (size) and [2](#) (power) reveal several key findings. In [Table 1](#), both $\hat{\mathcal{W}}$ and $\hat{\chi}^2$ tests, generally exhibit effective size control. Nevertheless, some (fairly small) size distortion persists, particularly depending on the partitioning methods and the number of cells, even for large sample sizes. This distortion arises from irregularities in the replication of the partitioning algorithm, which generates cells with too few observations, particularly when the number of observations in each cell of the initial two-cell split is close to the minimum limit set in the algorithm. Such irregularities disappear when n is large enough. Notably, the $KS1$ test experiences pronounced size distortions when $d_x = 10$, while the $KS2$ test demonstrates greater robustness in the face of higher covariate dimensions, as expected. Turning to [Table 2](#), our focus shifts to assessing the tests' ability to detect deviations from linearity when the null hypothesis is false ($b = 0.5$). The $\hat{\mathcal{W}}$ test generally exhibits higher power compared to the $\hat{\chi}^2$ test, especially in cases with larger values of c . However, power varies depending on the choice of partitioning method and other parameters. For example, the $\hat{\mathcal{W}}$ test, when employing model-based partitioning, consistently outperforms omnibus tests in most settings, while the J test performs better than the $KS2$ test primarily in scenarios involving high-frequency deviations ($c = 50$). As expected, both χ^2 tests using **PNP** partitions consistently exhibit the highest power across all scenarios, given their ability to leverage out-of-sample information. However, they exhibit low power when the partitioning is done with the **PSEB** method.

Table 1: Linear Model: Size ($b = 0$)

				PSEB		FSEB		FNP		PNP			
n	d_x	a	L	\hat{W}	$\hat{\chi}^2$	\hat{W}	$\hat{\chi}^2$	\hat{W}	$\hat{\chi}^2$	\hat{W}	$\hat{\chi}^2$	KS1	KS2
100	5	0	4	5.18	NA	4.86	NA	6.46	NA	4.86	NA	5.14	6.40
			6	4.88	5.70	4.32	5.84	4.88	6.06	4.32	5.84		
			8	3.83	5.44	4.48	5.86	4.24	5.92	4.48	5.86		
100	5	3	4	4.88	NA	5.09	NA	6.32	NA	5.09	NA	4.86	7.12
			6	4.42	5.72	3.98	5.08	4.38	5.72	3.98	5.08		
			8	3.62	4.90	3.42	4.62	3.92	5.26	3.42	4.62		
100	10	0	8	4.96	NA	4.74	NA	5.66	NA	4.74	NA	1.46	8.14
			10	4.92	NA	4.68	NA	5.28	NA	4.68	NA		
			12	4.72	6.24	4.54	6.24	4.60	6.38	4.54	6.24		
100	10	3	8	4.68	NA	4.94	NA	4.80	NA	4.94	NA	0.78	9.16
			10	4.06	NA	3.66	NA	4.40	NA	3.66	NA		
			12	3.48	5.88	3.56	6.75	3.86	6.26	3.56	6.75		
200	5	0	4	4.72	NA	4.56	NA	6.12	NA	4.56	NA	5.56	6.75
			6	4.16	5.08	4.51	3.64	5.04	5.56	4.51	3.64		
			8	4.72	4.44	4.92	5.68	5.16	6.44	4.92	5.68		
200	5	3	4	4.44	NA	4.56	NA	5.92	NA	4.56	NA	5.56	6.56
			6	4.60	5.40	4.00	5.20	5.40	6.64	4.00	5.20		
			8	4.48	5.52	4.08	5.48	4.44	5.60	4.08	5.48		
200	10	0	8	5.28	NA	4.08	NA	6.00	NA	4.08	NA	2.80	6.60
			10	4.08	NA	5.20	NA	5.20	NA	5.20	NA		
			12	4.68	5.64	4.32	6.24	4.76	6.12	4.32	6.24		
200	10	3	8	4.63	NA	4.96	NA	5.20	NA	4.96	NA	1.40	6.04
			10	5.20	NA	4.36	NA	4.56	NA	4.36	NA		
			12	4.68	6.28	4.00	5.76	4.36	5.96	4.00	5.76		
500	5	0	4	5.28	NA	5.28	NA	6.00	NA	5.28	NA	6.40	7.28
			6	5.04	5.20	5.52	6.56	5.28	5.68	5.52	6.56		
			8	5.04	4.72	5.12	6.24	6.08	7.20	5.12	6.24		
500	5	3	4	4.16	NA	4.96	NA	6.48	NA	4.96	NA	5.76	5.84
			6	5.04	4.32	4.32	3.92	5.52	5.44	4.32	3.92		
			8	4.24	4.08	5.20	6.32	5.68	6.00	5.20	6.32		
500	10	0	8	5.60	NA	4.72	NA	6.16	NA	4.72	NA	4.32	6.32
			10	4.72	NA	4.80	NA	5.12	NA	4.80	NA		
			12	5.04	5.20	5.44	5.52	5.76	5.68	5.44	5.52		
500	10	3	8	5.04	NA	4.32	NA	5.12	NA	4.32	NA	3.20	5.60
			10	3.92	NA	5.28	NA	4.08	NA	5.28	NA		
			12	4.72	5.12	4.16	4.08	4.56	5.76	4.16	4.08		
1000	5	0	4	4.74	NA	5.62	NA	5.62	NA	5.62	NA	3.40	5.92
			6	3.85	4.88	5.03	4.59	5.77	5.92	5.03	4.59		
			8	4.88	5.77	4.74	5.92	6.37	6.22	4.74	5.92		
1000	5	3	4	5.77	NA	5.77	NA	6.37	NA	5.77	NA	3.55	7.11
			6	4.59	5.77	6.51	6.51	6.07	5.92	6.51	6.51		
			8	6.51	6.07	3.70	4.00	4.44	4.59	3.70	4.00		
1000	10	0	8	5.03	NA	5.33	NA	6.07	NA	5.33	NA	6.37	4.88
			10	5.48	NA	5.77	NA	5.48	NA	5.77	NA		
			12	5.62	6.22	5.48	5.18	5.48	6.07	5.48	5.18		
1000	10	3	8	4.74	NA	5.48	NA	4.14	NA	5.48	NA	6.22	5.62
			10	5.03	NA	4.74	NA	3.40	NA	4.74	NA		
			12	4.44	5.92	4.44	5.03	7.11	6.07	4.44	5.03		

Percentage of rejections at the nominal level $\alpha = 0.05$ of the $\hat{\chi}^2$ and the Wald test under different partitioning methods. n , d_x , a , L , respectively denote the sample size, the number of covariates, the degree of heteroskedasticity, and the number of cells. The last two columns report the rejection rates of the Kolmogorov-Smirnov tests.

Table 2: Linear Model: Power ($b = 0.5$)

				PSEB		FSEB		FNP		PNP			
n	d_x	c	L	\hat{W}	$\hat{\chi}^2$	\hat{W}	$\hat{\chi}^2$	\hat{W}	$\hat{\chi}^2$	\hat{W}	$\hat{\chi}^2$	KS1	KS2
100	5	10	4	22.02	NA	39.12	NA	46.76	NA	44.60	NA	31.54	35.75
			6	18.32	6.74	37.01	15.34	40.34	16.36	40.08	30.12		
			8	15.16	6.02	30.71	19.60	33.18	19.54	33.28	35.64		
100	5	50	4	4.80	NA	6.36	NA	8.34	NA	72.14	NA	5.20	7.98
			6	4.48	5.38	7.04	7.32	7.66	8.16	63.72	42.62		
			8	3.70	4.51	7.60	8.18	7.43	8.10	55.18	53.12		
100	10	10	8	5.98	NA	27.80	NA	28.48	NA	40.64	NA	1.58	28.16
			10	4.97	NA	24.18	NA	24.32	NA	34.69	NA		
			12	4.76	6.34	20.08	16.08	19.44	15.16	28.04	32.80		
100	10	50	8	4.78	NA	4.97	NA	5.34	NA	53.24	NA	0.66	8.59
			10	4.24	NA	4.54	NA	4.84	NA	45.84	NA		
			12	3.48	5.70	4.26	6.48	4.26	6.30	39.58	35.42		
200	5	10	4	40.64	NA	73.36	NA	81.12	NA	80.80	NA	59.84	65.08
			6	35.44	7.12	75.44	25.92	79.04	24.76	79.96	49.64		
			8	31.92	6.20	72.44	36.12	75.28	34.12	75.12	62.92		
200	5	50	4	5.12	NA	12.28	NA	18.28	NA	96.52	NA	5.04	10.24
			6	4.40	5.44	19.91	10.96	20.92	12.04	94.96	59.28		
			8	4.56	4.68	22.88	17.96	23.40	16.44	92.96	82.44		
200	10	10	8	8.00	NA	72.40	NA	75.40	NA	82.60	NA	4.04	50.44
			10	7.24	NA	69.16	NA	71.04	NA	79.16	NA		
			12	6.68	6.68	66.04	30.24	65.68	29.48	74.96	53.08		
200	10	50	8	4.92	NA	7.88	NA	8.00	NA	91.08	NA	1.44	7.40
			10	4.56	NA	7.52	NA	7.12	NA	88.36	NA		
			12	4.16	5.76	7.16	6.84	6.80	7.04	86.04	56.72		
500	5	10	4	81.04	NA	99.04	NA	99.92	NA	99.84	NA	94.72	98.16
			6	75.36	8.80	99.20	40.24	99.44	40.48	99.60	68.72		
			8	71.28	9.27	99.60	69.19	99.68	64.72	99.52	92.40		
500	5	50	4	4.40	NA	35.04	NA	50.72	NA	100.00	NA	5.36	19.60
			6	5.68	5.44	73.44	33.43	70.88	30.24	100.00	75.68		
			8	5.28	6.88	82.64	57.92	81.28	55.92	100.00	96.64		
500	10	10	8	14.16	NA	99.84	NA	99.92	NA	100.00	NA	17.52	94.40
			10	11.76	NA	99.84	NA	100.00	NA	99.92	NA		
			12	10.00	5.76	99.84	59.84	99.84	56.72	99.84	81.44		
500	10	50	8	4.56	NA	42.64	NA	38.40	NA	100.00	NA	4.40	10.00
			10	5.28	NA	45.52	NA	45.20	NA	100.00	NA		
			12	5.28	5.12	48.96	23.12	49.28	21.92	100.00	80.32		
1000	5	10	4	97.48	NA	100.00	NA	100.00	NA	100.00	NA	99.85	100.00
			6	96.29	10.96	100.00	59.11	100.00	51.70	100.00	80.44		
			8	95.70	14.66	100.00	84.59	100.00	82.51	100.00	98.07		
1000	5	50	4	5.48	NA	68.29	NA	77.77	NA	100.00	NA	4.59	45.03
			6	5.62	5.77	99.55	57.33	96.44	49.48	100.00	85.33		
			8	6.96	6.07	99.40	87.11	99.85	84.29	100.00	99.11		
1000	10	10	8	22.07	NA	100.00	NA	100.00	NA	100.00	NA	37.92	100.00
			10	21.92	NA	100.00	NA	100.00	NA	100.00	NA		
			12	16.44	8.00	100.00	79.85	100.00	78.96	100.00	91.55		
1000	10	50	8	4.59	NA	93.48	NA	84.29	NA	100.00	NA	4.29	20.00
			10	5.03	NA	94.37	NA	95.55	NA	100.00	NA		
			12	4.29	5.62	96.88	50.07	94.66	48.74	100.00	91.11		

Percentage of rejections at the nominal level $\alpha = 0.05$ of the $\hat{\chi}^2$ and the Wald test under different partitioning methods. n , d_x , c , L , respectively denote the sample size, the number of covariates, the departure from the null, and the number of cells. The last two columns report the rejection rates of the Kolmogorov-Smirnov tests.

