

# Consistent Test for Conditional Moment Restriction Models in Reproducing Kernel Hilbert Spaces

**Yuhao Li**

*Economics and Management School  
Wuhan University*

LIYUHAO.ECON@WHU.EDU.CN

**Xiaojun Song**

*Guanghua School of Management  
Peking University*

SXJ@GSM.PKU.EDU.CN

**Draft Time:** November, 2022. [Click Here for Latest Version](#)

## Abstract

In this paper, we represent *Integrated Conditional Moment* (ICM) tests in *Reproducing Kernel Hilbert Spaces* (RKHS). There are several advantages to doing so. First, reproducing kernels embody dimension and integral measure, and hence, are effective dimension reduction tools. This phenomenon can be explained by the isometrically isomorphic relationship among infinite dimensional Hilbert spaces. Second, the test statistics, expressed in terms of kernels, have analytic closed forms, making them easy to compute in practice. Third, one can generate kernels easily and massively from existing kernels. Each kernel corresponds to an ICM test, thus, for certain models, one may obtain a more sensitive test than by using conventional ones. We further propose projection-based kernels to eliminate the estimation effect, leading to a simple multiplier bootstrap procedure to obtain critical values. A minimum distance estimator is developed as a byproduct. Monte Carlo exercises are performed to examine the finite-sample performance of the proposed test, and an empirical application is also provided.

**JEL Classification:** C12; C13; C15; C52

**Keywords:** Conditional Moment Restriction, Specification Tests, Reproducing Kernel Hilbert Spaces, Multiplier Bootstrap

## 1. Introduction

In this paper, we study the question of testing conditional moment restriction (CMR) models. Specifically, we develop a new framework for deriving integrated conditional moment (ICM) test statistics. This framework is based on the idea of embedding CMR in a reproducing kernel Hilbert space (RKHS). The test statistic is defined as the maximum

moment restriction (MMR) within the unit ball of the RKHS. Furthermore, we show that MMR corresponds to the RKHS norm of a Hilbert space embedding of conditional moments.

We contribute to the literature in the following aspects.

**Closed Form Expression.** ICM statistics, as its name suggests, is obtained after integrating the nuisance parameter of infinite many unconditional moment restrictions (UMR). This integration often leads to numerical challenges, and only a few weighting functions and integral measures are known in the literature to generate closed-form statistics, e.g., Bierens (1982); Escanciano (2006a). In our framework, the MMR is obtained directly from a user-chosen reproducing kernel without integration. Furthermore, we show that the MMR captures all information about the original CMR and it has a closed-form expression which eases practical implementation.

**Dimension Reduction.** The dimension of conditional variables often poses practical and theoretical challenges when conducting CMR specification tests. Most of the existing ICM tests depend on high-dimensional stochastic processes, e.g., Domínguez and Lobato (2015), and their power performance often drops significantly as the dimension  $d$  increases due to the data sparseness. One common solution is to project the original conditional covariates  $X$  onto the  $\beta^\top X$  for all  $\|\beta\|_2 = 1$ , see, e.g., Escanciano (2006a); Lavergne and Patilea (2012); Sant’Anna and Song (2019). Here,  $\|\cdot\|_2$  denotes the Euclidean norm. However, due to the involvement of infinite directions, projection-based tests are often computationally intensive and the powers of these tests are often low (Guo and Zhu, 2017). Reproducing kernels, on the other hand, embody both the dimension and the integral measure. As a result, the estimator of the MMR converges in the RKHS norm in a way that is independent of the dimension. This is an appealing property since tests based on this estimator are less sensitive to the curse of dimensionality.

**Massively Generate New Tests.** Existing literature has shown that the power of an ICM test is determined by the weighting function, the integral measure, the data-generating process (DGP) and the model itself. See, e.g., Escanciano (2009). Thus, an ICM test statistic might be powerful against one model and one DGP but could be powerless against another model or another DGP. Hence, it is desirable to have as many ICM test statistics as possible. We provide methods to construct new kernels from existing kernels. Since each kernel corresponds to an ICM test statistic, this means that one could generate infinitely many new ICM test statistics, all have closed-form expressions.

**Eliminate Estimation Effect.** We propose a projected kernel to cancel the estimation effect. The limiting null distribution, therefore, does not depend on how an estimator is obtained and does not require the estimator to be  $\sqrt{n}$ -asymptotically linear, with  $n$  the sample size. Without the estimation effect, critical values are obtained via a simple and fast multiplier bootstrap procedure, and perhaps most interestingly, the proposed test is capable of applying to certain “non-standard” estimators who have slower convergence rates.

**New Limit Distribution.** We derive the limit distribution of the proposed test statistic under the fixed alternative. This new result provides a framework to obtain a more powerful test by selecting an optimal kernel. Nevertheless, much work is needed to fully achieve this goal.

**A Minimum Distance Estimator.** We propose a minimum distance estimator based on the MMR. Compared to existing minimum distance estimators, e.g., [Domínguez and Lobato \(2004\)](#), it has the advantage of being less sensitive to the curse of dimensionality.

The rest of the paper is organized as follows. Section 2 presents our main idea of using the RKHS framework to develop the test statistic. We discuss some merits of doing so, as well as the challenges when unknown parameters are replaced by their estimators. Section 3 describes a method for projecting a kernel onto a tangent space of nuisance parameters so that the modified statistic is free from the estimation effect. We also establish the asymptotic properties of our test in this section. In section 4, we introduce a simple multiplier bootstrap procedure to obtain critical values and justify its asymptotic validity. Based on the test statistic, we propose a minimum distance estimator, its asymptotic properties are studied in Section 5. Section 6 conducts Monte Carlo experiments. Simulation results indicate that the proposed tests have an accurate empirical size and a good local power, even when the sample size is as small as  $n = 100$  and the dimension is as high as  $d = 20$ . Simulation results also suggest that the proposed tests have good power against high-frequency alternatives. One empirical application is studied in Section 7. Section 8 concludes. Some backgrounds on the RKHS are presented in Appendix A. A concise introduction of RKHS can be found in [Carrasco et al. \(2007\)](#), while for more comprehensive surveys on this subject, see, e.g., [Hofmann et al. \(2008\)](#); [Paulsen and Raghupathi \(2016\)](#).

## 2. Main Idea, Benefits and Challenges

### 2.1 Expressing the Conditional Moment Restrictions in RKHS

Let  $Z = (Y, X^\top)^\top$  be a random vector taking values in  $\mathcal{Z} \subseteq \mathbb{R}^{1+d}$  with distribution  $P_Z$ ,  $X$  a random vector taking values in  $\mathcal{X} \subseteq \mathbb{R}^d$  with distribution  $P_X$ , and  $\Theta \in \mathbb{R}^r$  a parameter space. Typically,  $Y$  represents the real-valued dependent (or response) variable, and  $X$  is the explanatory variable. Under  $\mathbb{E}|Y| < \infty$ , it is well-known that the regression function  $\mathbb{E}(Y|X)$  is well-defined and is the ‘best’ prediction of  $Y$  given  $X$ , in a mean squared sense. In empirical studies, it is common to consider the following expression:

$$\begin{aligned} Y &= \mathbb{E}_{\theta_0}(Y|X) + \varepsilon \\ &= \mathcal{M}(X; \theta_0) + \varepsilon \end{aligned}$$

We are interested in testing the moment restriction models where the only information about the unknown parameter  $\theta_0 \in \Theta$  is a set of conditional moment restrictions:

$$\mathcal{E}(X; \theta_0) = \mathbb{E}(\varepsilon(Z; \theta_0)|X) = \mathbb{E}(Y - \mathcal{M}(X; \theta_0) | X) = 0 \quad P_X\text{-a.s.}, \quad (1)$$

here,  $\varepsilon : \mathcal{Z} \times \Theta \rightarrow \mathbb{R}$  is the generalized residual function whose functional form is known up to the parameter  $\theta \in \Theta$ .

Given an i.i.d sample  $\{x_i, y_i; i = 1, \dots, n\}$  drawn from a distribution  $P_Z$ , our goal is to conduct specification testing:

$$\begin{aligned} H_0 : \mathcal{E}(X; \theta_0) &= 0 & P_X\text{-a.s.} \\ H_1 : \mathcal{E}(X; \theta_0) &\neq 0 & P_X\text{-a.s.} \quad \forall \theta \in \Theta \end{aligned} \tag{2}$$

where  $\theta_0$  has a consistent estimator  $\hat{\theta}$ . To do so, we follow the integrated conditional moment (ICM) approach, which converts the constraint on the conditional expectation to infinite and parametric unconditional orthogonality restrictions. Let  $\mathcal{H}$  be a set of measurable functions on  $\mathcal{X}$ , then

$$\mathcal{E}(X; \theta_0) = 0 \Leftrightarrow \mathbb{E}(\varepsilon(Z; \theta_0)h(X, t)) = 0, P_X\text{-a.s.} \quad \forall t \in \mathcal{T}, \quad h \in \mathcal{H} \tag{3}$$

where  $\mathcal{T}$  is some proper space. For sufficient conditions on the family  $\mathcal{H}$  to satisfy (3), see Bierens and Ploberger (1997); Escanciano (2006b). In the context of this work,  $\mathcal{H}$  must consist of infinitely many instruments for the conditional moment test to be consistent against all alternatives.

Equivalently, any  $\theta_0 \in \Theta$  that satisfies (3) must also satisfy the *maximum moment restriction* (MMR) (Muandet et al., 2020):

$$\sup_{h \in \mathcal{H}} \|\mathbb{E}(\varepsilon(Z; \theta_0)h(X, t))\|_2^2 = 0 \tag{4}$$

However, the sup operator makes it hard to optimize (4). We resolve this issue by restricting  $\mathcal{H}$  to be a unit ball in an RKHS. To express (4) using the RKHS, let  $h : \mathcal{X} \rightarrow \mathbb{R}$ , and  $\mathcal{H}(k)$  be the RKHS of functions on  $\mathcal{X}$  with reproducing kernel  $k$ . The subsequent analyses rely on the following assumptions:

- (A1) The random vector  $(X, Z)$  is a strictly stationary process with probability measure  $P_{XZ}$ .
- (A2) Some regularity conditions. (i) the function  $\varepsilon : \mathcal{Z} \times \Theta \rightarrow \mathbb{R}$  is continuous on  $\Theta$  for each  $z \in \mathcal{Z}$ ; (ii)  $\mathcal{E}(x; \theta)$  exists and is finite for every  $\theta \in \Theta$  and  $x \in \mathcal{X}$  for which  $P_X(x) > 0$ ; (iii)  $\mathcal{E}(x; \theta)$  is continuous on  $\Theta$  for all  $x \in \mathcal{X}$  for which  $P_X(x) > 0$ .
- (A3) There is a unique  $\theta_0 \in \Theta^\circ$  for which  $\mathcal{E}(X; \theta_0) = 0, a.s.$ , and  $P(\mathcal{E}(X; \Theta) = 0) < 1$  for all  $\theta \neq \theta_0$ , where  $\Theta^\circ$  is the interior of  $\Theta$ .
- (A4) The kernel  $k(\cdot, \cdot)$  is *integrally strictly positive definite* (ISPD), continuous and bounded, i.e.,  $\sup_{x \in \mathcal{X}} \sqrt{k(x, x)} < \infty$ .

Assumptions A1 and A2 are regular conditions that appeared in most literature. Assumption A3 is a global identification assumption, and Assumption A4 put restrictions on the kernel  $k$  and is essential for the identification of the model. An ISPD kernel satisfies

$$\int \int_{\mathcal{X}} f(x)k(x, x')f(x')dx dx' > 0, \quad \forall \|f\|_2 \neq 0$$

where  $\|\cdot\|_2$  denotes the  $L_2$  norm.

Define an operator  $\mathcal{C}_\theta : \mathcal{H}(k) \rightarrow \mathbb{R}$  that takes an instrument  $h \in \mathcal{H}(k)$  as input and returns the corresponding moment restrictions:

$$\mathcal{C}_\theta h = \mathbb{E}_{XZ} (\varepsilon(Z; \theta)h(X))$$

By the reproducing property of the RKHS, we have

$$h(x) = \langle h, \phi_x(\cdot) \rangle_{\mathcal{H}(k)}$$

where  $\phi_x(\cdot) = k(x, \cdot)$  is the feature map of this RKHS with  $k(x, x') = \langle \phi_x(\cdot), \phi_{x'}(\cdot) \rangle_{\mathcal{H}(k)}$ .

By Riesz's representation theorem, one can show that

$$\mathcal{C}_\theta h = \langle h, \boldsymbol{\mu}_\theta \rangle_{\mathcal{H}(k)} \tag{5}$$

where

$$\boldsymbol{\mu}_\theta = \mathbb{E}_X (\mathcal{E}(x; \theta)\phi_x(\cdot)) = \mathbb{E}_{XZ} (\varepsilon(Z; \theta)\phi_X(\cdot)) \tag{6}$$

is called the *Conditional Moment Embedding* (CME) (Muandet et al., 2017, 2020). Its verification can be found in Appendix C. The idea of CME is to extend the feature map  $\phi$  to the space of probability distribution  $P_X$  and the space of conditional moment restrictions  $\mathcal{E}(X; \theta)$  by representing both elements as a mean function. Through Equation (6), most RKHS methods can therefore be extended to conditional moment restrictions.

**Remark.** Since  $\phi_x(\cdot)$  takes values in the RKHS, the integral  $\int \varepsilon(z; \theta)\phi_x(\cdot)dP_{XZ}(x, z)$  should be interpreted as a Bochner integral (see Dinculeanu (2000) for the definition of the Bochner integral).

This representation is useful. Given ISPD kernels, the CME (or equivalently, the MMR) captures all information about the conditional moment restrictions. In other words,  $\boldsymbol{\mu}_\theta$  is injective, implying that for any  $\theta_1, \theta_2 \in \Theta$ ,  $\mathcal{E}(x; \theta_1) = \mathcal{E}(x; \theta_2)$  for  $P_X$ -almost surely if and only if  $\boldsymbol{\mu}_{\theta_1} = \boldsymbol{\mu}_{\theta_2}$ . An important consequence is  $\|\boldsymbol{\mu}_\theta\|_{\mathcal{H}(k)}^2 \geq 0$  and  $\|\boldsymbol{\mu}_\theta\|_{\mathcal{H}(k)}^2 = 0$  if and only if  $\theta = \theta_0$ . See, e.g., Muandet et al. (2020) for a detailed discussion.

To summarize so far, the MMR condition in (4) then can be written as

$$\sup_{\|h\|_{\mathcal{H}(k)} \leq 1} \|\mathbb{E} (\varepsilon(Z; \theta_0)h(X, t))\|_2^2 = \|\mathcal{C}_\theta\|^2 = \|\boldsymbol{\mu}_{\theta_0}\|_{\mathcal{H}(k)}^2$$

and the original null hypothesis is equivalent to

$$H_0 : \|\boldsymbol{\mu}_{\theta_0}\|_{\mathcal{H}(k)}^2 = 0, \quad P_x\text{-a.s}$$

To simplify notation, let  $\mathbb{M}^2(\theta_0) = \|\boldsymbol{\mu}_{\theta_0}\|_{\mathcal{H}(k)}^2$ , and further notice that

$$\begin{aligned}\mathbb{M}^2(\theta_0) &= \langle \mathbb{E}_{XZ}(\varepsilon(Z; \theta) \phi_X(\cdot)), \mathbb{E}_{XZ}(\varepsilon(Z; \theta) \phi_X(\cdot)) \rangle_{\mathcal{H}(k)} \\ &= \mathbb{E}_{XZ} \left( \langle \varepsilon(Z; \theta) \phi_X(\cdot), \varepsilon(Z; \theta) \phi_X(\cdot) \rangle_{\mathcal{H}(k)} \right) \\ &= \mathbb{E} \left( \langle \varepsilon(Z; \theta) \phi_X(\cdot), \varepsilon(Z'; \theta) \phi_{X'}(\cdot) \rangle_{\mathcal{H}(k)} \right) \\ &= \mathbb{E} (\varepsilon(Z; \theta_0) k(X, X') \varepsilon(Z'; \theta_0))\end{aligned}$$

Then, given a consistent estimator  $\hat{\theta}$ , we propose a simple test statistic  $n\widehat{\mathbb{M}}_n^2(\hat{\theta})$  as

$$n\widehat{\mathbb{M}}_n^2(\hat{\theta}) = \frac{n}{n(n-1)} \sum_{1 \leq i \neq j \leq n} \varepsilon(z_i; \hat{\theta}) k(x_i, x_j) \varepsilon(z_j; \hat{\theta})$$

The asymptotic distributions of  $n\widehat{\mathbb{M}}_n^2(\hat{\theta})$  under different hypotheses are complicated due to the presence of the estimator  $\hat{\theta}$ . In Section 3, we propose a projection-based test that has the ability to eliminate this estimation effect. Hence, we will not study the asymptotic distributions of  $n\widehat{\mathbb{M}}_n^2(\hat{\theta})$  in detail.

## 2.2 Benefits of using RKHS Techniques: Insensitive to Dimension and Massively Generate ICM Tests

### 2.2.1 DIMENSION REDUCTION

Let  $\widehat{\boldsymbol{\mu}}_{\hat{\theta}} = 1/n \sum_{i=1}^n \varepsilon(z_i; \hat{\theta}) \phi_{x_i}(\cdot) \in \mathcal{H}(k)$  be an estimator of  $\boldsymbol{\mu}_{\theta_0}$ , and suppose  $\hat{\theta} - \theta_0 = O_p(1/\sqrt{n})$ . Let  $g(z; \theta) = \nabla_{\theta} \varepsilon(z; \theta)$  be the first derivative of  $\varepsilon(z; \theta)$ , and  $\bar{\theta} = \gamma\theta_0 + (1-\gamma)\hat{\theta}$ ,  $\gamma \in (0, 1)$ . Notice that

$$\widehat{\boldsymbol{\mu}}_{\hat{\theta}} = \frac{1}{n} \sum_{i=1}^n \varepsilon(z_i; \theta_0) \phi_{x_i}(\cdot) + O_p(1/\sqrt{n})^{\top} \frac{1}{n} \sum_{i=1}^n g(z_i; \bar{\theta}) \phi_{x_i}(\cdot)$$

and

$$\begin{aligned}\|\widehat{\boldsymbol{\mu}}_{\hat{\theta}} - \boldsymbol{\mu}_{\theta_0}\|_{\mathcal{H}(k)} &= \|\widehat{\boldsymbol{\mu}}_{\theta_0} + O_p(1/\sqrt{n})^{\top} \frac{1}{n} \sum_{i=1}^n g(z_i; \bar{\theta}) \phi_{x_i}(\cdot) - \boldsymbol{\mu}_{\theta_0}\|_{\mathcal{H}(k)} \\ &\leq \|\widehat{\boldsymbol{\mu}}_{\theta_0} - \boldsymbol{\mu}_{\theta_0}\|_{\mathcal{H}(k)} + O_p(1/\sqrt{n}) \sqrt{\frac{1}{n^2} \sum_{i,j=1}^n g^{\top}(z_i; \bar{\theta}) k(x_i, x_j) g(z_j; \bar{\theta})} \\ &= \|\widehat{\boldsymbol{\mu}}_{\theta_0} - \boldsymbol{\mu}_{\theta_0}\|_{\mathcal{H}(k)} + O_p(1/\sqrt{n})\end{aligned}$$

Furthermore, [Muandet et al. \(2017\)](#); [Tolstikhin et al. \(2017\)](#) show that

$$\|\widehat{\boldsymbol{\mu}}_{\theta_0} - \boldsymbol{\mu}_{\theta_0}\|_{\mathcal{H}(k)} = O_p(1/\sqrt{n})$$

This observation states that  $\widehat{\boldsymbol{\mu}}_{\hat{\theta}}$  converges in the RKHS norm in a way that is independent of the dimension of  $(X, Z)$ . This is an appealing property since estimation and inference based on  $\widehat{\boldsymbol{\mu}}_{\hat{\theta}}$  is less sensitive to the curse of dimensionality.

The following isomorphic result helps us to investigate where the dimensionality is ‘hiding’. Let  $L_2(\mathbb{R}^{1+d}, \Pi) = \{f : \mathbb{R}^{1+d} \rightarrow \mathbb{R}, s.t. \|f\| = (\int |f|^2 d\Pi)^{1/2} < \infty\}$ .

Since both the RKHS and the  $L_2(\mathbb{R}^{1+d}, \Pi)$  are separable and infinite dimensional Hilbert spaces, these two spaces are isometrically isomorphic, i.e., there exists a one-to-one linear mapping  $J : \mathcal{H}(k) \rightarrow L_2(\mathbb{R}^{1+d}, \Pi)$  such that

$$\langle J(f), J(g) \rangle_{L_2(\mathbb{R}^{1+d}, \Pi)} = \langle f, g \rangle_{\mathcal{H}(k)}, \quad f, g \in \mathcal{H}(k)$$

See Carrasco et al. (2007) for details.