Probit Model

This second part considers the probit models denoted as DGP.6, DGP.7, DGP.8, DGP.9, and DGP.10 in Section 4 of SS. The data consists of random draws $\{(D_i, X_i')'\}_{i=1}^n$, where $D = \mathbb{I}\{D^* > 0\}$ and D^* is generated as,

$$D^* = -\frac{\sum_{j=1}^{10} X_j}{6} - \epsilon \quad (DGP.6),$$

under the null hypothesis. Here, $X_1 = Z_1$, $X_2 = (Z_1 + Z_2)/\sqrt{2}$, $X_j = Z_j$, for $j = 3, \dots, 10$, with $Z_j \sim N(0, 1)$, for all j , and $\epsilon \sim N(0, 1)$ independently of X . The alternative models for the latent variable are given by,

$$D^* = -1 - \frac{\sum_{j=1}^{10} X_j}{10} + \frac{X_1 X_2}{2} - \epsilon \quad (DGP.7),$$

$$D^* = -1 - \frac{\sum_{j=1}^{10} X_j}{10} + \frac{X_1 \sum_{j=2}^5 X_j}{4} - \epsilon \quad (DGP.8),$$

$$D^* = -1.5 - \frac{\sum_{j=1}^{10} X_j}{6} + \frac{\sum_{j=1}^{10} X_j^2}{10} - \epsilon \quad (DGP.9),$$

$$D^* = \frac{-0.1 + 0.1 \sum_{j=1}^5 X_j}{\exp(-0.2 \sum_{j=1}^{10} X_j)} - \epsilon \quad (DGP.10).$$

We generate $R = 1000$ samples of sizes $n \in \{200, 600, 1000\}$. The value of R is the same of SS, allowing fair comparisons with their results.

The null hypothesis to test is given by,

$$H_0 : p(X) = \Phi(X'\theta_0) \text{ a.s., for some } \theta_0 \in \Theta,$$

where, $p(X) = P(D = 1|X)$, $X = (1, X_1, \dots, X_{10})$, and $\Phi(x)$ denotes the CDF of the standard normal at x .

Since the model is nonlinear, the grouped GMM estimator does not have a closed-form solution, necessitating the use of numerical methods for estimation. In these simulations and in the empirical application discussed in the last section, the grouped GMM estimator is obtained with the "gmm" package in R ([Chaussé 2010](#)). This function requires the specification of the moment conditions and the weighting matrix, which are set equal

to $n^{-1/2}\hat{\Phi}_\gamma(\theta)$ and $\Sigma_\gamma(\tilde{\theta})$, respectively. Here, $\tilde{\theta}$ represents the (linear) probit maximum likelihood estimator (MLE):

$$\tilde{\theta} = \arg \max_{\theta \in \Theta} \sum_{i=1}^n \{D_i \log \Phi(X_i' \theta) + (1 - D_i) \log(1 - \Phi(X_i' \theta))\}.$$

The $\hat{\chi}^2$ criterion is minimized iteratively using the Nelder-Mead algorithm with 1000 iterations, the Conjugate Gradient algorithm with 1000 iterations, and the SANN algorithm with 200 iterations. At each step, if the criterion converges, the procedure stops; otherwise, it proceeds to minimize with the next algorithm. The tolerance level is set to e^{-15} . It's worth noting that this procedure does not guarantee the attainment of the global minimum, and slight size distortions may be encountered as a result.

The probit MLE estimator is also employed in the construction of the Wald test statistics. In this case, the covariance matrix is given by,

$$\text{Avar} \left(\Phi_\gamma(\tilde{\theta}) \right) = \Sigma_{\gamma,0} - \boldsymbol{\mu}_{\gamma,0}^* \mathcal{I}^{-1} \boldsymbol{\mu}_{\gamma,0}^{*'},$$

where $\boldsymbol{\mu}_{\gamma,0}^* = \mathbb{E}[\phi(X'\theta_0)\mathbf{I}_\gamma(X)X']$, $\phi(x)$ denoting the density of the standard normal at x , and,

$$\mathcal{I} = \mathbb{E} \left[\frac{\phi(X'\theta_0)^2 X X'}{\Phi(X'\theta_0)(1 - \Phi(X'\theta_0))} \right],$$

is the Fisher information matrix. Under the null, the MLE estimator is more efficient than the grouped GMM estimator and $\text{Avar} \left(\Phi_\gamma(\tilde{\theta}) \right)$ is positive definite. Therefore, it can be estimated using the plug-in estimator,

$$\hat{W}_\gamma(\tilde{\theta}) = \hat{\Sigma}_\gamma(\tilde{\theta}) - \hat{\boldsymbol{\mu}}_{\gamma,0}^* \left(\frac{1}{n} \sum_{i=1}^n \frac{\phi(X_i' \tilde{\theta})^2 X_i X_i'}{\Phi(X_i' \tilde{\theta})(1 - \Phi(X_i' \tilde{\theta}))} \right)^{-1} \hat{\boldsymbol{\mu}}_{\gamma,0}^{*'},$$

where $\hat{\sigma}_{\gamma_l}^2(\theta) = n^{-1} \sum_{i=1}^n \Phi(X_i' \tilde{\theta})(1 - \Phi(X_i' \tilde{\theta}))\mathbb{I}_{\gamma_l}(X_i)$, and $\hat{\boldsymbol{\mu}}_{\gamma,0}^* = n^{-1} \sum_{i=1}^n \phi(X_i' \tilde{\theta})\mathbf{I}_\gamma(X_i)X_i'$.

The comparison is done with the test proposed in SS, which is based on,

$$R_{3,n}(u) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \left(D_i - \Phi(X_i' \tilde{\theta}) \right) \left(\mathbb{I}\{\Phi(X_i' \tilde{\theta}) \leq u\} - g(X, \tilde{\theta})' \Delta^{-1}(\tilde{\theta}) G(u, \tilde{\theta}) \right),$$

where $g(X, \theta) = d/d\theta \Phi(X'\theta)$ is the score function, $\Delta(\theta) = n^{-1} \sum_{i=1}^n g(X, \theta)g(X, \theta)'$, and

$G(u, \theta) = n^{-1} \sum_{i=1}^n g(X, \theta) \mathbb{I}\{\Phi(X'\theta) \leq u\}$. The additional term in the right parenthesis is the projection of $\mathbb{I}\{\Phi(X'\tilde{\theta}) \leq u\}$ on $g(X, \theta)$ and has the purpose of making the limit null distribution of the test independent from the estimation method used for $\tilde{\theta}$. The test is given by the Kolmogorov-Smirnov functional,

$$KS3 = \sup_{u \in \mathbb{R}} |R_{3,n}(u)|,$$

and the limit null distribution is approximated using a multiplier bootstrap procedure, described in the paper, with $B = 999$ bootstrap repetitions.

The data is classified using the same partitioning methods as those used in the linear model simulations, considering L values from the set $\{10, 12, 14\}$ for the number of cells. For both the **FNP** and **PNP** algorithms, the minimum number of observations in each cell of the initial two-cell partition is set to $n_{\min} = n/5$. The **PNP** partitions under the null are constructed using the specification of DGP.6 as the alternative model. In the **FNP** method, we employ a polynomial expansion of the fitted values with an order of $q = 3$.

Table 3 below presents the simulation results for the probit models. Under the null hypothesis, i.e. under DGP.6, the Wald test exhibits rejection rates that closely match the nominal significance level of $\alpha = 0.05$, even for moderate sample sizes and across all partitioning methods. In contrast, the $\hat{\chi}^2$ test consistently displays size distortions, even with large sample sizes. The distortions stem from the approximation involved in the numerical methods used for the minimization of the χ^2 criterion. These approximations can lead the test to either over-reject or under-reject, depending on how far the approximation is from the minimum. Achieving greater accuracy in the approximation can be attained by increasing the number of iterations in the minimization algorithm and reducing the tolerance level. However, this enhanced accuracy comes at the expense of increased computational demands. Rejection rates of the χ^2 tests under the other DGPs demonstrate the tests' effectiveness in detecting deviations from the null hypothesis. Regardless of the DGP or the partitioning method employed, the Wald test outperforms the $\hat{\chi}^2$ test.

In contrast with the linear model simulations, χ^2 tests using the **PSEB** partitioning exhibit significantly superior performance compared to those under **FSEB** partitioning.

Table 3: Probit Model: Size and Power

			PSEB		FSEB		FNPF		PNP		
<i>DGP</i>	<i>n</i>	<i>L</i>	\hat{W}	$\hat{\chi}^2$	\hat{W}	$\hat{\chi}^2$	\hat{W}	$\hat{\chi}^2$	\hat{W}	$\hat{\chi}^2$	KS3
6	200	10	7.60	NA	7.00	NA	5.90	NA	6.50	NA	6.10
		12	6.40	5.92	7.30	8.37	6.80	7.05	6.00	8.56	
		14	7.60	2.51	6.90	1.87	6.50	2.33	7.00	1.98	
6	600	10	5.80	NA	5.20	NA	5.09	NA	4.39	NA	5.20
		12	4.39	6.90	3.80	6.20	4.80	8.20	4.00	7.12	
		14	5.60	2.00	5.30	2.60	5.30	1.70	5.50	2.81	
6	1000	10	4.59	NA	6.20	NA	5.20	NA	5.30	NA	6.00
		12	5.40	7.90	6.20	9.10	6.00	8.30	6.10	8.50	
		14	3.90	2.19	6.10	2.40	5.80	4.10	4.70	3.50	
7	200	10	89.70	NA	11.10	NA	89.30	NA	90.00	NA	14.90
		12	87.80	72.49	11.89	9.14	88.20	77.05	88.90	74.51	
		14	88.30	58.41	15.10	4.50	87.50	60.88	88.00	59.16	
7	600	10	100.00	NA	32.60	NA	100.00	NA	100.00	NA	33.00
		12	100.00	100.00	30.50	22.10	100.00	98.40	100.00	98.20	
		14	100.00	99.90	28.60	10.30	100.00	94.60	100.00	99.50	
7	1000	10	100.00	NA	51.50	NA	100.00	NA	100.00	NA	52.90
		12	100.00	100.00	49.90	37.70	100.00	99.20	100.00	98.70	
		14	100.00	100.00	45.40	20.70	100.00	98.80	100.00	99.00	
8	200	10	44.40	NA	19.40	NA	24.00	NA	82.69	NA	18.70
		12	44.40	23.31	20.90	13.37	25.40	16.07	82.19	67.65	
		14	39.20	11.20	20.10	6.08	21.20	4.59	80.50	47.01	
8	600	10	95.70	NA	70.00	NA	81.10	NA	100.00	NA	62.10
		12	94.10	73.30	66.80	46.70	78.20	60.00	100.00	98.10	
		14	95.00	54.10	67.00	32.30	78.00	43.70	100.00	96.09	
8	1000	10	100.00	NA	90.90	NA	96.40	NA	100.00	NA	82.10
		12	100.00	92.70	88.50	77.20	96.00	87.40	100.00	99.80	
		14	99.90	84.60	88.70	63.10	95.90	78.50	100.00	99.80	
9	200	10	27.30	NA	12.70	NA	25.80	NA	14.39	NA	15.10
		12	24.70	17.36	13.60	9.60	24.40	18.88	14.80	10.99	
		14	24.00	6.67	13.30	3.89	22.20	11.03	13.10	4.41	
9	600	10	71.30	NA	33.60	NA	64.20	NA	31.30	NA	36.10
		12	70.59	48.00	34.40	21.40	62.90	51.45	32.40	19.53	
		14	66.70	25.40	29.20	8.61	57.30	27.74	31.20	8.61	
9	1000	10	93.10	NA	53.20	NA	88.30	NA	53.40	NA	57.50
		12	91.80	72.70	52.80	38.60	86.30	76.00	49.90	32.50	
		14	93.40	53.10	49.10	16.30	84.70	54.55	49.80	16.90	
10	200	10	9.10	NA	8.69	NA	10.60	NA	10.40	NA	10.10
		12	8.60	9.62	9.90	11.13	11.30	11.79	10.50	9.54	
		14	10.50	2.52	7.50	2.03	9.19	2.88	8.10	2.64	
10	600	10	28.49	NA	25.30	NA	29.10	NA	30.10	NA	28.10
		12	27.80	13.60	22.80	21.40	26.40	23.82	28.19	20.20	
		14	23.40	3.80	21.70	9.10	24.20	11.70	25.60	9.60	
10	1000	10	56.39	NA	54.20	NA	54.70	NA	62.20	NA	51.70
		12	52.00	18.70	50.20	40.20	52.70	41.80	58.30	39.20	
		14	47.50	6.00	45.00	20.39	45.10	21.10	55.10	21.80	

Percentage of rejections at the nominal level $\alpha = 0.05$ of the $\hat{\chi}^2$ and the Wald test under different partitioning methods. *DGP*, *n*, and *L*, respectively denote the data generating process, the sample size, and the number of cells. The last column report the rejection rates of the [Sant'Anna and Song \(2019\)](#) Kolmogorov-Smirnov test.

Interestingly, under DGP.9, the tests utilizing **PSEB** partitioning even surpass those using

the **PNP** method. Part of this notable performance improvement can be attributed to the non-zero covariance between the first two regressors, which enhances the informativeness of the principal components reduction. Furthermore, with the exception of DGP.10, where the performance of $\hat{\mathcal{W}}$ is comparable to that of the KS3 test, the Wald test consistently outperforms the test proposed by SS under various partitioning schemes, including **PSEB**, **FNP**, and **PNP**, with slight enhancements observed in the case of **FSEB** partitioning. The comparison between KS3 and the $\hat{\chi}^2$ test exhibits greater nuance, depending on the specific DGPs and partitioning algorithms considered. For example, under DGP.7, the $\hat{\chi}^2$ test outperforms KS3 under **PSEB** partitioning, while the opposite is true under **FSEB** partitioning.

8 Empirical Application

I present an empirical application of the testing procedures by revisiting the analysis conducted by [Fryer and Greenstone \(2010\)](#) on the impact of attending Historically Black Colleges and Universities (HBCU), versus non-HBCU, on the labor market and educational outcomes of African American students. The question of whether HBCU are comparatively more effective than Traditional White Institutions (TWI) in advancing academic achievement and economic success for black students holds significant policy implications. As discussed by [Price and Viceisza \(2023\)](#), HBCU play a pivotal role in enhancing the identity, confidence, and self-esteem of students, particularly Black individuals. These institutions have a higher ratio of Black students, making it easier for students to foster a sense of belonging. HBCU also offer innovative programs, such as college preparation summer programs and Black-centered curricula, empowering students to engage in research experiences and actively contribute to their fields. However, HBCU have historically relied heavily on federal government financial support and have been consistently underfunded in comparison to TWI. This financial disparity can substantially influence the quality of education provided by HBCU, potentially limiting the opportunities and experiences accessible to their students. Thus, the suggestion of merging HBCU with other institutions or redirecting resources toward existing TWI could lead to a more efficient and cost-effective higher education system.

The analysis is based on the National Longitudinal Survey of the High School Class of 1972 (NLS) used by [Fryer and Greenstone \(2010\)](#), which is publicly available on OpenICPSR¹. The survey is a nationally representative sample of 23,451 high school seniors in 1972, with follow-up interviews conducted in 1973, 1974, 1976, 1979, and 1986. This dataset provides rich information concerning these students, including details about their demographics, academic performance, and labor market outcomes.

The sample comprises 624 African American students who pursued higher education after high school and were tracked throughout the 1972-1986 period. Among them, 260 attended an HBCU, while 364 attended a TWI. The considered outcomes include the logarithm of hourly wage in 1986, the probability of attaining a bachelor's degree, and the probability of obtaining a graduate degree.

Following [Fryer and Greenstone \(2010\)](#), I consider a set of covariates consisting of home environment variables and pre-college characteristics. The former include family income, mother and father's education, and whether or not a student lives in the south. The latter consists of SAT and ACT scores, and whether or not the student attended a private school.

Table 4: Summary Statistics of Black Students (NLS)

	HSBCU			TWI		
	Mean	St. deviation	NA	Mean	St. deviation	NA
ln(wage)	2.16	0.64	57	2.11	0.51	71
Bachelor's degree	0.67	0.46	0	0.62	0.48	2
Graduate degree	0.14	0.35	0	0.16	0.37	2
SAT	685.55	133.46	143	793.68	183.39	212
ACT	12.86	4.02	201	14.73	5.31	269
Family income	3.45	2.50	59	3.89	2.70	75
Father education	1.71	0.96	4	1.96	1.16	9
Mother education	1.95	1.08	3	2.08	1.11	4
Private high school	0.03	0.17	0	0.05	0.23	0
Female	0.66	0.47	0	0.64	0.47	0
South	0.87	0.33	0	0.48	0.50	0

The table report mean, standard deviation, and number of missing values of the outcomes and characteristics of African American who attended college. The columns HSBCU and TWI refer to the type of college attended. In the first three rows, the outcomes are the logarithm of hourly wage in 1986, the probability of receiving a bachelor's degree, and the probability of enrolling in graduate school. The remaining rows report the covariates. The total sample size is $n = 624$, with $n_1 = 260$ students enrolled in HBCU, and $n_2 = 364$ in TWI.

¹Ann Arbor, MI: Inter-university Consortium for Political and Social Research, 2019-10-12. <https://doi.org/10.3886/E113736V1>

The treatment of missing values consists of generating a dummy variable for each covariate indicating whether the value is missing or not. The relative NA is set to 0. Differently from [Fryer and Greenstone \(2010\)](#), however, I do not partition the set of covariates into categories.

The estimation of the effect of enrollment in an HBCU is conducted within the unconfoundedness framework (see, e.g., [Rosenbaum and Rubin 1983](#)). The individual outcomes are decomposed into two components: $Y_i = Y_i(1)D_i + Y_i(0)(1 - D_i)$, where $Y_i(1)$ and $Y_i(0)$ represent the potential outcomes under treatment and control, respectively, while D_i is a binary variable indicating whether the individual attended an HBCU. The following selection on observables and overlap assumptions allow to identify the average treatment effect (ATE) and the average treatment effect on the treated (ATT).

Assumption 6 $(Y_i(1), Y_i(0)) \perp\!\!\!\perp D_i | X_i$ and $0 < p(x) < 1 \quad \forall x \in \mathcal{X}$.

Here X_i denotes the vector of home-environment and pre-college characteristics, while $p(X) = P(D = 1|X)$ denotes the propensity score. Under Assumption 6, the ATE and ATT are identified ([Rosenbaum 1987](#)) by,

$$\text{ATE} = \mathbb{E} \left[\left(\frac{D}{p(X)} - \frac{1-D}{1-p(X)} \right) Y \right] \quad \text{and} \quad \text{ATT} = \frac{\mathbb{E} \left[\left(D - \frac{p(X)(1-D)}{1-p(X)} \right) Y \right]}{\mathbb{E}[D]}.$$

This motivates a two-step procedures where first we obtain an estimate of the propensity score, and then we estimate the ATE and ATT by using the analog principle.

It is worth remarking that the validity of the conditional independence in Assumption 6 relies on the comprehensive inclusion of all relevant factors in a student's decision-making process regarding enrollment in an HBCU institution. When uncertainties surround the validity of the independence assumption, the estimates of ATE and ATT presented below should be regarded as statistical associations rather than causal effects.

The estimation of the propensity score is done by fitting a probit model, i.e., $p(X) = \Phi(X'\theta_0)$, where $\Phi(\cdot)$ denotes the cumulative distribution function of the standard normal distribution. I consider two different specifications of the propensity score model,

- (a) This specification includes all the covariates listed in Table 4, along with the inclusion of the relative missing values dummy.

- (b) In addition to the elements in specification (a), specification (b) introduces interactions between the variables "south" and "SAT", "south" and "family income", as well as a dummy variable for family income above the upper quartile.

Specification (b) is more comprehensive and accommodates the consideration of varying effects of family income and SAT scores on individuals originating from the southern region, as well as the differential effect for families with higher income. This choice is influenced by the significant presence of HBCUs in the south, suggesting lower mobility costs for southern individuals attending these institutions. Proximity to HBCUs may diminish the perceived significance of standardized test scores, as students may find attending an HBCU more accessible and familiar, regardless of their SAT performance. Conversely, as noted by [Fryer and Greenstone \(2010\)](#), the presence of HBCUs in the south might intensify competition for students, potentially leading to higher enrollment of students with elevated SAT scores. It's also plausible that highly proficient students in the south, residing in predominantly black neighborhoods, may prefer enrolling in distant TWI institutions to circumvent the "acting white" stigma ([Fryer Jr and Torelli 2010](#)). The varied impact of high family income is driven by the inclination of students from affluent families to favor TWIs over HBCUs.

To ensure the stability and consistency of the estimated causal parameters using the procedure outlined above, it is crucial to check the correct specification of the propensity score estimates. To achieve this, a critical step involves verifying the correct specification of the propensity score estimates, i.e. testing the null hypothesis $H_0 : \mathbb{E}[D|X] = \Phi(X'\theta_0)$ a.s. for some $\theta_0 \in \Theta$. Table 5 below reports the $\hat{\chi}^2$, the Wald test, and the SS test for each of the two propensity score specifications, while Table 5 compare the ATE and ATT estimates obtained from these two specifications.

The parameter θ_0 is estimated through maximum likelihood estimation, and the χ^2 tests are constructed as described in Section 7. The J test, $\hat{\chi}^2$, is estimated via numerical approximation using the Nelder-Mead algorithm with 50,000 iterations, the Conjugate Gradient algorithm with 50,000 iterations, and the SANN algorithm with 10,000 iterations, with a tolerance level set to e^{-15} . The SS test is implemented using $B = 100,000$ bootstrap replications. The construction of **PNP** partition is done using the opposite specification as reference, meaning that the test for specification (a) assumes the alterna-

tive model to be specification (b), and vice-versa. The **FNP** partition uses a polynomial expansion of the fitted values of order $q = 3$, and the **PSEB** partition is build using only the first principal component. The comparison is done with the SS test using the Kolmogorov-Smirnov and the Cramer-von Mises distance functions. Following common practice, the ATE and ATT are estimated from the subsample of observations with estimated propensity score outside of the $[0.05, 0.95]$ interval.

Table 5 presents the specification tests results. In summary, the evidence suggests that model (a) is not correctly specified, whereas model (b) appears to offer a better fit to the data, although minor misalignments are still present. Over the 40 Wald tests (10 for each partitioning method), the null hypothesis of correct specification for model (a) is rejected 11 times at a 10% significance level, whereas for model (b), it is rejected only 2 times. Notably, most of these rejections occur when employing **PNP** partitioning, which is based on comparing specifications (a) and (b), and under **PSEB** partitioning. The effectiveness of the latter in detecting deviations is supported by the Monte Carlo study in Section 7. The $\hat{\chi}^2$ test, on the other hand, rejects the null hypothesis 3 times out of 12 tests for model (a) and 3 times out of 4 for model (b). However, it's important to exercise caution as these rejections are predominantly associated with the largest values of L , where cells have very few observations. The small number of within-cell observations in these cases can influence the test's size control. Regarding the KS test from [Sant'Anna and Song \(2019\)](#), it rejects the null hypothesis for model (a) but not for model (b). However, the Cramer-von Mises (CvM) functional does not detect misspecifications for either model.