Notice that the  $V$ -statistic version of  $n\widehat{\mathbb{M}}_n^2(\hat{\theta})$  can be thought as

$$\begin{aligned} \frac{n}{n^2} \sum_{i,j=1}^n \varepsilon(z_i; \hat{\theta}) k(x_i, x_j) \varepsilon(z_j; \hat{\theta}) &= \left\langle \frac{1}{\sqrt{n}} \sum_{i=1}^n \varepsilon(z_i; \hat{\theta}) \phi_{x_i}(\cdot), \frac{1}{\sqrt{n}} \sum_{j=1}^n \varepsilon(z_j; \hat{\theta}) \phi_{x_j}(\cdot) \right\rangle_{\mathcal{H}(k)} \\ &= \left\langle \frac{1}{\sqrt{n}} \sum_{i=1}^n \varepsilon(z_i; \hat{\theta}) J(\phi_{x_i}(\cdot)), \frac{1}{\sqrt{n}} \sum_{j=1}^n \varepsilon(z_j; \hat{\theta}) J(\phi_{x_j}(\cdot)) \right\rangle_{L_2(\mathbb{R}^{1+d}, \Pi)} \end{aligned} \quad (7)$$

where the first equality is a consequence of the reproducing property, and the last equality arises from the isometrically isomorphic relationship. Thus, the proposed test statistic  $n\widehat{\mathbb{M}}_n^2(\hat{\theta})$  is a  $U$ -statistic version of an ICM test with a weighting function  $J(\phi_x(\cdot))$ .

Furthermore, when the kernel is chosen to be ‘shift-invariant’, i.e., the kernel solely depends on the difference of its arguments,

$$k(x, x') = \psi(x - x')$$

a more specific ICM structure is revealed by the following characterization, which is due to Bochner (1933), see also Rudin (2017). We state it in the form given by Wendland (2004).

**Theorem 1** (Bochner). *Let  $k(x, x') = \psi(x - x')$  be a shift-invariant kernel for continuous function  $\psi : \mathbb{R}^d \rightarrow \mathbb{C}$ . Then  $\psi$  is positive definite if and only if it is the Fourier transform of a finite nonnegative Borel measure  $\Lambda$  on  $\mathbb{R}^d$ :*

$$\psi(t) = \int_{\mathbb{R}^d} \exp(-i\langle t, \omega \rangle) d\Lambda(\omega)$$

for  $t \in \mathbb{R}^d$ .

One may normalize  $\psi$  such that  $\psi(0) = 1$ , in which case  $\Lambda$  is a probability measure and  $\psi$  is its characteristic function. For example, if  $\Lambda$  is a normal distribution of the form  $(2\pi/\sigma^2)^{-d/2} e^{-\frac{\sigma^2 \|\omega\|^2}{2}} d\omega$ , then the corresponding ISPD kernel is the Gaussian  $\exp(-\|t\|^2/2\sigma^2)$ . By applying Bochner’s theorem, one can show that, see, e.g., Fan and Li (2000); Muandet et al. (2020)

$$\mathbb{M}^2(\theta) = (2\pi)^{-d/2} \int_{\mathbb{R}^d} \mathbb{E} \left( \varepsilon(Z; \theta) \exp(-i\omega^\top X) \right)^2 d\Lambda(\omega)$$

where  $\Lambda$  is a Fourier transform of the kernel  $k$ .

Several commonly studied ICM tests are indeed in a form of  $n\widehat{M}_n^2(\hat{\theta})$  (V-statistic version). For example, the ICM test of Bierens with exponential weighting function  $\exp(i\omega^\top x)$  and a uniform integral measure can be stated as

$$ICM_n = \frac{n}{n^2} \sum_{j,k=1}^n \varepsilon(z_j; \hat{\theta}) \varepsilon(z_k; \hat{\theta}) \exp\left(-\frac{1}{2}\|x_j - x_k\|^2\right)$$

where the kernel is chosen as the Gaussian RBF  $k(x, x') = \exp(-\|x - x'\|_2^2/2\sigma^2)$ ,  $\sigma = 1$ .

Another popular ICM test is Escanciano's  $PCvM_n$  (Escanciano, 2006a) with a weighting function of  $\mathbb{I}\{\omega^\top x \leq u\}$  and an empirical distribution as the integral measure, its analytic closed form is:

$$PCvM_n = \frac{n}{n^2} \sum_{j \neq k}^n \varepsilon(z_j; \hat{\theta}) \varepsilon(z_k; \hat{\theta}) \left( \frac{1}{n} \sum_{r=1}^n B_{jkq}^{(0)} \frac{\pi^{(d/2)-1}}{\Gamma(\frac{d}{2} + 1)} \right)$$

where  $\Gamma(\cdot)$  denotes the gamma function,  $B_{jkq}^{(0)}$  is the complementary angle between  $(x_j - x_q)$  and  $x_k - x_q$ , and is defined as

$$B_{jkq}^{(0)} = \left| \pi - \arccos \left( \frac{(x_j - x_q)^\top (x_k - x_q)}{|x_j - x_q| |x_k - x_q|} \right) \right|$$

Equation (7) and Bochner's theorem state that both weighting functions and integral measures are implicitly determined by the kernel.

### 2.2.2 GENERATE ICM TESTS, MASSIVELY AND CHEAPLY

Constructing ICM tests from weighting functions and integral measures is difficult, as one needs to perform the integration. Observe that each kernel corresponds to an ICM test, and the RKHS representation provides a cheap way to massively produce ICM tests by constructing kernels from existing kernels. To begin with, we first introduce some commonly used kernels listed in Table 1. For more examples, see, e.g., Muandet et al. (2017); Steinwart (2001); Steinwart and Christmann (2008).

Table 1: Various characterizations of known kernels

Kernel Function	$k(x, x')$	Domain $\mathcal{X}$	Characteristic	Shift Invariant	ISPD
Gaussian	$\exp(-\gamma\ x - x'\ _2^2), \gamma > 0$	$\mathbb{R}^d$	✓	✓	✓
Laplacian	$\exp(-\ x - x'\ _1/\sigma), \sigma > 0$	$\mathbb{R}^d$	✓	✓	✓
Inverse Multiquadric	$(c^2 + \ x - x'\ _2^2)^{-\gamma}, c, \gamma > 0$	$\mathbb{R}^d$	✓	✓	✓
Exponential	$\exp(\sigma\langle x, x' \rangle), \sigma > 0$	Compact sets of $\mathbb{R}^d$	✓	✗	✓
Matern	$2^{1-\nu}\Gamma^{-1}(\nu) (\sqrt{2\nu/\rho}\ x - x'\ _2)^\nu \kappa_\nu(\sqrt{2\nu/\rho}\ x - x'\ _2)$	$\mathbb{R}^d$	✓	✓	✓
Infinite Polynomial	$(1 - \langle x, x' \rangle)^{-\alpha}, \alpha > 0$	$\{x \in \mathbb{R}^d : \ x\ _2 < 1\}$	✓	✗	✓

**Notes:** In the Matern kernel,  $\Gamma(\cdot)$  is the Gamma function,  $\kappa_\nu$  is the modified Bessel function of the second type,  $\nu, \rho$  are non-negative parameters. When  $\nu \rightarrow \infty$ , it becomes equivalent to the Gaussian kernel, and when  $\nu = 1/2$ , it reduces to the Laplacian kernel.

The following lemmas describe ways of constructing new ISPD kernels.

**Lemma 2** *Let  $a \geq 0$ , and  $k$ ,  $k_1$  and  $k_2$  be ISPD kernels on  $X$ . Then  $ak$  and  $k_1 + k_2$  are also ISPD kernels on  $X$ .*

This lemma states that the set of ISPD kernels is a convex cone, its proof is trivial and will not be discussed here.

**Lemma 3** *A shift variant ISPD kernel,  $\tilde{k}$  can be obtained from a shift invariant ISPD kernel,  $k$ , as*

$$\tilde{k}(x, x') = f(x)k(x, x')f(x')$$

where  $f : \mathcal{X} \rightarrow \mathbb{R}$  is a bounded continuous function.

**Proof** See [Sriperumbudur et al. \(2010\)](#). ■

This lemma states that one can generate new ISPD kernels through conformal mapping, i.e., a transformation that preserves angles locally.

The next lemma is based on the fact that all bounded continuous shift-invariant kernels if they are characteristic, are also ISPD. A measurable and bounded kernel,  $k$  is said to be characteristic if

$$\mathbb{P} \rightarrow \int_{\mathcal{X}} k(\cdot, x)d\mathbb{P}(x)$$

is injective, that is,  $\mathbb{P}$  is embedded to a unique element in  $\mathcal{H}(k)$ . The above-mentioned shift-invariant kernels (i.e., the Gaussian, the Laplacian, the IMQ, and the Matern) are all characteristic kernels.

**Lemma 4** *Let  $k, k_1$  and  $k_2$  be a bounded continuous shift-invariant kernel on  $\mathbb{R}^d$ . Suppose  $k$  is characteristic and  $k_2 \neq 0$ , then  $k + k_1$  and  $k \times k_2$  are characteristic.*

**Proof** See [Sriperumbudur et al. \(2010\)](#). ■

[Escanciano \(2009\)](#) has shown that ICM tests only have substantial local power against alternatives in a finite-dimensional space, and there is only one direction with the highest asymptotic local power. This best direction depends on the weighting function, the integrated measure, the true model, and DGP. Since a kernel embodies the weighting function and the integral measure, it also affects the directions in which the corresponding ICM test has substantial power. Hence, different kernels would have different power properties, and it is desirable to have as many ICM tests as possible.

Choosing a kernel is important. Some kernels will gradually reduce to a constant function as the dimension  $d$  of  $X$  increases, making corresponding tests distorting size as well as losing power in almost all directions. Considering the Bierens' test, which corresponds to the Gaussian kernel with Euclidean distance and parameter  $\gamma = 1/2$ . When  $d$  is large, the corresponding kernel matrix  $K_{ij}$  (also known as the Gram matrix, see [Appendix A](#) for

details) becomes close to the identity matrix<sup>1</sup>. Nevertheless, one could “slow down” the decay rate by changing the parameter inside the Gaussian kernel, we will demonstrate this point in the simulation exercises.

### 2.3 Challenges when Using the Simple Statistic

We call the test statistic

$$n\widehat{\mathbb{M}}_n^2(\hat{\theta}) = \frac{n}{n(n-1)} \sum_{1 \leq i \neq j \leq n} \varepsilon(z_i; \hat{\theta})k(x_i, x_j)\varepsilon(z_j; \hat{\theta})$$

the simple statistic because it simply replaces unknown parts with their empirical counterparts. There are several potential drawbacks to using this simple statistic. First, if the estimator  $\hat{\theta}$  is standard in the sense that  $O_p(\|\hat{\theta} - \theta_0\|) = O_p(1/\sqrt{n})$ , then the distributions of the test statistic would depend on how  $\hat{\theta}$  is estimated. Furthermore, in most literature, to establish the limiting distribution of  $n\widehat{\mathbb{M}}_n^2(\hat{\theta})$  under the null, an asymptotic linear representation for  $\sqrt{n}(\hat{\theta} - \theta_0)$  is often required, see, e.g., [Delgado et al. \(2006\)](#); [Escanciano \(2006a\)](#).

Second, since the limiting distribution of  $n\widehat{\mathbb{M}}_n^2(\hat{\theta})$  is non-pivotal, a bootstrap procedure is needed to calculate critical values. However, the presence of  $\hat{\theta}$  often requires a case-by-case complicated parametric bootstrap procedure.

Finally, and perhaps more interestingly, certain ‘non-standard’ estimators  $\hat{\theta}$  with slower than  $1/\sqrt{n}$  rate of convergence are ruled out, as in these cases,  $\lim_{n \rightarrow \infty} O_p(\sqrt{n}\|\hat{\theta} - \theta_0\|) \rightarrow \infty$ .

To deal with this problem, one could try to find suitable transformations on the kernel to eliminate the “parameter estimation effect”. Most literature is based on the empirical process and are adopting two different transformation approaches. The first approach consists of martingale transformation of the empirical process, see, for instance, [Delgado and Stute \(2008\)](#); [Khmaladze \(1982, 1993\)](#); [Koul and Stute \(1999\)](#). However, the martingale transformation can be quite complicated even for some conventional econometrics models. More importantly, this transformation is based on a sequence of iterative regressions, the inversion of the projection matrix could be unstable, which will ultimately affect the sampling performance.

The second approach is based on the idea of projecting the weighting function  $h(X, t)$  onto a tangent space of nuisance parameters, see, for example, [Bickel et al. \(2006\)](#); [Escanciano and Goh \(2014\)](#); [Neyman \(1959\)](#); [Sant’Anna and Song \(2020\)](#); [Sant’Anna and Song \(2019\)](#). This approach is relatively easier to implement and requires weaker conditions than the Khmaladze transformation. In this study, we extend this projection idea to kernels.

---

1. Identity matrix for the V-statistic version of the test, and the kernel matrix is close to a zero matrix, where all elements are zeros, if one uses the U-statistic version.

### 3. A Projection Test Statistic and its Asymptotic Results

#### 3.1 The Projection Test Statistic

Let  $g(z; \theta)$  and  $\bar{\theta}$  are defined as before, the simple statistic can be expanded as

$$\begin{aligned} n\widehat{\mathbb{M}}_n^2(\hat{\theta}) &= \frac{n}{n(n-1)} \sum_{1 \leq i \neq j \leq n} \left( \varepsilon(z_i; \theta_0) + g^\top(z_i; \bar{\theta})(\hat{\theta} - \theta_0) \right) k(x_i, x_j) \left( \varepsilon(z_j; \theta_0) + g^\top(z_j; \bar{\theta})(\hat{\theta} - \theta_0) \right) \\ &= nA_{1,n}(k) + 2\sqrt{n}A_{2,n}(k)\sqrt{n}(\hat{\theta} - \theta_0) + \sqrt{n}(\hat{\theta} - \theta_0)^\top A_{3,n}(k)\sqrt{n}(\hat{\theta} - \theta_0) \end{aligned}$$

where

$$\begin{aligned} A_{1,n}(k) &= \frac{1}{n(n-1)} \sum_{1 \leq i \neq j \leq n} \varepsilon(z_i; \theta_0) k(x_i, x_j) \varepsilon(z_j; \theta_0) \\ A_{2,n}(k) &= \frac{1}{n(n-1)} \sum_{1 \leq i \neq j \leq n} \varepsilon(z_i; \theta_0) k(x_i, x_j) g^\top(z_j; \bar{\theta}) \\ A_{3,n}(k) &= \left( \frac{1}{n(n-1)} \sum_{1 \leq i \neq j \leq n} g(z_i; \bar{\theta}) k(x_i, x_j) g^\top(z_j; \bar{\theta}) \right) \end{aligned}$$

In order to eliminate the estimation effect  $\sqrt{n}(\hat{\theta} - \theta_0)$ , we need to find a  $k_p(\cdot, \cdot)$  such that

$$\mathbb{E}A_{2,n}(k_p) = \mathbb{E} \left( \varepsilon(Z; \theta_0) g(Z'; \theta_0) k_p(X, X') \right) = \mathbf{0}$$

and

$$\mathbb{E}A_{3,n}(k_p) = \mathbb{E} \left( g(Z; \theta_0) k_p(X, X') g(Z'; \theta_0)^\top \right) = 0$$

One possibility is

$$\begin{aligned} k_p(x, x') &= k(x, x') - g^\top(z; \theta_0) \Gamma_{\theta_0}^{-1} \mathbb{E}_{(X, Z)} \left( g(Z; \theta_0) k(X, x') \right) \\ &\quad - g^\top(z'; \theta_0) \Gamma_{\theta_0}^{-1} \mathbb{E}_{(X', Z')} \left( g(Z'; \theta_0) k(X', x) \right) \\ &\quad + g^\top(z; \theta_0) \Gamma_{\theta_0}^{-1} \mathbb{E} \left( g(Z; \theta_0) k(X, X') g^\top(Z'; \theta_0) \right) \Gamma_{\theta_0}^{-1} g(z'; \theta_0) \end{aligned} \quad (8)$$

where  $\Gamma_\theta = \mathbb{E} \left( g(Z; \theta) g^\top(Z; \theta) \right)$ , and

$$\mathbb{E}_{(X, Z)} \left( g(Z; \theta_0) k(X, x') \right) = \mathbb{E} \left( g(Z; \theta_0) k(X, x') | X' = x' \right)$$

The corresponding test statistic  $n\widehat{\mathbb{M}}_p^2(\hat{\theta})$  is specified as:

$$n\widehat{\mathbb{M}}_p^2(\hat{\theta}) = \frac{n}{n(n-1)} \sum_{1 \leq i \neq j \leq n} \varepsilon(z_i; \hat{\theta}) \hat{k}_p(x_i, x_j) \varepsilon(z_j; \hat{\theta}) \quad (9)$$

where  $\hat{k}_p(\cdot, \cdot)$  is the empirical counterpart of  $k_p(\cdot, \cdot)$ .

In the subsequent contents, we will explain (1) Where does  $k_p(\cdot, \cdot)$  come from? (2) What are the properties of  $k_p(\cdot, \cdot)$ , and (3) Does

$$\mathcal{E}(X; \theta_0) = 0 \Leftrightarrow \mathbb{E} \left( \varepsilon(Z; \theta_0) k_p(X, X') \varepsilon(Z'; \theta_0) \right) = 0, \quad P_x\text{-a.s.}$$

hold?

### A Projection Interpretation.

In the conventional ICM framework, the canonical way to “swipe out” the estimation effect is to project the weighting function onto the tangent space of nuisance parameters, see, for instance, [Escanciano and Goh \(2014\)](#); [Sant’Anna and Song \(2020\)](#); [Sant’Anna and Song \(2019\)](#). We adopt the same approach to the feature map  $\phi_x(\cdot)$  of  $k$ . Define a projection operator  $\mathcal{P}$  that takes a value in  $\mathcal{H}(k)$  and delivers the projected feature map  $\mathcal{P}\phi_x(\cdot)$ :

$$\mathcal{P}\phi_x(\cdot) = \phi_x(\cdot) - g^\top(z; \theta_0)\Gamma_{\theta_0}^{-1}\mathbb{E}_{XZ}(g(Z; \theta_0)\phi_X(\cdot)) \quad (10)$$

The intuition behind (10) is simple. First, note that  $\Gamma_{\theta_0}^{-1}\mathbb{E}_{XZ}(g(Z; \theta_0)\phi_X(\cdot))$  is the vector of linear projection coefficients of regressing  $\phi_x(\cdot)$  on the score function  $g(z; \theta_0)$ . Thus,  $g^\top(z; \theta_0)\Gamma_{\theta_0}^{-1}\mathbb{E}_{XZ}(g(Z; \theta_0)\phi_X(\cdot))$  is the best linear predictor of  $\phi_x(\cdot)$  given  $g(z; \theta_0)$ , and equation (10) is nothing more than the associated projection error, which, by definition, is orthogonal to  $g(z; \theta_0)$ .

Observe that  $\mathcal{P}$  is a linear operator, it follows the construction of RKHS that  $\mathcal{P}\phi_x(\cdot) \in \mathcal{H}(k)$ . This can be verified by checking its reproducing property (see Appendix C):

$$\begin{aligned} \langle \mathcal{P}\phi_x(\cdot), \phi_{x'}(\cdot) \rangle_{\mathcal{H}(k)} &= \mathcal{P}\phi_x(x') \\ &= k(x, x') - g^\top(z; \theta_0)\Gamma_{\theta_0}^{-1}\mathbb{E}_{XZ}(g(Z; \theta_0)k(X, x')) \end{aligned} \quad (11)$$

$k_p$  is then constructed from

$$k_p(x, x') = \langle \mathcal{P}\phi_x(\cdot), \mathcal{P}\phi_{x'}(\cdot) \rangle_{\mathcal{H}(k)} \quad (12)$$

It can be understood as the “residual square” of  $\phi_x(\cdot)$ .

### Properties of $K_p$ .

Note that the operator  $\mathcal{P}$  is an orthogonal projection operator in the Hilbert space of  $L_2(\mathbb{R}^{1+d}, P_{(X,Z)})$  but not necessarily in the RKHS  $\mathcal{H}(k)$ . Thus, in general

$$\langle \phi_x(\cdot), \mathcal{P}\phi_{x'}(\cdot) \rangle_{\mathcal{H}(k)} \neq \langle \mathcal{P}\phi_x(\cdot), \mathcal{P}\phi_{x'}(\cdot) \rangle_{\mathcal{H}(k)}$$

and

$$\langle \phi_x(\cdot), \mathcal{P}\phi_{x'}(\cdot) \rangle_{\mathcal{H}(k)} \neq \langle \mathcal{P}\phi_x(\cdot), \phi_{x'}(\cdot) \rangle_{\mathcal{H}(k)}$$

Nevertheless,  $\mathcal{P}$  is idempotent,

$$\mathcal{P}\mathcal{P}\phi_x(\cdot) = \mathcal{P}\phi_x(\cdot) - \mathcal{P}g^\top(z; \theta_0)\Gamma_{\theta_0}^{-1}\mathbb{E}_{XZ}(g(Z; \theta_0)\phi_X(\cdot)) = \mathcal{P}\phi_x(\cdot)$$

To investigate the positive definiteness of  $k_p(\cdot, \cdot)$ , it is tantamount to verifying the sign of

$$\int_{\mathcal{X} \times \mathcal{X}} f(x)k_p(x, x')f(x')dx dx', \quad f \in L_2$$

but the above equation can be written as

$$\begin{aligned} & \int_{\mathcal{X} \times \mathcal{X}} \langle f(x)\mathcal{P}\phi_x(\cdot), f(x')\mathcal{P}\phi_{x'}(\cdot) \rangle_{\mathcal{H}(k)} dx dx' \\ &= \left\langle \int_{\mathcal{X}} f(x)\mathcal{P}\phi_x(\cdot)dx, \int_{\mathcal{X}} f(x')\mathcal{P}\phi_{x'}(\cdot)dx' \right\rangle_{\mathcal{H}(k)} \\ &= \left\| \int_{\mathcal{X}} f(x)\mathcal{P}\phi_x(\cdot)dx \right\|_{\mathcal{H}(k)}^2 \geq 0 \end{aligned}$$

where the last equality comes from the independence between  $x$  and  $x'$ . The Moore–Aronszajn Theorem states that this positive definite kernel  $k_p(\cdot, \cdot)$  is associated with a unique RKHS  $\mathcal{H}(k_p)$ .