The estimated treatment effects under the two specifications present substantial differences. In specification (a), positive and statistically significant ATE estimates are observed for both log-wage and college degree outcomes, at the 10% nominal level. However, the ATT estimates for these outcomes display substantial noise, making it impossible to reject the null hypothesis of zero effect. When considering the distinct influence of family income and SAT scores on students from the South, and the influence of high-family-income students only the ATE estimate for the college degree outcome remains positive and statistically significant. The remaining estimates are not statistically significant at the 10% nominal level. Consequently, in the preferred model specification, attending HBCU institutions appears to have only a minimal (albeit positive) impact on the educational

Table 5: Tests Results

Spec.	L	PSEB		FSEB		FNP		PNP		KS3	CvM3
		\hat{W}	$\hat{\chi}^2$	\hat{W}	$\hat{\chi}^2$	\hat{W}	$\hat{\chi}^2$	\hat{W}	$\hat{\chi}^2$		
(a)	2	0.91	NA	0.13	NA	0.65	NA	0.00	NA	0.02	0.45
	4	0.05	NA	0.18	NA	0.88	NA	0.00	NA		
	6	0.11	NA	0.55	NA	0.02	NA	0.02	NA		
	8	0.06	NA	0.25	NA	0.80	NA	0.06	NA		
	10	0.03	NA	0.53	NA	0.64	NA	0.11	NA		
	12	0.26	NA	0.80	NA	0.08	NA	0.00	NA		
	14	0.15	NA	0.18	NA	0.45	NA	0.00	NA		
	16	0.20	0.02	0.41	0.18	0.85	0.30	0.39	0.02		
	18	0.12	0.62	0.14	0.34	0.15	0.45	0.14	0.08		
	20	0.10	0.68	0.45	0.59	0.14	0.15	0.14	0.41		
(b)	2	0.73	NA	0.63	NA	0.49	NA	0.63	NA	0.46	0.66
	4	0.37	NA	0.91	NA	0.33	NA	0.64	NA		
	6	0.42	NA	0.81	NA	0.37	NA	0.91	NA		
	8	0.65	NA	0.94	NA	0.15	NA	0.88	NA		
	10	0.10	NA	0.39	NA	0.19	NA	0.89	NA		
	12	0.27	NA	0.92	NA	0.56	NA	0.86	NA		
	14	0.52	NA	0.83	NA	0.79	NA	0.68	NA		
	16	0.14	NA	0.86	NA	0.43	NA	0.96	NA		
	18	0.42	NA	0.95	NA	0.17	NA	0.96	NA		
	20	0.07	0.00	0.39	0.10	0.49	0.12	0.86	0.02		

The table reports the p-values of the $\hat{\chi}^2$ and the Wald test, under different partitioning methods, for the null hypothesis of correct specification of model (a) and (b). L denotes the number of cells. The last two columns report the p-values of the Kolmogorov-Smirnov and the Cramer-von Mises functional of the [Sant'Anna and Song \(2019\)](#) test.

outcomes of African American students. In terms of long-term outcomes, the evidence indicates that HBCUs do not exhibit a higher level of effectiveness than TWIs in promoting the economic achievement of Black students. These findings contrast with the results reported by [Fryer and Greenstone \(2010\)](#), who found evidence of a positive effect of HBCU attendance on the wages of Black students in 1986 and their likelihood of attaining a bachelor's degree. Nonetheless, the results of [Fryer and Greenstone \(2010\)](#) align with our findings when considering the probability of obtaining a graduate degree. It's important to note that the estimation of HBCU attendance on wages is affected by selection bias, primarily because of the absence of wage outcomes for unemployed individuals. To mitigate this concern, one possible approach could involve employing the correction method proposed by [Heckman \(1979\)](#) and assessing the overall model fit through our testing procedures. Addressing this issue lies beyond the scope of the current section.

Table 6: ATE and ATT Estimates

Specification (a)				Specification (b)		
	ln(wage)	Col. Degree	Grad. Degree	ln(wage)	Col. Degree	Grad. Degree
ATE	14.44	8.51	2.72	6.78	7.80	0.92
	(6.28)	(4.53)	(3.77)	(5.43)	(4.56)	(3.30)
	[0.02]	[0.06]	[0.46]	[0.21]	[0.08]	[0.77]
ATT	2.87	3.59	0.64	1.03	3.43	0.63
	(6.43)	(4.94)	(3.11)	(6.89)	(5.04)	(3.15)
	[0.65]	[0.46]	[0.83]	[0.88]	[0.49]	[0.84]

The table reports estimates of the ATE and ATT using two different specifications of the propensity score for the following outcomes: the logarithm of hourly wage in 1986, the probability of obtaining a bachelor's degree, and the probability of obtaining a graduate degree. Estimates and standard errors (in parentheses) are multiplied by 100, and p-values are enclosed in brackets. The estimates are obtained by excluding observations with a propensity score outside the interval $[0.05, 0.95]$.

9 Conclusion

In conclusion, this article introduces a novel approach for validating the correct specification of conditional moment restriction (CMR) models, which are fundamental for identifying and estimating causal relationships. The proposed method harnesses the well-established chi-squared (χ^2) tests, commonly used in classical goodness-of-fit contexts, to goodness-of-fit checks for CMR specifications. Traditionally, CMR model checks have been adapted from tests for the cumulative density function (CDF) based on functionals of the standard empirical process (SEP), but these often exhibit limitations, especially when dealing with high-dimensional data and high-frequency alternatives. The introduced χ^2 tests, on the other hand, are distribution-free, do not necessitate complex bootstrapping or smoothing techniques, and offer flexibility in partitioning the data to favor specific alternative hypotheses. Thus, providing a valid complementary instruments to the existing model checking procedures. Monte Carlo simulations and empirical evidence highlight the effectiveness of these χ^2 tests, especially in scenarios involving a large number of covariates. The empirical analysis demonstrates the practical application of these tests in assessing the returns of attending historically black colleges and universities (HBCUs) for black students in the United States. The results suggest that, once we account for the effect of family income and SAT scores on students from the South, HBCUs are not more effective than traditional white institutions (TWIs) in advancing the economic success of

black students, although they do influence the probability of obtaining a college degree.

A Appendix A

A.1 Lemmas

I first state auxiliary lemmas for the propositions and theorems in the main text. Let \rightsquigarrow denote weak convergence on $l^\infty(\mathbb{D})$ (see definition 13.3 in [Van Der Vaart 1996](#), hereafter VW), where $l^\infty(\mathbb{D})$ is the space of all real-valued functions that are uniformly bounded on \mathbb{D} , and \rightarrow_d denote convergence of real-valued random variables. Troughout, to highlight the dependency on the partition, denote as $\hat{\Phi}_\theta(\gamma) = \hat{\Phi}_\gamma(\theta)$ and $\hat{\Phi}_0(\gamma) = \hat{\Phi}_\gamma(\theta_0)$.

Lemma 1 *Under the null H_0 ,*

- (a) *if Assumption 1, 2, and 2' hold, then $\Sigma_\gamma(\tilde{\theta}) = \Sigma_{\gamma,0} + o_p(1)$.*
- (b) *if Assumption 1, 2, 3 hold, and $\text{Avar}\left(\hat{\Phi}_\gamma(\tilde{\theta})\right)$ is p.d., then $\hat{W}_\gamma(\tilde{\theta}) = \text{Avar}\left(\hat{\Phi}_\gamma(\tilde{\theta})\right) + o_p(1)$*

Lemma 2 *Under the null hypothesis H_0 , Assumptions 1, and 5,*

$$\hat{\Phi}_0(\cdot) \rightsquigarrow \Phi_0(\cdot) \text{ as a process on } l^\infty(\mathbb{D}),$$

where $\Phi_0(\cdot)$ is an \mathbb{R}^L -valued Gaussian process with zero mean vector and covariance structure given by,

$$\mathbb{E}[\Phi_0(\gamma)\Phi_0(\tilde{\gamma})] = \mathbb{E}[\varepsilon_{\theta_0}(Z)^2 \mathbf{I}_\gamma(X)\mathbf{I}_{\tilde{\gamma}}(X)'] \quad \forall \gamma, \tilde{\gamma} \in \mathbb{D}.$$

Lemma 3 *Under the null hypothesis H_0 , and Assumptions 1-5, it holds that:*

- (a) $\sup_{\gamma \in \mathbb{D}} \left| \hat{\Phi}_\theta(\gamma) - (\hat{\Phi}_0(\gamma) - \hat{\mu}_\gamma^*(\theta_0)\sqrt{n}(\hat{\theta} - \theta_0)) \right| = o_p(1).$
- (b) $\hat{\mu}_\gamma^*(\theta_0) = \mu_{\gamma,0}^* + o_p(1).$

Lemma 4 *Under the null H_0 , and Assumptions 4,5,*

- (a) *if Assumption 1, 2, and 2' hold, then $\Sigma_{\tilde{\gamma}}(\tilde{\theta}) = \Sigma_{\gamma,0} + o_p(1)$.*
- (b) *if Assumption 1, 2, 3 hold, and $\text{Avar}\left(\hat{\Phi}_\gamma(\tilde{\theta})\right)$ is p.d., then $\hat{W}_{\tilde{\gamma}}(\tilde{\theta}) = \text{Avar}\left(\hat{\Phi}_\gamma(\tilde{\theta})\right) + o_p(1)$*

A.2 Proofs

For any class of functions \mathcal{F} , denote as $\{P_n f : f \in \mathcal{F}\}$ the empirical measure indexed by \mathcal{F} , such that $P_n f = n^{-1} \sum f(Z_i)$; alike, Pf denotes the population measure, $Pf = \int f(Z)dP$. We say that a class of functions is: i) Glivenko-Cantelli for P (hereafter, P -GC) whenever $\sup_{f \in \mathcal{F}} |P_n - P|f = o_p(1)$; ii) P -Donsker if $\{\sqrt{n}(P_n - P)f : f \in \mathcal{F}\}$ converge in distribution to a tight random element in the space $l^\infty(\mathcal{F})$. Throughout, we refer to both classes of sets with finite VC dimension and classes of functions with finite VC subgraph dimension as VC classes. These classes, having uniformly bounded covering numbers (Theorem 2.6.7 in VW), are Glivenko-Cantelli and Donsker (see Theorem 2.4.3 and 2.5.2 in VW) for any probability measure on the sample space, provided that they have integrable and square-integrable envelope function, respectively.

Proof of Lemma 1. For the first part of the Lemma, by the weak law of large numbers (WLLN) and a mean value theorem argument (MVT), suffices to show that

$$\frac{1}{n} \sum_{i=1}^n (\varepsilon_{\tilde{\theta}}^2(Z_i) - \varepsilon_{\theta_0}^2(Z_i)) \mathbb{I}_{\gamma_l}(X_i) = I + II + III = o_p(1)$$

for each $l \in 1, 2, \dots, L$, where,

$$\begin{aligned} I &= (\tilde{\theta} - \theta_0)' \frac{1}{n} \sum_{i=1}^n \nabla m_{\tilde{\theta}}(X_i) \nabla m_{\tilde{\theta}}(X_i)' \mathbb{I}_{\gamma_l}(X_i) (\hat{\theta} - \theta_0), \\ II &= \frac{2}{n} \sum_{i=1}^n m_{\theta_0}(X_i) \nabla m_{\tilde{\theta}}(X_i)' \mathbb{I}_{\gamma_l}(X_i) (\tilde{\theta} - \theta_0), \\ III &= \frac{2}{n} \sum_{i=1}^n Y_i (m_{\tilde{\theta}}(X_i) - m_{\theta_0}(X_i)) \mathbb{I}_{\gamma_l}(X_i), \end{aligned}$$

and $|\tilde{\theta} - \theta_0| \leq |\tilde{\theta} - \theta_0|$. The triangle inequality, Assumption 2, and the consistency of $\tilde{\theta}$ show that, $|I| \leq d_{\tilde{\theta}}^2 \left\| \tilde{\theta} - \theta_0 \right\|^2 n^{-1} \sum_{i=1}^n R(X_i)^2 = o_p(1)$, where $\|\cdot\|$ denotes the euclidean

norm. By a similar reasoning,

$$\begin{aligned} |II| &\leq d_\theta \left\| \tilde{\theta} - \theta_0 \right\| \frac{2}{n} \sum_{i=1}^n m_{\theta_0}(X_i) R(X_i) \\ &\leq d_\theta \left\| \tilde{\theta} - \theta_0 \right\| \left(\mathbb{E} [Y^2] \right)^{1/2} \left(\mathbb{E} [R(X)^2] \right)^{1/2} + o_p(1) = o_p(1) \end{aligned}$$

where the last inequality follows from the WLLN, the law of iterated expectation, and Cauchy-Schwarz inequality. Finally, after expanding again around θ_0 it is easy to see that $|III| \leq d_\theta \left\| \tilde{\theta} - \theta_0 \right\| 2n^{-1} \sum_{i=1}^n Y_i R(X_i) \mathbb{I}_{\gamma_l}(X_i) = o_p(1)$.