Let  $\mathcal{P}^1$  with  $\mathcal{P}^1\phi_x(\cdot) = \phi_x(\cdot) - \mathcal{P}\phi_x(\cdot)$  be another orthogonal projection operator, and by properties of  $\mathcal{P}^1$ , we have

$$\begin{aligned} \|\mathcal{P}^1\phi_x(\cdot)\|_{L_2(\mathbb{R}^{1+d}, P_{(X,Z)})}^2 &= \langle \mathcal{P}^1\phi_x(\cdot), \mathcal{P}^1\phi_x(\cdot) \rangle_{L_2(\mathbb{R}^{1+d}, P_{(X,Z)})} \\ &= \langle \mathcal{P}^1\phi_x(\cdot), \phi_x(\cdot) \rangle_{L_2(\mathbb{R}^{1+d}, P_{(X,Z)})} \\ &\leq \|\mathcal{P}^1\phi_x(\cdot)\|_{L_2(\mathbb{R}^{1+d}, P_{(X,Z)})} \|\phi_x(\cdot)\|_{L_2(\mathbb{R}^{1+d}, P_{(X,Z)})} \end{aligned}$$

Thus,

$$\|\mathcal{P}^1\phi_x(\cdot)\|_{L_2(\mathbb{R}^{1+d}, P_{(X,Z)})} \leq \|\phi_x(\cdot)\|_{L_2(\mathbb{R}^{1+d}, P_{(X,Z)})}$$

By the isometrically isomorphic relationship between  $\mathcal{H}(k)$  and  $L_2(\mathbb{R}^{1+d}, P_{(X,Z)})$ , we further have

$$\|\mathcal{P}^1\phi_x(\cdot)\|_{\mathcal{H}(k)} \leq \|\phi_x(\cdot)\|_{\mathcal{H}(k)} \tag{13}$$

with equality holds if  $\phi_x(\cdot) \in \text{span}\{g(z; \theta_0) : z \in \mathcal{Z}\}$ .

Thus, as long as  $\phi_x(\cdot) \notin \text{span}\{g(z; \theta_0) : z \in \mathcal{Z}\}$ ,

$$\|\mathcal{P}\phi_x(\cdot)\|_{\mathcal{H}(k)} = \|\phi_x(\cdot) - \mathcal{P}^1\phi_x(\cdot)\|_{\mathcal{H}(k)} > 0, \quad \forall x \in \mathcal{X}$$

$$\left\| \int_{\mathcal{X}} f(x)\mathcal{P}\phi_x(\cdot)dx \right\|_{\mathcal{H}(k)}^2 > \left\| \int_{\mathcal{X}} f(x)dx \inf_{x \in \mathcal{X}} \mathcal{P}\phi_x(\cdot) \right\|_{\mathcal{H}(k)}^2 \geq 0$$

and  $k_p(\cdot, \cdot)$  is an ISPD kernel.

**(Almost) Equivalence between  $\mathcal{E}(X; \theta_0)$  and  $\mathbb{E}(\varepsilon(Z; \theta_0)k_p(X, X')\varepsilon(Z'; \theta_0))$**

Similar to the story presented in Section 2, we show  $\mathcal{E}(X; \theta)$  is almost injective to a conditional moment embedding  $\mu_\theta^{(p)} \in \mathcal{H}(k_p)$  and  $\mathbb{E}(\varepsilon(Z; \theta_0)k_p(X, X')\varepsilon(Z'; \theta_0)) = \|\mu_{\theta_0}^{(p)}\|_{\mathcal{H}(k_p)}^2$ .

Redefine the operator  $\mathcal{C}_\theta$  as  $\mathcal{C}_\theta^{(p)} : \mathcal{H}(k_p) \rightarrow \mathbb{R}$ :

$$\mathcal{C}_\theta^{(p)} h = \mathbb{E}_{XZ}(\varepsilon(Z; \theta)h(X)), \quad h \in \mathcal{H}(k_p)$$

Let  $\phi_x^{(p)}(\cdot)$  be the feature map associated with  $k_p(\cdot, \cdot)$ . By the reproducing property, we have  $h(x) = \langle h, \phi_x^{(p)}(\cdot) \rangle_{\mathcal{H}(k_p)}$  and

$$\mathcal{C}_\theta^{(p)} h = \left\langle h, \mathbb{E}_{XZ} \left( \varepsilon(Z; \theta) \phi_x^{(p)}(\cdot) \right) \right\rangle_{\mathcal{H}(k_p)} = \left\langle h, \mu_\theta^{(p)} \right\rangle_{\mathcal{H}(k_p)}$$

By Riesz's representer theorem,

$$|\mathcal{C}_\theta^{(p)}| = \|\mu_\theta^{(p)}\|_{\mathcal{H}(k_p)}$$

The following theorem states the almost injectivity between  $\mathcal{E}(X; \theta)$  and  $\mu_\theta^{(p)}$ :

**Theorem 5** *For any  $\theta_1, \theta_2 \in \Theta$ , assume  $\mathcal{E}(x; \theta)$  is not collinear with  $g(Z; \theta_0)$ , then we have  $\mathcal{E}(X; \theta_1) = \mathcal{E}(X; \theta_2)$  if and only if  $\mu_{\theta_1}^{(p)} = \mu_{\theta_2}^{(p)}$ . Consequently,*

$$\mathcal{E}(X; \theta_0) = 0 \Leftrightarrow \|\mu_{\theta_0}^{(p)}\|_{\mathcal{H}(k_p)}^2 = 0 \quad P_x\text{-a.s.}$$

**Proof** See Appendix C. ■

Finally, by the construction of  $\mu_\theta^{(p)}$ , it is easy to check that

$$\|\mu_{\theta_0}^{(p)}\|_{\mathcal{H}(k_p)}^2 = \mathbb{E}(\varepsilon(Z; \theta_0)k_p(X, X')\varepsilon(Z'; \theta_0))$$

### 3.2 Asymptotic Null Distribution

One can estimate  $\mathbb{M}_p^2(\theta_0) = \|\mu_{\theta_0}^{(p)}\|_{\mathcal{H}(k_p)}^2$  by

$$\widehat{\mathbb{M}}_p^2(\hat{\theta}) = \frac{1}{n(n-1)} \sum_{1 \leq i \neq j \leq n} \varepsilon(z_i; \hat{\theta}) \hat{k}_p(x_i, x_j) \varepsilon(z_j; \hat{\theta})$$

Consequently, define the test statistic as  $n\widehat{\mathbb{M}}_p^2(\hat{\theta})$ . In this subsection, we study the asymptotic properties of  $n\widehat{\mathbb{M}}_p^2(\hat{\theta})$  under the null hypothesis.

To derive our theoretical results, we impose additional assumptions:

- (A5). (i)  $\mathbb{E}\|\varepsilon(Z; \theta_0)\|^2 < \infty$ . Let  $\nabla_\theta g(z; \theta)$  exist almost surely in an open neighborhood  $\mathcal{N}(\theta_0)$  of  $\theta_0$ . (ii)  $\mathbb{E}\|g(Z; \theta_0)\| < \infty$  and  $\sup_{\theta \in \mathcal{N}(\theta_0)} \|\nabla_\theta g(\cdot, \theta)\| < S(\cdot)$ , with  $\mathbb{E}S(Z) < \infty$ , where  $\|\cdot\|$  denotes either a vector or matrix norm. (iii)  $\nabla_\theta g(z, \theta)$  are continuous in  $\theta$  for  $\theta \in \mathcal{N}(\theta_0)$  and uniformly in  $Z$  almost everywhere.

- (A6). (i)  $\|\hat{\theta} - \theta_0\| = o_p(n^{-1/4})$ ; (ii)  $\Gamma_\theta$  is nonsingular uniformly in  $\theta \in \Theta^\circ$ .

**Remark.** Assumption A5 contains regularity conditions for  $\varepsilon(z; \theta)$ , and these conditions are similar in, e.g., [Delgado et al. \(2006\)](#); [Newey \(1985\)](#); [Robinson \(1991\)](#). A sufficient condition for Assumption A6 (i) is  $n^\varsigma \|\hat{\theta} - \theta_0\| = O_p(1)$  for some  $\varsigma > 1/4$ . In addition, we do not require  $\hat{\theta}$  to have an asymptotically linear representation.

**Lemma 6** *Under the null, we have*

$$\hat{k}_p(\cdot, \cdot) = k_p(\cdot, \cdot) + O_p(1/\sqrt{n}) + O_p(\|\hat{\theta} - \theta_0\|) \quad (14)$$

while expanding  $n\widehat{\mathbb{M}}_p^2(\hat{\theta})$  around  $\theta_0$  yields

$$\begin{aligned} n\widehat{\mathbb{M}}_p^2(\hat{\theta}) &= \frac{n}{n(n-1)} \sum_{i \neq j} \varepsilon(z_i; \theta_0) k_p(x_i, x_j) \varepsilon(z_j, \theta_0) + O_p(\|\hat{\theta} - \theta_0\|) + O_p(\|\hat{\theta} - \theta_0\|^2) \\ &\quad + O_p(1/\sqrt{n}) \end{aligned} \quad (15)$$

**Theorem 7** *Assume that  $\mathbb{M}_p^2(\theta) < \infty$  for all  $\theta \in \Theta$  and Assumption A6 (i) hold, under the null, we have*

$$n\widehat{\mathbb{M}}_p^2(\hat{\theta}) \xrightarrow{d} \sum_{k=1}^{\infty} \tau_k^{(p)} (W_k^2 - 1) \quad (16)$$

where  $W_k \sim N(0, 1)$ ,  $\{\tau_k^{(p)}\}$  are eigenvalues of the operator  $A$  defined as  $(A\psi)(v) = \int f(v, v') \psi(v') dP_v(v')$  for non-zero  $\psi$ ,  $v = (x, y)$ , and  $f(v, v') = \varepsilon(z; \theta_0) k_p(x, x') \varepsilon(z'; \theta_0)$

**Proof** See, [Serfling \(1980\)](#). ■

### 3.3 Asymptotic Power

We now study the asymptotic distribution of  $n\widehat{\mathbb{M}}_p^2(\hat{\theta})$  under fixed alternative and a sequence of local alternatives converging to null at a parametric rate  $n^{-1/2}$ .

We first consider the fixed alternative hypothesis. Observe that<sup>2</sup>

$$\begin{aligned} \sqrt{n}\widehat{\mathbb{M}}_p^2(\hat{\theta}) &= \frac{\sqrt{n}}{n(n-1)} \sum_{i \neq j} \varepsilon(z_i; \theta_0) k_p(x_i, x_j) \varepsilon(z_j, \theta_0) + O_p(\|\hat{\theta} - \theta_0\|) + O_p(\|\hat{\theta} - \theta_0\|^2) + O_p(1/\sqrt{n}) \\ &= \frac{\sqrt{n}}{n(n-1)} \sum_{i \neq j} \varepsilon(z_i; \theta_0) k_p(x_i, x_j) \varepsilon(z_j, \theta_0) + O_p(\|\hat{\theta} - \theta_0\|) + O_p(1/\sqrt{n}) \end{aligned}$$

---

2. Since both  $A_{1,n}^{(p)}$  and  $A_{2,n}^{(p)}$  are non-degenerate, and  $A_{1,n}^{(p)}, A_{2,n}^{(p)} = O_P(1/\sqrt{n})$ .  $A_{1,n}^{(p)}$  and  $A_{2,n}^{(p)}$  are defined in the proof of Lemma 6.

**Theorem 8** Assume that  $M_p^2(\theta) < \infty$  for all  $\theta \in \Theta$ , and  $\hat{\theta} \xrightarrow{P} \theta_1 \in \Theta$ . Under the fixed alternative, we have

$$\sqrt{n} \left( \widehat{M}_p^2(\hat{\theta}) - M_p^2(\theta_1) \right) \xrightarrow{d} N(0, \sigma_{\theta_1, p}^2) \quad (17)$$

where

$$\sigma_{\theta_1, p}^2 = 4\text{Var}_{(X, Z)} \left( \mathbb{E}_{(X', Z')} \left( \varepsilon(Z; \theta_1) k_p(X, X') \varepsilon(Z'; \theta_1) \right) \right)$$

**Proof** See [Serfling \(1980\)](#). ■

It is readily to see that for large  $n$  and fixed critical value  $c_\alpha$ , the test power can be approximated by

$$P_{H_1}(n\widehat{M}_p^2(\hat{\theta}) > c_\alpha) \approx \Phi \left( \frac{\sqrt{n}M_p^2(\theta_1)}{\sigma_{\theta_1, p}} - \frac{c_\alpha}{\sqrt{n}\sigma_{\theta_1, p}} \right)$$

where  $\Phi$  denotes the cumulative distribution function of the standard normal distribution.

Assume that  $n$  is sufficiently large, in  $\sqrt{n}M_p^2(\theta_1)/\sigma_{\theta_1, p} - c_\alpha/\sqrt{n}\sigma_{\theta_1, p}$ , we observe that the second term  $c_\alpha/\sqrt{n}\sigma_{\theta_1, p} = O(n^{-1/2})$  going to 0 as  $n \rightarrow \infty$ , while the first term  $\sqrt{n}M_p^2(\theta_1)/\sigma_{\theta_1, p} = O(n^{1/2})$ , dominating the second. Thus, the best kernel that maximizes the test power is given by

$$k^* = \arg \sup_{k \in \mathcal{K}} \frac{\sqrt{n}M_p^2(\theta_1)}{\sigma_{\theta_1, p}}$$

where  $\mathcal{K}$  is a proper kernel space, e.g.,  $\mathcal{K} = \{\exp(-\gamma\|x - x'\|), \gamma > 0\}$ .

A heuristic way to estimate  $k^*$  is to divide the sample  $\{(x_i, z_i), i = 1, \dots, n\}$  into two disjoint training and test sets, and use the training set to compute  $(\widehat{M}_p(\hat{\theta})/\hat{\sigma}_{\theta_1, p})(k)$ , which can be maximized by choosing the kernel parameter (e.g., in Gaussian kernel, the kernel parameter is  $\gamma$ ). We denote the kernel that maximize  $(\widehat{M}_p(\hat{\theta})/\hat{\sigma}_{\theta_1, p})(k)$  as  $\hat{k}^*$ . We then, use  $\hat{k}^*$  and perform testing in the test set.

A similar idea has been discussed in the machine learning literature, where the test of interest is the equality of two samples, see, e.g., [Gretton et al. \(2012\)](#). Nevertheless, extending the idea to our content is not trivial. Specifically, there are several key questions needed to be answered for validating this heuristic procedure:

- Does  $\hat{k}^* \xrightarrow{P} k^*$ , and if so, what is the convergence rate.
- Does  $(\widehat{M}_p(\hat{\theta})/\hat{\sigma}_{\theta_1, p})(\hat{k}^*) \xrightarrow{P} (M_p^2(\theta_1)/\sigma_{\theta_1, p})(k^*)$ , and if so, what is the convergence rate.

We leave these questions in future research.

**Remark.** On the other hand, Theorem 8 indicates that our test statistic might not be consistent against all fixed alternative hypotheses if  $\varepsilon(Z; \theta_0)$  is collinear to the function  $g(Z; \theta_0)$ . However, given the nonlinear nature of our model, we do not think of this type of alternative being relevant.

We now proceed to consider the asymptotic local power properties. To this end, we study the asymptotic distribution of  $n\widehat{\mathbb{M}}_p^2(\hat{\theta})$  under a certain sequence of Pitman-type local alternatives converging to null at a parametric rate:

$$H_{1,n} : \mathbb{E}(Y|X = x) = \mathcal{M}(x; \theta_0) + \frac{R(x)}{\sqrt{n}} \quad (18)$$

where the random variable  $R(X)$  is  $P_X$ -integrable and satisfies  $P(R(X) = 0) < 1$ .

**Theorem 9** *Assume that  $\mathbb{M}_p^2(\theta) < \infty$  for all  $\theta \in \Theta$ . Under  $H_{1,n}$ , we have*

$$\begin{aligned} n\widehat{\mathbb{M}}_n^2(\hat{\theta}) &\xrightarrow{d} \sum_{k=1}^{\infty} \tau_k^{(p)} (W_k^2 - 1) + 2N(0, 4\text{Var}_{X,Z}(\mathbb{E}_{X',Z'}\varepsilon(Z; \theta_0)k_p(X, X')R(X'))) \\ &\quad + \mathbb{E}(R(X)k_p(X, X')R(X')) \end{aligned} \quad (19)$$

where  $\sum_{k=1}^{\infty} \tau_k^{(p)} (W_k^2 - 1)$  is defined in Theorem 7.

**Proof** See Appendix C. ■

**Remark.** A pathological situation in which our test will only have trivial local power against such alternatives is when  $R(X)$  is a linear combination of  $g(Z; \theta_0)$ , i.e.,  $R(x) = \nu^\top g(Z; \theta_0)$  a.s. for some nonzero vector  $\nu$ . In such a case, the limiting distribution of  $n\widehat{\mathbb{M}}_p^2(\hat{\theta})$  under  $H_0$  and  $H_{1,n}$  is the same so that  $H_{1,n}$  can not be detected. However, such a specific class of local alternatives is of very limited practical interest.

The following lemma states that the proposed test only has non-trivial local power in a finite-dimensional space, and there is only one direction with the highest asymptotic local power. Although this lemma is essentially Theorem 1 of Escanciano (2009), it provides a clear viewpoint that highlights the importance of a kernel.

To begin with, let  $T_k$  be an integral operator defined as

$$T_k f(x) = \int_{\mathcal{X}} k(x, x') f(x') dP_X(x')$$

Mercer's theorem states that one can characterize a kernel  $k_p(\cdot, \cdot)$  as:

$$k_p(x, x') = \sum_{j \geq 1} \lambda_j e_j(x) e_j(x')$$

where the convergence is absolute and uniform, and  $\{\lambda_j\}_{j \geq 1}, \{e_j(\cdot)\}_{j \geq 1}$  are eigenvalues and eigenfunctions of the operator  $T_k$ , respectively.  $\lambda_1 > \lambda_2 > \dots$  and  $\lambda_j \rightarrow 0$  as  $j \rightarrow \infty$ . Since  $\{e_j(\cdot)\}_{j \geq 1}$  are also basis of the space  $L_2(\mathbb{R}^d, P_X)$ , one can write  $R(x) = \sum_{s \geq 1} \alpha_s e_s(x)$ ,  $\alpha_s = \langle R, e_s \rangle_{L_2(\mathbb{R}^d, P_X)} \in \mathbb{R}$ .

**Lemma 10** *Under local alternatives, let*

$$\mathbb{M}_p^2(\theta_0) = \mathbb{E} \left( \left( \varepsilon(Z; \theta_0) + \frac{R(X)}{\sqrt{n}} \right) k_p(X, X') \left( \varepsilon(Z'; \theta_0) + \frac{R(X')}{\sqrt{n}} \right) \right)$$

we have

$$\mathbb{M}_p^2(\theta_0) = \mathbb{E} \left( \varepsilon(Z; \theta_0) k_p(X, X') \varepsilon(Z'; \theta_0) \right) + \lambda_j \frac{\alpha_j^2}{n} + 2\lambda_j \mathbb{E} \left( \varepsilon(Z; \theta_0) e_j(X) \right) \frac{\alpha_j}{\sqrt{n}} \quad (20)$$

**Proof** See Appendix C. ■

Immediately, we conclude that if  $\alpha_j \neq 0$  but  $\{\alpha_s\}_{s \neq j} = 0$ , then when  $j = 1$ , i.e.,  $R(x) = \alpha_1 e_1(x)$ , the proposed test have highest asymptotic local power. The local power decreases when  $j$  increases, and when  $j \rightarrow \infty$  one can only have trivial power. Nevertheless, for any fixed direction, e.g.,  $R(x) = \alpha_s e_s(x)$ , we can change the value of  $\lambda_s$  (equivalently, change the kernel) to increase the local power.

#### 4. A Multiplier Bootstrap Procedure

Our test statistic  $n\widehat{\mathbb{M}}_p^2(\hat{\theta})$  is non-pivotal, in this section, we propose a simple-to-use multiplier bootstrap procedure to approximate the null distribution. Its implementation is listed below:

1. Generate a sequence of i.i.d random variables  $\{v_i : i = 1, 2, \dots, n\}$  with mean zero and variance one; e.g., Rademacher random variable, standard normal random variable, or Bernoulli random variable with  $P(v = 1 - \kappa) = \kappa/\sqrt{5}$  and  $P(v = \kappa) = 1 - \kappa/\sqrt{5}$ , where  $\kappa = (\sqrt{5} + 1)/2$  (Mammen, 1993).

2. Compute

$$\left( n\widehat{\mathbb{M}}_p^{2,*}(\hat{\theta}) \right)_b = \frac{1}{n-1} \sum_{i \neq j} \varepsilon(z_i; \hat{\theta}) v_i \hat{k}_p(x_i, x_j) \varepsilon(z_j; \hat{\theta}) v_j$$

3. Repeat steps 1 and 2  $B$  times, and collect  $\left\{ \left( n\widehat{\mathbb{M}}_p^{2,*}(\hat{\theta}) \right)_b, b = 1, 2, \dots, B \right\}$
4. Define a confidence level  $\alpha$ , obtain the  $(1-\alpha)$ -th quantile of  $\left\{ \left( n\widehat{\mathbb{M}}_p^{2,*}(\hat{\theta}) \right)_b, b = 1, 2, \dots, B \right\}$ ,  $c_{n,\alpha}^*$ .
5. Reject the null if  $n\widehat{\mathbb{M}}_p^2(\hat{\theta}) > c_{n,\alpha}^*$ , and fail to reject otherwise.