For the second part of the lemma, we need to show that $\hat{L}(\tilde{\theta}) = L_0 + o_p(1)$, $\hat{C}_\gamma(\tilde{\theta}) = C_{\gamma,0} + o_p(1)$, and $\hat{\mu}_\gamma^*(\tilde{\theta}) = \mu_{\gamma,0}^* + o_p(1)$. By the usual MVT argument and the law of large numbers,

$$L_n = L_0 + I + II + II' + o_p(1)$$

with,

$$\|I\| = \left\| \frac{1}{n} \sum_{i=1}^n \nabla l_{\tilde{\theta}}(Z_i) (\hat{\theta} - \theta_0) (\hat{\theta} - \theta_0)' \nabla l_{\tilde{\theta}}(Z_i)' \right\| \leq d_\theta^4 \left\| \hat{\theta} - \theta_0 \right\|^2 \frac{1}{n} \sum_{i=1}^n R_2^2(Z_i) = o_p(1),$$

$$\begin{aligned} \|II\| &= \left\| \frac{1}{n} \sum_{i=1}^n \nabla l_{\tilde{\theta}}(Z_i) (\hat{\theta} - \theta_0) l'_{\theta_0}(Z_i) \right\| \leq d_\theta^2 \left\| \hat{\theta} - \theta_0 \right\| \frac{1}{n} \sum_{i=1}^n \|l_{\theta_0}(Z_i)\| R_2(Z_i) \\ &\leq d_\theta^2 \left\| \hat{\theta} - \theta_0 \right\| \left(\frac{1}{n} \sum_{i=1}^n \|l_{\theta_0}(Z_i)\|^2 \right)^{1/2} \left(\frac{1}{n} \sum_{i=1}^n R_2^2(Z_i) \right)^{1/2} = o_p(1). \end{aligned}$$

Alike, we write $\hat{C}_\gamma(\tilde{\theta})$ as,

$$\hat{C}_\gamma(\tilde{\theta}) = I - II - III + C_{\gamma,0} + o_p(1),$$

where,

$$\begin{aligned}
I &= \frac{1}{n} \sum_{i=1}^n l_{\theta_0}(Z_i) \nabla m_{\bar{\theta}}(X_i)' (\hat{\theta} - \theta_0) \mathbf{I}_{\gamma}(X_i)' \\
II &= \frac{1}{n} \sum_{i=1}^n \nabla l_{\bar{\theta}}(Z_i)' (\hat{\theta} - \theta_0) \varepsilon_{\theta_0}(Z_i) \mathbf{I}_{\gamma}(X_i)' \\
III &= \frac{1}{n} \sum_{i=1}^n \nabla l_{\bar{\theta}}(Z_i) (\hat{\theta} - \theta_0) (\hat{\theta} - \theta_0)' \nabla m_{\bar{\theta}}(X_i) \mathbf{I}_{\gamma}(X_i)'
\end{aligned}$$

By Assumptions 2, and 3,

$$\begin{aligned}
\|I\| &\leq \left\| \hat{\theta} - \theta_0 \right\| \frac{1}{n} \sum_{i=1}^n \|l_{\theta_0}(Z_i)\| \|\nabla m_{\bar{\theta}}(X_i)\| \|\mathbf{I}_{\gamma}(X_i)\| \\
&\leq \sqrt{L} \left\| \hat{\theta} - \theta_0 \right\| \left(\frac{1}{n} \sum_{i=1}^n \|l_{\theta_0}(Z_i)\|^2 \right)^{1/2} \left(\frac{1}{n} \sum_{i=1}^n R^2(X_i) \right)^{1/2} = o_p(1),
\end{aligned}$$

An analogous reasoning shows that $\|II\| = o_p(1)$, and,

$$\begin{aligned}
\|III\| &\leq \left\| \hat{\theta} - \theta_0 \right\|^2 \frac{1}{n} \sum_{i=1}^n \|\nabla l_{\bar{\theta}}(Z_i)\| \|\nabla m_{\bar{\theta}}(X_i)\| \|\mathbf{I}_{\gamma}(X_i)\| \\
&\leq \sqrt{L} d_{\theta}^3 \left\| \hat{\theta} - \theta_0 \right\|^2 \left(\frac{1}{n} \sum_{i=1}^n R^2(X_i) \right)^{1/2} \left(\frac{1}{n} \sum_{i=1}^n R_2^2(Z_i) \right)^{1/2} = o_p(1).
\end{aligned}$$

Finally, $\hat{\boldsymbol{\mu}}_{\gamma}^*(\tilde{\theta}) = \boldsymbol{\mu}_{\gamma,0}^* + o_p(1)$ follows from the proof of Lemma 3 below. ■

Proof of Lemma 2. By Lemma 2.6.17 in VW and Assumption 5, both \mathbb{D} and $\{\mathbf{I}_{\gamma}(x) : \gamma \in \mathbb{D}\}$ are VC classes. Therefore, $\mathcal{F} = \{\varepsilon_{\theta_0}(z) \mathbf{I}_{\gamma}(x) : \gamma \in \mathbb{D}\}$ is a VC class (Lemma 2.6.18 in VW), with square integrable envelope function $F = |\varepsilon_{\theta_0}|$, and, hence, is P -Donsker. The convergence of the finite-dimensional distributions (fidis) of $\hat{\Phi}_0(\cdot)$ to those of $\Phi_0(\cdot)$, by the multivariate central limit theorem, characterize the limit process. ■

Proof of Lemma 3. By an MVT argument,

$$\hat{\Phi}_{\hat{\theta}}(\gamma) = \hat{\Phi}_0(\gamma) - I' \sqrt{n}(\hat{\theta} - \theta_0) - \hat{\boldsymbol{\mu}}_{\gamma}'(\theta_0) \sqrt{n}(\hat{\theta} - \theta_0)$$

where,

$$I = \hat{\mu}_{\gamma}^*(\bar{\theta}) - \hat{\mu}_{\gamma}^*(\theta_0) = \frac{1}{n} \sum_{i=1}^n (\nabla m_{\bar{\theta}}(X_i) - \nabla m_{\theta_0}(X_i)) \mathbf{L}_{\gamma}(X_i)'$$

and $|\bar{\theta} - \theta_0| \leq |\tilde{\theta} - \theta_0|$. The class $\{\nabla m_{\theta}(x) : \theta \in \Theta\}$ is a collection of continuous mapping, $\theta \rightarrow \nabla m_{\theta}$, over the compact metric space Θ with integrable envelope function $R(\cdot)$ and, therefore, is P -GC (e.g., Example 19.8 in [Van der Vaart 2000](#)). Thus,

$$\begin{aligned} \sup_{\gamma \in \mathbb{D}} \|I\| &\leq \sup_{\gamma \in \mathbb{D}} \frac{1}{n} \sum_{i=1}^n \|\nabla m_{\bar{\theta}}(X_i) - \nabla m_{\theta_0}(X_i)\| \|\mathbf{L}_{\gamma}(X_i)\| \\ &\leq \sqrt{L} \frac{1}{n} \sum_{i=1}^n \|\nabla m_{\bar{\theta}}(X_i) - \nabla m_{\theta_0}(X_i)\| = o_p(1). \end{aligned}$$

where the last equality follows from an application of the uniform law of large numbers (e.g., [Davidson 1994](#), Theorem 21.6). For the second part of the Lemma is sufficient to prove that,

$$|II| = \left| \frac{1}{n} \sum_{i=1}^n \nabla^{(j)} m_{\theta_0}(X_i) (\mathbf{L}_{\tilde{\gamma}}(X_i) - \mathbf{L}_{\gamma}(X_i)) \right| = o_p(1),$$

for each $j \in \{1, \dots, d_{\theta}\}$. To see this is true, notice that by Assumption 5, $\mathbb{D}\Delta\mathbb{D} = \{\gamma_1 \Delta \gamma_2 : \gamma_1, \gamma_2 \in \mathbb{D}\}$ is a class of subsets of unions of VC classes, and hence is VC. Therefore, $\{|\nabla^{(j)} m_{\theta_0}(x)| \mathbf{L}_{\tilde{\gamma}}(x) : \tilde{\gamma} \in \mathbb{D}\Delta\mathbb{D}\}$ is also VC with integrable envelope $R(\cdot)$ and, hence, P -GC. Thus, for each $j \in \{1, \dots, d_{\theta}\}$,

$$\begin{aligned} |II| &\leq \frac{1}{n} \sum_{i=1}^n |\nabla^{(j)} m_{\theta_0}(X_i)| \mathbf{L}_{\hat{\gamma} \Delta \gamma}(X_i) \\ &\leq \sup_{\tilde{\gamma} \in \mathbb{D}\Delta\mathbb{D}} (P_n - P) |\nabla^{(j)} m_{\theta_0}| \mathbf{L}_{\tilde{\gamma} \Delta \gamma} + \mathbb{E} [|\nabla^{(j)} m_{\theta_0}(X)| \mathbf{L}_{\hat{\gamma} \Delta \gamma}(X)] \\ &\leq o_p(1) + \mu_R(\hat{\gamma} \Delta \gamma) = o_p(1) \end{aligned}$$

where $\mu_R(\hat{\gamma} \Delta \gamma) = (\mu_R(\hat{\gamma}_1 \Delta \gamma_1), \dots, \mu_R(\hat{\gamma}_L \Delta \gamma_L))'$, and $\mu_R(A) = \int_A R(X) dP$ is a (signed) measure absolutely continuous with respect to P . The last equality follows from Assumption 4. ■

Proof of Lemma 4. For each element on the main diagonal of $\hat{\Sigma}_{\tilde{\gamma}}(\tilde{\theta}) - \hat{\Sigma}_{\gamma}(\tilde{\theta})$ write,

$$\frac{1}{n} \sum_{i=1}^n \varepsilon_{\hat{\theta}}(Z_i)^2 (\mathbb{I}_{\hat{\gamma}_l}(X_i) - \mathbb{I}_{\gamma_l}(X_i)) = I + II + III$$

where

$$\begin{aligned} I &= \frac{1}{n} \sum_{i=1}^n \varepsilon_{\theta_0}(Z_i)^2 (\mathbb{I}_{\hat{\gamma}_l}(X_i) - \mathbb{I}_{\gamma_l}(X_i)) \\ II &= (\hat{\theta} - \theta_0)' \frac{1}{n} \sum_{i=1}^n \nabla m_{\hat{\theta}}(X_i) \nabla m_{\hat{\theta}}(X_i)' (\mathbb{I}_{\hat{\gamma}_l}(X_i) - \mathbb{I}_{\gamma_l}(X_i)) (\hat{\theta} - \theta_0) \\ III &= -(\hat{\theta} - \theta_0)' \frac{2}{n} \sum_{i=1}^n \varepsilon_{\theta_0}(Z_i) \nabla m_{\hat{\theta}}(X_i) (\mathbb{I}_{\hat{\gamma}_l}(X_i) - \mathbb{I}_{\gamma_l}(X_i)) \end{aligned}$$

The class $\{\varepsilon_{\theta_0}(z)^2 \mathbb{I}_{\tilde{\gamma}_l} : \tilde{\gamma} \in \mathbb{C}\Delta\mathbb{C}\}$ is VC with integrable envelope function $\varepsilon_{\theta_0}^2$ and, hence, is P -GC. Therefore, $|I| \leq \mu_{\sigma}(\hat{\gamma}_l \Delta \gamma_l) + o_p(1) = o_p(1)$, by Assumption 4. Also, $|II| \leq \sqrt{L} d_{\hat{\theta}}^2 \left\| \hat{\theta} - \theta_0 \right\|^2 n^{-1} \sum_{i=1}^n R(X_i)^2 = o_p(1)$, by Assumptions 1-2 and the consistency of $\hat{\theta}$, and $|III| \leq \sqrt{L} d_{\hat{\theta}} \left\| \hat{\theta} - \theta_0 \right\| n^{-1} \sum_{i=1}^n \varepsilon_{\theta_0} R(X_i) = o_p(1)$ by Cauchy-Schwarz inequality. Thus, $\hat{\Sigma}_{\tilde{\gamma}}(\tilde{\theta}) = \hat{\Sigma}_{\gamma}(\tilde{\theta}) + o_p(1)$, and the first part of the lemma follows from Lemma 1(a).