The multiplier bootstrapped test statistic  $n\widehat{\mathbb{M}}_p^{2,*}(\hat{\theta})$  has several attractive properties. First, it does not require computing new parameter estimates at each bootstrap draw, reducing the computational intensity of the proposed procedure. Second, due to the employment of the projection, its implementation does not require using estimators that admit

an asymptotic linear representation. These computational conveniences are important when the dimension  $d$  is high.

The next theorem establishes the asymptotic validity of the proposed multiplier bootstrap procedure.

**Theorem 11** *Assume that  $\mathbb{M}_p^2(\theta) < \infty$  for all  $\theta \in \Theta$ . Then, we have  $n\widehat{\mathbb{M}}_p^{2,*}(\hat{\theta}) \xrightarrow{d,*} \sum_{k=1}^{\infty} \tau_k^{(p)}(W_k^2 - 1)$ , with probability one under the bootstrap law. Here  $\sum_{k=1}^{\infty} \tau_k^{(p)}(W_k^2 - 1)$  is defined as the same in Theorem 7, and  $\xrightarrow{d,*}$  denotes weak convergence under the bootstrap law, i.e., conditional on the original sample  $\{z_i, x_i : i = 1, 2, \dots, n\}$ .*

**Proof** See Appendix C. ■

Theorem 11 states that the bootstrap statistic  $n\widehat{\mathbb{M}}_p^{2,*}(\hat{\theta})$  converges to the null distribution of  $n\widehat{\mathbb{M}}_p^2(\hat{\theta})$  conditional on the original sample under  $H_0, H_1$  and  $H_{1,n}$ . This fact is what allows the proposed procedure to work.

## 5. A Minimum Distance Estimator

Based on the U-statistic expression derived in Section 2, we present a minimum distance estimator in this section. It is known that when the number of the arbitrarily chosen instruments is finite, the GMM estimation procedure could render inconsistent estimates due to an identification problem, see, e.g., Domínguez and Lobato (2004) for various examples. Integrated conditional moment, on the other hand, introduces infinitely many instruments, and therefore, does not arise the identification issue. Domínguez and Lobato (2004) is the first in the literature to introduce a consistent estimation procedure based on the ICM framework, their estimator reads as,

$$\hat{\theta} = \arg \min_{\theta \in \Theta} \frac{1}{n^3} \sum_{j=1}^n \left( \sum_{i=1}^n \varepsilon(z_i; \theta) \mathbb{I}\{x_i \leq x_j\} \right)^2$$

This estimator corresponds to an indicator weighting function  $\mathbb{I}\{x \leq u\}$ , and suffers from the curse of dimensionality due to data sparseness.

The representation of ICM statistics in the RKHS provides a natural channel to develop a minimum distance type estimator:

$$\hat{\theta} = \arg \min_{\theta \in \Theta} \frac{1}{n(n-1)} \sum_{i \neq j} \varepsilon(z_i; \theta) k(x_i, x_j) \varepsilon(z_j; \theta) \quad (21)$$

The objective function  $\widehat{R}_U(\theta)$  can be rewritten as

$$\widehat{R}_U(\theta) = \varepsilon(z; \theta)^\top W_U \varepsilon(z; \theta)$$

where  $W_U \in \mathbb{R}^{n \times n}$  is a symmetric weight matrix that depends on the kernel matrix  $K$  with  $K_{i,j} = k(x_i, x_j)$ . Here  $W_U = (K - \text{diag}(K_{11}, \dots, K_{nn})) / (n(n-1))$ , where  $\text{diag}(a_1, \dots, a_n)$  denotes an  $n \times n$  diagonal matrix whose diagonal elements are  $a_1, \dots, a_n$ .

Although using the objective function  $\widehat{R}_U(\theta)$ , one could obtain a minimum-variance unbiased estimator, the weight matrix  $W_U$ , unfortunately, is indefinite, since  $\text{trace}(W_U) = \sum_{i=1}^n \varpi_i = 0$ , where  $\{\varpi_i; i = 1, \dots, n\}$  are the eigenvalues of  $W_U$ . Thus, we conclude that there exist both positive and negative eigenvalues.

We, therefore, focus on the V-statistic version of this estimator:

$$\hat{\theta} = \arg \min_{\theta \in \Theta} \frac{1}{n^2} \sum_{i,j=1}^n \varepsilon(z_i; \theta) k(x_i, x_j) \varepsilon(z_j; \theta) \quad (22)$$

whose objective function  $\widehat{R}_V(\theta)$  can be written as

$$\widehat{R}_V(\theta) = \varepsilon(z; \theta)^\top W_V \varepsilon(z; \theta)$$

with  $W_V = K/n^2$ .

Let  $R_k(\theta) = \mathbb{E}(f_\theta(V, V'))$ , where  $f_\theta(v, v') = \varepsilon(z; \theta) k(x, x') \varepsilon(z'; \theta)$  with  $v = (x, z)$ , and denote  $\|\cdot\|_F$  as the Frobenius norm. The following theorems establish the asymptotic properties of this estimator.

**Theorem 12** *Assume that  $\mathbb{E}(|Y|^2 < \infty)$ ,  $\mathbb{E}(\sup_{\theta \in \Theta} |\mathcal{M}(X; \theta)|^2) < \infty$ ,  $\Theta$  is compact and convex,  $R_k(\theta)$  is uniquely minimized at  $\theta_0 \in \Theta^\circ$ , and Assumption A4 holds, then*

$$\hat{\theta} \xrightarrow{P} \theta_0$$

**Proof** See Appendix C. ■

**Theorem 13** *Suppose that  $\mathcal{M}(X; \theta)$  is twice continuously differentiable about  $\theta$ ,  $\Theta$  is compact,  $H = \mathbb{E}(\nabla_\theta^2 f_{\theta_0}(V, V'))$  is non-singular,  $\mathbb{E}(|Y|^2 < \infty)$ ,  $\mathbb{E}(\sup_{\theta \in \Theta} |\mathcal{M}(X; \theta)|^2) < \infty$ ,  $\mathbb{E}(\sup_{\theta \in \Theta} \|\nabla_\theta \mathcal{M}(X; \theta)\|_2^2) < \infty$ ,  $\mathbb{E}(\sup_{\theta \in \Theta} \|\nabla_\theta^2 \mathcal{M}(X; \theta)\|_F^2) < \infty$ ,  $R_k(\theta)$  is uniquely minimized at  $\theta_0 \in \Theta^\circ$ , and Assumption A4 holds, then*

$$\sqrt{n}(\hat{\theta} - \theta_0) \xrightarrow{d} N(0, \Sigma_V)$$

where

$$\Sigma_V = 4H^{-1} \text{Var}_V(\mathbb{E}_{V'}^2(\nabla_\theta f_{\theta_0}(V, V'))) H^{-1}$$

**Proof** See Appendix C. ■

In general, the estimator given by (22) is not efficient. An efficient estimator based on the infinite number of moment conditions can be constructed following the ideas of Carrasco and Florens (2000). For regularized and infinite dimensional estimators based on  $\widehat{R}_{U(V)}(\theta)$ , see Zhang et al. (2020).

## 6. Monte Carlo Studies

This section conducts a sequence of Monte Carlo simulations to evaluate the finite sample performance of the kernel-based tests. The running data-generating processes (DGPs) are linear models, for additional simulation exercises, we refer readers to Appendix D. Nevertheless, before specifying DGPs, first observe that for any  $\theta \in \Theta$ , we have

$$\mathbb{E} [\varepsilon(Z; \theta)k_p(X, X')\varepsilon(Z'; \theta)] = \mathbb{E} [\varepsilon_p(Z; \theta)k(X, X')\varepsilon_p(Z'; \theta)] \quad (23)$$

where

$$\varepsilon_p(z; \theta) = \varepsilon(z; \theta) - g^\top(z; \theta)\Gamma_{\theta_0}^{-1}\mathbb{E}(g(Z; \theta)\varepsilon(Z; \theta))$$

is the projection residual of  $\varepsilon(z; \theta)$ . Its verification can be found in Appendix C. This representation greatly simplifies computation, as

$$n\widehat{\mathbb{M}}_p^2(\hat{\theta}) = \frac{1}{n-1} \sum_{i \neq j} \hat{\varepsilon}_p(z_i; \hat{\theta})k(x_i, x_j)\hat{\varepsilon}_p(z_j; \hat{\theta})$$

and

$$n\widehat{\mathbb{M}}_p^{2,*}(\hat{\theta}) = \frac{1}{n-1} \sum_{i \neq j} \hat{\varepsilon}_p(z_i; \hat{\theta})v_i k(x_i, x_j)\hat{\varepsilon}_p(z_j; \hat{\theta})v_j$$

where

$$\hat{\varepsilon}_p(z_i; \hat{\theta}) = \varepsilon(z_i; \hat{\theta}) - g^\top(z_i)\Gamma_{n,\hat{\theta}}^{-1} \left( \frac{1}{n} \sum_{s=1}^n g(z_s; \hat{\theta})\varepsilon(z_s; \hat{\theta}) \right)$$

and  $\Gamma_{n,\hat{\theta}}$  is the empirical counterpart of  $\Gamma_{\theta_0}$ .

### 6.1 Data Generating Processes

We consider the following DGPs:

- DGP(m):  $Y_i = \beta_0 + \sum_{j=1}^m \beta_j X_{ji} + \sigma_i^{(m)} \varepsilon_i$ .
- DGP-LOCAL(m):  $Y_i = \beta_0 + \sum_{j=1}^m \beta_j X_{ji} + n^{-1/2} \sum_{j=1}^m \beta_j X_{ji}^2 + \sigma_i^{(m)} \varepsilon_i$ .

DGP(m) specifies  $m$  covariates and is used to evaluate the size performance of proposed tests. DGP-LOCAL(m) is used to evaluate the local powers of the corresponding null DGPs.

We allow for conditional heteroskedasticity in all models and generate the covariates and heteroskedasticity as follows.

In DGP(m) and DGP-LOCAL(m),

- When  $m = 2$ ,  $X_1, X_2 \sim N(0, 1)$ , and  $\sigma^{(2)} = (0.1 + X_1^2 + X_2^2)^{1/2}$ .
- When  $m = 5$ ,  $X_j \sim U(0, j)$  for  $j = 1, 2, 3$ ,  $X_j \sim N(0, (j-3)^2)$  for  $j = 4, 5$ .  $\sigma^{(5)} = (0.1 + \sum_{j=1}^3 X_j + \sum_{j=4}^5 X_j^2)^{1/2}$ .

- When  $m = 10$ ,  $X_j \sim U(0, j)$  for  $j = 1, \dots, 5$ ,  $X_j \sim N(0, (j - 5)^2)$  for  $j = 6, \dots, 10$ .  
 $\sigma^{(10)} = \left(0.1 + \sum_{j=1}^5 X_j + \sum_{j=5}^{10} X_j^2\right)^{1/2}$ .
- When  $m = 20$ ,  $X_j \sim U(0, j)$  for  $j = 1, \dots, 10$ ,  $X_j \sim N(0, (j - 10)^2)$  for  $j = 11, \dots, 20$ .  
 $\sigma^{(20)} = \left(0.1 + \sum_{j=1}^{10} X_j + \sum_{j=11}^{20} X_j^2\right)^{1/2}$ .

In all cases, we set  $\varepsilon_i \sim N(0, 1)$ , and set  $\beta_j$ 's to be 1.

## 6.2 Test Statistics and Simulation Results

From the discussion of asymptotic power properties, it should be clear now that the choice of a kernel is important for ICM testing. Nevertheless, finding a case-dependent optimal kernel is challenging. In this subsection, we provide a heuristic algorithm for tuning the parameter of a Gaussian kernel  $k(x, x') = \exp(-\gamma\|x - x'\|_2^2)$ ,  $\gamma > 0$ . Note that we are not claiming such an algorithm would lead to an optimal testing statistic. Rather, we believe that this heuristic algorithm would not deliver a bad test (in the sense of small power) with greater probability.

The Gaussian kernel is the default kernel in many kernel-based algorithms. Conventional wisdom suggests that ‘Gaussian kernels tend to yield good performance under general smoothness assumptions and should be considered especially if no additional knowledge of the data is available’ (Smola et al., 1998). A Gaussian kernel takes the form of a normal distribution and is smooth. The tuning parameter  $\gamma$  determines how well this kernel fits the data  $\mathbb{X}$ , here  $\mathbb{X}$  is a  $n \times d$  matrix consisting of conditional variables. Fixing an input data  $x'$ , a large  $\gamma$  would lead to an over-fitting scenario since large weight would be concentrated around  $x'$ , while points that are far away from  $x'$  would have a kernel value that decay to zero exponentially. A small  $\gamma$  corresponds to an under-fitting scenario, as points would have kernel values close to one. These two cases are essentially the same thing: the resulting kernel values concentrate around one point (zero or one), making a test powerless.

The desired parameter  $\gamma$  would ‘spread’ the kernel values in the range  $(0, 1]$ , one nature method is to normalize the data using the second-moment information of the input matrix  $\mathbb{X}$ . We propose to perform a principal component analysis (PCA) for  $\mathbb{X}$ , and set  $\gamma = 1/(2\zeta_1)$ , where  $\zeta_1$  is the largest principal value. The idea is simple: In high-dimensional cases, the norm  $\|x - x'\|_2^2 = \sum_{s=1}^d |x_s - x'_s|^2$  are more likely to make kernel values concentrate around zero than one, and the higher the variance of the data, the higher the probability of occurring such concentration. By setting  $\gamma = 1/(2\zeta_1)$ , one could avoid such a phenomenon.

We consider five kernels in simulation studies:

- Gaussian Kernel,  $k_1(x, x') = \exp(-(1/\zeta_1)\|x - x'\|_2^2)$ .
- Inverse Multiquadric (IMQ) Kernel,  $k_2(x, x') = (1 + \|x - x'\|_2^2)^{-1.5}$ .
- Gaussian+IMQ Kernel,  $k_3(x, x') = k_1(x, x') + (1 + \|x - x'\|_2^2)^{-0.5}$ .

- Shift Variant Kernel,  $k_4(x, x') = (2 + \sin(4\|x\|_2))k_1(x, x')(2 + \sin(4\|x'\|_2))$
- Local Periodic Kernel,  $k_5(x, x') = (2 + \sin(0.1\|x - x'\|_2^2))k_1(x, x')$ .

The Gaussian+IMQ kernel, the Shift Variant kernel, and the Local Periodic kernel are constructed as results of Lemmas 2, 3 and 4, respectively.

We report the simulation results in Table 2. The nominal significance levels are given by 0.01, 0.05, and 0.1, while the sample sizes range from  $N = 100, N = 200$  to  $N = 400$ . For each experiment, i.e., each DGP and sample size, we run 1000 simulations. For each round of the simulation, the bootstrap procedure repeats 500 times to estimate the critical values. The parameters  $\beta_j$ 's are estimated by the ordinary least squares.

Table 2: Simulation Results, Five Kernels

N=100	0.1					0.05					0.01				
	Gaussian	IMQ	Gaussian+IMQ	Shift Variant	Local Periodic	Gaussian	IMQ	Gaussian+IMQ	Shift Variant	Local Periodic	Gaussian	IMQ	Gaussian+IMQ	Shift Variant	Local Periodic
DGP	0.113	0.124	0.113	0.116	0.118	0.068	0.063	0.065	0.065	0.065	0.011	0.016	0.017	0.015	0.016
DGP(2)	0.112	0.109	0.107	0.101	0.102	0.054	0.057	0.059	0.051	0.059	0.009	0.013	0.017	0.01	0.014
DGP(5)	0.109	0.117	0.116	0.109	0.109	0.055	0.057	0.062	0.054	0.059	0.015	0.007	0.02	0.02	0.018
DGP(10)	0.081	0.081	0.082	0.102	0.11	0.053	0.047	0.047	0.064	0.07	0.016	0.011	0.02	0.02	0.019
DGP(20)	0.247	0.205	0.213	0.255	0.238	0.17	0.129	0.132	0.179	0.155	0.057	0.045	0.04	0.052	0.045
DGP-LOCAL(2)	0.329	0.293	0.322	0.297	0.348	0.247	0.205	0.24	0.21	0.233	0.098	0.086	0.102	0.064	0.101
DGP-LOCAL(5)	0.963	0.907	0.965	0.939	0.948	0.938	0.86	0.946	0.888	0.923	0.86	0.718	0.849	0.796	0.821
DGP-LOCAL(10)	0.999	1	1	1	1	0.999	1	1	0.998	1	0.999	1	0.998	0.992	0.999
DGP-LOCAL(20)															
N=200	0.1					0.05					0.01				
	Gaussian	IMQ	Gaussian+IMQ	Shift Variant	Local Periodic	Gaussian	IMQ	Gaussian+IMQ	Shift Variant	Local Periodic	Gaussian	IMQ	Gaussian+IMQ	Shift Variant	Local Periodic
DGP	0.111	0.113	0.111	0.109	0.123	0.053	0.059	0.056	0.046	0.065	0.007	0.013	0.009	0.011	0.013
DGP(2)	0.101	0.127	0.103	0.105	0.111	0.05	0.062	0.046	0.058	0.064	0.01	0.01	0.01	0.012	0.009
DGP(5)	0.099	0.112	0.097	0.119	0.154	0.054	0.058	0.044	0.069	0.091	0.012	0.016	0.011	0.019	0.023
DGP(10)	0.108	0.093	0.109	0.116	0.112	0.055	0.05	0.059	0.061	0.062	0.013	0.012	0.015	0.015	0.015
DGP(20)	0.234	0.228	0.212	0.234	0.239	0.149	0.134	0.137	0.153	0.146	0.045	0.045	0.047	0.052	0.038
DGP-LOCAL(2)	0.347	0.284	0.352	0.314	0.353	0.261	0.178	0.272	0.218	0.264	0.108	0.065	0.119	0.085	0.115
DGP-LOCAL(5)	0.972	0.931	0.977	0.962	0.977	0.958	0.903	0.958	0.933	0.953	0.891	0.778	0.878	0.81	0.855
DGP-LOCAL(10)	1	1	1	1	1	1	1	1	1	1	1	1	1	0.999	1
DGP-LOCAL(20)															
N=400	0.1					0.05					0.01				
	Gaussian	IMQ	Gaussian+IMQ	Shift Variant	Local Periodic	Gaussian	IMQ	Gaussian+IMQ	Shift Variant	Local Periodic	Gaussian	IMQ	Gaussian+IMQ	Shift Variant	Local Periodic
DGP	0.101	0.119	0.099	0.098	0.126	0.051	0.066	0.044	0.047	0.066	0.013	0.007	0.009	0.01	0.022
DGP(2)	0.1	0.098	0.102	0.105	0.119	0.051	0.057	0.05	0.051	0.061	0.011	0.008	0.016	0.013	0.011
DGP(5)	0.102	0.124	0.106	0.118	0.094	0.049	0.062	0.061	0.064	0.051	0.014	0.013	0.015	0.019	0.012
DGP(10)	0.12	0.095	0.118	0.102	0.109	0.068	0.051	0.054	0.059	0.058	0.013	0.011	0.009	0.015	0.011
DGP(20)	0.214	0.193	0.205	0.227	0.225	0.137	0.119	0.117	0.152	0.145	0.05	0.04	0.04	0.044	0.057
DGP-LOCAL(2)	0.338	0.294	0.353	0.3	0.365	0.239	0.192	0.236	0.208	0.278	0.103	0.068	0.091	0.072	0.114
DGP-LOCAL(5)	0.981	0.955	0.988	0.959	0.977	0.97	0.926	0.974	0.933	0.959	0.906	0.816	0.905	0.824	0.892
DGP-LOCAL(10)	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
DGP-LOCAL(20)															

We first analyze the size of the proposed tests. From the results of  $DGP(m), m = 2, 5, 10, 20$ , we find that the actual empirical sizes of all the proposed tests are close to their nominal sizes, even when the sample size is as small as 100 and the dimension is as high as 20. Results of  $DGP-LOCAL(m), m = 2, 5, 10, 20$  confirm that all proposed tests have non-trivial local power, the power is especially high when tests are facing high-dimensional models. We want to emphasize that the size, as well as the power of conventional CMR tests, often diminish rapidly to zero as dimension increases, and the degeneracy of levels or powers does not improve when sample size increases. Lastly, we want to emphasize that different kernels (different test statistics) have different power against different models. Each test necessarily exploits certain features of the data-generating process at the expense of others, and complementarities between tests can easily arise.

To further illustrate our intuition on  $\gamma$ , we propose a parallel set of DGPs and focus on Bierens' test statistic, which is a  $V$ -statistic associated with a Gaussian kernel with tuning

parameter  $\gamma = 1/2$ :

$$T_{n,v} = \frac{1}{n} \sum_{j,k=1}^n \varepsilon(z_j; \hat{\theta}) \varepsilon(z_k; \hat{\theta}) \exp\left(-\frac{1}{2}\|x_j - x_k\|^2\right)$$

its U-statistic version is

$$T_{n,u} = \frac{1}{n-1} \sum_{j \neq k} \varepsilon(z_j; \hat{\theta}) \varepsilon(z_k; \hat{\theta}) \exp\left(-\frac{1}{2}\|x_j - x_k\|^2\right)$$

The set of parallel DGPs are identical to previous ones (and denoted as DGP(m)\* and DGP-LOCAL(m)\*), except in the following areas,

- When  $m = 2$ ,  $X_1, X_2 \sim N(0, 10)$ , and  $\sigma^{(2)} = (0.1 + X_1^2 + X_2^2)^{1/2}$ .
- When  $m = 5$ ,  $X_j \sim U(0, 10 * j)$  for  $j = 1, 2, 3$ ,  $X_j \sim N(0, (10 * (j - 3))^2)$  for  $j = 4, 5$ .  
 $\sigma^{(5)} = \left(0.1 + \sum_{j=1}^3 X_j + \sum_{j=4}^5 X_j^2\right)^{1/2}$ .
- When  $m = 10$ ,  $X_j \sim U(0, 1 + 0.1 * (j - 1))$  for  $j = 1, \dots, 5$ ,  $X_j \sim N(0, (1 + 0.1 * (j - 5))^2)$  for  $j = 6, \dots, 10$ .  $\sigma^{(10)} = \left(0.1 + \sum_{j=1}^5 X_j + \sum_{j=6}^{10} X_j^2\right)^{1/2}$ .
- When  $m = 20$ ,  $X_j \sim U(0, 1 + 0.1 * (j - 1))$  for  $j = 1, \dots, 10$ ,  $X_j \sim N(0, (1 + 0.1 * (j - 11))^2)$  for  $j = 11, \dots, 15$ ,  $X_j \sim N(1, (1 + 0.1 * (j - 15))^2)$  for  $j = 15, \dots, 20$ .  
 $\sigma^{(20)} = \left(0.1 + \sum_{j=1}^{10} X_j + \sum_{j=11}^{20} X_j^2\right)^{1/2}$ .

In a nutshell, this set of parallel DGPs only differs in the variances of conditional variables such that in low-dimensional cases (i.e., DGP(2)\*, DGP(5)\*, DGP-LOCAL(2)\*, DGP-LOCAL(5)\*), the largest principal value  $\zeta_1 > 2$ , and in high-dimensional cases (i.e., rest of the DGPs), the largest principal value  $\zeta_1 \in (1, 2.25)$ . While in contrast, in the original DGPs, we have  $\zeta_1 \in (1, 2.25)$  for low-dimensional cases and  $\zeta_1 > 2$  for high-dimensional cases.

Tables 3-5 present the results. We also study the performance of a Gaussian kernel with tuning parameter  $\gamma = 1/(2\zeta_1)$  under these parallel DGPs, the results are shown in Table 6. We draw the following remarks:

- Overall,  $V$ -Statistics are inferior to  $U$ -statistics, this is especially true when conditional variables have high variance, i.e., the parallel DGPs. In extreme cases, the  $V$ -Statistic test sizes are completely wrong, losing almost all the power against alternatives. In the contrast,  $U$ -Statistic performs well when the variance of conditional variables matches the tuning parameter (e.g., low-dimensional cases in the original DGPs and high-dimensional cases in the parallel DGPs). We conjecture that this is because the  $V$ -statistic is a biased statistic, and a high dimension increases the bias level.

- Tuning parameter works as conjectured. High-dimensional cases in the parallel DGPs have an accurate size and good power against local alternatives, but these cases perform poorly under the original DGPs. A reverse pattern holds true for low-dimensional cases under parallel and original cases.
- Model complexity (i.e., dimension) is important. Observe that even though in the parallel DGPs, low-dimensional cases are mismatched by its tuning parameter, DGP(2)\* has a more accurate test size compared to DGP(5)\*, and DGP-LOCAL(2)\* has more power than DGP-LOCAL(5)\*. These patterns suggest model complexity plays an important role in a test’s power properties.

Table 3: Bierens’ Test with significant level  $\alpha = 0.1$

$\alpha = 0.1$	U-Statistic			V-Statistic		
	N=100	N=200	N=400	N=100	N=200	N=400
DGP(2)	0.116	0.126	0.105	0.148	0.111	0.124
DGP(5)	0.12	0.11	0.095	0.1	0.081	0.091
DGP(10)	0.127	0.127	0.117	0	0	0
DGP(20)	0.188	0.175	0.141	0.699	0	0
DGP-LOCAL(2)	0.226	0.222	0.2	0.224	0.209	0.212
DGP-LOCAL(5)	0.289	0.312	0.277	0.227	0.264	0.294
DGP-LOCAL(10)	0.378	0.398	0.371	0	0	0
DGP-LOCAL(20)	0.232	0.207	0.199	0.694	0	0

  

$\alpha = 0.1$	U-Statistic			V-Statistic		
	N=100	N=200	N=400	N=100	N=200	N=400
DGP(2)*	0.104	0.11	0.097	0.001	0	0
DGP(5)*	0.092	0.126	0.158	0	0	0
DGP(10)*	0.134	0.116	0.116	0.027	0.01	0.022
DGP(20)*	0.129	0.117	0.113	0.731	0	0
DGP-LOCAL(2)*	0.996	1	1	0.121	0.616	0.958
DGP-LOCAL(5)*	0.41	0.581	0.713	0	0	0
DGP-LOCAL(10)*	0.919	0.954	0.973	0.692	0.795	0.881
DGP-LOCAL(20)*	0.505	0.589	0.611	0.812	0	0

## 7. An Empirical Illustration

In this section, we use the proposed method to examine the validity of instrument variables used in Angrist and Krueger (1991). This influential paper investigates does compulsory school attendance affect schooling and earnings. The authors exploits the variation induced by compulsory school laws in the US, and show that these variations affect a student’s schooling attainment. When investigating how schooling would affect the earnings, the author use quarter of birth (QoB) as instrument variables. They argue that QoB should not affect income directly, nor does it correlate with ability, motivation or family incomes, etc. Furthermore, QoB is correlated with educational attainment via the compulsory school laws. To support such claim, the authors present several tabulations to demonstrate that

Table 4: Bierens' Test with significant level  $\alpha = 0.05$ 

$\alpha = 0.05$	U-Statistic			V-Statistic		
	N=100	N=200	N=400	N=100	N=200	N=400
DGP(2)	0.062	0.065	0.056	0.071	0.049	0.065
DGP(5)	0.068	0.056	0.043	0.034	0.037	0.04
DGP(10)	0.034	0.033	0.036	0	0	0
DGP(20)	0.137	0.085	0.035	0.01	0	0
DGP-LOCAL(2)	0.144	0.142	0.12	0.122	0.121	0.15
DGP-LOCAL(5)	0.196	0.215	0.185	0.113	0.166	0.168
DGP-LOCAL(10)	0.168	0.214	0.224	0	0	0
DGP-LOCAL(20)	0.152	0.083	0.041	0.014	0	0
$\alpha = 0.05$	U-Statistic			V-Statistic		
	N=100	N=200	N=400	N=100	N=200	N=400
DGP(2)*	0.041	0.057	0.044	0	0	0
DGP(5)*	0.002	0.001	0.009	0	0	0
DGP(10)*	0.056	0.058	0.069	0.001	0	0.003
DGP(20)*	0.074	0.056	0.044	0.009	0	0
DGP-LOCAL(2)*	0.988	0.999	1	0.021	0.285	0.855
DGP-LOCAL(5)*	0.016	0.039	0.183	0	0	0
DGP-LOCAL(10)*	0.88	0.922	0.955	0.335	0.583	0.749
DGP-LOCAL(20)*	0.403	0.453	0.492	0.018	0	0

Table 5: Bierens' Test with significant level  $\alpha = 0.01$ 

$\alpha = 0.01$	U-Statistic			V-Statistic		
	N=100	N=200	N=400	N=100	N=200	N=400
DGP(2)	0.011	0.014	0.01	0.012	0.011	0.011
DGP(5)	0.016	0.013	0.007	0.005	0.007	0.006
DGP(10)	0	0	0.006	0	0	0
DGP(20)	0.044	0.018	0.011	0	0	0
DGP-LOCAL(2)	0.048	0.05	0.033	0.028	0.034	0.047
DGP-LOCAL(5)	0.074	0.075	0.076	0.025	0.039	0.047
DGP-LOCAL(10)	0.012	0.036	0.055	0	0	0
DGP-LOCAL(20)	0.042	0.021	0.01	0	0	0
$\alpha = 0.01$	U-Statistic			V-Statistic		
	N=100	N=200	N=400	N=100	N=200	N=400
DGP(2)*	0.003	0.007	0.006	0	0	0
DGP(5)*	0	0	0	0	0	0
DGP(10)*	0.012	0.01	0.012	0	0	0
DGP(20)*	0.011	0.007	0.006	0	0	0
DGP-LOCAL(2)*	0.874	0.984	1	0	0.018	0.428
DGP-LOCAL(5)*	0	0	0.005	0	0	0
DGP-LOCAL(10)*	0.706	0.836	0.881	0.039	0.181	0.417
DGP-LOCAL(20)*	0.193	0.205	0.263	0	0	0

compulsory attendance laws are part of the mechanism generating a relationship between QoB and educational attainment. The exclusion restriction is examined by performing an over-identification test.

Table 6: Gaussian Kernel with  $\gamma = 1/(2\zeta_1)$  Under Parallel DGPs

	$\alpha = 0.1$			$\alpha = 0.05$			$\alpha = 0.01$		
	N=100	N=200	N=400	N=100	N=200	N=400	N=100	N=200	N=400
DGP(2)*	0.141	0.106	0.104	0.067	0.054	0.054	0.022	0.013	0.015
DGP(5)*	0.129	0.094	0.116	0.072	0.044	0.06	0.016	0.015	0.015
DGP(10)*	0.107	0.101	0.098	0.052	0.058	0.052	0.016	0.014	0.009
DGP(20)*	0.089	0.093	0.116	0.062	0.043	0.055	0.019	0.012	0.013
DGP-LOCAL(2)*	1	1	1	1	1	1	1	1	1
DGP-LOCAL(5)*	1	1	1	1	1	1	1	1	1
DGP-LOCAL(10)*	0.979	0.992	0.992	0.965	0.977	0.988	0.903	0.935	0.958
DGP-LOCAL(20)*	0.788	0.871	0.891	0.731	0.802	0.836	0.592	0.651	0.695

However, there are other channels (rather than the compulsory schooling laws) that QoB could correlate with school attainment, and hence undermining the validity of QoB as instruments:

- QoB could affect a student’s performance in school;
- There are differences in the physical and mental health of individuals bore at different times of the year;
- Regional patterns in QoB;
- Redshirting, i.e., parents voluntarily delaying their children’s enrollment.

See [Aliprantis \(2007\)](#); [Bound et al. \(1995\)](#) for detailed discussions.

We examine one specification of [Angrist and Krueger \(1991\)](#). Specifically, column (2) of Table IV in their paper. The 2SLS model investigates how education attainment affects earnings of men who were born in 1920-1929. The model reads as:

$$\begin{aligned}
 \ln(W_i) &= \sum_{c=1}^9 Y_{ic}\xi_c + \rho E_i + \mu_i \\
 E_i &= \sum_{c=1}^9 Y_{ic}\delta_c + \sum_{c=1}^{10} \sum_{j=1}^3 Y_{ic}Q_{ij}\theta_{ij} + e_i
 \end{aligned} \tag{24}$$

where  $W_i, E_i$  are weakly wage and the education of the  $i$ th individual, respectively.  $Q_{ij}$  is a dummy variable indicating whether the individual was born in quarter  $j, j = 1, 2, 3$ , and  $Y_{ic}$  is a dummy variable indicating whether the individual was born in year  $c, c = 1, \dots, 10$ . The coefficient  $\rho$  is the return to education. In this specification, the dimension of exogenous variables is  $d = 39$ . The first row of Table 7 reports  $p$ -values of the specification tests against model (24) using five kernels mentioned in the previous section. All the tests overwhelmingly reject the null that the QoB are valid IVs.

We further investigate a low dimensional specification, where only individuals born in a fixed year is considered. The specification reads as

$$\begin{aligned} \ln(W_i) &= \rho E_i + \mu_i \\ E_i &= \sum_{j=1}^3 Q_{ij} \theta_{ij} + e_i \end{aligned} \tag{25}$$

In this case, the dimension of exogenous variables is  $d = 3$ . Row 2-6 of Table 7 report the results. In most cases (each test and each fixed year is considered as a case), our tests suggest us to reject the null, one exception happens for the case Year 1923+ Shift Variant kernel, where the  $p$ -value is close to the 5% threshold.

Table 7:  $p$ -values under Different Specifications

	Gaussian	IMQ	Gaussian+IMQ	Shift Variant	Local Periodic
FULL	0.584	0.54	0.554	0.522	0.612
1921	0.488	0.466	0.484	0.456	0.494
1923	0.12	0.186	0.174	0.052	0.138
1925	0.684	0.692	0.67	0.704	0.668
1927	0.69	0.71	0.7	0.624	0.7
1929	0.164	0.19	0.194	0.188	0.17

## 8. Conclusion

In this paper, we propose to represent ICM tests in the RKHS. There are several motivations behind this representation. First, conventional ICM tests are based on empirical processes and require integration to obtain the Cramer–Von Mises statistics. This integration is often unable to present closed-form test statistics. Applications of ICM tests are then forced to focus on Birens’s or Escanciano’s ICM tests, which enjoy closed-form presentations. Second, existing literature has well documented that when conditional variables are of high dimensional, ICM tests typically have power-loss issues. Existing dimension-reduction tools rely on projecting the covariates onto a one-dimensional space, and integrate projected statistics from all directions. This procedure leads to a kernel that is hard to compute (its algorithm complexity is  $O(n^3)$ ). Third, although ICM tests are admissible, i.e., there exists no test that is uniformly more powerful than ICM tests, they do not have non-trivial power in all directions. In fact, one ICM test only has substantial power in a finite-dimensional space (Escanciano, 2009). Thus, it is desired to have as many ICM tests as possible.

Once we represent ICM tests in RKHS, we found that (i) after specifying a kernel, the CvM statistics is a closed-form  $U$ -statistic; (ii) a kernel embodies both the dimension and integral measure, and hence, is a valid dimension reduction tool; (iii) with some assumptions, new kernels could be constructed using existing kernels by addition or multiplication.

The main idea behind this representation consists of several steps: (i) the conditional moment restriction is transformed into an unconditional moment restriction with an infinite

number of moment conditions. This transformation is over a function space, and under the null, the supremum (over this function space) value of the squared unconditional moment condition is zero, i.e., the maximum moment restriction; (ii) To obtain a closed form of the maximum moment restriction, one could restrict the function space to a unit ball of an RKHS  $\mathcal{H}(k)$ . The maximum moment restriction corresponds to an integral operator, and Riesz’s representation theorem states that there exists a unique element in the corresponding  $\mathcal{H}(k)$  that represents such an integral operator. We call the such element in  $\mathcal{H}(k)$  the conditional moment embedding, and its  $\mathcal{H}(k)$ -norm is the norm of the integral operation. Under the null, this norm should be zero; (iii) Using the kernel trick, the squared norm has a closed-form presentation, and we can estimate it and further build a test statistic using a U-statistic.

Kernels are essential in this framework, only ISPD kernels could lead to a conditional moment embedding that is injective to the original null hypothesis. Commonly used ISPD kernels are the Gaussian, the Laplacian, the Inverse multiquadric, and the Matern kernels. One could also construct new ISPD kernels from existing ones with additional assumptions imposed.

We further propose a projected kernel to eliminate estimation effects. The advantages of using such a kernel are (i) the limiting null distribution of the test statistic does not depend on how an estimator is obtained; (ii) We do not need to require the estimator to be  $\sqrt{n}$ -asymptotically linear; (iii) The corresponding tests could include certain ‘non-standard’ estimators whose convergence rate is slower than  $1/\sqrt{n}$ . We propose a simple multiplier bootstrap to find the critical value. This is particularly appreciated if the underlying model is non-linear and estimation is time-consuming.

A minimum distance estimator based on conditional moment embedding is developed as a byproduct. This estimator inherits the merits of the corresponding test statistic, e.g., dimension-reduction properties.

Monte Carlo experiments are conducted, and simulation results indicate that the proposed tests have an accurate empirical size and admirably good local power even when the sample size is as small as  $n = 100$  and the dimension is as high as  $d = 20$ . In addition, simulation results also suggest that the proposed tests have good power against high-frequency alternatives. Lastly, a simple empirical application is studied.

## References

- ALIPRANTIS, D. (2007): “A note on why quarter of birth is not a valid instrument for educational attainment,” .
- ANGRIST, J. D. AND A. B. KRUEGER (1991): “Does compulsory school attendance affect schooling and earnings?” *The Quarterly Journal of Economics*, 106, 979–1014.
- ARONSZAJN, N. (1950): “Theory of reproducing kernels,” *Transactions of the American mathematical society*, 68, 337–404.
- BICKEL, P. J., Y. RITOV, AND T. M. STOKER (2006): “Tailor-made tests for goodness of fit to semiparametric hypotheses,” *The Annals of Statistics*, 34, 721–741.
- BIERENS, H. J. (1982): “Consistent model specification tests,” *Journal of Econometrics*, 20, 105–134.
- BIERENS, H. J. AND W. PLOBERGER (1997): “Asymptotic theory of integrated conditional moment tests,” *Econometrica: Journal of the Econometric Society*, 1129–1151.
- BOCHNER, S. (1933): “Monotone funktionen, stieltjessche integrale und harmonische analyse,” *Mathematische Annalen*, 108, 378–410.
- BOUND, J., D. A. JAEGER, AND R. M. BAKER (1995): “Problems with instrumental variables estimation when the correlation between the instruments and the endogenous explanatory variable is weak,” *Journal of the American statistical association*, 90, 443–450.
- CARRASCO, M. AND J.-P. FLORENS (2000): “Generalization of GMM to a continuum of moment conditions,” *Econometric Theory*, 16, 797–834.
- CARRASCO, M., J.-P. FLORENS, AND E. RENAULT (2007): “Linear inverse problems in structural econometrics estimation based on spectral decomposition and regularization,” *Handbook of econometrics*, 6, 5633–5751.
- DELGADO, M. A., M. A. DOMÍNGUEZ, AND P. LAVERGNE (2006): “Consistent tests of conditional moment restrictions,” *Annales d’Économie et de Statistique*, 33–67.
- DELGADO, M. A. AND W. STUTE (2008): “Distribution-free specification tests of conditional models,” *Journal of Econometrics*, 143, 37–55.
- DINCULEANU, N. (2000): *Vector integration and stochastic integration in Banach spaces*, vol. 48, John Wiley & Sons.
- DOMÍNGUEZ, M. A. AND I. N. LOBATO (2004): “Consistent estimation of models defined by conditional moment restrictions,” *Econometrica*, 72, 1601–1615.

- (2015): “A simple omnibus overidentification specification test for time series econometric models,” *Econometric Theory*, 31, 891–910.
- ESCANCIANO, J. C. (2006a): “A consistent diagnostic test for regression models using projections,” *Econometric Theory*, 22, 1030–1051.
- (2006b): “Goodness-of-fit tests for linear and nonlinear time series models,” *Journal of the American Statistical Association*, 101, 531–541.
- (2009): “On the lack of power of omnibus specification tests,” *Econometric Theory*, 25, 162–194.
- ESCANCIANO, J. C. AND S.-C. GOH (2014): “Specification analysis of linear quantile models,” *Journal of Econometrics*, 178, 495–507.
- FAN, Y. AND Q. LI (2000): “Consistent model specification tests: Kernel-based tests versus Bierens’ ICM tests,” *Econometric Theory*, 16, 1016–1041.
- GRETTON, A., D. SEJDINOVIC, H. STRATHMANN, S. BALAKRISHNAN, M. PONTIL, K. FUKUMIZU, AND B. K. SRIPERUMBUDUR (2012): “Optimal kernel choice for large-scale two-sample tests,” *Advances in neural information processing systems*, 25.
- GUO, X. AND L. ZHU (2017): “A review on dimension-reduction based tests for regressions,” *From statistics to mathematical finance: Festschrift in Honour of Winfried Stute*, 105–125.
- HOFMANN, T., B. SCHÖLKOPF, AND A. J. SMOLA (2008): “Kernel methods in machine learning,” *The annals of statistics*, 36, 1171–1220.
- KHMALADZE, E. V. (1982): “Martingale approach in the theory of goodness-of-fit tests,” *Theory of Probability & Its Applications*, 26, 240–257.
- (1993): “Goodness of fit problem and scanning innovation martingales,” *The Annals of Statistics*, 798–829.
- KIMELDORF, G. AND G. WAHBA (1971): “Some results on Tchebycheffian spline functions,” *Journal of mathematical analysis and applications*, 33, 82–95.
- KOUL, H. L. AND W. STUTE (1999): “Nonparametric model checks for time series,” *The Annals of Statistics*, 27, 204–236.
- LAVERGNE, P. AND V. PATILEA (2012): “One for all and all for one: regression checks with many regressors,” *Journal of business & economic statistics*, 30, 41–52.
- MAMMEN, E. (1993): “Bootstrap and wild bootstrap for high dimensional linear models,” *The annals of statistics*, 21, 255–285.