For the second part of the lemma, notice that by Lemma 3(b) (and the proof of the first part of Lemma 3), $\hat{\mu}_{\tilde{\gamma}}^*(\hat{\theta}) = \mu_{\gamma,0}^* + o_p(1)$, and for each element of $C_n(\hat{\gamma}) - C_n(\gamma_0)$ it holds that,

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \varepsilon_{\hat{\theta}}(Z_i) l_{\hat{\theta},j} (\mathbb{I}_{\hat{\gamma}_l}(X_i) - \mathbb{I}_{\gamma_l}(X_i)) &\leq \left(\frac{1}{n} \sum_{i=1}^n \varepsilon_{\hat{\theta}}(Z_i)^2 \mathbb{I}_{\hat{\gamma}_l \Delta \gamma_l}(X_i) \right)^{1/2} \left(\frac{1}{n} \sum_{i=1}^n l_{\hat{\theta},j}^2 \right)^{1/2} \\ &= o_p(1) O_p(1), \end{aligned}$$

where $l_{\theta,j}$ denotes the j -th component of l_{θ} and the last equality follows from the first part of this proof and Lemma 1. ■

Proof of Theorem 1. The consistency of the grouped GMM estimator follows from $\sup_{\theta \in \Theta} |Q_n(\theta) - Q_0(\theta)| = o_p(1)$, where $Q_0(\theta) = \mathbb{E} [\varepsilon_{\theta}(Z) \mathbf{I}_{\gamma}(X)]' (\Sigma_{\gamma,0})^{-1} \mathbb{E} [\varepsilon_{\theta}(Z) \mathbf{I}_{\gamma}(X)]$ and $Q_n(\theta) = n^{-1} \hat{\chi}_{\tilde{\gamma},\hat{\theta}}^2(\theta)$ (see Theorem 2.1 in [Newey and McFadden 1994](#), for instance).

To see that the condition holds, notice that,

$$\sup_{\theta \in \Theta} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_{\theta}(Z_i) \mathbf{L}_{\hat{\gamma}}(X_i) - \mathbb{E} [\varepsilon_{\theta}(Z) \mathbf{L}_{\gamma}(X)] \right| \leq \sup_{\theta \in \Theta} |I| + \sup_{\theta \in \Theta} |II| + \sup_{\theta \in \Theta} |III|$$

where

$$\begin{aligned} I &= (P_n - P) \varepsilon_{\theta} \mathbf{L}_{\gamma}, \\ II &= \frac{1}{n} \sum_{i=1}^n \varepsilon_{\theta_0}(Z_i) (\mathbf{L}_{\hat{\gamma}}(X_i) - \mathbf{L}_{\gamma}(X_i)), \\ III &= \frac{1}{n} \sum_{i=1}^n [\nabla m_{\bar{\theta}}(Z_i) (\mathbf{L}_{\hat{\gamma}}(X_i) - \mathbf{L}_{\gamma}(X_i))]' (\theta - \theta_0). \end{aligned}$$

The mapping $\theta \rightarrow \varepsilon_{\theta}$ is continuous over the compact Θ , with

$$\mathbb{E} \left[\sup_{\theta \in \Theta} \varepsilon_{\theta}(Z) \right] \leq \mathbb{E} \left[\sup_{\bar{\theta}, \theta \in \Theta} \varepsilon_{\theta_0}(Z) - \nabla m_{\bar{\theta}}(X)' (\theta - \theta_0) \right] \leq d_{\theta} \mathbb{E} [R(Z)] D < \infty,$$

by Assumption 2 and the compactness of Θ , where D denotes the diameter of Θ . Therefore, both $\{\varepsilon_{\theta}(z) : \theta \in \Theta\}$ and $\{\varepsilon_{\theta}(z) \mathbf{L}_{\gamma}(X) : \theta \in \Theta\}$ are P -GC classes (see Corollary 8.6 in [Giné and Zinn 1984](#), for instance) and, hence, $\sup_{\theta \in \Theta} |I| = o_p(1)$. Then, by Cauchy-Schwarz inequality and Assumptions 1 and 4,

$$\sup_{\theta \in \Theta} |II| \leq \left(\frac{1}{n} \sum_{i=1}^n \varepsilon_{\theta_0}^2(Z_i) \right)^{1/2} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{L}_{\hat{\gamma} \Delta \gamma}(X_i) \right)^{1/2} = o_p(1).$$

Finally,

$$\begin{aligned} \sup_{\theta \in \Theta} \|III\| &\leq \sup_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \|\mathbf{L}_{\hat{\gamma} \Delta \gamma}(X_i)\| \|\nabla m_{\bar{\theta}}(X_i)\| \|\theta - \theta_0\| \\ &\leq d_{\theta} D \frac{1}{n} \sum_{i=1}^n \|\mathbf{L}_{\hat{\gamma} \Delta \gamma}(X_i)\| R(X_i) = o_p(1), \end{aligned}$$

where the second inequality follows from Assumption 2 and the compactness of Θ . This result, together with the consistency of $\hat{\Sigma}_{\hat{\gamma}}(\tilde{\theta})$ implies that $\sup_{\theta \in \Theta} |Q_n(\theta) - Q_0(\theta)| = o_p(1)$ and, therefore, $\hat{\theta}_{\hat{\gamma}} = \theta_0 + o_p(1)$. For the asymptotic normality: by Assumptions 2 and

1(c), the first-order conditions of the minimization problem are satisfied with probability approaching one, $\hat{\boldsymbol{\mu}}_{\hat{\gamma}}^*(\hat{\theta}_{\gamma})\hat{\Sigma}_{\hat{\gamma}}(\tilde{\theta})^{-1}n^{-1/2}\hat{\boldsymbol{\Phi}}_{\hat{\gamma}}(\hat{\theta}_{\gamma}) = 0$. Expanding $\hat{\boldsymbol{\Phi}}_{\hat{\gamma}}(\hat{\theta}_{\gamma})$ around θ_0 and solving gives the Bahadur representation,

$$\sqrt{n}(\hat{\theta}_{\gamma} - \theta_0) = - \left[\hat{\boldsymbol{\mu}}_{\hat{\gamma}}^*(\hat{\theta}_{\gamma})\hat{\Sigma}_{\hat{\gamma}}(\tilde{\theta})^{-1}\hat{\boldsymbol{\mu}}_{\hat{\gamma}}^{*'}(\tilde{\theta}) \right]^{-1} \hat{\boldsymbol{\mu}}_{\hat{\gamma}}^*(\hat{\theta}_{\gamma})\hat{\Sigma}_{\hat{\gamma}}(\tilde{\theta})^{-1}\hat{\boldsymbol{\Phi}}_0(\hat{\gamma})$$

where $|\tilde{\theta} - \theta_0| \leq |\hat{\theta}_{\gamma} - \theta_0|$.

From the proof of Lemma 3, it follows that $\hat{\boldsymbol{\mu}}_{\hat{\gamma}}^*(\theta) - \hat{\boldsymbol{\mu}}_{\hat{\gamma}}^*(\theta_0) = o_p(1)$ for any $\theta \xrightarrow{p} \theta_0$. Thus by consistency of $\hat{\theta}_{\hat{\gamma}}$ and Lemma 1(b), it follows that both $\hat{\boldsymbol{\mu}}_{\hat{\gamma}}^*(\hat{\theta}_{\gamma})$ and $\hat{\boldsymbol{\mu}}_{\hat{\gamma}}^*(\tilde{\theta})$ converge in probability to $\boldsymbol{\mu}_{\gamma,0}^*$. Finally, by Assumptions 4-5, and the uniform continuity of the sample paths of $\Phi_0(\cdot)$,

$$\Phi_0(\hat{\gamma}) \xrightarrow{d} N(0, \Sigma_{\gamma,0}).$$

The result follows by Slutsky's lemma. ■

Proof of Theorem 2. By Lemma 3, 4, and Assumption 3 (or Theorem 1 and Assumption 2' for the $\hat{\chi}^2$ tests) both test statistics are asymptotically equivalent to the following quadratic form,

$$q(\hat{\theta}, W, \hat{\gamma}) = \left(\Phi_0(\hat{\gamma}) - \sqrt{n}\boldsymbol{\mu}_{\gamma,0}^{*'}\bar{l}_{\theta_0} \right)' W^{-1} \left(\Phi_0(\hat{\gamma}) - \sqrt{n}\boldsymbol{\mu}_{\gamma,0}^*\bar{l}_{\theta_0} \right), \quad (\text{A1})$$

where $\bar{l}_{\theta_0} = n^{-1} \sum_{i=1}^n l_{\theta_0}(Z_i)$ and W^{-1} is the probability limit of the weighting matrix W_n^{-1} . In particular, the couple $(\hat{\theta}, W_n^{-1})$ is equal to $(\hat{\theta}_{\hat{\gamma}}, \hat{\Sigma}_{\hat{\gamma}}(\tilde{\theta}))$ in the $\hat{\chi}^2$ test and $(\tilde{\theta}, \widehat{\text{Avar}}(\hat{\boldsymbol{\Phi}}_{\gamma}(\tilde{\theta})))$ in the Wald test. The functional

$$\phi(z, w, \gamma) = (z(\gamma) - \boldsymbol{\mu}_{\gamma,0}^{*'}w)' W^{-1} (z(\gamma) - \boldsymbol{\mu}_{\gamma,0}^*w),$$

mapping $(\hat{\boldsymbol{\Phi}}_0(\cdot), \sqrt{n}\bar{l}_{\theta_0}, \hat{\gamma})$ into $q(\hat{\theta}, W, \hat{\gamma})$ is continuous with respect to the product topology on $l^\infty(\mathbb{D}) \times \mathbb{R}^L \times \mathbb{D}$ (see Lemma 4 in [Andrews 1988a](#)). Thus, Theorem 2 follows by establishing the limit null distribution of $(\Phi_0(\hat{\gamma}) - \boldsymbol{\mu}_{\gamma,0}^{*'}\sqrt{n}\bar{l}_{\theta_0})$ and an application of the continuous mapping theorem (e.g., Theorem 1.3.6 in VW). Lemma 2, Assumptions 1, 4, 5, and the central limit theorem imply that $(\Phi_0(\cdot), \sqrt{n}\bar{l}_{\theta_0}, \hat{\gamma})$ is a uniformly tight process

on \mathbb{D} with fidis converging weakly to those of $(\Phi_0(\cdot), l_0, \gamma_0)$, where $l_0 \stackrel{d}{=} N(0, L_0)$ and $\mathbb{E}[l_0 \Phi_0(\gamma)] = \mathbb{E}[l_{\theta_0}(Z) \varepsilon_{\theta_0}(Z) \mathbf{I}_\gamma(X)']$. Thus,

$$(\hat{\Phi}_0(\cdot), \sqrt{n} \bar{l}, \hat{\gamma}) \rightsquigarrow (\Phi_0(\cdot), l_0, \gamma_0) \text{ on } l^\infty(\mathbb{D}),$$

and by the continuous mapping theorem,

$$q(\hat{\theta}, W, \hat{\gamma}) \xrightarrow{d} \tilde{Y}' W^{-1} \tilde{Y},$$

where $\tilde{Y} \stackrel{d}{=} N(0, \Sigma_{\tilde{Y}})$ and $\Sigma_{\tilde{Y}} = \text{Avar}(\hat{\Phi}_\gamma(\hat{\theta}))$. In the Wald test, where W^{-1} is a generalized inverse of $\Sigma_{\tilde{Y}}$, $\tilde{Y}' W^{-1} \tilde{Y} \stackrel{d}{=} \chi_{r(\text{Avar}(\hat{\Phi}_\gamma(\hat{\theta})))}^2$ by Theorem 7.3(i) in [Rao and Mitra \(1972\)](#). In the $\hat{\chi}^2$ test, $W = \Sigma_{\gamma,0}$, and $\Sigma_{\tilde{Y}} = \text{Avar}(\hat{\Phi}_\gamma(\hat{\theta}_\gamma)) = \Sigma_{\gamma,0} - \boldsymbol{\mu}_\gamma^{*'} (\boldsymbol{\mu}_\gamma^* \Sigma_{\gamma,0}^{-1} \boldsymbol{\mu}_\gamma^{*'}) \boldsymbol{\mu}_\gamma^*$. Since $\text{Var}(\Sigma_{\gamma,0}^{-1/2} \tilde{Y}) = \Sigma_{\gamma,0}^{-1/2} (\Sigma_{\tilde{Y}}) \Sigma_{\gamma,0}^{-1/2}$ is idempotent with rank equal to $L - d_\theta$, it follows that $\tilde{Y}' W^{-1} \tilde{Y} = (\Sigma_{\gamma,0}^{-1/2} \tilde{Y})' (\Sigma_{\gamma,0}^{-1/2} \tilde{Y}) \stackrel{d}{=} \chi_{L-d_\theta}^2$, proving the theorem. ■

Proof of Proposition 1. By an MVT argument, for fixed x , one can rewrite the difference between the two models fit as,

$$m_{\theta_0}(x) + \Delta m_{\bar{\theta}}(x)(\tilde{\theta} - \theta_0) = m_{1,\theta_1^*}(x) + \Delta m_{1,\bar{\theta}^*}(x)(\tilde{\theta}^* - \theta_1^*).$$

where $|\bar{\theta} - \theta_0| \leq |\tilde{\theta} - \theta_0|$ and $|\bar{\theta}^* - \theta_1^*| \leq |\tilde{\theta}^* - \theta_1^*|$. Then using conditions (a) and (b) of Proposition 2, we obtain, that the splitting points, x_0 , solve,

$$\Delta m_{\bar{\theta}}(x_0) C n^{-1/2} \sum_{i=1}^n g(X_i, \varepsilon_{\theta_0}) - \Delta m_{1,\bar{\theta}^*}(x_0) D n^{-1/2} \sum_{i=1}^n g(X_i, \varepsilon_{\theta_0}) + o_p(1) = 0,$$

and, thus,

$$\Delta m_{\bar{\theta}}(x_0) C - \Delta m_{1,\bar{\theta}^*}(x_0) D = o_p(1).$$

■

Proof of Proposition 2. Let

$$d_\gamma(I_L) = \delta(\gamma)' \delta(\gamma) = \sum_{l=1}^L \delta_l^2 = \sum_{l=1}^L \mathbb{E}[h(X) \mathbb{I}_\gamma(X)]^2.$$

If $h(x) > 0$ ($h(x) \leq 0$) for all $x \in \mathcal{X}$, then $\delta_l \geq 0$ ($\delta_l \leq 0$) for all $l = 1, 2, \dots, L$, and, since $(\delta_l + \delta_f)^2 \geq (\delta_l^2 + \delta_f^2)$ for each δ_l, δ_f such that $\text{sign}(\delta_l) = \text{sign}(\delta_f)$, it follows that,

$$\sum_{l=1}^L \delta_l^2 \leq \left(\sum_{l=1}^L \delta_l \right)^2.$$

Thus, a unique cells maximize the drifts norm.

Consider, instead, the case $h(x) = 0$ for some $x \in \mathcal{X}$, then for any finite split of \mathcal{X} with $L > 1$, we have that

$$\delta_l^2 + \delta_f^2 = (\delta_l^+ + \delta_l^-)^2 + (\delta_f^+ + \delta_f^-)^2 \leq (\delta_l^+ + \delta_f^+)^2 + (\delta_l^- + \delta_f^-)^2 = \delta_l^{*2} + \delta_f^{*2},$$

where δ_l^+ and δ_l^- denote the positive and negative part of δ_l , respectively. Therefore, any partition is dominated by a partition with the same number of cells obtained by merging the positive and negative part of each δ_l into a new cell. Since the δ_l^* are either non-negative or non-positive, by merging them in two cells containing only the positive and negative δ_l^* we obtain a partition with two cells that dominates any partition with $L > 1$ cells. Of course, since

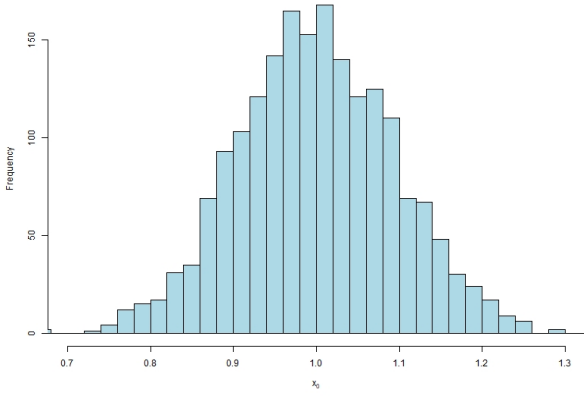
$$(\delta_l^+ + \delta_l^-)^2 \leq (\delta_l^+)^2 + (\delta_l^-)^2,$$

the two-cells partition dominates the one-cell partition as well. Thus, the two-cells partition into positive and negative values maximizes the drifts norm. ■

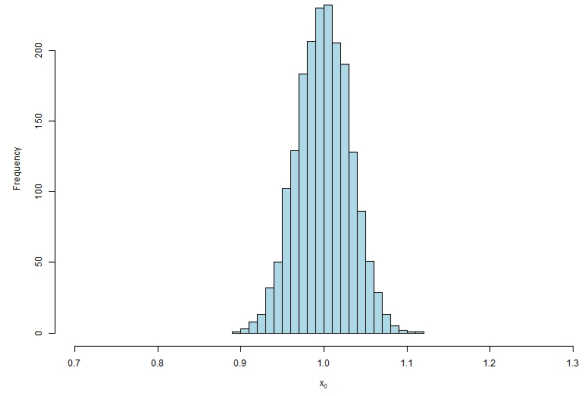
B Appendix B

I present Monte Carlo evidence demonstrating NP partitions converging to fixed cells in \mathbb{D} . In Figure 1, we generate data as $Y = 1 + \varepsilon$, where $\varepsilon \sim N(0, 1)$. Under the null hypothesis, $H_0 : m(X) = c$ a.s., with c as a constant. The alternative model, $H_1 : m(X) = c + \beta X$ a.s., includes a constant c and slope β , where X is generated as $X \sim N(1, 1)$. In Figure 2, we create a linear model, $Y = 1 + X + \varepsilon$, with $\varepsilon \sim N(0, 1)$. The null hypothesis is $H_0 : m(X) = c + \beta X$ a.s., with c and β as constants. The alternative model, $H_1 : m(X) = c + \beta X + \gamma X^2$ a.s., introduces a curvature term γ . I consider sample sizes of $n = 100$, $n = 1000$, $n = 10000$, and $n = 10000$, to compute the splitting points

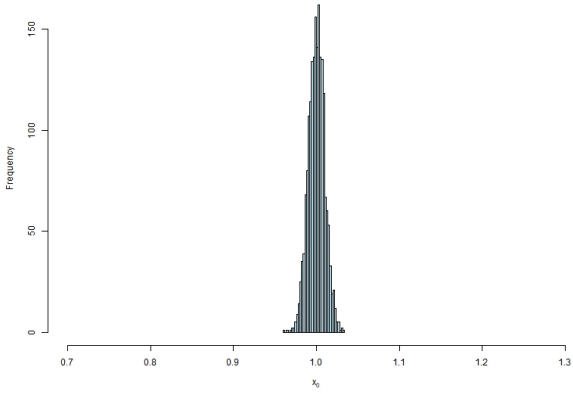
Figure 5: Convergence of x_0 for the constant model under linear alternative.



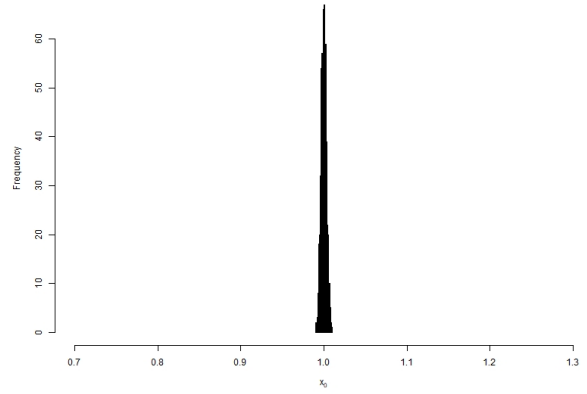
(a) $n = 100$



(b) $n = 1000$



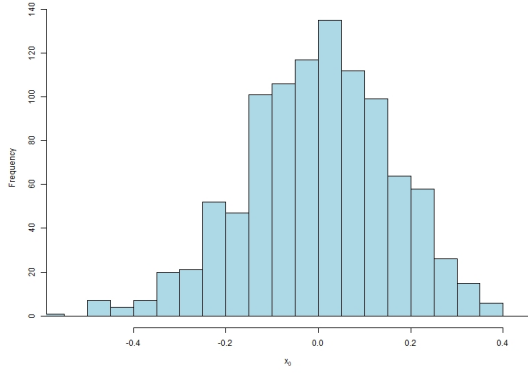
(c) $n = 10,000$



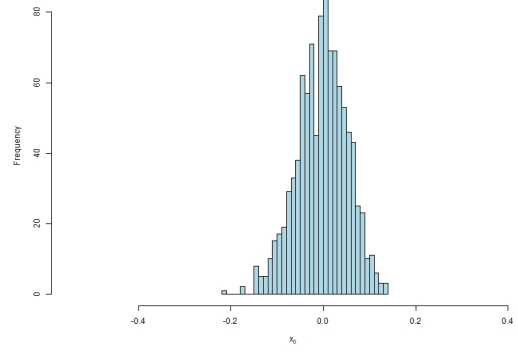
(d) $n = 100,000$

using the NP partition algorithm over $R = 1000$ replications. The results, depicted in the figures below, align with theoretical findings, illustrating the convergence of splitting points to fixed points in \mathcal{X} .

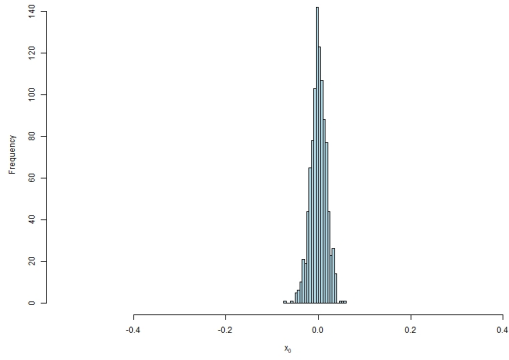
Figure 6: Convergence of x_0 for the linear model under quadratic alternative. The first four graphs and the last four refer to the two roots of the quadratic equation.