- MINH, H. Q., P. NIYOGI, AND Y. YAO (2006): “Mercer’s theorem, feature maps, and smoothing,” in *International Conference on Computational Learning Theory*, Springer, 154–168.
- MUANDET, K., K. FUKUMIZU, B. SRIPERUMBUDUR, B. SCHÖLKOPF, ET AL. (2017): “Kernel mean embedding of distributions: A review and beyond,” *Foundations and Trends® in Machine Learning*, 10, 1–141.
- MUANDET, K., W. JITKRITTUM, AND J. KÜBLER (2020): “Kernel conditional moment test via maximum moment restriction,” in *Conference on Uncertainty in Artificial Intelligence*, PMLR, 41–50.
- NEWKEY, W. K. (1985): “Generalized method of moments specification testing,” *Journal of econometrics*, 29, 229–256.
- NEWKEY, W. K. AND D. MCFADDEN (1994): “Large sample estimation and hypothesis testing,” *Handbook of econometrics*, 4, 2111–2245.
- NEYMAN, J. (1959): “Optimal tests of composite statistical hypotheses,” *The Harald Cramér Volume*, 213–234.
- PAULSEN, V. I. AND M. RAGHUPATHI (2016): *An introduction to the theory of reproducing kernel Hilbert spaces*, vol. 152, Cambridge university press.
- ROBINSON, P. M. (1991): “Best nonlinear three-stage least squares estimation of certain econometric models,” *Econometrica: Journal of the Econometric Society*, 755–786.
- 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, 212–218.
- RUDIN, W. (2017): *Fourier analysis on groups*, Courier Dover Publications.
- SANT’ANNA, P. H. AND X. SONG (2020): “Specification tests for generalized propensity scores using double projections,” *arXiv preprint arXiv:2003.13803*.
- SANT’ANNA, P. H. AND X. SONG (2019): “Specification tests for the propensity score,” *Journal of Econometrics*, 210, 379–404.
- SERFLING, R. J. (1980): *Approximation theorems of mathematical statistics*, John Wiley & Sons.
- SMOLA, A. J., B. SCHÖLKOPF, AND K.-R. MÜLLER (1998): “The connection between regularization operators and support vector kernels,” *Neural networks*, 11, 637–649.

- SRIPERUMBUDUR, B. K., A. GRETTON, K. FUKUMIZU, B. SCHÖLKOPF, AND G. R. LANCKRIET (2010): “Hilbert space embeddings and metrics on probability measures,” *The Journal of Machine Learning Research*, 11, 1517–1561.
- STEINWART, I. (2001): “On the influence of the kernel on the consistency of support vector machines,” *Journal of machine learning research*, 2, 67–93.
- STEINWART, I. AND A. CHRISTMANN (2008): *Support vector machines*, Springer Science & Business Media.
- TOLSTIKHIN, I., B. K. SRIPERUMBUDUR, AND K. MUANDET (2017): “Minimax estimation of kernel mean embeddings,” *The Journal of Machine Learning Research*, 18, 3002–3048.
- WENDLAND, H. (2004): *Scattered data approximation*, vol. 17, Cambridge university press.
- ZHANG, R., M. IMAIZUMI, B. SCHÖLKOPF, AND K. MUANDET (2020): “Maximum moment restriction for instrumental variable regression,” *arXiv preprint arXiv:2010.07684*.

## A. Backgrounds on the Reproducing Kernel Hilbert Space

Reproducing Kernel Hilbert Space (RKHS), first proposed in Aronszajn (1950), is a special Hilbert space with some properties. It is a Hilbert space of functions with reproducing kernels. Formally, RKHS is defined as

**Definition 14** *A Reproducing Kernel Hilbert Space is a Hilbert space  $\mathcal{H}(k)$  of functions  $f : \mathcal{X} \rightarrow \mathbb{R}$  with a reproducing kernel  $k : \mathcal{X}^2 \rightarrow \mathbb{R}$ , where (i)  $k(x, \cdot) \in \mathcal{H}(k)$ , and (ii)  $\langle f, k(x, \cdot) \rangle_{\mathcal{H}(k)} = f(x)$*

**Remark.** Property (ii) is called the reproducing property of  $\mathcal{H}(k)$ . By Aronszajn (1950), every positive definite kernel  $k$  uniquely determines the RKHS for which  $k$  is a reproducing kernel.

To gain an understanding of the kernel  $k$ , some backgrounds are needed.

**Definition 15** *Given a kernel  $k : \mathcal{X}^2 \rightarrow \mathbb{R}$  and inputs  $x_1, \dots, x_n \in \mathcal{X}$ , the  $n \times n$  matrix*

$$K_{ij} = k(x_i, x_j), \forall i, j \in \{1, \dots, n\}$$

*is called the Gram Matrix (also known as the Kernel Matrix) of  $k$  with respect to  $x_1, \dots, x_n$ .*

**Definition 16** *A real  $n \times n$  symmetric matrix  $K_{ij}$  satisfying*

$$\sum_{i,j=1}^n c_i c_j K_{ij} \geq 0$$

*for all  $c_i \in \mathbb{R}$  is called positive definite. If equality in the above equation only occurs  $c_1 = \dots = c_n = 0$ , then we shall call the matrix strictly positive definite.*

**Definition 17** *Let  $\mathcal{X}$  be a nonempty set. A function  $k : \mathcal{X}^2 \rightarrow \mathbb{R}$  which for all  $n \in \mathbb{N}$ ,  $x_i \in \mathcal{X}$  give rise to a positive definite Gram matrix is called a positive definite kernel. Similarly, a function  $k : \mathcal{X}^2 \rightarrow \mathbb{R}$  which for all  $n \in \mathbb{N}$  and distinct  $x_i \in \mathcal{X}$  gives rise to a strictly positive definite Gram matrix is called a strictly positive definite kernel.*

The RKHS is better explained in the following way. Define a map from  $\mathcal{X}$  into the space of functions mapping  $\mathcal{X}$  to  $\mathcal{H}(k)$ , via

$$\begin{aligned} \phi : \mathcal{X} &\rightarrow \mathcal{H}(k) \\ x &\rightarrow k(x, \cdot) \end{aligned}$$

Here,  $\phi_x(\cdot) = k(x, \cdot)$  denotes the function that assigns the value  $k(x, x')$  to  $x' \in \mathcal{X}$ .

We next construct a dot product space containing the images of the inputs under  $\phi$ . To this end, considering the kernel function  $k(x, x')$ , suppose for  $n$  points, we fix one of the variables to have  $k(x_1, x'), \dots, k(x_n, x')$ . There are all functions of the variable  $x'$ .

$$\mathcal{H}(k) = \left\{ f(\cdot) = \sum_{i=1}^n \alpha_i k(x_i, \cdot) \right\} \quad (26)$$

Here,  $n \in \mathbb{N}$ ,  $\alpha_i \in \mathbb{R}$  and  $x_i \in \mathcal{X}$  are arbitrary. RKHS is a function space that is the set of all possible linear combinations of these functions (Kimeldorf and Wahba, 1971). This equation shows that the bases of an RKHS are kernels, hence every function in the RKHS can be written as a linear combination.

Consider two functions in this space represented as  $f = \sum_{i=1}^n \alpha_i k(x_i, \cdot)$ , and  $g = \sum_{j=1}^n \beta_j k(x_j, \cdot)$ , the inner product in RKHS is defined as

$$\begin{aligned} \langle f, g \rangle_{\mathcal{H}(k)} &= \left\langle \sum_{i=1}^n \alpha_i k(x_i, \cdot), \sum_{j=1}^n \beta_j k(x_j, \cdot) \right\rangle_{\mathcal{H}(k)} \\ &= \sum_{i,j=1}^n \alpha_i \beta_j k(x_i, x_j) \end{aligned}$$

The feature map  $\phi_x(\cdot)$  is a (possibly infinite-dimensional) vector whose elements are

$$\phi_x(\cdot) = \left( \sqrt{\lambda_1} \psi_1(x), \sqrt{\lambda_2} \psi_2(x), \dots \right)^\top = \Phi(x)$$

where  $\{\lambda_j; j = 1, 2, \dots\}$  and  $\{\psi_j; j = 1, 2, \dots\}$  are eigenvalues and eigenfunctions of the following eigen problem:

$$\int k(x, x') \psi_j(x') dx' = \lambda_j \psi_j(x)$$

See Minh et al. (2006) for details. One can understand the kernel as a similarity measurement since the kernel can be expressed as an inner product, which is a measure of similarity in terms of angles of vectors:

$$\begin{aligned} k(x, x') &= \langle \phi_x(\cdot), \phi_{x'}(\cdot) \rangle_{\mathcal{H}(k)} \\ &= \Phi(x)^\top \Phi(x') \end{aligned}$$

Hence, the relative similarity of inputs is known by the kernel. However, in most of the kernels, we cannot find an explicit expression for the feature map. Therefore, the exact location of inputs to RKHS is not necessarily known but the relative similarity, which is the kernel, is known.

## B. Auxiliary Lemmas

**Lemma 18** *Suppose Assumption A5 holds, we have*

$$\left\| \frac{1}{n} \sum_{s=1}^n g(z_s; \hat{\theta}) k(x_s, x_j) - \mathbb{E}_{XZ} (g(Z; \theta_0) k(X, x_j)) \right\| = O_p(1/\sqrt{n}) + O_p(\|\hat{\theta} - \theta_0\|)$$

**Proof**

$$\begin{aligned} \frac{1}{n} \sum_{s=1}^n g(z_s; \hat{\theta}) k(x_s, x_j) &= \frac{1}{n} \sum_{s=1}^n g(z_s; \theta_0) k(x_s, x_j) + \frac{1}{n} \sum_{s=1}^n G(z_s; \bar{\theta}) k(x_s, x_j) (\hat{\theta} - \theta_0) \\ &= I_{1,n} + I_{2,n}(\hat{\theta} - \theta_0) \end{aligned}$$

Observe that

$$\begin{aligned}\|I_{1,n} - \mathbb{E}g(Z; \theta_0)k(X, x_j)\| &= O_p(1/\sqrt{n}) \\ \|I_{2,n} - \mathbb{E}G(Z; \theta_0)k(X, x_j)\| &= O_p(1/\sqrt{n})\end{aligned}$$

and

$$\begin{aligned}I_{2,n}(\hat{\theta} - \theta_0) &= \mathbb{E}G(Z; \theta_0)k(X, x_j)O_p(\|\hat{\theta} - \theta_0\|) + O_p(1/\sqrt{n}\|\hat{\theta} - \theta_0\|) \\ &= O_p(\|\hat{\theta} - \theta_0\|)\end{aligned}$$

■

**Lemma 19** *Suppose Assumption 5 holds, we have*

$$\begin{aligned}\left\| \frac{1}{n(n-1)} \sum_{s \neq k} g(z_s; \hat{\theta})k(x_s, x_k)g^\top(z_k; \hat{\theta}) - \mathbb{E} \left( g(Z; \theta)k(X, X')g^\top(Z', \theta_0) \right) \right\| &= O_p(1/\sqrt{n}) \\ &+ O_p(\|\hat{\theta} - \theta_0\|)\end{aligned}$$

**Proof**

$$\begin{aligned}\frac{1}{n(n-1)} \sum_{s \neq k} g(z_s; \hat{\theta})k(x_s, x_k)g^\top(z_k; \hat{\theta}) &= \frac{1}{n(n-1)} \sum_{s \neq k} g(z_s; \theta_0)k(x_s, x_k)g^\top(z_k; \theta_0) \\ &+ \frac{1}{n(n-1)} \sum_{s \neq k} g(z_s; \theta_0)k(x_s, x_k) \left( \nabla_{\theta} g(z_k; \bar{\theta})(\hat{\theta} - \theta_0) \right)^\top \\ &+ \frac{1}{n(n-1)} \sum_{s \neq k} g(z_k; \theta_0)k(x_s, x_k) \left( \nabla_{\theta} g(z_s; \bar{\theta})(\hat{\theta} - \theta_0) \right)^\top \\ &+ \frac{1}{n(n-1)} \sum_{s \neq k} \nabla_{\theta} g(z_s; \bar{\theta})(\hat{\theta} - \theta_0)k(x_s, x_k) \left( \nabla_{\theta} g(z_k; \bar{\theta})(\hat{\theta} - \theta_0) \right)^\top \\ &= I_{1,n} + I_{21,n} + I_{22,n} + I_{3,n}\end{aligned}$$

It is clear that

$$\begin{aligned}I_{1,n} &= \mathbb{E} \left( g(Z; \theta)k(X, X')g^\top(Z', \theta_0) \right) + O_p(1/\sqrt{n}) \\ I_{21,n} = I_{22,n} &= \mathbb{E} \left( \nabla_{\theta} g(Z; \theta_0)(\hat{\theta} - \theta_0)k(X, X')g^\top(Z'; \theta_0) \right) + O_p(1/\sqrt{n}) = O_p(\|\hat{\theta} - \theta_0\|) + O_p(1/\sqrt{n})\end{aligned}$$

$$\begin{aligned}I_{3,n} &= \mathbb{E} \left( \nabla_{\theta} g(Z; \theta_0)(\hat{\theta} - \theta_0)k(X, X') \left( \nabla_{\theta} g(Z'; \theta_0)(\hat{\theta} - \theta_0) \right)^\top \right) + O_p(1/\sqrt{n}) \\ &= O_p(\|\hat{\theta} - \theta_0\|^2) + O_p(1/\sqrt{n})\end{aligned}$$

Putting all pieces together, we yield what is asserted. ■

**Lemma 20** *Suppose Assumption 5 holds, we have*

$$\|\Gamma_{n,\hat{\theta}}^{-1} - \Gamma_{\theta_0}^{-1}\| = O_p(1/\sqrt{n}) + O_p(\|\hat{\theta} - \theta_0\|)$$

**Proof** Observe that

$$\begin{aligned} \Gamma_{n,\hat{\theta}} - \Gamma_{\theta_0} &= \frac{1}{n} \sum_{s=1}^n g(z_s; \theta_0) g^\top(z_s; \theta_0) - \mathbb{E}g(Z; \theta_0) g^\top(Z; \theta_0) \\ &= \frac{2}{n} \sum_{s=1}^n (g(z_s; \hat{\theta}) - g(z_s; \theta_0)) g^\top(z_s; \theta_0) \\ &= \frac{1}{n} \sum_{s=1}^n (g(z_s; \hat{\theta}) - g(z_s; \theta_0))(g(z_s; \hat{\theta}) - g(z_s; \theta_0))^\top \\ &= O_p(1/\sqrt{n}) + O_p(\|\hat{\theta} - \theta_0\|) + O_p(\|\hat{\theta} - \theta_0\|^2) \\ &= O_p(1/\sqrt{n}) + O_p(\|\hat{\theta} - \theta_0\|) \end{aligned}$$

By the continuous mapping theorem, the above fact yields:

$$\Gamma_{n,\hat{\theta}}^{-1} - \Gamma_{\theta_0}^{-1} = o_p(1)$$

Furthermore, we have the decomposition:

$$\Gamma_{n,\hat{\theta}}^{-1} - \Gamma_{\theta_0}^{-1} = -\Gamma_{n,\hat{\theta}}^{-1} (\Gamma_{n,\hat{\theta}} - \Gamma_{\theta_0}) \Gamma_{\theta_0}^{-1}$$

Putting everything together, we have the desired result. ■

## C. Verifications and Proofs

### Verify the Riesz Represor of (5)

**Proof** Note that

$$\begin{aligned} \mathcal{C}_\theta h &= \int \varepsilon(z; \theta) h(x) dP_{XZ}(x, z) \\ &= \int \varepsilon(z; \theta) \langle h, \phi_x(\cdot) \rangle_{\mathcal{H}(k)} dP_{XZ}(x, z) \\ &= \int \langle h, \varepsilon(z; \theta) \phi_x(\cdot) \rangle_{\mathcal{H}(k)} dP_{XZ}(x, z) \\ &= \langle h, \mathbb{E}_{XZ}(\varepsilon(Z; \theta) \phi_X(\cdot)) \rangle_{\mathcal{H}(k)} \\ &= \langle h, \boldsymbol{\mu}_\theta \rangle_{\mathcal{H}(k)} \end{aligned}$$

To use Riesz's theorem, we need to show that the operator is bounded. Let  $\|\mathcal{C}_\theta\|$  be the operator norm,

$$\|\mathcal{C}_\theta\| = \sup_{\|h\|_{\mathcal{H}(k)} \leq 1} \mathcal{C}_\theta h = \sup_{\|h\|_{\mathcal{H}(k)} \leq 1} \langle h, \boldsymbol{\mu}_\theta \rangle_{\mathcal{H}(k)} = \left\langle \frac{\boldsymbol{\mu}_\theta}{\|\boldsymbol{\mu}_\theta\|_{\mathcal{H}(k)}}, \boldsymbol{\mu}_\theta \right\rangle_{\mathcal{H}(k)} = \|\boldsymbol{\mu}_\theta\|_{\mathcal{H}(k)}$$

with

$$\|\boldsymbol{\mu}_\theta\|_{\mathcal{H}(k)}^2 = \mathbb{E}(\varepsilon(Z; \theta)k(X, X')\varepsilon(Z'; \theta)) < \infty$$

by Assumption (A4). Here,  $(X', Z')$  is an independent copy of  $(X, Z)$ . Furthermore, by Assumption (A2),

$$|\mathcal{C}_\theta h| \leq \|h\|_{\mathcal{H}(k)} \|\mathcal{C}_\theta\|_{\mathcal{H}(k)} < \infty$$

Thus, the  $\mathcal{C}_\theta$  is a bounded linear operator. By Riesz's representation theorem,  $\boldsymbol{\mu}_\theta$  is the unique representer of  $\mathcal{C}_\theta$  in  $\mathcal{H}(k)$ .  $\blacksquare$

### Verify the Reproducing Property of (11)

**Proof**

$$\begin{aligned} \langle \mathcal{P}\phi_x(\cdot), \phi'_x(\cdot) \rangle_{\mathcal{H}(k)} &= k(x, x') - \left\langle g^\top(z; \theta_0)\Gamma_{\theta_0}^{-1}\mathbb{E}_{XZ}(g(Z; \theta_0)\phi_X(\cdot)), \phi'_x(\cdot) \right\rangle_{\mathcal{H}(k)} \\ &= k(x, x') - g^\top(z; \theta_0)\Gamma_{\theta_0}^{-1}\mathbb{E}_{XZ}\left(g(Z; \theta_0)\langle \phi_X(\cdot), \phi'_x(\cdot) \rangle_{\mathcal{H}(k)}\right) \\ &= k(x, x') - g^\top(z; \theta_0)\Gamma_{\theta_0}^{-1}\mathbb{E}_{XZ}(g(Z; \theta_0)k(X, x')) \end{aligned}$$

$\blacksquare$

### Proof of Lemma (6)

**Proof**

$$\begin{aligned} \hat{k}_p(x_i, x_j) &= k(x_i, x_j) - g^\top(z_i; \hat{\theta})\Gamma_{n, \hat{\theta}}^{-1}\left(\frac{1}{n}\sum_{s=1}^n g(z_s; \hat{\theta})k(x_s, x_j)\right) \\ &\quad - g^\top(z_j; \hat{\theta})\Gamma_{n, \hat{\theta}}^{-1}\left(\frac{1}{n}\sum_{s=1}^n g(z_s; \hat{\theta})k(x_s, x_i)\right) \\ &\quad + g^\top(z_i; \hat{\theta})\Gamma_{n, \hat{\theta}}^{-1}\left(\frac{1}{n(n-1)}\sum_{s \neq k}^n g(z_s; \hat{\theta})k(x_s, x_k)g^\top(z_k; \hat{\theta})\right)\Gamma_{n, \hat{\theta}}^{-1}g(z_j; \hat{\theta}) \end{aligned}$$

where

$$\Gamma_{n,\hat{\theta}} = \frac{1}{n} \sum_{s=1}^n g(z_s; \hat{\theta}) g^\top(z_s; \hat{\theta})$$

Hence,

$$\hat{k}_p(x_i, x_j) = k(x_i, x_j) - g^\top(z_i; \hat{\theta}) \Gamma_{n,\hat{\theta}}^{-1} I_{11,n} - g^\top(z_j; \hat{\theta}) \Gamma_{n,\hat{\theta}}^{-1} I_{12,n} + g^\top(z_i; \hat{\theta}) \Gamma_{n,\hat{\theta}}^{-1} I_{2,n} \Gamma_{n,\hat{\theta}}^{-1} g(z_j; \hat{\theta})$$

By Lemmas 18, 19 and 20, we have

- $g^\top(z_i; \hat{\theta}) = g(z_i; \theta_0) + O_p(\|\hat{\theta} - \theta_0\|)$
- $\Gamma_{n,\hat{\theta}}^{-1} = \Gamma_{\theta_0}^{-1} + O_p(1/\sqrt{n}) + O_p(\|\hat{\theta} - \theta_0\|)$
- $I_{11,n} = \mathbb{E}g(Z; \theta_0)k(X, x_j) + O_p(1/\sqrt{n}) + O_p(\|\hat{\theta} - \theta_0\|)$
- $I_{12,n} = \mathbb{E}g(Z; \theta_0)k(X, x_i) + O_p(1/\sqrt{n}) + O_p(\|\hat{\theta} - \theta_0\|)$
- $I_{2,n} = \mathbb{E}(g(Z; \theta)k(X, X')g^\top(Z', \theta_0)) + O_p(1/\sqrt{n}) + O_p(\|\hat{\theta} - \theta_0\|)$

Putting all pieces together, we have

$$\hat{k}_p(x_i, x_j) = k_p(x_i, x_j) + O_p(1/\sqrt{n}) + O_p(\|\hat{\theta} - \theta_0\|)$$

Now, we are ready to show the expanding result.

$$\begin{aligned} n\widehat{\mathbb{M}}_p^2(\hat{\theta}) &= \frac{n}{n(n-1)} \sum_{i \neq j} \varepsilon(z_i; \theta_0) k_p(x_i, x_j) \varepsilon(z_j, \theta_0) \\ &\quad + \frac{2n}{n(n-1)} \sum_{i \neq j} \varepsilon(z_i; \theta_0) k_p(x_i, x_j) g^\top(z_j; \bar{\theta}) (\hat{\theta} - \theta_0) \\ &\quad + \sqrt{n} (\hat{\theta} - \theta_0)^\top \frac{1}{n(n-1)} \sum_{i \neq j} g(z_i; \bar{\theta}) k_p(x_i, x_j) g^\top(z_j; \bar{\theta}) \sqrt{n} (\hat{\theta} - \theta_0) + O_p(1/\sqrt{n}) \\ &= nA_{1,n}^{(p)} + 2nA_{2,n}^{(p)} (\hat{\theta} - \theta_0) + \sqrt{n} (\hat{\theta} - \theta_0)^\top A_{3,n}^{(p)} \sqrt{n} (\hat{\theta} - \theta_0) \\ &\quad + O_p(1/\sqrt{n}) + O_p(\|\hat{\theta} - \theta_0\|) \end{aligned}$$

where  $\bar{\theta} = \gamma \hat{\theta} + (1 - \gamma) \theta_0$ ,  $\gamma \in (0, 1)$ , and the last term  $(O_p(1/\sqrt{n}) + O_p(\|\hat{\theta} - \theta_0\|))$  comes from the fact that  $1/(n(n-1)) \sum_{i \neq j} \varepsilon(z_i; \theta_0) \varepsilon(z_j, \theta_0)$  is a degenerate U-statistic, hence,

$$(O_p(1/\sqrt{n}) + O_p(\|\hat{\theta} - \theta_0\|)) \frac{1}{n-1} \sum_{i \neq j} \varepsilon(z_i; \theta_0) \varepsilon(z_j, \theta_0) = O_p(1/\sqrt{n}) + O_p(\|\hat{\theta} - \theta_0\|)$$

One can easily check that  $A_{1,n}^{(p)}$ ,  $A_{2,n}^{(p)}$  and  $A_{3,n}^{(p)}$  are degenerate U-statistic, and hence  $A_{1,n}^{(p)}$ ,  $A_{2,n}^{(p)}$ ,  $A_{3,n}^{(p)} = O_p(1/n)$ . Thus, we have

$$\begin{aligned} n\widehat{\mathbb{M}}_p^2(\hat{\theta}) &= \frac{n}{n(n-1)} \sum_{i \neq j} \varepsilon(z_i; \theta_0) k_p(x_i, x_j) \varepsilon(z_j, \theta_0) + O_p(\|\hat{\theta} - \theta_0\|) + O_p(\|\hat{\theta} - \theta_0\|^2) \\ &\quad + O_p(1/\sqrt{n}) \end{aligned}$$

**Proof of Theorem 5.**

**Proof** Suppose that  $\mu_{\theta_1}^{(p)} = \mu_{\theta_2}^{(p)}$  and let  $\delta(x) = \mathcal{E}(x; \theta_1) - \mathcal{E}(x; \theta_2)$ . Then we have

$$\begin{aligned}
\|\mu_{\theta_1}^{(p)} - \mu_{\theta_2}^{(p)}\|_{\mathcal{H}(k_p)}^2 &= \left\| \int \xi_{\theta_1}^{(p)}(x, z) dP_{XZ}(x, z) - \int \xi_{\theta_2}^{(p)}(x, z) dP_{XZ}(x, z) \right\|_{\mathcal{H}(k_p)}^2 \\
&= \left\| \int \mathcal{E}(x; \theta_1) \phi_x^{(p)}(\cdot) dP_X(x) - \int \mathcal{E}(x; \theta_2) \phi_x^{(p)}(\cdot) dP_X(x) \right\|_{\mathcal{H}(k_p)}^2 \\
&= \left\| \int (\mathcal{E}(x; \theta_1) - \mathcal{E}(x; \theta_2)) \phi_x^{(p)}(\cdot) dP_X(x) \right\|_{\mathcal{H}(k_p)}^2 \\
&= \int \int \delta(x) k_p(x, x') \delta(x') dP_X(x) dP_{X'}(x') = 0
\end{aligned}$$

where  $X'$  is an independent copy of  $X$ . Since  $k_p(\cdot, \cdot)$  is ISPD kernel and the assumption that  $\delta(x)$  is not colinear with  $g(Z; \theta_0)$ , it follows that the function  $\varphi(x) = \delta(x)p_X(x)$  has zero L2-norm, i.e.,  $\|\varphi\|_2^2 = 0$  where  $p_X$  denotes the density of  $P_X$ . As a result,  $\delta(x) = 0$  a.s  $P_X$  implying that  $P_X(B_0) = 1$  where  $B_0 = \{x \in \mathcal{X} : \mathcal{E}(x; \theta_1) - \mathcal{E}(x; \theta_2) = 0\}$ . Therefore,  $\mathcal{E}(x; \theta_1) = \mathcal{E}(x; \theta_2)$  for  $P_X$  almost surely.  $\blacksquare$

**Proof of Theorem 9**

**Proof**

Under local alternatives,  $\varepsilon(z; \hat{\theta}) = \varepsilon(z; \theta_0) + (\hat{\theta} - \theta_0)^\top g(z; \bar{\theta}) + R(x)/\sqrt{n}$ . Hence,

$$\begin{aligned}
n\widehat{\mathbb{M}}_n^2(\hat{\theta}) &= \frac{n}{n(n-1)} \sum_{1 \leq i \neq j \leq n} \varepsilon(z_i; \theta_0) \hat{k}_p(x_i, x_j) \varepsilon(z_j; \theta_0) \\
&+ \frac{2n}{n(n-1)} \sum_{1 \leq i \neq j \leq n} \varepsilon(z_i; \theta_0) \hat{k}_p(x_i, x_j) g^\top(z_j; \bar{\theta}) (\hat{\theta} - \theta_0) \\
&+ \sqrt{n} (\hat{\theta} - \theta_0)^\top \left( \frac{1}{n(n-1)} \sum_{1 \leq i \neq j \leq n} g(z_i; \bar{\theta}) \hat{k}_p(x_i, x_j) g^\top(z_j; \bar{\theta}) \right) \sqrt{n} (\hat{\theta} - \theta_0) \\
&+ \frac{2n}{n(n-1)} \sum_{1 \leq i \neq j \leq n} \varepsilon(z_i; \theta_0) \hat{k}_p(x_i, x_j) \frac{R(x_j)}{\sqrt{n}} \\
&+ \frac{2n}{n(n-1)} \sum_{1 \leq i \neq j \leq n} (\hat{\theta} - \theta_0)^\top g(z_i; \bar{\theta}) \hat{k}_p(x_i, x_j) \frac{R(x_j)}{\sqrt{n}} \\
&+ \frac{n}{n(n-1)} \sum_{1 \leq i \neq j \leq n} \frac{R(x_i)}{\sqrt{n}} \hat{k}_p(x_i, x_j) \frac{R(x_j)}{\sqrt{n}} \\
&= nA_{1,n} + 2nA_{2,n}(\hat{\theta} - \theta_0) + \sqrt{n}(\hat{\theta} - \theta_0)^\top A_{3,n} \sqrt{n}(\hat{\theta} - \theta_0) + 2\sqrt{n}A_{4,n} \\
&+ 2\sqrt{n}(\hat{\theta} - \theta_0)^\top A_{5,n} + A_{6,n}
\end{aligned}$$

Note that

$$nA_{1,n} \xrightarrow{d} \sum_{k=1}^{\infty} \tau_k^{(p)} (W_k^2 - 1)$$

by Theorem 7.  $A_{2,n}$  is a degenerate  $U$ -statistic by the orthogonality argument, hence  $2nA_{2,n}(\hat{\theta} - \theta_0) = O_p(\|\hat{\theta} - \theta_0\|) = o_p(1)$ .

Furthermore,

$$A_{3,n} = O_p(1/n)$$

,

$$A_{l,n} = O_p(1/\sqrt{n}), \quad l = 4, 5$$

by the orthogonality between  $k_p(\cdot, \cdot)$  and  $g(z; \theta)$  and the fact that  $\hat{k}_p(x, x') = k_p(x, x') + O_p(1/\sqrt{n}) + O_p(\|\hat{\theta} - \theta_0\|)$ . Thus, by [Serfling \(1980\)](#)

$$\begin{aligned} & \sqrt{n}(\hat{\theta} - \theta_0)^\top A_{3,n} \sqrt{n}(\hat{\theta} - \theta_0) + 2\sqrt{n}A_{4,n} + 2\sqrt{n}(\hat{\theta} - \theta_0)^\top A_{5,n} \\ & \xrightarrow{d} 2N(0, 4\text{Var}_{X,Z}(\mathbb{E}_{X',Z'}\varepsilon(Z; \theta_0)k_p(X, X')R(X'))) \end{aligned}$$

Lastly,

$$A_{6,n} = \mathbb{E}(R(X)k_p(X, X')R(X')) + O_p(1/\sqrt{n})$$

Putting these pieces together, we yield what is asserted. ■

### Proof of Lemma 10

**Proof**

$$\begin{aligned} \mathbb{M}_p^2(\theta_0) &= \mathbb{E} \left( \varepsilon(Z; \theta_0)k_p(X, X')\varepsilon(Z'; \theta_0) + \frac{R(X)R(X')}{n}k_p(X, X') + 2\varepsilon(Z; \theta_0)\frac{R(X')}{\sqrt{n}}k_p(X, X') \right) \\ &= \sum_{j \geq 1} \lambda_j (\mathbb{E}\varepsilon(Z; \theta_0)e_j(X))^2 + \sum_{j \geq 1} \lambda_j \left( \mathbb{E}\frac{R(X)}{\sqrt{n}}e_j(X) \right)^2 \\ &\quad + 2 \sum_{j \geq 1} \lambda_j \mathbb{E}\varepsilon(Z; \theta_0)e_j(X) \mathbb{E} \left( \frac{R(X)}{\sqrt{n}}e_j(X) \right) \end{aligned}$$

Recall  $R(x) = \sum_{s \geq 1} \alpha_s e_s(x)$ , we can further conclude that

$$\begin{aligned} \mathbb{M}_p^2(\theta_0) &= \sum_{j \geq 1} \lambda_j (\mathbb{E}\varepsilon(Z; \theta_0)e_j(X))^2 + \lambda_j \frac{\alpha_j^2}{n} \mathbb{E}e_j^2(X) + 2\lambda_j \mathbb{E}\varepsilon(Z; \theta_0)e_j(X) \frac{\alpha_j}{\sqrt{n}} \mathbb{E}e_j^2(X) \\ &= \mathbb{E}(\varepsilon(Z; \theta_0)k_p(X, X')\varepsilon(Z'; \theta_0)) + \lambda_j \frac{\alpha_j^2}{n} + 2\lambda_j \mathbb{E}(\varepsilon(Z; \theta_0)e_j(X)) \frac{\alpha_j}{\sqrt{n}} \end{aligned}$$

■

**Proof of Theorem 11.**

**Proof** To streamline the presentation, let

$$\hat{f}_\theta(u_i, u_j) = \varepsilon(z_i; \theta) v_i \hat{k}_p(x_i, x_j) \varepsilon(z_j; \theta) v_j$$

$$f_\theta(u_i, u_j) = \varepsilon(z_i; \theta) v_i k_p(x_i, x_j) \varepsilon(z_j; \theta) v_j$$

where  $u_i = (v_i, x_i, z_i)$ . Furthermore, let  $(x, z)^{(n)} = \{(x_i, z_i); i = 1, \dots, n\}$ ,  $\hat{\theta} - \theta^* = O_p(1/\sqrt{n})$  under different hypotheses, and

$$k_p(x, x') = k(x, x') - g^\top(z; \theta^*) \Gamma_{\theta^*}^{-1} \mathbb{E}_{(X, Z)}(g(Z; \theta^*) k(X, x'))$$

By the fact that  $\hat{k}_p(\cdot, \cdot) = k_p(\cdot, \cdot) + O_p(\|\hat{\theta} - \theta^*\|) + O_p(1/\sqrt{n})$ , we have

$$\begin{aligned} n\widehat{\mathbb{M}}_p^{2,*}(\hat{\theta}) &= \frac{n}{n(n-1)} \sum_{i \neq j} \hat{f}_\theta(u_i, u_j) \\ &= \frac{n}{n(n-1)} \sum_{i \neq j} f_{\theta_0}(u_i, u_j) + o_p(1) \\ &= n\mathbb{M}_p^{2,*}(\theta_0) + o_p(1) \end{aligned}$$

Let

$$T_n = \frac{1}{n} \sum_{i \neq j} f_\theta(u_i, u_j)$$

we have  $n\widehat{\mathbb{M}}_p^{2,*}(\theta_0) = \frac{n}{n-1} T_n$ , the goal is to show that

$$T_n \xrightarrow{d,*} Y = \sum_{k=1}^{\infty} \tau_k^{(p)} (W_k^2 - 1)$$

We shall carry this out by the method of characteristic functions, i.e., to show that

$$\mathbb{E}(e^{i\omega T_n} | (x, z)^{(n)}) \rightarrow \mathbb{E}(e^{i\omega Y}), \quad n \rightarrow \infty, \forall \omega$$

Denote  $\{\rho_k(\cdot)\}$  as the orthonormal eigenfunctions corresponding to the eigenvalues  $\{\tau_k^{(p)}\}$  defined in connection with  $\varepsilon(z; \theta_0) k_p(x, x') \varepsilon(z'; \theta_0)$ . Thus,

$$f_{\theta_0}(u_1, u_2) = \sum_{k \geq 1} \tau_k^{(p)} v_1 v_2 \rho_k(y_1) \rho_k(y_2)$$

with  $y_1 = (x_1, z_1)$

Thus,  $T_n$  might be expressed as

$$T_n = \frac{1}{n} \sum_{i \neq j} \sum_{k \geq 1} \tau_k^{(p)} v_i v_j \rho_k(y_i) \rho_k(y_j)$$

Now put

$$T_{nK} = \frac{1}{n} \sum_{i \neq j} \sum_{k=1}^K \tau_k^{(p)} v_i v_j \rho_k(y_i) \rho_k(y_j)$$

Using the equality  $|e^{iz} - 1| < |z|$ , we have for an arbitrary  $\delta > 0$ , there exists a  $K$  such that

$$\begin{aligned} |\mathbb{E}(e^{i\omega T_n} | (x, z)^{(n)}) - \mathbb{E}(e^{i\omega T_{nK}} | (x, z)^{(n)})| &\leq \mathbb{E}(|e^{i\omega T_n} - e^{i\omega T_{nK}}| | (x, z)^{(n)}) \\ &\leq |\omega| \mathbb{E}(|T_n - T_{nK}| | (x, z)^{(n)}) \\ &\leq |\omega| \left( \mathbb{E}(T_n - T_{nK})^2 | (x, z)^{(n)} \right)^{1/2} \end{aligned}$$

Observe that  $T_n - T_{nK}$  is in the form of a  $U$ -statistic, that is,

$$T_n - T_{nK} = \frac{2}{n} \binom{n}{2} U_{nK}$$

where

$$U_{nK} = \binom{n}{2}^{-1} \sum_{i \neq j} g_K(u_i, u_j)$$

with

$$g_K(u_1, u_2) = \sum_{k=K+1}^{\infty} \tau_k^{(p)} \rho_k(y_1) \rho_k(y_2) v_1 v_2$$

Note that

$$\begin{aligned} \mathbb{E} \left( U_{nK}^2 | (x, z)^{(n)} \right) &= \left[ \sum_{k=K+1}^{\infty} \left( \binom{n}{2}^{-1} \sum_{i \neq j} \tau_k^{(p)} \rho_k(y_i) \rho_k(y_j) \right) \right]^2 \\ &= \left( \sum_{k=K+1}^{\infty} U_{nk}^* \right)^2 \end{aligned}$$

Thus,

$$\begin{aligned} \mathbb{E} \left( (T_n - T_{nK})^2 | (x, z)^{(n)} \right) &= (n-1)^2 \binom{n}{2}^{-1} \left( \sum_{k=K+1}^{\infty} U_{nk}^* \right)^2 \\ &\leq 2 \left( \sum_{k=K+1}^{\infty} U_{nk}^* \right)^2 \end{aligned}$$

Since

$$\left( \sum_{k=1}^{\infty} U_{nk}^* \right)^2 = \left( \binom{n}{2}^{-1} \sum_{i \neq j} \varepsilon(z_i; \theta^*) k_p(x_i, x_j) \varepsilon(z_j; \theta^*) \right)^2 < \infty$$

One can fix  $\omega$  and let  $\delta > 0$  be given, then choose and fix  $K$  large enough that

$$|\omega| \left( 2 \left( \sum_{k=K+1}^{\infty} U_{nk}^* \right)^2 \right)^{1/2} < \delta$$

Thus we have

$$|\mathbb{E}(e^{i\omega T_n} | (x, z)^{(n)}) - \mathbb{E}(e^{i\omega T_{nK}} | (x, z)^{(n)})| < \delta \quad (27)$$

Next we show that  $T_{nK} | (x, z)^{(n)} \xrightarrow{d} Y_k = \sum_{k=1}^K \tau_k^{(p)} (W_k^2 - 1)$ . Let

$$W_{kn} = \frac{1}{\sqrt{n}} \sum_{i=1}^n v_i \rho_k(y_i); \quad Z_{kn} = \frac{1}{n} \sum_{i=1}^n v_i^2 \rho_k^2(y_i)$$

then,

$$T_{nK} = \sum_{k=1}^K \tau_k^{(p)} (W_{nk}^2 - Z_{nk})$$

Notice that

$$\mathbb{E}(W_{nk} | (x, z)^{(n)}) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \rho_k(y_i) (\mathbb{E}V) = 0$$

and

$$\text{Cov}(W_{jn}, W_{kn} | (x, z)^{(n)}) = \frac{1}{n} \sum_{i=1}^n \rho_k(y_i) \rho_j(y_i) \xrightarrow{p} \begin{cases} 1, & j = k \\ 0, & j \neq k \end{cases}$$

Therefore, by Lindeberg–Levy CLT,

$$(W_{1n}, \dots, W_{Kn}) | (x, z)^{(n)} \xrightarrow{d} N(0, I_K)$$

Furthermore, by SLLN,

$$(Z_{1n}, \dots, Z_{Kn}) | (x, z)^{(n)} \xrightarrow{p} (1, \dots, 1)$$

Thus,

$$T_{nK} | (x, z)^{(n)} \xrightarrow{d} Y_k = \sum_{k=1}^K \tau_k^{(p)} (W_k^2 - 1)$$

and for all  $n$  sufficiently large

$$|\mathbb{E}(e^{i\omega T_{nK}} | (x, z)^{(n)}) - \mathbb{E}e^{i\omega Y_K}| < \delta \quad (28)$$

Finally, denote  $Y$  as the limit in the mean square of  $Y_k$  as  $K \rightarrow \infty$ . Then

$$\begin{aligned} |\mathbb{E}e^{i\omega Y_K} - \mathbb{E}e^{i\omega Y}| &\leq |\omega| [\mathbb{E}(Y - Y_K)^2]^{1/2} \\ &\leq |\omega| [\mathbb{E}(W_1^2 - 1)^2]^{1/2} \left[ \sum_{k=K+1}^{\infty} (\tau_k^{(p)})^2 \right]^{1/2} \\ &< \delta [\mathbb{E}(W_1^2 - 1)^2]^{1/2} \end{aligned} \quad (29)$$

Combining inequality equations (27), (28) and (29), we have, for any  $\omega$  and any  $\delta > 0$ , and for all  $n$  sufficiently large,

$$|\mathbb{E}(e^{i\omega T_n} | (x, z)^{(n)}) - \mathbb{E}(e^{i\omega Y})| \leq \delta \left( 2 + [\mathbb{E}(W_1^2 - 1)^2]^{1/2} \right)$$

Thus,

$$T_n \xrightarrow{d,*} Y$$

■

**Proof of Theorem 12.**

**Proof** Essentially, we need to show the uniform convergence of  $\widehat{R}_V(\theta)$ , the rest of the consistency proof follows immediately from Theorem 2.1 of [Newey and McFadden \(1994\)](#).

To prove that  $\sup_{\theta \in \Theta} |\widehat{R}_V(\theta) - R_k(\theta)| \xrightarrow{P} 0$ , we need to show that (1)  $f_\theta(v, v')$  is continuous at each  $\theta$  with probability one; and (2)  $\mathbb{E}_{V, V'}(\sup_{\theta \in \Theta} |f_\theta(V, V')|) < \infty$ , and  $\mathbb{E}_{V, V}(\sup_{\theta \in \Theta} |f_\theta(V, V)|) < \infty$  (Lemma 8.5 of [Newey and McFadden \(1994\)](#)).

To this end, we can check that

$$\begin{aligned} |f_\theta(v, v')| &= |\varepsilon(z; \theta)k(x, x')\varepsilon(z'; \theta)| \\ &\leq |\varepsilon(z; \theta)||\varepsilon(z'; \theta)||k(x, x')| \\ &\leq |\varepsilon(z; \theta)||\varepsilon(z'; \theta)|\sqrt{k(x, x)k(x', x')} \end{aligned}$$

Since  $\Theta$  is compact,  $\mathbb{E}(|Y|) < \infty$ , and  $\mathbb{E}_\theta(Y|X) < \infty$ , we have  $|\varepsilon(z; \theta)| < \infty$  for all  $\theta \in \Theta$ . Furthermore,  $k(\cdot, \cdot)$  is bounded by Assumption A4, we have  $f_\theta(v, v') < \infty$  and thus  $f_\theta(v, v')$  is continuous at each  $\theta$ .