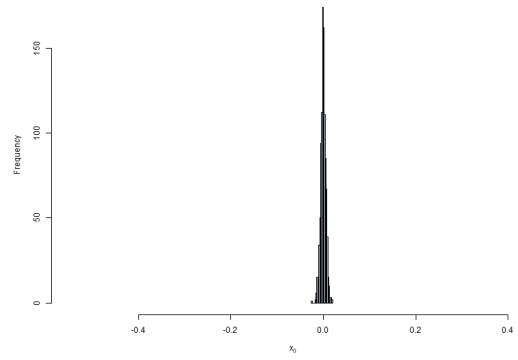
(a) $n = 100$



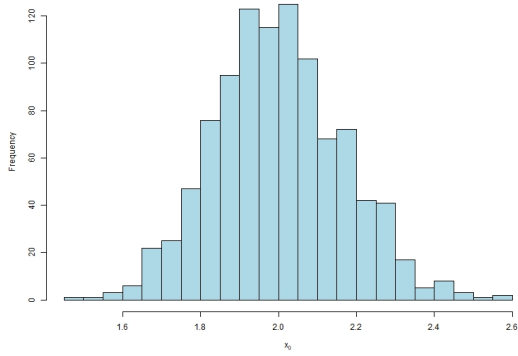
(b) $n = 1000$



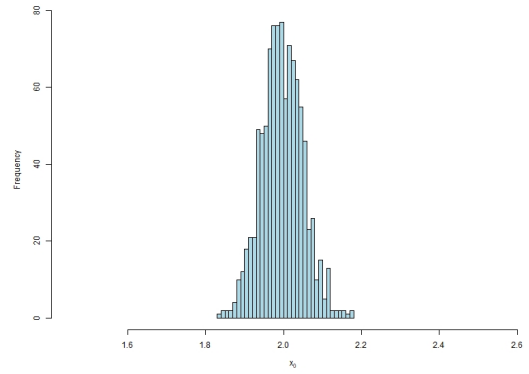
(c) $n = 10,000$



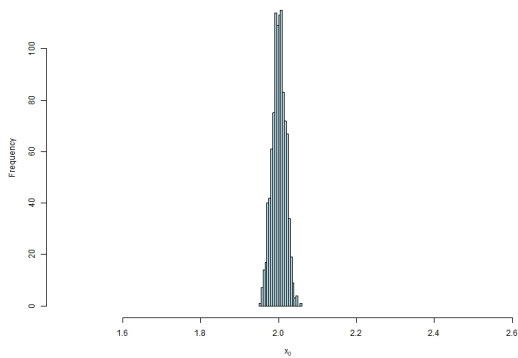
(d) $n = 100,000$



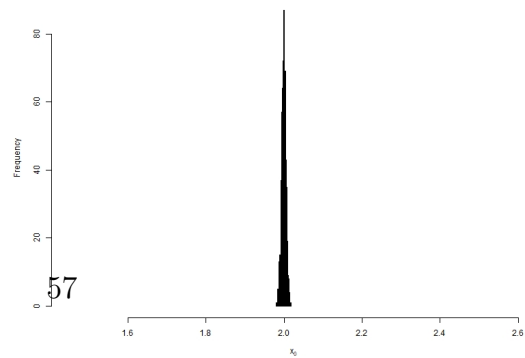
(e) $n = 100$



(f) $n = 1000$



(g) $n = 10,000$



(h) $n = 100,000$

References

- ANDREWS, D. W. (1987): “Asymptotic results for generalized Wald tests,” *Econometric Theory*, 3(3), 348–358.
- (1988a): “Chi-square diagnostic tests for econometric models: theory,” *Econometrica: Journal of the Econometric Society*, pp. 1419–1453.
- (1988b): “Chi-square diagnostic tests for econometric models: Introduction and applications,” *Journal of Econometrics*, 37(1), 135–156.
- (1997): “A conditional Kolmogorov test,” *Econometrica: Journal of the Econometric Society*, pp. 1097–1128.
- BALAKRISHNAN, N., V. VOINOV, AND M. S. NIKULIN (2013): *Chi-squared goodness of fit tests with applications*. Academic Press.
- BICKEL, P. J., AND M. ROSENBLATT (1973): “On some global measures of the deviations of density function estimates,” *The Annals of Statistics*, pp. 1071–1095.
- BIERENS, H. J. (1982): “Consistent model specification tests,” *Journal of Econometrics*, 20(1), 105–134.
- BIERENS, H. J., AND W. PLOBERGER (1997): “Asymptotic theory of integrated conditional moment tests,” *Econometrica: Journal of the Econometric Society*, pp. 1129–1151.
- BILLINGSLEY, P. (2017): *Probability and measure*. John Wiley & Sons.
- CHAUSSÉ, P. (2010): “Computing generalized method of moments and generalized empirical likelihood with R,” *Journal of Statistical Software*, 34, 1–35.
- COOK, R. D. (1994): “On the interpretation of regression plots,” *Journal of the American Statistical Association*, 89(425), 177–189.
- COX, D. R. (1972): “Regression models and life-tables,” *Journal of the Royal Statistical Society: Series B (Methodological)*, 34(2), 187–202.

- CRAMÉR, H. (1946): *Mathematical methods of statistics*, vol. 26. Princeton university press.
- CRISTOBAL, J. C., P. F. ROCA, AND W. G. MANTEIGA (1987): “A class of linear regression parameter estimators constructed by nonparametric estimation,” *The Annals of Statistics*, 15(2), 603–609.
- DAVIDSON, J. (1994): *Stochastic limit theory: An introduction for econometricians*. OUP Oxford.
- DAVIDSON, R., J. G. MACKINNON, ET AL. (2004): *Econometric theory and methods*, vol. 5. Oxford University Press New York.
- DELGADO, M. A., AND J. VAINORA (2023): “Conditional Distribution Model Specification Testing Using Chi-Square Goodness-of-Fit Tests,” .
- DOMÍNGUEZ, M. A., AND I. N. LOBATO (2004): “Consistent estimation of models defined by conditional moment restrictions,” *Econometrica*, 72(5), 1601–1615.
- DURBIN, J., AND M. KNOTT (1972): “Components of Cramér–von Mises statistics. I,” *Journal of the Royal Statistical Society: Series B (Methodological)*, 34(2), 290–307.
- FRYER, R. G., AND M. GREENSTONE (2010): “The changing consequences of attending historically black colleges and universities,” *American economic journal: Applied economics*, 2(1), 116–148.
- FRYER JR, R. G., AND P. TORELLI (2010): “An empirical analysis of ‘acting white’,” *Journal of Public Economics*, 94(5-6), 380–396.
- GESSAMAN, M. (1970): “A consistent nonparametric multivariate density estimator based on statistically equivalent blocks,” *The Annals of Mathematical Statistics*, 41(4), 1344–1346.
- GINÉ, E., AND J. ZINN (1984): “Some limit theorems for empirical processes,” *The Annals of Probability*, pp. 929–989.

- GREENWOOD, P. E., AND M. S. NIKULIN (1996): *A guide to chi-squared testing*, vol. 280. John Wiley & Sons.
- HAN, Y., P. MA, H. REN, AND Z. WANG (2023): “Model checking in large-scale data set via structure-adaptive-sampling,” *Statistica Sinica*, 33, 1–27.
- HARDLE, W., AND E. MAMMEN (1993): “Comparing nonparametric versus parametric regression fits,” *The Annals of Statistics*, pp. 1926–1947.
- HASTIE, T., R. TIBSHIRANI, AND J. H. FRIEDMAN (2009): *The elements of statistical learning: data mining, inference, and prediction*, vol. 2. Springer.
- HECKMAN, J. J. (1979): “Sample selection bias as a specification error,” *Econometrica: Journal of the econometric society*, pp. 153–161.
- KALLENBERG, W. C. M., J. OOSTERHOFF, AND B. SCHRIEVER (1985): “The number of classes in chi-squared goodness-of-fit tests,” *Journal of the American Statistical Association*, 80(392), 959–968.
- KHMALADZE, E. V. (1982): “Martingale approach in the theory of goodness-of-fit tests,” *Theory of Probability & Its Applications*, 26(2), 240–257.
- KHMALADZE, E. V., AND H. L. KOUL (2004): “Martingale transforms goodness-of-fit tests in regression models,” .
- KIEFER, N. M., T. J. VOGELSANG, AND H. BUNZEL (2000): “Simple robust testing of regression hypotheses,” *Econometrica*, 68(3), 695–714.
- KOUL, H. L., AND P. NI (2004): “Minimum distance regression model checking,” *Journal of Statistical Planning and Inference*, 119(1), 109–141.
- KOUTROUVELIS, I. A., AND J. KELLERMEIER (1981): “A goodness-of-fit test based on the empirical characteristic function when parameters must be estimated,” *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 43(2), 173–176.
- KUAN, C.-M., AND W.-M. LEE (2006): “Robust M tests without consistent estimation of the asymptotic covariance matrix,” *Journal of the American Statistical Association*, 101(475), 1264–1275.

- LEE, W.-M., C.-M. KUAN, AND Y.-C. HSU (2014): “Testing over-identifying restrictions without consistent estimation of the asymptotic covariance matrix,” *Journal of Econometrics*, 181(2), 181–193.
- LÜTKEPOHL, H., AND M. M. BURDA (1997): “Modified Wald tests under nonregular conditions,” *Journal of Econometrics*, 78(2), 315–332.
- MOORE, D. S. (1977): “Generalized inverses, Wald’s method, and the construction of chi-squared tests of fit,” *Journal of the American Statistical Association*, 72(357), 131–137.
- MOORE, D. S., AND M. C. SPRUILL (1975): “Unified large-sample theory of general chi-squared statistics for tests of fit,” *The Annals of Statistics*, pp. 599–616.
- NEWKEY, W. K. (1985): “Generalized method of moments specification testing,” *Journal of econometrics*, 29(3), 229–256.
- (1990): “Efficient instrumental variables estimation of nonlinear models,” *Econometrica: Journal of the Econometric Society*, pp. 809–837.
- NEWKEY, W. K., AND D. MCFADDEN (1994): “Large sample estimation and hypothesis testing,” *Handbook of econometrics*, 4, 2111–2245.
- NIKULIN, M. S. (1973): “Chi-square test for continuous distributions with location and scale parameters,” *Teoriya Veroyatnostei i ee Primeneniya*, 18(3), 583–591.
- PEARSON, K. (1900): “X. On the criterion that a given system of deviations from the probable in the case of a correlated system of variables is such that it can be reasonably supposed to have arisen from random sampling,” *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, 50(302), 157–175.
- POLACHEK, S. W., ET AL. (2008): “Earnings over the life cycle: The mincer earnings function and its applications,” *Foundations and Trends® in Microeconomics*, 4(3), 165–272.

- POLLARD, D. (1979): “General chi-square goodness-of-fit tests with data-dependent cells,” *Zeitschrift für Wahrscheinlichkeitstheorie und verwandte Gebiete*, 50(3), 317–331.
- (1984): *Convergence of stochastic processes*. Springer Science & Business Media.
- PRICE, G. N., AND A. C. VICEISZA (2023): “What Can Historically Black Colleges and Universities Teach about Improving Higher Education Outcomes for Black Students?,” *Journal of Economic Perspectives*, 37(3), 213–232.
- RAO, C. R., AND S. K. MITRA (1972): “Generalized inverse of a matrix and its applications,” in *Proceedings of the Sixth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Theory of Statistics*, vol. 6, pp. 601–621. University of California Press.
- RAO, K. C., AND B. ROBSON (1974): “A chi-square statistic for goodness-of-fit tests within the exponential family,” *Communications in Statistics-Theory and Methods*, 3(12), 1139–1153.
- ROSENBAUM, P. R. (1987): “Model-based direct adjustment,” *Journal of the American statistical Association*, 82(398), 387–394.
- ROSENBAUM, P. R., AND D. B. RUBIN (1983): “Assessing sensitivity to an unobserved binary covariate in an observational study with binary outcome,” *Journal of the Royal Statistical Society: Series B (Methodological)*, 45(2), 212–218.
- SANT’ANNA, P. H., AND X. SONG (2019): “Specification tests for the propensity score,” *Journal of Econometrics*, 210(2), 379–404.
- SCHOTT, J. R. (2016): *Matrix analysis for statistics*. John Wiley & Sons.
- SEETHARAMAN, P., AND P. K. CHINTAGUNTA (2003): “The proportional hazard model for purchase timing: A comparison of alternative specifications,” *Journal of Business & Economic Statistics*, 21(3), 368–382.
- STUTE, W. (1997): “Nonparametric model checks for regression,” *The Annals of Statistics*, pp. 613–641.

- STUTE, W., W. G. MANTEIGA, AND M. P. QUINDIMIL (1998): “Bootstrap approximations in model checks for regression,” *Journal of the American Statistical Association*, 93(441), 141–149.
- STUTE, W., AND L.-X. ZHU (2002): “Model checks for generalized linear models,” *Scandinavian Journal of Statistics*, 29(3), 535–545.
- TAUCHEN, G. (1985): “Diagnostic testing and evaluation of maximum likelihood models,” *Journal of Econometrics*, 30(1-2), 415–443.
- VAN DER VAART, A. W. (2000): *Asymptotic statistics*, vol. 3. Cambridge university press.
- VAN DER VAART, J. W. (1996): *Weak convergence*. Springer.
- WATSON, G. S. (1959): “Some recent results in chi-square goodness-of-fit tests,” *Biometrics*, pp. 440–468.
- WOOLDRIDGE, J. M. (1990): “A unified approach to robust, regression-based specification tests,” *Econometric Theory*, 6(1), 17–43.