Next, observe that

$$\begin{aligned} \mathbb{E}_{V, V'} \left( \sup_{\theta \in \Theta} |f_\theta(V, V')| \right) &\leq \mathbb{E} \left( \sup_{\theta \in \Theta} |\varepsilon(Z; \theta)||\varepsilon(Z'; \theta)|\sqrt{k(X, X)k(X', X')} \right) \\ &\leq \mathbb{E} \left( \sup_{\theta \in \Theta} |\varepsilon(Z; \theta)||\varepsilon(Z'; \theta)| \right) \sup_x k(x, x) \\ &= \mathbb{E}^2 \left( \sup_{\theta \in \Theta} |\varepsilon(Z; \theta)| \right) \sup_x k(x, x) < \infty \end{aligned}$$

and

$$\begin{aligned} \mathbb{E}_{V, V} \left( \sup_{\theta \in \Theta} |f_\theta(V, V)| \right) &\leq \mathbb{E} \left( \sup_{\theta \in \Theta} (|\varepsilon(Z; \theta)|)^2 \right) \sup_x k(x, x) \\ &= \mathbb{E} \left( \sup_{\theta \in \Theta} (|Y - \mathbb{E}_\theta(Y|X)|)^2 \right) \sup_x k(x, x) \\ &\leq 2 \left( \mathbb{E}(|Y|^2) + \mathbb{E} \left( \sup_{\theta \in \Theta} |\mathbb{E}_\theta(Y|X)|^2 \right) \right) \sup_x k(x, x) < \infty \end{aligned}$$

■

**Proof of Theorem 13.**

**Proof** By Theorem 3.1 of [Newey and McFadden \(1994\)](#), we need to show

- (a)  $\hat{\theta} - \theta_0 \xrightarrow{p} 0$ ;
- (b)  $\widehat{R}_V(\theta)$  is twice continuously differentiable;
- (c)  $\sqrt{n}\widehat{R}_V(\theta) \xrightarrow{d} N(0, 4\text{Var}_V(\mathbb{E}_{V'}^2(\nabla_\theta f_{\theta_0}(V, V'))))$ ;
- (d) there exist  $H(\theta)$  that is continuous at  $\theta_0$  and  $\sup_{\theta \in \Theta} \|\nabla_\theta^2 \widehat{R}_V(\theta) - \mathbb{E}(\nabla_\theta^2 f_\theta(V, V'))\|_F \xrightarrow{p} 0$ ;
- (e)  $H(\theta_0)$  is non-singular.

With the consistency result and assumptions made, we only need to check conditions (c) and (d).

Since  $\sqrt{n}(\nabla_\theta \widehat{R}_V(\theta) - \nabla_\theta \widehat{R}_U(\theta)) \xrightarrow{p} 0$  (see, Section 5.7.3 of [Serfling \(1980\)](#)), we can check the asymptotic properties of  $\sqrt{n}\nabla_\theta \widehat{R}_U(\theta)$  instead. Note that

$$\nabla_\theta \widehat{R}_U(\theta) = \frac{1}{n(n-1)} \sum_{i \neq j} \nabla_\theta f_\theta(v_i, v_j)$$

$$\nabla_\theta f_\theta(v_i, v_j) = (g(z_i; \theta)\varepsilon(z_j; \theta) + g(z_j; \theta)\varepsilon(z_i; \theta))k(x_i, x_j)$$

First, we show that

$$\sqrt{n}\nabla_\theta \widehat{R}_U(\theta) \xrightarrow{d} N(0, 4\text{Var}_V(\mathbb{E}_{V'}^2(\nabla_\theta f_{\theta_0}(V, V'))))$$

The proof follows from Section 5.5.1 and 5.5.2 of [Serfling \(1980\)](#). We need to show (i)  $\nabla_\theta \widehat{R}_U(\theta_0) \xrightarrow{p} 0$ ; and (ii) whether  $\text{Var}_V(\mathbb{E}_{V'}^2(f_{\theta_0}(V, V')))$   $> 0$  or not. (i) can be easily obtained by L.L.N.

To verify (ii), note that

$$\text{Var}_V(\mathbb{E}_{V'}^2(\nabla_\theta f_{\theta_0}(V, V'))) = \mathbb{E}_V(\mathbb{E}_{V'}^2(\nabla_\theta f_{\theta_0}(V, V'))) \geq 0$$

where the equality hold if for any  $V$ , there is  $\mathbb{E}_{V'}(\nabla_\theta f_{\theta_0}(V, V')) = 0$ , i.e.,

$$\mathbb{E}_{V'}(\nabla_\theta f_{\theta_0}(V, V')) = \mathbb{E}_{V'}(g(Z'; \theta_0)k(X, X'))\varepsilon(Z; \theta_0) + \mathbb{E}_{V'}(\varepsilon(Z'; \theta_0)k(X, X'))g(Z; \theta_0) = 0$$

Since

$$\mathbb{E}_{V'}(\varepsilon(Z'; \theta_0)k(X, X')) = \mathbb{E}_{X'}(\mathbb{E}(\varepsilon(Z'; \theta_0)|X')k(X, X')) = 0$$

Equality holds if

$$\mathbb{E}_{V'}(g(Z'; \theta_0)k(X, X')) = \mathbb{E}_{X'}(\mathbb{E}(g(Z'; \theta_0)|X')k(X, X')) = 0$$

However, the assertion that  $\mathbb{E}(g(Z'; \theta_0)|X') = 0$  contradicts with the condition that  $H$  is non-singular. Thus, we conclude

$$\text{Var}_V(\mathbb{E}_{V'}^2(\nabla_\theta f_{\theta_0}(V, V'))) > 0$$

Next, we show the uniform consistency of  $\nabla_{\theta}^2 \widehat{R}_U(\theta)$ . Note that

$$\nabla_{\theta}^2 \widehat{R}_U(\theta) = \frac{1}{n(n-1)} \sum_{i \neq j} \nabla_{\theta}^2 f_{\theta}(v_i, v_j)$$

$$\nabla_{\theta}^2 f_{\theta}(v_i, v_j) = \left( g(z_i; \theta) g^{\top}(z_j; \theta) + g(z_j; \theta) g^{\top}(z_i; \theta) + \nabla_{\theta} g(z_i; \theta) \varepsilon(z_j; \theta) + \nabla_{\theta} g(z_j; \theta) \varepsilon(z_i; \theta) \right) k(x_i, x_j)$$

We need to show that (i)  $\nabla_{\theta}^2 \widehat{R}_U(\theta)$  is continuous at each  $\theta$  with probability one, and (ii) there exists  $\mathbb{E}(\sup_{\theta \in \Theta} \|\nabla_{\theta}^2 f_{\theta}(V, V')\|_F^2) < \infty$  and  $\mathbb{E}(\sup_{\theta \in \Theta} \|\nabla_{\theta}^2 f_{\theta}(V, V)\|_F^2) < \infty$ .

To prove (i), we exploit the triangle inequality of the Frobenius norm,

$$\begin{aligned} \|\nabla_{\theta}^2 \widehat{R}_U(\theta)\|_F &\leq \left( 2\|g(z; \theta) g^{\top}(z'; \theta)\|_F + |\varepsilon(z; \theta)| \|\nabla_{\theta} g(z'; \theta)\|_F + |\varepsilon(z'; \theta)| \|\nabla_{\theta} g(z; \theta)\|_F \right) k(x, x') \\ &= d(v, v') \end{aligned}$$

Since  $\mathbb{E}_{\theta}(Y|X)$  is twice continuously differentiable about  $\theta$  and  $\Theta$  is compact, we have  $\mathbb{E}_{\theta}(Y|X)$  bounded as well as each entry of  $g(z; \theta)$  and  $\nabla_{\theta} g(z; \theta)$  for  $\|z\| < \infty$ . Furthermore, since  $k(\cdot, \cdot)$  is also bounded, thus,  $d(v, v') < \infty$  if  $v, v'$  are bounded. We conclude then (i) must hold.

To show (ii), note that

$$\begin{aligned} \mathbb{E} \left( \sup_{\theta \in \Theta} \|\nabla_{\theta}^2 f_{\theta}(V, V')\|_F^2 \right) &\leq 2\mathbb{E} \left( \sup_{\theta \in \Theta} \|g(Z; \theta) g^{\top}(Z'; \theta)\|_F + |\varepsilon(Z; \theta)| \|\nabla_{\theta} g(Z'; \theta)\|_F \right) \sup_x k(x, x) \\ &= 2 \left( \left( \mathbb{E} \sup_{\theta \in \Theta} \|g(Z; \theta)\|_F \right)^2 + \mathbb{E} \left( \sup_{\theta \in \Theta} |\varepsilon(Z; \theta)| \right) \mathbb{E} \left( \sup_{\theta \in \Theta} \|\nabla_{\theta} g(Z'; \theta)\|_F \right) \right) \\ &\quad \times \sup_x k(x, x) \\ &< \infty \end{aligned}$$

$$\begin{aligned} \mathbb{E} \left( \sup_{\theta \in \Theta} \|\nabla_{\theta}^2 f_{\theta}(V, V)\|_F^2 \right) &\leq 2\mathbb{E} \left( \|g(Z; \theta) g^{\top}(Z; \theta)\|_F + |\varepsilon(Z; \theta)| \|\nabla_{\theta} g(Z; \theta)\|_F \right) \sup_x k(x, x) \\ &\leq \left( 2\mathbb{E} \left( \sup_{\theta \in \Theta} \|g(Z; \theta)\|_F^2 \right) + 2\mathbb{E} \left( \sup_{\theta \in \Theta} |\varepsilon(Z; \theta)| \right) \mathbb{E} \left( \sup_{\theta \in \Theta} \|\nabla_{\theta} g(Z; \theta)\|_F \right) \right) \\ &\quad \times \sup_x k(x, x) \\ &< \infty \end{aligned}$$

Therefore, by Lemma 8.5 of [Newey and McFadden \(1994\)](#), we have

$$\sup_{\theta \in \Theta} \|\nabla_{\theta}^2 \widehat{R}_V(\theta) - \mathbb{E}(\nabla_{\theta}^2 f_{\theta}(V, V'))\|_F \xrightarrow{P} 0$$

The rest of the asymptotic normality proof follows from Theorem 3.1 of [Newey and McFadden \(1994\)](#). ■

### Verification of Equation (23).

#### Proof

$$\begin{aligned}\mathbb{E} [\varepsilon(Z; \theta)k_p(X, X')\varepsilon(Z'; \theta)] &= \mathbb{E} \langle \varepsilon(Z; \theta)\mathcal{P}\phi_X(\cdot), \varepsilon(Z'; \theta)\mathcal{P}\phi_{X'}(\cdot) \rangle_{\mathcal{H}(k)} \\ &= \langle \mathbb{E}_{(X,Z)}\varepsilon(Z; \theta)\mathcal{P}\phi_X(\cdot), \mathbb{E}_{(X',Z')}\varepsilon(Z'; \theta)\mathcal{P}\phi_{X'}(\cdot) \rangle_{\mathcal{H}(k)}\end{aligned}$$

Observe that

$$\begin{aligned}\mathbb{E}_{(X,Z)}\varepsilon(Z; \theta)\mathcal{P}\phi_X(\cdot) &= \int_{\mathcal{X}, \mathcal{Z}} \varepsilon(z; \theta) \left( \phi_x(\cdot) - g^\top(z; \theta)\Gamma_{\theta_0}^{-1}\mathbb{E}_{(X,Z)}(g(Z; \theta)\phi_X(\cdot)) \right) dP_{(X,Z)}(x, z) \\ &= \int_{\mathcal{X}, \mathcal{Z}} \varepsilon(z; \theta)\phi_x(\cdot) dP_{(X,Z)}(x, z) \\ &\quad - \int_{\mathcal{X}, \mathcal{Z}} \varepsilon(z; \theta)g^\top(z; \theta)\Gamma_{\theta_0}^{-1}\mathbb{E}_{(X,Z)}(g(Z; \theta)\phi_X(\cdot)) dP_{(X,Z)}(x, z)\end{aligned}$$

Further analysis of the second part, we have

$$\begin{aligned}&\int_{\mathcal{X}, \mathcal{Z}} \varepsilon(z; \theta)g^\top(z; \theta)\Gamma_{\theta_0}^{-1}\mathbb{E}_{(X,Z)}(g(Z; \theta)\phi_X(\cdot)) dP_{(X,Z)}(x, z) \\ &= \int_{\mathcal{X}, \mathcal{Z}} \varepsilon(z; \theta)g^\top(z; \theta)\Gamma_{\theta_0}^{-1} \int_{\mathcal{X}, \mathcal{Z}} g(z; \theta)\phi_x(\cdot) dP_{(X,Z)}(x, z) dP_{(X,Z)}(x, z) \\ &= \int_{\mathcal{X}, \mathcal{Z}} \int_{\mathcal{X}, \mathcal{Z}} \phi_x(\cdot)g^\top(z; \theta)\Gamma_{\theta_0}^{-1}g(z; \theta)\varepsilon(z; \theta) dP_{(X,Z)}(x, z) dP_{(X,Z)}(x, z) \\ &= \int_{\mathcal{X}, \mathcal{Z}} \phi_x(\cdot)g^\top(z; \theta)\Gamma_{\theta_0}^{-1}\mathbb{E}_{(X,Z)}(g(Z; \theta)\varepsilon(Z; \theta)) dP_{(X,Z)}(x, z)\end{aligned}$$

Thus,

$$\begin{aligned}\mathbb{E}_{(X,Z)}\varepsilon(Z; \theta)\mathcal{P}\phi_X(\cdot) &= \int_{\mathcal{X}, \mathcal{Z}} \varepsilon(z; \theta)\phi_x(\cdot) dP_{(X,Z)}(x, z) \\ &\quad - \int_{\mathcal{X}, \mathcal{Z}} \phi_x(\cdot)g^\top(z; \theta)\Gamma_{\theta_0}^{-1}\mathbb{E}_{(X,Z)}(g(Z; \theta)\varepsilon(Z; \theta)) dP_{(X,Z)}(x, z) \\ &= \int_{\mathcal{X}, \mathcal{Z}} \left( \varepsilon(z; \theta) - g^\top(z; \theta)\Gamma_{\theta_0}^{-1}\mathbb{E}_{(X,Z)}(g(Z; \theta)\varepsilon(Z; \theta)) \right) \phi_x(\cdot) dP_{(X,Z)}(x, z) \\ &= \int_{\mathcal{X}, \mathcal{Z}} \varepsilon_p(z; \theta)\phi_x(\cdot) dP_{(X,Z)}(x, z) = \mathbb{E}_{(X,Z)}\varepsilon_p(Z; \theta)\phi_X(\cdot)\end{aligned}$$

Hence,

$$\mathbb{E} [\varepsilon(Z; \theta)k_p(X, X')\varepsilon(Z'; \theta)] = \mathbb{E} [\varepsilon_p(Z; \theta)k(X, X')\varepsilon_p(Z'; \theta)]$$

■

## D. Additional Simulations

We provide more simulation results here. We focus on two sets of DGPs. The first batch of DGPs focuses on the high-frequency alternatives:

- DGP-Freq( $m$ ):  $Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + 2\sin(mX_{1i})\sin(mX_{2i}) + \sigma_i \varepsilon_i$

where  $X_1, X_2 \sim N(0, 3)$ ,  $\sigma_i = (0.1 + X_1^2 + X_2^2)^{1/2}$ , and we set  $\varepsilon_i \sim N(0, 1)$ ,  $\beta_j = 1; j = 0, 1, 2$ . We specify  $m = 0.5, 1.0, 2.0$ , corresponding to low-, moderate-, and high-frequency alternatives, respectively. We use the Gaussian kernel with tuning parameter set as described in Section 6, and the IMQ kernel with parameters  $c = 1, \gamma = 5$ . The following simulation results further illustrate the point that different kernels would have different power properties. Specifically, for moderate- and high-frequency alternatives, the Gaussian kernel only has trivial power, while for the IMQ kernel, we obtain admirably well power properties.

Table 8: Simulation Results, Frequency Alternatives

N=100	$\alpha = 0.1$		$\alpha = 0.05$		$\alpha = 0.01$	
DGP	Gaussian	IMQ(1,5)	Gaussian	IMQ(1,5)	Gaussian	IMQ(1,5)
DGP-Freq(0.5)	0.35	0.28	0.238	0.149	0.093	0.012
DGP-Freq(1.0)	0.11	0.283	0.062	0.167	0.013	0.034
DGP-Freq(2.0)	0.214	0.967	0.066	0.1	0.007	0.016
N=200	0.1		0.05		0.01	
DGP	Gaussian	IMQ(1,5)	Gaussian	IMQ(1,5)	Gaussian	IMQ(1,5)
DGP-Freq(0.5)	0.497	0.384	0.373	0.268	0.168	0.07
DGP-Freq(1.0)	0.128	0.462	0.072	0.325	0.013	0.113
DGP-Freq(2.0)	0.12	0.316	0.062	0.196	0.013	0.039
N=400	0.1		0.05		0.01	
DGP	Gaussian	IMQ(1,5)	Gaussian	IMQ(1,5)	Gaussian	IMQ(1,5)
DGP-Freq(0.5)	0.793	0.725	0.692	0.613	0.409	0.339
DGP-Freq(1.0)	0.105	0.751	0.054	0.641	0.012	0.376
DGP-Freq(2.0)	0.099	0.565	0.039	0.397	0.005	0.166

The second batch of DGPs is about non-linear model specification tests. Specifically, we focus on testing the propensity score models. The propensity score was initially introduced by Rosenbaum and Rubin (1983) to adjust for observable differences between the treatment and control groups when treatment is binary. It is defined as the conditional probability of receiving treatment given a vector of pre-treatment covariates. It is well understood that one can use propensity scores to estimate causal effects through matching, weighting, regression, subclassification, or their combinations.

Given the high dimensionality of available pre-treatment covariates and limited sample size, researchers are coerced to adopt a parametric model for the propensity score to bypass the “curse of dimensionality”.

We consider the following DGPs, which are similar to Sant’Anna and Song (2020); Sant’Anna and Song (2019):

- DGP1:  $D^* = -\frac{\sum_{j=1}^{10} X_j}{6} - \varepsilon;$
- DGP2:  $D^* = -1 - \frac{\sum_{j=1}^{10} X_j}{10} + \frac{X_1 X_2}{2} - \varepsilon;$
- DGP3:  $D^* = -1 - \frac{\sum_{j=1}^{10} X_j}{10} + \frac{X_1 \sum_{k=2}^5 X_k}{4} - \varepsilon;$
- DGP4:  $D^* = -1 - \frac{\sum_{j=1}^{10} X_j}{10} + \frac{\sum_{k=1}^{10} X_k^2}{10} - \varepsilon;$
- DGP5:  $D^* = \frac{-0.1+0.1\sum_{j=1}^5 X_j}{\exp(-0.2\sum_{k=1}^{10} X_k)} - \varepsilon;$

For each DGPs,  $D = \mathbb{I}\{D^* > 0\}$ ,  $\varepsilon \perp X$ , with  $X = (1, X_1, \dots, X_{10})^\top$ , where  $X_1 = Z_1$ ,  $X_2 = (Z_1 + Z_2)/\sqrt{2}$ ,  $X_k = Z_k; k = 3, \dots, 10$ , and  $\{Z_k; k = 1, \dots, 10\}$  and  $\varepsilon$  are i.i.d standard normal random variables.

For DGP1-DGP5, the null  $H_0$  considered is

$$H_0 : \exists \theta^* = (\theta_0, \theta_1, \dots, \theta_{10}) \in \Theta : \mathbb{E}(D|X) = \Phi(X^\top \theta^*) \quad P_X - a.s$$

where  $\Phi(\cdot)$  is the CDF of the standard normal distribution. We estimate  $\theta^*$  using the probit maximum likelihood. Under the null, the expectation of the general residual is now

$$\mathcal{E}(X; \theta^*) = \mathbb{E}\left(D - \Phi(X^\top \theta^*)|X\right) = 0$$

Clearly, DGP1 falls under  $H_0$ , whereas DGP2-DGP5 fall under  $H_1$ , i.e., the negation of the null.

We use the Gaussian kernel with the tuning parameter described in Section 6 to perform the test. The simulation results are presented in Table 9. From the results of DGP1, we find that the actual finite sample size of the proposed test is close to its nominal size, even when the sample size is as small as 100. The proposed test performs admirably well in most alternatives (DGP2-DGP4), however, in DGP5, the proposed test has a weak power.

Table 9: Simulation Results, Propensity Scores

	$\alpha = 0.1$			$\alpha = 0.05$			$\alpha = 0.01$		
	N=100	N=200	N=400	N=100	N=200	N=400	N=100	N=200	N=400
DGP1	0.099	0.104	0.1	0.054	0.057	0.051	0.014	0.011	0.01
DGP2	0.436	0.767	0.978	0.341	0.669	0.957	0.191	0.445	0.883
DGP3	0.256	0.479	0.816	0.177	0.368	0.721	0.073	0.177	0.49
DGP4	0.557	0.858	0.992	0.483	0.802	0.985	0.295	0.669	0.967
DGP5	0.11	0.167	0.221	0.072	0.096	0.125	0.02	0.023	0.032